

Noise-Resilient Designs for Optical Neural Networks

Gianluca Kosmella^{a,b,*}, Ripalta Stabile^a, Jaron Sanders^b

^aDepartment of Electrical Engineering, Eindhoven University of Technology, PO Box 513, 5600 MB Eindhoven, The Netherlands

^bDepartment of Mathematics and Computer Science, Eindhoven University of Technology, PO Box 513, 5600 MB Eindhoven, The Netherlands

Abstract

All analog signal processing is fundamentally subject to noise, and this is also the case in modern implementations of **Optical Neural Networks (ONNs)**. Therefore, to mitigate noise in **ONNs**, we propose two designs that are constructed from a given, possibly trained, **Neural Network (NN)** that one wishes to implement. Both designs have the capability that the resulting **ONNs** gives outputs close to the desired **NN**.

To establish the latter, we analyze the designs mathematically. Specifically, we investigate a probabilistic framework for the first design that establishes that the design is correct, i.e., for any feed-forward **NN** with Lipschitz continuous activation functions, an **ONN** can be constructed that produces output arbitrarily close to the original. **ONNs** constructed with the first design thus also inherit the universal approximation property of **NNs**. For the second design, we restrict the analysis to **NNs** with linear activation functions and characterize the **ONNs'** output distribution using exact formulas.

Finally, we report on numerical experiments with LeNet **ONNs** that give insight into the number of components required in these designs for certain accuracy gains. We specifically study the effect of noise as a function of the depth of an **ONN**. The results indicate that in practice, adding just a few components in the manner of the first or the second design can already be expected to increase the accuracy of **ONNs** considerably.

Keywords: Optical Neural Networks, Law of Large Numbers, Universal Approximation

1. Introduction

Machine Learning (ML) is a computing paradigm in which problems that are traditionally challenging for programmers to explicitly write algorithms for, are solved by learning algorithms that improve automatically through experience. That is, they “learn” structure in data. Prominent examples include image recognition [1], semantic segmentation [2], human-level control in video games [3], visual tracking [4], and language translation [5].

Classical computers are designed and best suited for serialized operations (they have a central processing unit and separated memory), while the data-driven **ML** approach requires decentralized and parallel calculations at high bandwidth as well as continuous processing of parallel data. To illustrate how **ML** can benefit from a different architecture, we can consider performance relative to the number of executed operations, also indicated as **Multiply–Accumulate Operation (MAC)** rates, and the energy efficiency, i.e., the amount of energy spent to execute one single operation. Computational efficiency in classical computers levels off below 10 **GMAC/s/W** [6].

An alternative computing architecture with a more distributed interconnectivity and memory would allow for greater energy efficiency and computational speed. An inspiring example would be an architecture such as the brain. The brain

is able to perform about 10^{18} **MAC/s** using only 20 W of power [6], and operates approximately 10^{11} neurons with an average number of inputs for each of about 10^4 synapses. This leads to an estimated total of 10^{15} synaptic connections, all conveying signals up to 1 kHz bandwidth. The brain’s computational efficiency (being less than 1 aJ per **MAC**) is then about 8 orders of magnitude higher than the one of current supercomputers, which operate instead at 100 pJ per **MAC** [6].

Connecting software to hardware through computing architecture tailored to **ML** tasks is the endeavor of research within the field of neuromorphic computing. The electronics community is now busy developing non-von Neumann computing architectures to enable information processing with an energy efficiency down to a few pJ per operation. Aiming to replicate fundamentals of biological neural circuits in dedicated hardware, important advances have been made in neuromorphic accelerators [7]. These advances are based on the spiking architectural models, which are still not fully understood. **Deep Learning (DL)**-focused approaches, on the other hand, aim to construct hardware that efficiently realizes **DL** architectures, while eliminating as much of the complexity of biological neural networks as possible. Among the most powerful **DL** hardware we can name the GPU-based **DL** accelerators hardware [8, 9, 10, 11, 12], as well as emerging analogue electronic Artificial Intelligence chipsets that tend to collocate processing and memory to minimize the memory–processor communication energy costs (e.g. the analogue

*Corresponding author

Email address: g.k.kosmella@tue.nl (Gianluca Kosmella)

crossbar approaches [13]). The *Mythic*'s architecture, for example, can yield high accuracy in inference applications within a remarkable energy efficiency of just half a pJ per MAC. Even if the implementation of neuromorphic approaches is visibly bringing outstanding record energy efficiencies and computation speeds, neuromorphic electronics is already struggling to offer the desired data throughput at the neuron level. Neuromorphic processing for high-bandwidth applications requires GHz operation per neuron, which calls for a fundamentally different technology approach.

1.1. Optical Neural Networks

A major concern with neuromorphic electronics is that the distributed hardware needed for parallel interconnections is impractical to realize with classical metal wiring: a trade-off applies between interconnectivity and bandwidth, limiting these engine's utilization to applications in the kHz and sub-GHz regime. When sending information not through electrical signals but via optical signals, the optical interconnections do not undergo interference and the optical bandwidth is virtually unlimited. This can for example be achieved when exploiting the color and/or the space and/or the polarization and/or the time domain, thus allowing for applications in the GHz regime. It has been theorized that photonic neuromorphic processors could operate ten thousand times faster while using less energy per computation [14, 15, 16, 17]. Photonics therefore seems to be a promising platform for advances in neuromorphic computing.

Implementations of weighted addition for *Optical Neural Networks (ONNs)* include Mach–Zehnder Interferometer-based Optical Interference Units [18], time-multiplexed and coherent detection [19], free space systems using spatial light modulators [20] and Micro–Ring–Resonator-based weighting bank on silicon [21]. Furthermore, Indium phosphide-integrated optical cross-connect using Semiconductor Optical Amplifiers as single stage weight elements, as well as Semiconductor Optical Amplifier-based wavelength converters [22, 23, 24] have been demonstrated for allowing *All-Optical (AO) Neural Networks (NNs)*. A comprehensive review of all the approaches used in integrated photonics can be found in [25].

Next to these promises, aspects like implementation of nonlinearities, access and storage of weights in on-chip memory, and noise sources in analog photonic implementations, all pose challenges in devising scalable photonic neuromorphic processors and accelerators. These challenges also occur when they are embedded within end-to-end systems. Fortunately, arbitrary scalability of these networks has been demonstrated, with a certain noise and accuracy. However, it would be useful to envision new architectures to reduce noise even more.

1.2. Noise in ONNs

The types of noise in *ONNs* include thermal crosstalk [26], cumulative noise in optical communication links [27, 28] and noise deriving from applying an activation function [29].

In all these studies, the noise is considered to be approximated well by *Additive White Gaussian Noise (AWGN)*.

For example, taking the studies [26, 28, 27, 29, 30] as starting point, the authors of [31] model an *ONN* as a communication channel with *AWGN*. We follow this assumption and will model an *ONN* as having been built up from interconnected nodes with noise in between them. This generic approach does not restrict us to any specific device that may be used in practice.

The model also applies to the two alternative designs of an *AO* implementation of a *NN* (see for example [32]) and the case of an *optical/electrical/optical (O/E/O) NN* [22]. In an *AO NN*, the activation function is applied by manipulating an incoming electromagnetic wave. Modulation (and the *AWGN* it causes) only occurs prior to entering an *AO NN* (or equivalently, in the first layer). For the remainder of the network the signal remains in the optical domain. Here, when applying the optical activation function a new source of noise is introduced as *AWGN* at the end of each layer. Using the *O/E/O* network architecture, the weighted addition is performed in the optical realm, but the light is captured soon after each layer, where it is converted into an electrical and digital signal and the activation function is applied via software on a computer. The operation on the computer can be assumed to be noiseless. However, since the result again needs to be modulated (to be able to act as input to the next layer), modulation noise is added. We can further abstract from the specifics of the *AO* and *O/E/O* design and see that in either implementation noise occurs at the same locations within the mathematical modeling, namely *AWGN* for weighted addition and afterwards *AWGN* from an optical activation function or from modulation, respectively. This means that we do not need to distinguish between the two design choices in our modeling; we only need to choose the corresponding *AWGN* term after activation.

The operation of a layer of a feed-forward *NN* can be modeled by multiplying a matrix W with an input vector x (a bias term b can be absorbed into the matrix–vector product and will therefore be suppressed in notation here) and then applying an activation function $f : \mathbb{R} \rightarrow \mathbb{R}$ element-wise to the result. Symbolically,

$$x \mapsto f(Wx).$$

Now, concretely, the noise model that we study is described by

$$x \mapsto f(Wx + \text{Normal}(0, \Sigma_w)) + \text{Normal}(0, \Sigma_a), \quad (1)$$

for each hidden layer of the *ONN*. Here $\text{Normal}(0, \Sigma)$ denotes the multivariate normal distribution with mean vector 0 and covariance matrix Σ . More specifically, Σ_w , Σ_a and Σ_m are the covariance matrices associated with weighted addition, application of the activation function, and modulation, respectively. Figure 1 gives a schematic representation of the noise model under study. As we have seen above, in the *O/E/O* case we have $\Sigma_a = \Sigma_m$, otherwise Σ_a is due to the specific structure of the photonic activation function. The first layer, regardless of an *AO* or *O/E/O* network, sees a modulated input x , i.e., $x + \text{Normal}(0, \Sigma_m)$, and afterwards the same

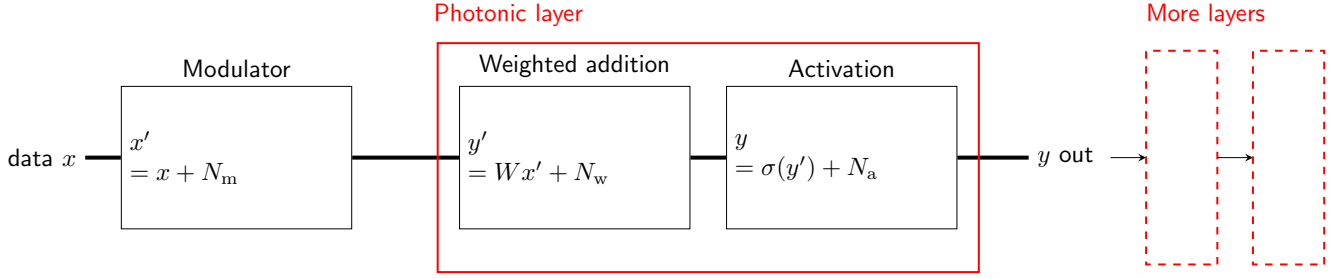


Figure 1: Schematic depiction of the noise model of ONNs that we study. First, data x is modulated onto light. This step adds an AWGN term N_m . This light enters the Photonic Layer, in which a weighted addition takes place, adding AWGN N_w . The activation function is then applied, adding AWGN N_a . The activation function may be applied by photo-detecting the signal of the weighted addition, turning it to a digital signal and applying the activation function on a computer. The result of that action would then be modulated again, to produce the optical output of the photonic neuron. The modulator is thus only required in the first layer, as each photonic neuron takes in light and outputs light.

steps of weighing and applying an activation function, that is (1). Arguably the hidden layers and their noise structure are the most important parts, especially in deep NNs. Therefore, the main equation governing the behavior of the noise propagation in an ONN will remain (1).

1.3. Noise-resistant designs for ONNs

The main contribution of this paper lies in analyzing two noise reduction mechanisms for feed-forward ONNs. The mechanisms are derived from the insight that noise can be mitigated through averaging because of the law of large numbers, and they are aimed at using the enormous bandwidth that photonics offer. The first design (Design A) and its analysis are inspired by recent advancements for NNs with random edges in [33]; the second design (Design B) is new and simpler to implement, but comes without a theoretical guarantee of correctness for nonlinear ONNs, specifically.

Both designs—illustrated in Figure 2—are built from a given NN for which an optical implementation is desired. Each design proposes a larger ONN by taking parts of the original NN, and duplicating and arranging them in a certain way. If noise is absent, then this larger ONN produces the same output as the original NN; and, if noise is present, then this ONN produces an output closer to the desired NN than the direct implementation of the NN as an ONN without modifications would give.

The first mechanism to construct a larger ONN suppressing inherent noise of analog systems starts with a certain number of copies N of the input data. The copies are all processed independently by (in parallel arranged copies of) the layers. Each copy of a layer takes in multiple input copies to produce the result of weighted addition, to which the activation mechanism is applied. The copies that are transmitted to each layer (or set of parallel arrayed layers) are independent of each other. The independent outputs function as inputs to the upcoming (copies of the) layers, and so on and so forth.

The idea of the second design is to use multiple copies of the input, on which the weighted addition is performed. The noisy products of the weighted addition are averaged to a single number/light beam. This average is then copied and

the multiple copies are fed through the activation function, creating multiple noisy activations to be used as the next layer's input, and so on.

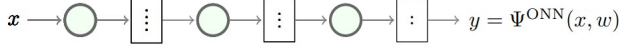
1.4. Summary of results

Using Design A, we are able to establish that ONNs possess the same theoretical properties as NNs. Specifically, we can prove that any NN can be approximated arbitrarily well by an ONN built using Design A (Theorem 1). Similar considerations for NNs with random edges can be found in [33], but the noise model and proof method are different. Here, we first bound the deviation of an ONN and a noiseless NN. To this bound Hoeffding's inequality is then applied.

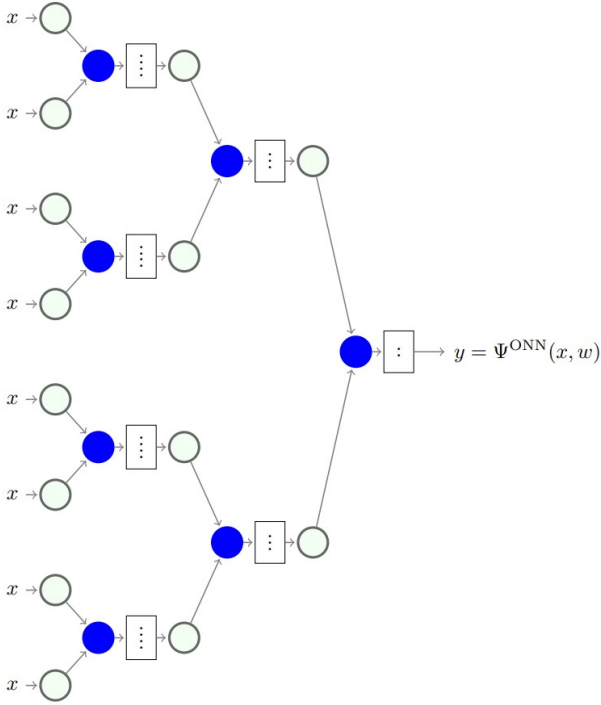
Establishing this theoretical guarantee, however, is done by increasing the number of components exponentially as the depth of the network increases. The current proof shows that for an ONN with Design A meant to approximate a NN with L layers arbitrarily well (and thus reduce the noise to negligible levels), a sufficient number of components is $\omega(K^{L(L+1)}L^L)$ for some constant $K > 0$. This is however not to say that such a large number is necessary: it is merely sufficient.

From a practical viewpoint, however, having to use as few components as possible would be more attractive. We therefore also investigate Design B, in which the number of components increases only linearly with the depth of the network. Because Design A already allows us to establish the approximation property of ONNs, we limit our analysis of Design B to linear NNs for simplicity. We specifically establish in Theorem 2 for any linear NN the exact output distribution of an ONN built using Design B. Similar to the guarantee for Design A in Theorem 1, but more restrictively, this implies that any linear NN can be approximated arbitrarily well by some ONN built using Design B. Strictly speaking, Design B now has no guarantee of correctness for nonlinear NNs, but this should practically not withhold us (especially when activations, for instance, are close to linear).

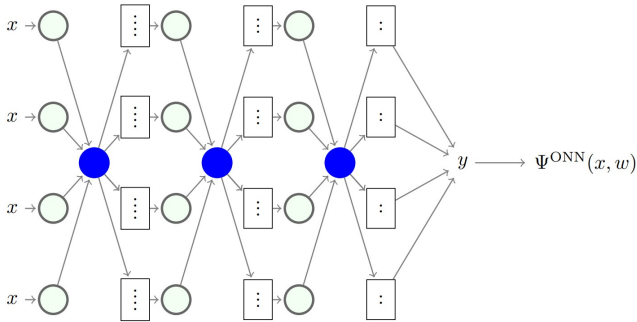
We conduct numerical experiments with Designs A and B by constructing to LeNet ONNs. The numerical results indicate that in practice, adding some components for noise negation is already sufficient to increase the accuracy of an



(a) The original NN.



(b) Design A.



(c) Design B.

Figure 2: (a) Base 4 – 3 – 2 network, light circles indicate activations, boxes indicate post-activations. (b) Example for Design A with 2 layers as input copies to each subsequent layer. The light circles indicate the linear operations/matrix-vector products. The results of the linear operation is averaged (single solid-blue circle) and fed through the activation function, producing the multiple version of the layers output (boxes). (c) Example of Design B.

ONN; an exponential number does not appear not to be necessary (see Figures 3 to 4).

Finally, we want to remark that the high bandwidth of photonic circuits can be exploited to implement the designs as efficiently as possible.

1.5. Outline of the paper

We introduce the AWGN model formally in Section 2. This model is the basis for the analysis of the proposed noise reduction schemes that are next discussed in Sections 3 and 4. There, we specifically define Designs A and B, and each design is followed by a mathematical analysis. The main results are Theorems 1 and 2. Section 5 contain numerical simulations on LeNet ONNs to which we apply Designs A and B. Section 6 concludes; technical details are deferred to the Appendix.

2. Model

We consider general feed-forward NNs implemented on analog optical devices. Noise occurs due to various reasons in those optical devices. Reasons include quantum noise in modulation, chip imperfections, and crosstalk [26, 28, 27, 29, 30].

The noise profiles and levels of different devices differ, but we can, to good approximation, expect AWGN to occur at three separate instances [31]: when modulating, when weighting, and when applying an activation function. The thus proposed AWGN model is formalized next in Section 2.1.

2.1. Feed-forward nonlinear ONNs

We assume that our aim is to implement a feed-forward nonlinear NN with domain \mathbb{R}^{d_0} and range \mathbb{R}^{d_L} , that can be represented by a parameterized function $\Psi^{\text{NN}} : \mathbb{R}^{d_0} \times \mathbb{R}^n \rightarrow \mathbb{R}^{d_L}$ as follows. For $\ell = 1, \dots, L \in \mathbb{N}_+$, Ψ^{NN} must be the composition of the functions

$$\Psi_\ell^{\text{NN}} : \mathbb{R}^{d_{\ell-1}} \rightarrow \mathbb{R}^{d_\ell}, \quad x \mapsto \sigma^{(\ell)}(W^{(\ell)}x + b^{(\ell)}).$$

Here, $W^{(\ell)} \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$ denotes the weight matrix in the ℓ -th layer, $b^{(\ell)} \in \mathbb{R}^{d_\ell \times 1}$ the bias vector in the ℓ -th layer, and $\sigma^{(\ell)} : \mathbb{R}^{d_\ell \times 1} \rightarrow \mathbb{R}^{d_\ell \times 1}$ the activation function in the ℓ -th layer. Specifically, the NN satisfies

$$\Psi^{\text{NN}}(\cdot, w) = \Psi_L^{\text{NN}}(\cdot, w^{(L)}) \circ \dots \circ \Psi_1^{\text{NN}}(\cdot, w^{(1)}), \quad (2)$$

where $w^{(\ell)} = (W^{(\ell)}, b^{(\ell)})$ represents the parameters in the ℓ -th layer. Note that we do not necessarily assume that the activation function is applied component-wise (it could be any high-dimensional function). Such cases are simply contained within the model.

Suppose now that the NN in (2) is implemented as an ONN, but *without* amending its design. AWGN will then disrupt the output of each layer. Specifically, for depths $L \in \mathbb{N}_+$, the ONN will be representable by a function Ψ^{ONN} that is the composition of the noisy functions

$$\Psi_\ell^{\text{ONN}} : \mathbb{R}^{d_{\ell-1}} \rightarrow \mathbb{R}^{d_\ell}, \quad x \mapsto \sigma^{(\ell)}(W^{(\ell)}x + b^{(\ell)} + N_w^{(\ell)}) + N_a^{(\ell)} \quad (3)$$

for $\ell = 1, \dots, L \in \mathbb{N}_+$. Here,

$$N_w^{(\ell)} \stackrel{(d)}{=} \text{Normal}(0, \Sigma_w^{(\ell)}) \text{ and } N_a^{(\ell)} \stackrel{(d)}{=} \text{Normal}(0, \Sigma_a^{(\ell)})$$

denote multivariate normal distributions that describe the **AWGN** within the **ONN**. In other words, the **ONN** will satisfy

$$\Psi^{\text{ONN}}(\cdot, w) = \Psi_L^{\text{ONN}}(\cdot, w^{(L)}) \circ \dots \circ \Psi_1^{\text{ONN}}(\cdot, w^{(1)}) \quad (4)$$

instead of (2). Observe that (4) is a random **NN**; its outcome is uncertain, but hopefully close to that of (2).

2.2. Feed-forward linear ONNs

Let us briefly examine the special case of a feed-forward **linear ONN** in more detail. That is, we now assume additionally that for $\ell = 1, \dots, L$, there exist $e^{(\ell)} \in \mathbb{R}^{d_\ell}$ such that $\sigma^{(\ell)}(y) = D^{(\ell)}y$ where $D^{(\ell)} = \text{diag}(e^{(\ell)})$. In other words, each activation function $\sigma^{(\ell)}$ does element-wise multiplications by constants.

If each activation function is linear, then the output distribution of each layer will remain multivariate normal distributed due to the so-called linear transformation theorem [34, Theorem 1.2.6]. The mean and covariance matrix of the underlying multivariate normal distribution will however be transformed in each layer.

Let us illustrate how the covariance matrix transforms by discussing the first layer in detail. Each layer in (3) can be interpreted as a random function that takes the noisy vector $\mathbf{A}^{(\ell-1)} = (\mathbf{A}_1^{(\ell-1)}, \dots, \mathbf{A}_{d_{\ell-1}}^{(\ell-1)})$ say as input, and produces the even noisier vector $\mathbf{A}^{(\ell)} = (\mathbf{A}_1^{(\ell)}, \dots, \mathbf{A}_{d_\ell}^{(\ell)})$ say as output. Specifically, the noisy input to the first layer is modeled by

$$\mathbf{A}^{(0)} \mid x \stackrel{(d)}{=} x + \mathcal{N}(0, \Sigma_m) \quad (5)$$

because of the modulation error within the first layer. Here $\cdot \mid \cdot$ indicates a conditional random variable. This input next experiences weighted addition and more noise is introduced: the noisy preactivation of the first layer satisfies

$$\mathbf{U}^{(1)} \mid \mathbf{A}^{(0)} \stackrel{(d)}{=} W^{(1)}\mathbf{A}^{(0)} + b^{(1)} + \mathcal{N}(0, \Sigma_w^{(1)}). \quad (6)$$

Combining (5) and (6) with the linear transformation theorem for the multivariate normal distribution as well as the fact that sums of independent multivariate normal random variables are again multivariate normally distributed [34, Theorem 1.2.14], we find that

$$\begin{aligned} \mathbf{U}^{(1)} \mid x &\stackrel{(d)}{=} W^{(1)}x + b^{(1)} + W^{(1)}\mathcal{N}(0, \Sigma_m) + \mathcal{N}(0, \Sigma_w^{(1)}) \\ &\stackrel{(d)}{=} W^{(1)}x + b^{(1)} + \mathcal{N}(0, W^{(1)}\Sigma_m(W^{(1)})^\top + \Sigma_w^{(1)}). \end{aligned}$$

After applying the linear activation function, we obtain

$$\begin{aligned} \mathbf{A}^{(1)} \mid x &\stackrel{(d)}{=} \sigma_1(\mathbf{U}^{(1)}) + \mathcal{N}(0, \Sigma_a^{(1)}) \mid x \\ &\stackrel{(d)}{=} D^{(1)}(W^{(1)}x + b^{(1)}) \\ &\quad + \mathcal{N}\left(0, \Sigma_a^{(1)} + D^{(1)}(W^{(1)}\Sigma_m(W^{(1)})^\top + \Sigma_w^{(1)})(D^{(1)})^\top\right) \end{aligned}$$

$$= \mathcal{N}(\Psi_1^{\text{NN}}(x, w), \Sigma_{\text{ONN}}^{(1)})$$

say. Observe that the unperturbed network's output remains intact, and is accompanied by a centered normal distribution with an increasingly involved covariance matrix:

$$\begin{aligned} \Sigma_{\text{ONN}}^{(1)} &= D^{(1)}(W^{(1)}\Sigma_m(W^{(1)})^\top + \Sigma_w^{(1)})(D^{(1)})^\top + \Sigma_a^{(1)} \\ &= D^{(1)}W^{(1)}\Sigma_m(W^{(1)})^\top(D^{(1)})^\top + D^{(1)}\Sigma_w^{(1)}(D^{(1)})^\top \\ &\quad + \Sigma_a^{(1)}. \end{aligned} \quad (7)$$

Observe furthermore that the covariance matrix in (7) is independent of the bias $b^{(1)}$.

The calculations in eqs. (5) to (7) can readily be extended into a recursive proof that establishes the covariance matrix of the entire linear **ONN**. Specifically, for $\ell = 1, \dots, L$, define the maps

$$\begin{aligned} T^{(\ell)}(\Sigma) &:= D^{(\ell)}W^{(\ell)}\Sigma(W^{(\ell)})^\top(D^{(\ell)})^\top \\ &\quad + D^{(\ell)}\Sigma_w^{(\ell)}(D^{(\ell)})^\top + \Sigma_a^{(\ell)}. \end{aligned} \quad (8)$$

We then have the following:

Proposition 1 (Distribution of linear ONNs) *Assume that there exist vectors $e^{(\ell)} \in \mathbb{R}^{d_\ell}$ such that $\sigma^{(\ell)}(y) = \text{diag}(e^{(\ell)})y$. The feed-forward linear **ONN** in (4) then satisfies*

$$\Psi^{\text{ONN}}(\cdot, w) \stackrel{(d)}{=} \mathcal{N}(\Psi^{\text{NN}}(\cdot, w), \Sigma_{\text{ONN}}^{(L)}),$$

where for $\ell = L, L-1, \dots, 1$,

$$\Sigma_{\text{ONN}}^{(\ell)} = T^{(\ell)}(\Sigma_{\text{ONN}}^{(\ell-1)}); \quad \text{and} \quad \Sigma_{\text{ONN}}^{(0)} = \Sigma_m.$$

In linear **ONNs** with symmetric noise (that is, the **AWGN** of each layer's noise sources has the same covariance matrix), Proposition 1's recursion simplifies. Introduce $P^{(\ell)} := \prod_{i=\ell+1}^L D^{(i)}W^{(i)}$ for notational convenience. The following is proved in Appendix A.1.1:

Corollary 1 (Symmetric noise case) *Within the setting of Proposition 1, assume additionally that for all $\ell \in \mathbb{N}_+$, $\Sigma_a^{(\ell)} = \Sigma_a$ and $\Sigma_w^{(\ell)} = \Sigma_w$. Then,*

$$\begin{aligned} \Sigma_{\text{ONN}}^{(L)} &= P^{(0)}\Sigma_m(P^{(0)})^\top + \sum_{\ell=1}^L P^{(\ell)}\Sigma_a(P^{(\ell)})^\top \\ &\quad + \sum_{\ell=1}^L P^{(\ell)}D^{(\ell)}\Sigma_w(D^{(\ell)})^\top(P^{(\ell)})^\top. \end{aligned}$$

If moreover for all $\ell \in \mathbb{N}_+$, $W^{(\ell)} = W$, $D^{(\ell)} = D$, and $\|D\|_F\|W\|_F < 1$, then

$$\lim_{L \rightarrow \infty} \Sigma_{\text{ONN}}^{(L)} = \sum_{n=0}^{\infty} (DW)^n (D\Sigma_w D^\top + \Sigma_a) ((DW)^n)^\top.$$

Proposition 1 and Corollary 1 describe the output distribution of linear **ONNs** completely.

2.3. Discussion

One way to think of the **AWGN** model in Section 2.1 is to take a step back from the microscopic analysis of individual devices, and consider an **ONN** as a series of black box devices (recall also Figure 1). Each black box device performs their designated task and acts as communication channels with **AWGN**. This way of modeling in order to analyze the impact of noise can also be seen in [31]; and other papers modeling optical channels include [28, 27]. Further papers considering noise in optical systems with similar noise assumptions are [35, 36], where furthermore multiplicative noise is considered when an amplifier is present within the circuit [35]. Qualitatively the results for Design A also apply for multiplicative noise, the scaling however may differ.

Limitations of the model. We note firstly that modeling the noise in **ONNs** as **AWGN** is warranted only in an operating regime with many photons, and is thus unlikely to be a good model for **ONNs** that operate in a regime with just a few photons.

Secondly, due to physical device features and operation conditions, weights, activations, and outputs can only be realized in **ONNs** if their values lie in certain ranges. Such constraints are no part of the model in Section 2. Fortunately, however, the implied range restrictions are usually not a problem in practice. For example, activation functions like sigmoid and tanh map into $[0, 1]$ and $[-1, 1]$, respectively. Additional regularization rules like weight decay also move the entries of weight matrices in **NNs** towards smaller values. In case physical constraints were met one can increase the weight decay parameter to further penalize large weights during training, leading to smaller weights so that the **ONN** is again applicable.

3. Results—Design A

3.1. Reducing the noise in feed-forward **ONNs** (Design A)

Recall that an example of Design A is presented in Figure 2(b). Algorithm 1 constructs this tree-like network, given the desired number of copies n_0, \dots, n_L per layer.

Observe that in Design A, the number of copies utilized in each layer, the n_ℓ , are fixed. There is however only a single copy in the last layer. Its output is the unique output of the **ONN**. Each other layer receives multiple independent inputs. With each of the independent copies weighted addition is performed, and the results are averaged to produce the layer's single output. Having independent incoming copies is achieved by having multiple independent branches of the prior partial networks incoming into a given layer. This means that the single layer L receives n_{L-1} independent inputs of n_{L-1} independent layers $L-1$. Each of the n_{L-1} copies of layer $L-1$ receives n_{L-2} inputs from independent copies of layer $L-2$. Generally, let $n_{\ell-1}$ be the number of copies of layer $\ell-1$ that act as inputs to layer ℓ .

Observe that all copies are created upfront. That means there are $\prod_{\ell=0}^{L-1} n_\ell$ copies of the data. By Algorithm 1,

Algorithm 1 Algorithm to construct a noise reducing network

Require: Input $\mathbf{n} = (n_\ell)_{\ell=0, \dots, L}$

Require: $\prod_{\ell=0}^L n_\ell$ copies of input $x^{(0)}$, named ${}_1x^{(0)}, \dots, (\prod_{\ell=0}^L n_\ell) x^{(0)}$

for $\ell = 0, \dots, L-1$ **do**

for $\alpha = 1, \dots, \prod_{i=\ell}^{L-1} n_i$ **do**

$\alpha \xi^{(\ell)} \leftarrow W^{(\ell+1)} \alpha x^{(\ell)} + b^{(\ell+1)} + \text{Normal}(0, \Sigma_w)$

end for

for $\alpha = 0, \dots, (\prod_{i=\ell}^{L-1} n_i) - 1$ **do**

$\alpha y^{(\ell)} \xleftarrow{\text{averaging}} n_\ell^{-1} (\alpha \cdot n_\ell + 1 \xi^{(\ell)} + \dots + \alpha \cdot n_\ell + n_\ell \xi^{(\ell)})$

$\alpha x^{(\ell+1)} \leftarrow \sigma^{(\ell)}(\alpha y^{(\ell)}) + \text{Normal}(0, \Sigma_a)$

end for

end for

return ${}_1x^{(L)}$

$\prod_{\ell=1}^{L-1} n_\ell$ copies of the first layer are arrayed in parallel to each other, and each of them processes n_0 copies of the data. The outputs of the $\prod_{\ell=1}^{L-1} n_\ell$ arrayed copies of the first layer are the input to the $\prod_{\ell=2}^{L-1} n_\ell$ arrayed copies of the second layer, and so on.

Notice that noise stemming from applying the activation function is subject to a linear transformation in the next layer. The activation function noise can therefore be considered as weight-noise by inserting an identity layer with $\sigma = \text{id}$, $W = I$ and $b = 0$.

We want to verify that a Design A **ONN** yields outputs that are with high probability close to the original noiseless **NN**. Let $\tilde{\Psi}^{\text{ONN}}(x, w)$ the Design A **ONN** and then let

$$\mathbb{P} \left[\sup_{x \in \mathbb{R}^d} \|\Psi^{\text{NN}}(x, w) - \tilde{\Psi}^{\text{ONN}}(x, w)\|_2 < D_L \right] > 1 - C_L, \quad (9)$$

be the desired property. The main result of this section is the following:

Theorem 1 For any $C_L \in (0, 1)$, any $D_L \in (0, \infty)$, and any nonlinear **NN** Ψ^{NN} , with Lipschitz-continuous activations functions with Lipschitz-constants $a^{(i)}$ and weight-matrices $W^{(i)}$, Algorithm 1 is able to construct an **ONN** $\tilde{\Psi}^{\text{ONN}}$ that satisfies (9).

Let the covariance matrices of the occurring **AWGN** be diagonal matrices and let each of the values of the covariance matrices be upper bounded by $\sigma^2 \geq 0$. For any set of $(\kappa_i)_{i=1, \dots, L}$, $(\delta_i)_{i=1, \dots, L}$ such that $\prod (1 - \kappa_\ell) > 1 - C_L$ and $\sum \delta_\ell \leq D_L$, a sufficient number of copies to construct an **ONN** $\tilde{\Psi}^{\text{ONN}}$ that satisfies (9) is given by

$$n_L = 1$$

$$n_\ell \geq \frac{\sigma^2 \left(\prod_{i=\ell+1}^L a^{(i)} \prod_{i=\ell+2}^L \|W^{(i)}\|_{\text{op}} \right)^2}{\delta_{\ell+1}^2} \times \left(\sqrt{2} \frac{\Gamma((d_{\ell+1} + 1)/2)}{\Gamma(d_{\ell+1}/2)} \right)$$

$$+ \sqrt{\frac{C^2}{c} \frac{4^{d_{\ell+1}\sqrt{4}}}{2^{d_{\ell+1}\sqrt{4}-2}} (-\ln(\kappa_{\ell+1}/2)) \frac{1}{\prod_{i=\ell+1}^L n_i}}^2, \\ \ell = L-1, \dots, 0.$$

Here Γ is the gamma function and $C, c > 0$ are absolute constants.

This result is proven in Section 3.3. A consideration on the asymptotic total amount of copies in deep ONNs is relegated to Appendix A.2.1.

3.2. Idea behind Design A

Having the law of large numbers in mind it seems reasonable that the average of multiple experiments would help in achieving a more precise output in the presence of noise. However, it would typically not be correct to just input n identical, deterministic copies of x into n independent ONNs—thus producing n noisy realizations $\Psi^{\text{ONN},1}(x, w), \dots, \Psi^{\text{ONN},n}(x, w)$ say—and then calculate their average in the hope to recover $\Psi^{\text{NN}}(x, w)$. This is because while by the law of large numbers it is true that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \Psi^{\text{ONN},i}(x, w) = \mathbb{E} [\Psi^{\text{ONN}}(x, w)],$$

it is not necessarily true that the expectation $\mathbb{E} [\Psi^{\text{ONN}}(x, w)]$ equals $\Psi^{\text{NN}}(x, w)$. The reason is that activation functions in NNs are typically nonlinear.

We can circumvent the issue by modifying the approach and instead exploit the law of large numbers layer-wise. Recall that in the noiseless NN, layer ℓ maps a fixed input

$$x \mapsto \sigma^{(\ell)}(W^{(\ell)}x + b^{(\ell)}), \quad (10)$$

and that the same layer in the ONN maps the same fixed input

$$x \mapsto \sigma^{(\ell)}(W^{(\ell)}x + b^{(\ell)} + N_w^{(\ell)})$$

instead. If we let $(N^{(i)})_{i \in \{1, \dots, n\}}$ be independent realizations of the distribution of $N_w^{(\ell)}$ (which has mean zero), we can expect by the law of large numbers that for sufficiently large n , the realized quantities

$$\frac{1}{n} \left(\sum_{i=1}^n W^{(\ell)}x + b^{(\ell)} + N^{(i)} \right) \quad \text{and} \quad W^{(\ell)}x + b^{(\ell)}$$

are close to each other. If $\sigma^{(\ell)}$ is moreover sufficiently regular, then we may expect that the realized quantity

$$\sigma^{(\ell)} \left(\frac{1}{n} \left(\sum_{i=1}^n W^{(\ell)}x + b^{(\ell)} + N^{(i)} \right) \right) \quad (11)$$

is close to (10) for sufficiently large n , i.e., close to the unperturbed output of the original layer.

The implementation in (11) can be realized by using n times as many nodes in the hidden layer; thus to essentially

create n copies of the original hidden layer. These independent copies are then averaged. Furthermore, one can allow for different inputs $(x^i)_{i \in \{1, \dots, n\}}$, assuming some statistical properties of their distribution. This will be formalized next in the proof of Theorem 1 in Section 3.3.

3.3. Proof of Theorem 1

For the proof we will first upper bound the deviation between an ONN constructed with Design A and the noiseless NN (Section 3.3.1) and then we find a probabilistic bound on the deviations bound (Section 3.3.2).

3.3.1. An upper-bound for the ONN — NN deviation

The output of the Design A network is

$$\tilde{x} = \sigma^{(L)} \left(\frac{1}{n_{L-1}} \sum_{i=1}^{n_{L-1}} (W^{(L)}\tilde{x}^i + b^{(L)} + N^{(i)}) \right), \quad (12)$$

where each \tilde{x}^i is recursively calculated as

$$\tilde{x}^i = \sigma^{(L-1)} \left(\frac{1}{n_{L-2}} \sum_{j_i=1}^{n_{L-2}} (W^{(L-1)}\tilde{x}^{j_i} + b^{(L-1)} + N^{(j_i)}) \right),$$

the \tilde{x}^{j_i} are calculated as

$$\tilde{x}^{j_i} = \sigma^{(L-2)} \left(\frac{1}{n_{L-3}} \sum_{k_{j_i}=1}^{n_{L-3}} (W^{(L-2)}\tilde{x}^{k_{j_i}} + b^{(L-2)} + N^{(k_{j_i})}) \right),$$

and so on and so forth. The difference in L_2 -norm of (12) and the noiseless NN

$$\sigma^{(L)}(W^{(L)}\sigma^{(L-1)}(W^{(L-1)}(\dots) + b^{(L-1)}) + b^{(L)})$$

can iteratively be bounded by using the Lipschitz property of the activation functions, triangle inequality, and submultiplicativity of the norms.

We start the iteration by bounding

$$\begin{aligned} & \left\| \sigma^{(L)} \left(\frac{1}{n_{L-1}} \sum_{i=1}^{n_{L-1}} (W^{(L)}\tilde{x}^i + b^{(L)} + N^{(i)}) \right) \right. \\ & \quad \left. - \sigma^{(L)}(W^{(L)}\sigma^{(L-1)}(W^{(L-1)}(\dots) + b^{(L-1)}) + b^{(L)}) \right\|_2 \\ & \leq a^{(L)} \left\| \frac{1}{n_{L-1}} \sum_{i=1}^{n_{L-1}} \right. \\ & \quad \left. (W^{(L)}(\tilde{x}^i - \sigma^{(L-1)}(W^{(L-1)}(\dots) + b^{(L-1)})) + N^{(i)}) \right\|_2 \\ & \leq \frac{a^{(L)} \|W^{(L)}\|_{\text{op}}}{n_{L-1}} \\ & \quad \times \left\| \sum_{i=1}^{n_{L-1}} (\tilde{x}^i - \sigma^{(L-1)}(W^{(L-1)}(\dots) + b^{(L-1)})) \right\|_2 \\ & \quad + a^{(L)} \left\| \frac{1}{n_{L-1}} \sum_{i=1}^{n_{L-1}} N^{(i)} \right\|_2. \end{aligned}$$

In the next iteration step the term

$$\left\| \sum_{i=1}^{n_{L-1}} \left(\tilde{x}^i - \sigma^{(L-1)}(W^{(L-1)}(\dots) + b^{(L-1)}) \right) \right\|_2$$

is further bounded by first using the triangle inequality and thereafter bounding in the same way as we did in the first layer:

$$\begin{aligned} & \left\| \sum_{i=1}^{n_{L-1}} \left(\sigma^{(L-1)} \left(\frac{1}{n_{L-2}} \sum_{j_i=1}^{n_{L-2}} \left(W^{(L-1)} \tilde{x}^{j_i} + b^{(L-1)} + N^{(j_i)} \right) \right) \right. \right. \\ & \left. \left. - \sigma^{(L-1)} \left(W^{(L-1)} \left(\sigma^{(L-2)} \left(W^{(L-2)}(\dots) + b^{(L-2)} \right) + b^{(L-1)} \right) \right) \right) \right\|_2 \\ & \leq \frac{a^{(L-1)} \|W^{(L-1)}\|_{\text{op}}}{n_{L-2}} \\ & \quad \times \sum_{i=1}^{n_{L-1}} \left\| \sum_{j_i=1}^{n_{L-2}} \left(\tilde{x}^{j_i} - \sigma^{(L-2)} \left(W^{(L-2)}(\dots) + b^{(L-2)} \right) \right) \right\|_2 \\ & \quad + a^{(L-1)} \sum_{i=1}^{n_{L-1}} \left\| \frac{1}{n_{L-2}} \sum_{j_i=1}^{n_{L-2}} N^{(j_i)} \right\|_2. \end{aligned}$$

Here,

$$\sum_{i=1}^{n_{L-1}} \left\| \sum_{j_i=1}^{n_{L-2}} \left(\tilde{x}^{j_i} - \sigma^{(L-2)} \left(W^{(L-2)}(\dots) + b^{(L-2)} \right) \right) \right\|_2$$

may again be bounded in the same fashion. This leads to the following recursive argument.

Let $\mathcal{F}^{(\ell)}$ be the sum of the differences between—loosely speaking—the ends of the remaining Design A “subtrees” and noiseless NNs “subtrees” at layer ℓ . More specifically, let

$$\begin{aligned} \mathcal{F}^{(L)} & := \left\| \tilde{x} - \sigma^{(L)} \left(W^{(L)}(\dots) + b^{(L)} \right) \right\|; \\ \mathcal{F}^{(\ell)} & := \sum_{i_L=1}^{n_{L-1}} \sum_{i_{L-1}=1}^{n_{L-2}} \dots \sum_{i_{\ell+2}, \dots, i_{L-1}=1}^{n_{\ell+1}} \left\| \sum_{i_{\ell+1}, \ell+2, \dots, L-1=1}^{n_{\ell}} \right. \\ & \quad \left. \left(\tilde{x}^{i_{\ell+1}, \ell+2, \dots, L-1} - \sigma^{(\ell)} \left(W^{(\ell)}(\dots) + b^{(\ell)} \right) \right) \right\|_2, \\ & \quad \forall \ell = 1, \dots, L-1, \end{aligned}$$

the special case of $\mathcal{F}^{(0)}$ will be considered in detail later. For simplicity, we join the sums outside the norm into one. Notice that because $n_L = 1$, we have $\prod_{k=\ell+1}^{L-1} n_k = \prod_{k=\ell+1}^L n_k$, and we can write

$$\mathcal{F}^{(\ell)} = \sum_{i=1}^{\prod_{k=\ell+1}^{L-1} n_k} \left\| \sum_{j_i=1}^{n_{\ell}} \left(\tilde{x}^{j_i} - \sigma^{(\ell)} \left(W^{(\ell)}(\dots) + b^{(\ell)} \right) \right) \right\|_2,$$

where specifically

$$\tilde{x}^{j_i} = \sigma^{(\ell)} \left(\frac{1}{n_{\ell-1}} \sum_{k_{j_i}=1}^{n_{\ell-1}} \left(W^{(\ell)} \tilde{x}^{k_{j_i}} + b^{(\ell)} + N^{(k_{j_i})} \right) \right),$$

and the j_i and k_{j_i} are nothing more than relabelings.

Bounding $\mathcal{F}^{(\ell)}$ using the triangle inequality, Lipschitz property, and submultiplicativity yields

$$\begin{aligned} \mathcal{F}^{(\ell)} & \leq a^{(\ell)} \sum_{i=1}^{\prod_{k=\ell+1}^{L-1} n_k} \sum_{j_i=1}^{n_{\ell}} \left\| \frac{1}{n_{\ell-1}} \sum_{k_{j_i}=1}^{n_{\ell-1}} \left(N^{(k_{j_i})} \right. \right. \\ & \quad \left. \left. + W^{(\ell)} \left(\tilde{x}^{k_{j_i}} - \sigma^{(\ell-1)} \left(W^{(\ell-1)}(\dots) + b^{(\ell-1)} \right) \right) \right) \right\|_2 \\ & \leq \frac{a^{(\ell)} \|W^{(\ell)}\|_{\text{op}}}{n_{\ell-1}} \mathcal{F}^{(\ell-1)} \\ & \quad + a^{(\ell)} \sum_{i=1}^{\prod_{k=\ell+1}^{L-1} n_k} \sum_{j_i=1}^{n_{\ell}} \left\| \frac{1}{n_{\ell-1}} \sum_{k_{j_i}=1}^{n_{\ell-1}} N^{(k_{j_i})} \right\|_2. \end{aligned} \quad (13)$$

We thus found a recursive formula for the bound.

The recursion ends at $\mathcal{F}^{(0)}$. The noiseless NN receives x as input, while the ONN receives modulated input $x + N^{(j_i)}$, where $N^{(j_i)}$ is the modulation noise, i.e., AWGN. Therefore,

$$\begin{aligned} \mathcal{F}^{(0)} & = \sum_{i=1}^{\prod_{k=1}^L n_k} \left\| \sum_{j_i=1}^{n_0} \left((x + N^{(j_i)}) - x \right) \right\|_2 \\ & = \sum_{i=1}^{\prod_{k=1}^L n_k} \left\| \sum_{j_i=1}^{n_0} N^{(j_i)} \right\|_2. \end{aligned} \quad (14)$$

Observe that the x -dependence disappeared.

Readily iterating (13) leads to the bound

$$\begin{aligned} \mathcal{F}^{(L)} & \leq \sum_{\ell=L, L-1, \dots, 1} \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}} \\ & \quad \times \frac{1}{\prod_{k=\ell}^L n_k} \frac{1}{n_{\ell-1}} \sum_{i=1}^{\prod_{k=\ell}^L n_k} \left\| \sum_{j_i=1}^{n_{\ell-1}} N^{(j_i)} \right\|_2. \end{aligned}$$

Therefore, if all the L_2 -norms of the sums of the Gaussians are small at the same time, the network is close to the noiseless NN. Let

$$\begin{aligned} \mathcal{S}_{\ell} & := \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}} \\ & \quad \times \frac{1}{\prod_{k=\ell}^L n_k} \frac{1}{n_{\ell-1}} \sum_{i=1}^{\prod_{k=\ell}^L n_k} \left\| \sum_{j_i=1}^{n_{\ell-1}} N^{(j_i)} \right\|_2. \end{aligned}$$

If for all ℓ

$$\mathbb{P}[\mathcal{S}_{\ell} \leq \delta_{\ell}] > 1 - \kappa_{\ell}, \quad (15)$$

and moreover $\sum \delta_{\ell} \leq D_L$ as well as $\prod(1 - \kappa_{\ell}) > 1 - C_L$, then (9) holds. This can be seen by bounding

$$\mathbb{P} \left[\sup_{x \in \mathbb{R}^d} \left\| \Psi^{\text{NN}}(x, w) - \tilde{\Psi}^{\text{ONN}}(x, w) \right\|_2 < D_L \right]$$

$$\begin{aligned} &\geq \mathbb{P}\left[\sum_{\ell} \mathcal{S}_{\ell} < D_L\right] \geq \mathbb{P}\left[\bigcap_{\ell} \{\mathcal{S}_{\ell} < \delta_{\ell}\}\right] \\ &= \prod_{\ell} \mathbb{P}\left[\mathcal{S}_{\ell} < \delta_{\ell}\right] > \prod_{\ell} (1 - \kappa_{\ell}) > 1 - C_L. \end{aligned}$$

Here, in the first inequality the dependence on x disappears due to (14).

3.3.2. Bound for deviations

We next consider the \mathcal{S}_{ℓ} for which we want to guarantee that

$$\mathbb{P}[\mathcal{S}_{\ell} < \delta_{\ell}] > 1 - \kappa_{\ell}.$$

Let $m_{\ell} = \prod_{k=\ell}^L n_k$. By assumption the $N_k^{(j_i)}$ are independent and identically $\text{Normal}(0, \sigma_k^2)$ distributed, where $\sigma_k^2 \leq \sigma^2$, for some common σ^2 . We are lower bounding the number of copies required, therefore using AWGN with higher variance only increases the lower bound, as the calculations below show. We calculate the bound exemplary for $N^{(j_i)}$ distributed according to $\text{Normal}(0, \sigma^2)$, re-substituting σ_k^2 below in (18) (which is the bound given in Theorem 1) thus covers the case of $N_k^{(j_i)} \stackrel{(d)}{=} \text{Normal}(0, \sigma_k^2)$.

Each component of the vector

$$\sum_{j_i=1}^{n_{\ell-1}} N^{(j_i)} = \left(\sum_{j_i=1}^{n_{\ell-1}} N_1^{(j_i)}, \dots, \sum_{j_i=1}^{n_{\ell-1}} N_{d_{\ell}}^{(j_i)} \right)^{\top}$$

is assumed to be $\text{Normal}(0, n_{\ell-1}\sigma^2) = \sqrt{n_{\ell-1}}\sigma \text{Normal}(0, 1)$ distributed. It then holds that

$$\sum_{i=1}^{m_{\ell}} \left\| \sum_{j_i=1}^{n_{\ell-1}} N^{(j_i)} \right\|_2 \stackrel{(d)}{=} \sum_{i=1}^{m_{\ell}} \sqrt{n_{\ell-1}}\sigma \|\text{Normal}(0, I_d)\|_2.$$

This is a sum of independent chi-distributed random variables, which means they are sub-gaussian (see below that we can calculate the sub-gaussian norm and it is indeed finite). Thus Hoeffding's inequality applies, according to which, for X_1, \dots, X_n independent, mean zero, sub-gaussian random variables, for every $t \geq 0$

$$\mathbb{P}\left[\left|\sum_{i=1}^N X_i\right| < t\right] > 1 - 2 \exp\left(-\frac{ct^2}{\sum_{i=1}^N \|X_i\|_{\psi_2}^2}\right) \quad (16)$$

holds; see e.g. [37, Theorem 2.6.2]. Here $c > 0$ is an absolute constant (see [37, Theorem 2.6.2]) and

$$\|X\|_{\psi_2} := \inf\{t > 0 : \mathbb{E}[\exp(X^2/t^2)] \leq 2\}.$$

To apply Hoeffding's inequality in our setting, we need to center the occurring random variables. For $N^{(i)} \sim \text{Normal}(0, I_d)$, the term $\|N^{(i)}\|_2$ is chi distributed with mean

$$\mu_d = \sqrt{2} \frac{\Gamma((d+1)/2)}{\Gamma(d/2)}, \quad (17)$$

where Γ is the gamma function, see e.g. [38, p.238].

Consider

$$\begin{aligned} &\mathbb{P}\left[\frac{\prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}}{m_{\ell} \sqrt{n_{\ell-1}}} \sigma \sum_{i=1}^{m_{\ell}} \|N^{(i)}\|_2 < \delta_{\ell}\right] \\ &= \mathbb{P}\left[\frac{\prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}}{m_{\ell} \sqrt{n_{\ell-1}}} \sigma \sum_{i=1}^{m_{\ell}} \left(\|N^{(i)}\|_2 - \mu_{d_{\ell}}\right) \right. \\ &\quad \left. < \delta_{\ell} - \frac{\prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}}{m_{\ell} \sqrt{n_{\ell-1}}} \sigma m_{\ell} \mu_{d_{\ell}}\right] \end{aligned}$$

which equals

$$\begin{aligned} &\mathbb{P}\left[\sum_{i=1}^{m_{\ell}} \left(\|N^{(i)}\|_2 - \mu\right) \right. \\ &\quad \left. < \frac{m_{\ell} \sqrt{n_{\ell-1}} \delta_{\ell}}{\sigma \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}} - m_{\ell} \mu\right] \end{aligned}$$

and is lower bounded (compare to (16)) by

$$1 - 2 \exp\left(\frac{-c \left(\frac{m_{\ell} \sqrt{n_{\ell-1}} \delta_{\ell}}{\sigma \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}} - m_{\ell} \mu_{d_{\ell}}\right)^2}{\sum_{i=1}^{m_{\ell}} \left\| \|N^{(i)}\|_2 - \mu\right\|_{\psi_2}^2}\right),$$

which in turn is lower bounded by

$$1 - 2 \exp\left(\frac{-c \left(\frac{m_{\ell} \sqrt{n_{\ell-1}} \delta_{\ell}}{\sigma \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}} - m_{\ell} \mu_{d_{\ell}}\right)^2}{\sum_{i=1}^{m_{\ell}} C^2 \left\| \|N^{(i)}\|_2 \right\|_{\psi_2}^2}\right),$$

where $C > 0$ is an absolute constant (see [37, Lemma 2.6.8]). For a chi distributed random variable \mathbf{X} it holds that

$$\mathbb{E}[\exp(\mathbf{X}^2/t^2)] = M_{\mathbf{X}^2}(1/t^2)$$

where $M_{\mathbf{X}^2}(s)$ is the moment generating function of \mathbf{X}^2 —a chi-squared distributed random variable. It is known (see e.g. [39, Appendix 13]) that

$$M_{\mathbf{X}^2}(s) = (1 - 2s)^{-d_{\ell}/2}$$

for $s < \frac{1}{2}$. Accordingly for $2 < t^2$, the property in the definition of the sub-gaussian norm

$$\mathbb{E}[\exp(\mathbf{X}^2/t^2)] = \left(1 - 2\frac{1}{t^2}\right)^{-d_{\ell}/2} \leq 2$$

is satisfied for all t for which

$$t \geq \max\left\{\sqrt{\frac{4 \sqrt[d_{\ell}]{4}}{2 \sqrt[d_{\ell}]{4} - 2}}, \sqrt{2}\right\} = \sqrt{\frac{4 \sqrt[d_{\ell}]{4}}{2 \sqrt[d_{\ell}]{4} - 2}}$$

holds. The square of the sub-gaussian norm of the chi distributed random variables is thus

$$\left\| \|N^{(i)}\|_2 \right\|_{\psi_2}^2 = \frac{4 \sqrt[d_{\ell}]{4}}{2 \sqrt[d_{\ell}]{4} - 2}.$$

Substituting the norm into the lower bound yields

$$1 - 2 \exp \left(\frac{-c \left(\frac{m_\ell \sqrt{n_{\ell-1}} \delta_\ell}{\sigma \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}} - m_\ell \mu_{d_\ell} \right)^2}{C^2 m_\ell \frac{4}{2} \frac{d_\ell^{d_\ell/4}}{d_\ell/4 - 2}} \right).$$

In order to achieve (15), a sufficient criterion is

$$\frac{\kappa_\ell}{2} \geq \exp \left(\frac{-cm_\ell \left(\frac{\sqrt{n_{\ell-1}} \delta_\ell}{\sigma \prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}}} - \mu_{d_\ell} \right)^2}{C^2 \frac{4}{2} \frac{d_\ell^{d_\ell/4}}{d_\ell/4 - 2}} \right).$$

Solving for $n_{\ell-1}$ leads to

$$n_{\ell-1} \geq \frac{\sigma^2 \left(\prod_{i=\ell}^L a^{(i)} \prod_{i=\ell+1}^L \|W^{(i)}\|_{\text{op}} \right)^2}{\delta_\ell^2} \times \left(\sqrt{C^2 \frac{4}{2} \frac{d_\ell^{d_\ell/4}}{d_\ell/4 - 2}} (-\ln(\kappa_\ell/2)) \frac{1}{cm_\ell} + \mu_{d_\ell} \right)^2. \quad (18)$$

If we substitute the expression in (17) for μ_{d_ℓ} , (18) becomes the bound as seen in Theorem 1. \square

3.4. Conclusion

Within the context of the model described in Section 2, we have established that any feed-forward NN can be approximated arbitrarily well by ONNs constructed using Design A. This is Theorem 1 in essence.

This result has two consequences when it comes to the physical implementation of ONNs. On the one hand, it is guaranteed that the theoretical expressiveness of NNs can be retained in practice. On the other hand, Design A allows one to improve the accuracy of a noisy ONN to a desired level, and in fact bring the accuracy arbitrarily close to that of any state-of-the-art feed-forward noiseless NNs. Let us finally remark that the high bandwidth of photonic circuits may be of use when implementing Design A.

4. Results—Design B

4.1. Reducing noise in feed-forward linear ONNs (Design B)

Recall that an example of Design B is presented in Figure 2(c). Algorithm 2 constructs this network, given a desired number of copies m in each layer.

Calculating the output of a NN by using Design B first requires to fix a number m . The input data $x^{(0)}$ is then modulated m times, creating m noisy realizations of the input $(\alpha x^{(0)})_{\alpha=1, \dots, m}$. The weighted addition step and the activation function of each layer are singled out and copied m times. Both the copies of the weighted addition step and of the activation function of each layer are arrayed parallel to each other and performed on the m inputs, resulting in m outputs. The m parallel outputs of the weighted addition are

Algorithm 2 Algorithm to construct a noise reducing network

Require: Fix number $m \in \mathbb{N}$

Require: m copies of input $^1x^{(0)}, \dots, ^m x^{(0)}$

for $\ell = 1, \dots, L$ **do**

for $\alpha = 1, \dots, m$ **do**

$\alpha \xi^{(\ell)} \leftarrow W^{(\ell)} \alpha x^{(\ell-1)} + b^{(\ell)} + \text{Normal}(0, \Sigma_w)$

end for

$y^{(\ell)} \xleftarrow{\text{combining}} \alpha \xi^{(\ell)} + \dots + m \xi^{(\ell)}$

$(^1y^{(\ell)}, \dots, ^m y^{(\ell)}) \xleftarrow{\text{splitting}} m^{-1} y^{(\ell)}$

for $\alpha = 1, \dots, m$ **do**

$\alpha x^{(\ell)} \leftarrow \sigma^{(\ell)} (\alpha y^{(\ell)}) + \text{Normal}(0, \Sigma_{\text{act}})$

end for

end for

return $m^{-1} \sum_{\alpha=1}^m \alpha x^{(L)}$

merged to a single output, and afterwards split into m pieces. The m pieces are each sent to one of the m activation function mechanisms for processing. The resulting m activation values are the output of the layer. If it is the last layer, the m activation values are merged to produce the final output. These steps are formally described in Algorithm 2. A schematic representation of Design B can be seen in Figure 2(c).

4.2. Analysis of Design B

We now consider the physical and mathematical consequences of Design B.

Observe that in Design B, the m weighted additions of the ℓ -th layer's input $x^{(\ell-1)}$ result in realizations $(\alpha \xi^{(\ell)})_{\alpha=1, \dots, m}$ of $W^{(\ell)} x^{(\ell-1)} + b^{(\ell)} + \text{Normal}(0, \Sigma_w)$. These realizations are then combined resulting in

$$mW^{(\ell)} x^{(\ell-1)} + mb^{(\ell)} + \text{Normal}(0, m\Sigma_w) + \text{Normal}(0, \Sigma_{\text{sum}}).$$

Splitting the signal again into m parts, each signal carries information following the distribution

$$W^{(\ell)} x^{(\ell-1)} + b^{(\ell)} + \text{Normal}(0, m^{-1}\Sigma_w + m^{-2}\Sigma_{\text{sum}}) + \text{Normal}(0, \Sigma_{\text{spl}}).$$

The mean of the normal distribution therefore is the original networks pre-activation obtained from this input (that is without perturbations). The covariance matrix of the normal distribution is $m^{-1}\Sigma_w + m^{-2}\Sigma_{\text{sum}} + \Sigma_{\text{spl}}$. Each of those signals is fed through the mechanism applying the activation function, yielding m noisy versions of the output, distributed according to

$$\begin{aligned} & x^{(\ell)} | x^{(\ell-1)} \\ & \stackrel{(d)}{=} \sigma^{(\ell)} (W^{(\ell)} x^{(\ell-1)} + b^{(\ell)} \\ & \quad + \text{Normal}(0, m^{-1}\Sigma_w + m^{-2}\Sigma_{\text{sum}} + \Sigma_{\text{spl}})) \\ & \quad + \text{Normal}(0, \Sigma_{\text{act}}). \end{aligned}$$

The effect of Design B is thus that $T^{(\ell)}(\Sigma)$ in (8) is replaced by

$$T_m^{(\ell)}(\Sigma) := \frac{1}{m} D^{(\ell)} W^{(\ell)} \Sigma (D^{(\ell)} W^{(\ell)})^\top + \frac{1}{m} D^{(\ell)} \Sigma_w (D^{(\ell)})^\top$$

$$+ \frac{1}{m^2} D^{(\ell)} \Sigma_{\text{sum}} (D^{(\ell)})^\top + D^{(\ell)} \Sigma_{\text{spl}} (D^{(\ell)})^\top + \Sigma_a;$$

see Appendix A.2.2. Observe also that $\Sigma_a^{(\ell)}$ can be written as $(1/m)m\Sigma_a^{(\ell)}$. Therefore, if we substitute the matrix $\bar{\Sigma}_a^{(\ell)} = m\Sigma_a^{(\ell)}$ for $\Sigma_a^{(\ell)}$ in $T^{(\ell)}(\Sigma)$, we can write

$$T_m^{(\ell)}(\Sigma) = m^{-1}T^{(\ell)}(\Sigma) + m^{-2}D^{(\ell)}\Sigma_{\text{sum}}(D^{(\ell)})^\top + D^{(\ell)}\Sigma_{\text{spl}}(D^{(\ell)})^\top.$$

We have the following analogs to Proposition 1 and Corollary 1:

Theorem 2 (Distribution of Design B) *Assume that there exist vectors $a^{(\ell)} \in \mathbb{R}^{d_\ell}$ such that $\sigma^{(\ell)}(y) = \text{diag}(a^{(\ell)})y$. The feed-forward linear ONN constructed using Design B with m copies then satisfies*

$$\Psi_m^{\text{ONN}}(\cdot, w) \stackrel{(d)}{=} \text{Normal}(\Psi^{\text{NN}}(\cdot, w), \Sigma_{\text{ONN},m}^{(L)}),$$

where for $\ell = L, L-1, \dots, 1$,

$$\Sigma_{\text{ONN},m}^{(\ell)} = T_m^{(\ell)}(\Sigma_{\text{ONN},m}^{(\ell-1)}); \quad \text{and} \quad \Sigma_{\text{ONN},m}^{(0)} = \Sigma_m.$$

Under the assumption of symmetric noise, a similar simplification of the recursion in Theorem 2, similar to that in Proposition 1, is possible. Assume $\Sigma_{\text{sum}} = \Sigma_{\text{spl}} = 0$. Introduce again $P^{(\ell)} := \prod_{i=\ell+1}^L D^{(i)}W^{(i)}$ for notational convenience. The following is proved in Appendix A.2:

Corollary 2 (Symmetric noise case) *Assume that for all $\ell \in \mathbb{N}_+$, $\Sigma_a^{(\ell)} = \Sigma_a$ and $\Sigma_w^{(\ell)} = \Sigma_w$. Then,*

$$\begin{aligned} \Sigma_{\text{ONN},m}^{(L)} &= \sum_{\ell=1}^L (m^{-1})^{L-\ell} P^{(\ell)} m \Sigma_a (P^{(\ell)})^\top \\ &+ \sum_{\ell=1}^L (m^{-1})^{L-\ell} P^{(\ell)} D^{(\ell)} \Sigma_w (D^{(\ell)})^\top (P^{(\ell)})^\top \\ &+ (m^{-L}) P^{(0)} \Sigma_m (P^{(0)})^\top. \end{aligned}$$

We will next consider the limit of the covariance matrix in a large, symmetric linear ONN with Design B, that we can grow infinitely deep. Algorithm 2 is namely able to guarantee boundedness of the covariance matrix in such deep ONN if the parameter m is chosen appropriately:

Corollary 3 *Consider a linear ONN with Design B and parameter m , that has L layers, and that satisfies the following symmetry properties: for all $\ell \in \{1, \dots, L\}$, $W^{(\ell)} = W$, $D^{(\ell)} = D$, $\Sigma_a^{(\ell)} = \Sigma_a$ and $\Sigma_w^{(\ell)} = \Sigma_w$. Then, if $\|D\|_F \|W\|_F < \sqrt{m}$, the limit $\lim_{L \rightarrow \infty} \Sigma_{\text{ONN}}^{(L)}$ exists.*

Moreover,

$$\begin{aligned} &\lim_{L \rightarrow \infty} \Sigma_{\text{ONN},m}^{(L)} \\ &= \sum_{n=0}^{\infty} m^{-(n+1)} (DW)^n (D\Sigma_w D^\top + m\Sigma_a) ((DW)^n)^\top. \end{aligned}$$

Notice that the bound on the number of copies needed for the covariance matrix of an ONN to converge to a limit is independent of e.g. the Frobenius norms of the covariance matrices that describe the noise distributions. This is because, here, we are not interested in bounding the covariance matrix to a specific level; instead, we are merely interested in the existence of a limit.

4.3. Discussion & Conclusion

Compared to Theorem 2's recursive description of the covariance matrix in any linear ONN with Design B, Corollary 2 provides a series that describes the covariance matrix in any linear, symmetric ONN with Design B. While the result holds more restrictively, it is more insightful. For example, it allows us to consider the limit of the covariance matrix in an extremely deep ONNs (see Corollary 3). Corollary 3 suggests that in deep ONNs with Design B, one should choose $m \approx \lceil (\|D\|_F \|W\|_F)^2 \rceil$ in order to control the noise and not be too inefficient with the number of copies.

These results essentially mean that in a physical implementation of an increasingly deep and linear ONN, the covariance matrix can be reduced (and thus remain bounded) by applying Design B with multiple copies. The quality of the ONN's output increases as the number of copies in Design B (or Design A for that matter) is increased. Finally, it is worth mentioning that Design B could potentially be implemented such that it leverages the enormous bandwidth of optics.

5. Simulations

We investigate the improvements of output quality achieved by Designs A and B on a benchmark example: the convolutional neural network LeNet [40]. As measure for quality we consider the Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \left(\Psi^{\text{ONN}}(x^{(i)}, w) - \Psi^{\text{NN}}(x^{(i)}, w) \right)^2$$

and the prediction accuracy

$$\frac{\#\{\text{correctly classified images}\}}{\#\{\text{images}\}}.$$

5.1. Empirical variance

We extracted plausible values for Σ_w and Σ_a from the ONN implementation [41] of a 2-layer NN for classification of the Modified National Institute of Standards and Technology (MNIST) database [42]. In [41], the authors trained the NN on a classical computer and implemented the trained weights afterwards in an ONN. They then tuned noise (with the same noise model as in Section 2 of this paper) into the noiseless computer model, assuming that $\Sigma_w = \text{diag}(\sigma_w^2)$ and $\Sigma_a = \text{diag}(\sigma_a^2)$. They found $\sigma_w \in [0.08, 0.1] \cdot d$ and $\sigma_a \in [0.1, 0.15] \cdot d$ to reach the same accuracy levels as the ONN, where d denotes the diameter of the range.

5.2. LeNet ONN: Performance when implemented via Design A and B

Convolutional NNs can be regarded as feedforward NNs by stacking the (2D or 3D) images into column vectors and arranging the filters to a weight matrix. Thus Design A and B are well-defined for Convolutional NNs. We apply the designs to LeNet5 [40], which is trained for classifying the handwritten digits in the MNIST dataset [42]. The layers are:

1. 2D convolutional layer with kernel size 5, stride 1 and 2-padding. Output has 6 channels of 28x28 pixel representations, with the activation function being tanh;
2. average pooling layer, pooling 2x2 block, the output therefore is 14x14;
3. 2D convolutional layer with kernel size 5, stride 1 and no padding. The output has 16 channels of 10x10 pixel representations and the activation function is tanh;
4. average pooling layer, pooling 2x2 block, the output therefore is 5x5;
5. 2D convolutional layer with kernel size 5, stride 1 and no padding. The output has 120 channels of 1 pixel representations and the activation function used is tanh;
- (5.) flattening layer, which turns the 120 one-dimensional channels into one 120-dimensional vector;
6. dense layer with 84 neurons and tanh activation function;
7. dense layer with 10 neurons and softmax activation function.

Figures 3 and 4 show the MSE and the prediction accuracy of Design A and B for an increasing number of copies, respectively.

For simplicity we set all individual copies n_i per layer i in Design A to equal m , that is $n_i = m$ for all i . The total number of copies that Design A starts with then is m^L . Here L is equal to 7. In Design B the number of copies is m per layer and the total number of copies is mL . In the case of one copy the designs A and B are identical to the original network, while we focus on the effect once the designs deviate from the original network ($m \geq 2$).

The axis in Figures 3 and 4 denote the number of copies per layer. Here, we scale the copies per layer for Design A linearly, because the total amount of copies for Design A grows exponentially and we scale the copies per layer for Design B exponentially, because the total number of copies for Design B grows linearly. This way the comparison is on equal terms.

Figure 3 displays the MSE seen for LeNet, depending on the amount of copies for each design. In the trade-off between additional resources needed for the additional copies against the diminishing benefits of adding further copies, we see that, for both measures MSE (Figure 3) and relative accuracy (Figure 4), already 2 to 5 copies per layer yield good results. The relative accuracy in Figure 4 is scaled such

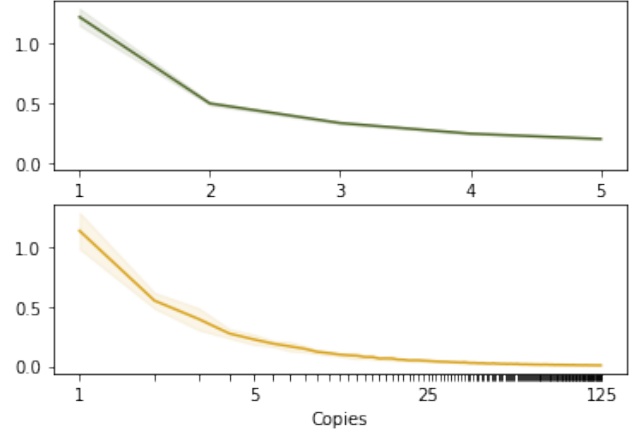


Figure 3: $MSE(\cdot 10^2)$ for Design A (top) and Design B (bottom) as function of copies on LeNet5 trained for MNIST classification. The pale area contains the 95%-confidence intervals.

that 0 corresponds to the accuracy of the original NN with noise profile (i.e., the ONN without modifications, we call this the original ONN) and 1 to the accuracy of the original NN without noise. The designs do not alter the fundamental operation of the original NN, therefore there should be no performance gain and the original NN's accuracy should be considered the highest achievable, thus constituting the upper bound in relative accuracy of 1. Likewise the lowest accuracy should be given by the original ONN, as there is no noise reduction involved.

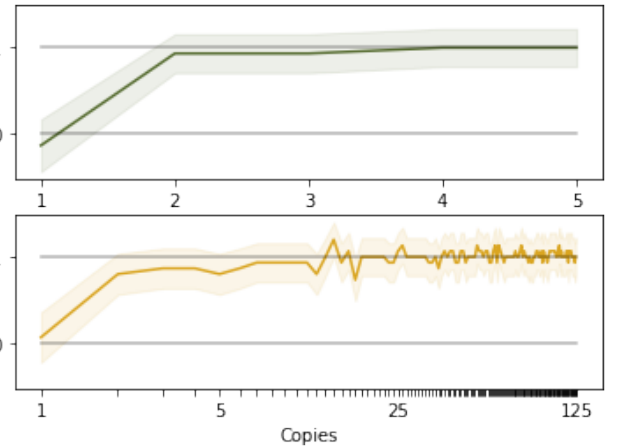


Figure 4: Relative accuracy for Design A (top) and Design B (bottom) as function of copies on LeNet5 trained for MNIST classification. The pale area contains the 56.5%-confidence intervals.

5.3. Effect of additional layers in LeNet

In order to investigate how the depth affects the noise at the output, while keeping the operation of the network the same to ensure the results are commensurable, we insert additional layers with identity matrix and identity activation function (we will call them identity layers) into a network.

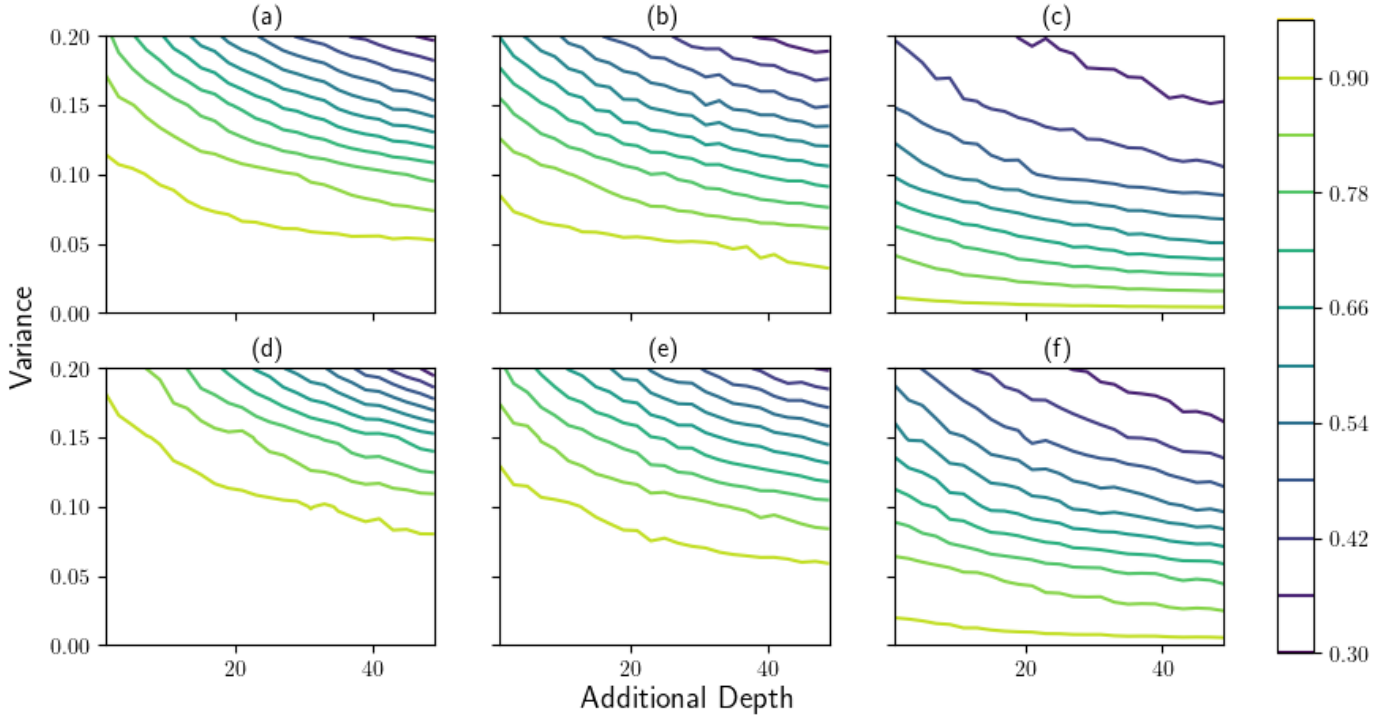


Figure 5: Accuracy of LeNet ONNs, depending on the amount of inserted identity layers and the variance level of the ONN, for (a) a network with tanh activation function and one copy, (b) a network with ReLU activation function and one copy, (c) a network with linear activation function and one copy, (d) a network with tanh activation function and two copies, (e) a network with ReLU activation function and two copies, (f) a network with linear activation function and two copies.

Specifically, we take networks with the LeNet architecture as in Section 5.2, using different activation functions, while fixing the output layer to be softmax. We then insert identity layers between layers 1 and 2, 3 and 4, 5 and 6, as well as between layers 6 and 7. For a fixed total of additional layers, the layers are inserted in the four spots between layers 1&2, 3&4, 5&6, and 7&8 according to the tuple

$$n \mapsto \left(\left\lfloor \frac{n+3}{4} \right\rfloor, \left\lfloor \frac{n+2}{4} \right\rfloor, \left\lfloor \frac{n+1}{4} \right\rfloor, \left\lfloor \frac{n}{4} \right\rfloor \right).$$

The insertion pattern is illustrated in Table 1:

# of additional layers	1&2	3&4	5&6	7&8
1	1	0	0	0
2	1	1	0	0
3	1	1	1	0
4	1	1	1	1
5	2	1	1	1
6	2	2	1	1
...

Table 1: Insertion pattern.

Finally, we tune the variance terms of the covariance matrix in our noise model. The results are displayed in Figure 5.

In Figure 5, we observe that the tanh and the ReLU networks perform as expected. Additional noisy layers decrease the accuracy and thus the same level of performance can only be achieved if the variance is lower. This trend can also be seen in the linear network, but to a lesser extent.

5.4. Simulations on effective values for Design B

According to Corollary 3, the covariance matrix of a linear ONN constructed by Design B is bounded if $m > (\|D\|_F \|W\|_F)^2$, and therefore $m = \lceil (\|D\|_F \|W\|_F)^2 \rceil$, is sufficient to ensure that the covariance matrix of the output distribution $\Psi_m^{\text{ONN}}(\cdot, w) \stackrel{(d)}{=} \text{Normal}(\Psi^{\text{NN}}(\cdot, w), \Sigma_{\text{ONN}, m}^{(L)})$ in Theorem 2 is bounded in linear NNs. This is derived by using submultiplicativity of the norm (see (A.8)) and is therefore possibly a loose bound. We use the exact relation given by Corollary 2 for the covariance matrix in Theorem 2 to investigate the lowest values for m for which the covariance matrix starts being bounded. In Figure 6 we depict a linear NN with constant width 4. We vary the values for $\|D\|_F$ and $\|W\|_F$. Upon close inspection we see that the lowest value for m seems to be $g(x, y) \simeq \lceil (xy)^2 / \|I_d\|_F^4 \rceil$, where I_d is the identity matrix of dimension d , see Figure 6. Because $\|I_d\|_F = \sqrt{d}$, the value for m found numerically is $m \simeq \lceil (\|D\|_F \|W\|_F / d)^2 \rceil$.

6. Discussion & Conclusion

Design A, introduced in Section 3, guarantees an approximation property (Theorem 1). This is achieved through technical machinery to control the noise, even though there are nonlinear activation functions involved. This method is powerful enough to yield the universal approximation property, as NNs can be approximated arbitrarily well with ONNs that

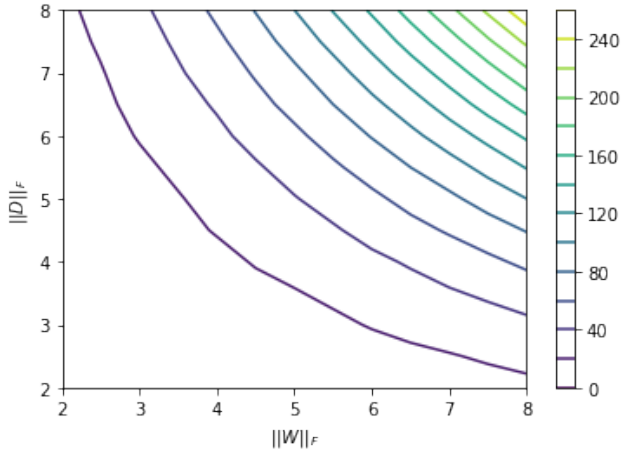



Figure 6: The contour lines denote the lowest m for which $\|\Sigma_{\text{ONN},m}^{(L)}\|_F$ stops growing exponentially, as a function of $\|W\|_F$ and $\|D\|_F$.

are constructed through the first design, and NNs themselves can approximate any continuous function arbitrarily well [43, Theorem 1]. Our mathematical guarantee however, only states a sufficient number of copies required, and this number grows exponentially as the number of layers increases.

We then introduced Design B in Section 4, in which the growth of number of copies is much more benign. However, the analysis of Design B was restricted to linear NNs, and Design B might therefore not be expressive enough to have the universal approximation property. Linear NNs, or NNs with algebraic polynomials as activation functions for that matter, namely do not possess the universal approximation property. The assumption of linear activation functions did allow us to characterize the distribution of the output exactly on the flipside (Theorem 2).

In short, in this paper, we have discussed the noise present in ONNs and described a mathematical model for the noise. We also investigate the numerical implications of the mathematical model, with a specific focus on the effects of depth (Figure 5). The proposed noise reduction schemes yield greater accuracy and the theoretical results (Theorem 1 and Corollary 3) guarantee that ONNs work just as noiseless NNs in the many copies limit. With the designs and findings of Sections 3 to 4 we have a framework to exploit known NN wisdom, as no new training is required. Further research should address optimization algorithms that take the noise of ONNs into account to investigate the regularization, generalization and minimization properties of trained ONNs.

Acknowledgments

This research was supported by the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement no. 945045, and by the NWO Gravitation project NETWORKS under grant no. 024.002.003. 

We would finally like to thank Bin Shi for advise on the noise level parameters of ONNs for our simulations. Furthermore we want to thank Albert Senen-Cerda, Sanne van Kempen, and Alexander Van Werde for feedback to a draft version of this document.

References

- [1] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.
- [2] J. Long, E. Shelhamer, T. Darrell, Fully Convolutional Networks for Semantic Segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015.
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al., Human-level control through deep reinforcement learning, *Nature* (2015).
- [4] H. Nam, B. Han, Learning Multi-Domain Convolutional Neural Networks for Visual Tracking, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.
- [5] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, et al., Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation, arXiv preprint arXiv:1609.08144 (2016).
- [6] J. Hasler, H. Marr, Finding a roadmap to achieve large neuromorphic hardware systems, *Frontiers in Neuroscience* (2013).
- [7] Z. Du, D. Rubin, Y. Chen, L. Hel, T. Chen, L. Zhang, C. Wu, O. Temam, Neuromorphic Accelerators: A Comparison Between Neuroscience and Machine-Learning Approaches, in: MICRO-48: Proceedings of the 48th International Symposium on Microarchitecture, 2015.
- [8] F. Akopyan, J. Sawada, A. Cassidy, R. Alvarez-Icaza, J. Arthur, P. Merolla, N. Imam, Y. Nakamura, P. Datta, G.-J. Nam, B. Taba, M. Beakes, B. Brezzo, J. B. Kuang, R. Manohar, W. P. Risk, B. Jackson, D. S. Modha, TrueNorth: Design and Tool Flow of a 65 mW 1 Million Neuron Programmable Neurosynaptic Chip, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (2015).
- [9] B. V. Benjamin, P. Gao, E. McQuinn, S. Choudhary, A. R. Chandrasekaran, J.-M. Bussat, R. Alvarez Icaza, J. V. Arthur, P. A. Merolla, K. Boahen, Neurogrid: A Mixed-Analog-Digital Multichip System for Large-Scale Neural Simulations, *Proceedings of the IEEE* (2014).
- [10] S. Furber, F. Galluppi, S. Temple, L. Plana, The SpiNNaker Project, *Proceedings of the IEEE* (2014).
- [11] P. A. Merolla, J. V. Arthur, R. Alvarez-Icaza, A. S. Cassidy, J. Sawada, F. Akopyan, B. L. Jackson, N. Imam, C. Guo, Y. Nakamura, et al., A million spiking-neuron integrated circuit with a scalable communication network and interface, *Science* (2014).
- [12] J. Schemmel, D. Brüderle, A. Grübl, M. Hock, K. Meier, S. Millner, A wafer-scale neuromorphic hardware system for large-scale neural modeling, in: 2010 IEEE International Symposium on Circuits and Systems (ISCAS), 2010.
- [13] S. A. Siddiqui, S. Dutta, A. Tang, L. Liu, C. A. Ross, M. A. Baldo, Magnetic Domain Wall Based Synaptic and Activation Function Generator for Neuromorphic Accelerators, *Nano Letters* (2019).
- [14] T. De Lima, B. Shastri, A. Tait, M. Nahmias, P. Prucnal, Progress in neuromorphic photonics, *Nanophotonics* (2017).
- [15] K. Kitayama, M. Notomi, M. Naruse, K. Inoue, S. Kawakami, A. Uchida, Novel frontier of photonics for data processing—Photonic accelerator, *APL Photonics* (2019).
- [16] M. Miscuglio, J. Meng, O. Yesiliurt, Y. Zhang, L. J. Prokopena, A. Mehrabian, J. Hu, A. V. Kildishev, V. J. Sorger, Artificial Synapse with Mnemonic Functionality using GSST-based Photonic Integrated Memory, in: 2020 International Applied Computational Electromagnetics Society Symposium (ACES), 2020.
- [17] B. J. Shastri, A. N. Tait, T. F. de Lima, M. A. Nahmias, H.-T. Peng, P. R. Prucnal, Principles of Neuromorphic Photonics, arXiv preprint arXiv:1801.00016 (2017).

- [18] Y. Shen, N. C. Harris, S. Skirlo, M. Prabhu, T. Baehr-Jones, M. Hochberg, X. Sun, S. Zhao, H. Larochelle, D. Englund, M. Soljačić, Deep Learning with Coherent Nanophotonic Circuits, *Nature Photonics* (2017).
- [19] R. Hamerly, L. Bernstein, A. Sludds, M. Soljačić, D. Englund, Large-Scale Optical Neural Networks Based on Photoelectric Multiplication, *Physical Review X* (2019).
- [20] L. Bernstein, A. Sludds, R. Hamerly, V. Sze, J. Emer, D. Englund, Freely scalable and reconfigurable optical hardware for deep learning, *scientific reports* (2021).
- [21] C. Huang, S. Fujisawa, T. F. De Lima, A. N. Tait, E. Blow, Y. Tian, S. Bilodeau, A. Jha, F. Yaman, H. G. Batshon, et al., Demonstration of photonic neural network for fiber nonlinearity compensation in long-haul transmission systems, in: 2020 Optical Fiber Communications (OFC) Conference and Exhibition, 2020.
- [22] B. Shi, N. Calabretta, R. Stabile, Deep Neural Network through an InP SOA-Based Photonic Integrated Cross-Connect, *IEEE Journal of Selected Topics in Quantum Electronics* (2019).
- [23] B. Shi, K. Prifti, E. Magalhães, N. Calabretta, R. Stabile, Lossless Monolithically Integrated Photonic InP Neuron for All-Optical Computation, in: *Optical Fiber Communication Conference*, 2020.
- [24] B. Shi, N. Calabretta, R. Stabile, First Demonstration of a Two-Layer All-Optical Neural Network by Using Photonic Integrated Chips and SOAs, in: 45th European Conference on Optical Communication (ECOC 2019), 2019.
- [25] B. J. Shastri, A. N. Tait, T. Ferreira de Lima, W. H. Pernice, H. Bhaskaran, C. Wright, P. R. Prucnal, Photonics for artificial intelligence and neuromorphic computing, *Nature Photonics* (2021).
- [26] A. N. Tait, T. F. De Lima, E. Zhou, A. X. Wu, M. A. Nahmias, B. J. Shastri, P. R. Prucnal, Neuromorphic photonic networks using silicon photonic weight banks, *scientific reports* (2017).
- [27] R.-J. Essiambre, G. Kramer, P. Winzer, G. Foschini, B. Goebel, Capacity Limits of Optical Fiber Networks, *Journal of Lightwave Technology* (2010).
- [28] X. Li, R. Mardling, J. Armstrong, Channel Capacity of IM/DD Optical Communication Systems and of ACO-OFDM, in: 2007 IEEE International Conference on Communications, 2007.
- [29] T. de Lima, A. Tait, H. Saeidi, M. Nahmias, H. Peng, S. Abbaslou, B. Shastri, P. Prucnal, Noise Analysis of Photonic Modulator Neurons, *IEEE Journal of Selected Topics in Quantum Electronics* (2019).
- [30] I. Chakraborty, G. Saha, A. Sengupta, K. Roy, Toward Fast Neural Computing using All-Photonic Phase Change Spiking Neurons, *Scientific Reports* (2018).
- [31] N. Passalis, M. Kirtas, G. Mourgias-Alexandris, G. Dabos, N. Pleros, A. Tefas, Training Noise-Resilient Recurrent Photonic Networks for Financial Time Series Analysis, in: 2020 28th European Signal Processing Conference (EUSIPCO), 2021.
- [32] G. Mourgias-Alexandris, A. Tsakyridis, N. Passalis, A. Tefas, K. Vyrsokinos, N. Pleros, An all-optical neuron with sigmoid activation function, *Optics Express* (2019).
- [33] O. A. Manita, M. A. Peletier, J. W. Portegies, J. Sanders, A. Senen-Cerda, Universal approximation in dropout neural networks, *Journal of Machine Learning Research* (2022).
- [34] R. J. Muirhead, *Aspects of Multivariate Statistical Theory*, 2009.
- [35] N. Semenova, L. Larger, D. Brunner, Understanding and mitigating noise in trained deep neural networks, *Neural Networks* (2022).
- [36] N. Semenova, X. Porte, L. Andreoli, M. Jacquot, L. Larger, D. Brunner, Fundamental aspects of noise in analog-hardware neural networks, *Chaos: An Interdisciplinary Journal of Nonlinear Science* (2019).
- [37] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*, Cambridge University Press, 2018.
- [38] M. L. Abell, J. P. Braselton, J. A. Rafter, *Statistics with mathematics*, Academic Press, 1999.
- [39] C. Clapham, J. Nicholson, J. R. Nicholson, *The concise Oxford dictionary of mathematics*, Oxford University Press, 2014.
- [40] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-Based Learning Applied to Document Recognition, *Proceedings of the IEEE* (1998).
- [41] B. Shi, N. Calabretta, R. Stabile, InP photonic integrated multi-layer neural networks: Architecture and performance analysis, *APL Photonics* (2022).
- [42] L. Deng, The mnist database of handwritten digit images for machine learning research, *IEEE Signal Processing Magazine* (2012).
- [43] M. Leshno, V. Y. Lin, A. Pinkus, S. Schocken, Multilayer Feed-forward Networks with Non-Polynomial Activation Function Can Approximate Any Function, *Neural Networks* (1993).
- [44] S. Banach, Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales, *Fundamenta mathematicae* (1922).

A.

A.1. Proofs of Section 2

A.1.1. Proof of Corollary 1

The first conclusion of Corollary 1 follows immediately from expanding the recursion.

Assume that for all $\ell \in \mathbb{N}$ $\text{diag}(a^{(\ell)}) = a$ and $W^{(\ell)} = W$ as well as $\Sigma_a^{(\ell)} = \Sigma_a$ and $\Sigma_w^{(\ell)} = \Sigma_w$. Let $T := T^{(\ell)}$ be the common map under those conditions.

To prove the second conclusion of Corollary 1, observe that for any two matrices X and Y of the same dimension as a and W ,

$$\begin{aligned} \|T(X) - T(Y)\|_F &= \|aW(X - Y)(aW)^\top\|_F \\ &\leq \|a\|_F^2 \|W\|_F^2 \|X - Y\|_F \end{aligned} \quad (\text{A.1})$$

by submultiplicativity of the Frobenius norm. Let us now consider the setting of Proposition 1 for a moment, that is, we initialize $X_1 = \Sigma_m$ and calculate $X_{\ell+1} = T(X_\ell)$ recursively. The sequence X_1, X_2, \dots converges if $\|a\|_F \|W\|_F < 1$, as a consequence of the Banach fixed point theorem [44] combined with (A.1).

We may therefore consider the unique fixed point $X_* := \lim_{\ell \rightarrow \infty} X_\ell$. It must satisfy the fixed point equation $T(X_*) = X_*$, which reads $X_* - (aW)X_*(aW)^\top = a\Sigma_w a^\top + \Sigma_a$. Equivalently,

$$\begin{aligned} \text{vec}(X_*) - ((aW) \otimes (aW))\text{vec}(X_*) \\ = \text{vec}(a\Sigma_w a^\top) + \text{vec}(\Sigma_a). \end{aligned}$$

Here, \otimes denotes the Kronecker product and vec the vectorization of a matrix (effectively, we stack the columns of the matrix A on top of one another). This vectorization trick allows us to write the solution to the fixed point equation as

$$\text{vec}(X_*) = (I - ((aW) \otimes (aW)))^{-1} (\text{vec}(a\Sigma_w a^\top) + \text{vec}(\Sigma_a)). \quad (\text{A.2})$$

Here, $(aW) \otimes (aW)$ denotes $(aW) \otimes (aW)$.

Formally we rewrite the inverse in (A.2) in terms of a von Neumann series

$$(I - ((aW) \otimes (aW)))^{-1} = \sum_{n=0}^{\infty} ((aW) \otimes (aW))^n. \quad (\text{A.3})$$

This is however justified only if

$$\|((aW) \otimes (aW))^n\|_F \rightarrow 0, \quad (\text{A.4})$$

which we verify next.

For the Kronecker product it holds that $\text{tr}(A \otimes B) = \text{tr}(A)\text{tr}(B)$. Therefore, $\|(aW) \otimes (aW)\|_F = \text{tr}((aW)^\top (aW)) = \|aW\|_F^2$ by definition of the Frobenius norm. Furthermore, by submultiplicativity, $\|((aW) \otimes (aW))^n\|_F \leq \|aW\|_F^{2n}$. Thus, by the assumption that $\|a\|_F \|W\|_F < 1$, condition (A.4) holds and (A.3)'s expression is proper. This leads to the representation of X_* as

$$\text{vec}(X_*) = \sum_{n=0}^{\infty} ((aW) \otimes (aW))^n \text{vec}(a\Sigma_w a^\top + \Sigma_a)$$

$$= \text{vec}(a\Sigma_w a^\top + \Sigma_a) + \sum_{n=1}^{\infty} ((aW) \otimes (aW))^n \text{vec}(a\Sigma_w a^\top + \Sigma_a).$$

Returning to matrix notation we have

$$X_* = a\Sigma_w a^\top + \Sigma_a + \sum_{n=1}^{\infty} (aW)^n (a\Sigma_w a^\top + \Sigma_a) ((aW)^n)^\top.$$

This proves the second conclusion of Corollary 1.

A.2. Additional material Section 4

A.2.1. Additional considerations on Design A

To consider the total number of copies in Design A to guarantee (9), we need to multiply all the n_ℓ in Theorem 1. To simplify the terms we upper bound

$$\sqrt{2} \frac{\Gamma((d_{\ell+1} + 1)/2)}{\Gamma(d_{\ell+1}/2)} \quad \ell \in \{0, \dots, L-1\}$$

by a constant D (assuming the sequence of d_ℓ is bounded). We also replace

$$\sqrt{\frac{C^2}{c} \frac{4^{d_{\ell+1} + \sqrt{4}}}{2^{d_{\ell+1} + \sqrt{4} - 2}} (-\ln(\kappa_{\ell+1}/2)) \frac{1}{\prod_{i=\ell+1}^L n_i}}$$

by a constant E . If the total number of copies satisfies

$$\begin{aligned} \prod_{\ell=0}^L n_\ell &\geq \frac{\left(\prod_{\ell=1}^L (a^{(\ell)})^{2\ell}\right) \left(\prod_{\ell=1}^L (\|W^{(\ell)}\|_{\text{op}})^{2(\ell-1)}\right)}{\prod_{\ell=1}^L \delta_\ell^2} \\ &\quad \times \sigma^{2L} (D + E)^{2L}, \end{aligned}$$

then we are able to construct an ONN $\tilde{\Psi}^{\text{ONN}}$ that satisfies (9). The product $\prod_{\ell=1}^L \delta_\ell^2$ is maximized if all $\delta_\ell = \mathcal{D}_L/L$. We furthermore upper-bound $\prod_{\ell=1}^L (a^{(\ell)})^{2\ell}$ by a^{2L^2} and $\prod_{\ell=1}^L (\|W^{(\ell)}\|_{\text{op}})^{2(\ell-1)}$ by W^{2L^2} . We then have

$$N = \prod_{\ell=0}^L n_\ell = \omega(K^{2L+2L^2} L^L) = \omega(K^{2L(L+1)} L^L).$$

A.2.2. Deriving the covariance matrix for Design B

We now derive $T_m^{(\ell)}(\Sigma)$ —the transformation of the covariance matrix which an input undergoes as it becomes the output of layer ℓ . Recall that this input is distributed as $\text{Normal}(x^{(\ell-1)}, \Sigma)$. Denote the random variable for the pre-activation (from which the realizations are drawn) after joining and splitting beams by $\mathbf{P}^{(\ell)}$. Then

$$\begin{aligned} \mathbf{P}^{(\ell)} &| x^{(\ell-1)} \\ &\stackrel{(d)}{=} m^{-1} \left[\left(\sum_{i=1}^m W^{(\ell)} \text{Normal}(x^{(\ell-1)}, \Sigma) + \text{Normal}(0, \Sigma_w) \right) \right. \\ &\quad \left. + \text{Normal}(0, \Sigma_{\text{sum}}) \right] + \text{Normal}(0, \Sigma_{\text{spl}}) \\ &\stackrel{(d)}{=} m^{-1} \left(\text{Normal}(mW^{(\ell)} x^{(\ell-1)}, mW^{(\ell)} \Sigma (W^{(\ell)})^\top) \right) \end{aligned}$$

$$\begin{aligned}
& + \text{Normal}(0, m\Sigma_w) + \text{Normal}(0, \Sigma_{\text{sum}}) \Big) + \text{Normal}(0, \Sigma_{\text{spl}}) \\
& \stackrel{(d)}{=} \text{Normal}\left(W^{(\ell)}x^{(\ell-1)}, \right. \\
& \quad \left. m^{-1}(W^{(\ell)}\Sigma(W^{(\ell)})^\top + \Sigma_w + m^{-1}\Sigma_{\text{sum}}) + \Sigma_{\text{spl}}\right).
\end{aligned}$$

The random variable $\mathbf{P}^{(\ell)}$ is then channeled through the activation function, which subsequently adds another noise term. The resulting activation is the random variable $\mathbf{A}^{(\ell)}$

$$\begin{aligned}
& \mathbf{A}^{(\ell)} \mid x^{(\ell-1)} \\
& \stackrel{(d)}{=} \sigma^{(\ell)}(\mathbf{P}^{(\ell)}) + \text{Normal}(0, \Sigma_a) \\
& \stackrel{(d)}{=} D^{(\ell)}W^{(\ell)}x^{(\ell-1)} \\
& + \text{Normal}\left(0, \frac{D^{(\ell)}W^{(\ell)}\Sigma(D^{(\ell)}W^{(\ell)})^\top}{m} + \frac{D^{(\ell)}\Sigma_w(D^{(\ell)})^\top}{m} \right. \\
& \quad \left. + \frac{D^{(\ell)}\Sigma_{\text{sum}}(D^{(\ell)})^\top}{m^2} + D^{(\ell)}\Sigma_{\text{spl}}(D^{(\ell)})^\top + \Sigma_a\right) \\
& \stackrel{(d)}{=} \text{Normal}(\Phi_\ell^{\text{NN}}(x^{(\ell-1)}), T_m^{(\ell)}(\Sigma))
\end{aligned}$$

As we can see, instead of

$$T^{(\ell)}(\Sigma) := D^{(\ell)}W^{(\ell)}\Sigma(D^{(\ell)}W^{(\ell)})^\top + D^{(\ell)}\Sigma_w^{(\ell)}(D^{(\ell)})^\top + \Sigma_a^{(\ell)}$$

we have

$$\begin{aligned}
T_m^{(\ell)}(\Sigma) & := m^{-1}D^{(\ell)}W^{(\ell)}\Sigma(D^{(\ell)}W^{(\ell)})^\top \\
& + m^{-1}D^{(\ell)}\Sigma_w(D^{(\ell)})^\top + m^{-2}D^{(\ell)}\Sigma_{\text{sum}}(D^{(\ell)})^\top \\
& + D^{(\ell)}\Sigma_{\text{spl}}(D^{(\ell)})^\top + \Sigma_a.
\end{aligned}$$

A.2.3. Proof of Corollary 2

As Corollary 2 is the Design B analog of Corollary 1, the proofs are similar. The first expression in Corollary 2 is again immediate from expansion. For the limit we use the same Banach fixpoint argument, where only the variables have to be exchanged. The following executes these steps.

Assume again that for all $\ell \in \mathbb{N}$ $D^{(\ell)} = D$ and $W^{(\ell)} = W$ as well as $\Sigma_a^{(\ell)} = \Sigma_a$ and $\Sigma_w^{(\ell)} = \Sigma_w$. Let $T_m := T_m^{(\ell)}$ be the common map under those conditions.

Recall (A.1). In the setting of Theorem 2, that is $X_1 = \Sigma_m$ and $X_{\ell+1} = T_m(X_\ell)$, the so-defined sequence converges if $\|D\|_F \|W\|_F < \sqrt{m}$ (see also below (A.8)) due to (A.1) and the Banach fixed point theorem [44]. We therefore let the unique fixed point be

$$\lim_{\ell \rightarrow \infty} X_\ell = X_\star.$$

We can write the fixed point equation $T_m(X_\star) = X_\star$ as

$$\begin{aligned}
& X_\star - m^{-1}(DW)X_\star(DW)^\top \\
& = m^{-1}D\Sigma_w D^\top + m^{-1}\Sigma_a + m^{-2}D\Sigma_{\text{sum}}D^\top + D\Sigma_{\text{spl}}D^\top,
\end{aligned}$$

and further write it as

$$\text{vec}(X_\star) - m^{-1}((DW) \otimes (DW))\text{vec}(X_\star)$$

$$\begin{aligned}
& = \text{vec}(m^{-1}D\Sigma_w D^\top + m^{-1}\Sigma_a + m^{-2}D\Sigma_{\text{sum}}D^\top \\
& \quad + D\Sigma_{\text{spl}}D^\top).
\end{aligned}$$

Here, \otimes denotes again the Kronecker product and vec the vectorization of a matrix. Applying the vectorization trick as in the proof of Corollary 1 allows us to write the solution to the fixed point equation as

$$\begin{aligned}
& \text{vec}(X_\star) \\
& = (I - (m^{-1}(DW)^{\otimes 2}))^{-1} \\
& \quad \text{vec}(m^{-1}D\Sigma_w D^\top + m^{-1}\Sigma_a + m^{-2}D\Sigma_{\text{sum}}D^\top \\
& \quad + D\Sigma_{\text{spl}}D^\top).
\end{aligned} \tag{A.5}$$

Again, $(DW)^{\otimes 2}$ denotes $(DW) \otimes (DW)$.

Formally we rewrite the inverse in (A.5) in terms of a von Neumann series

$$(I - (m^{-1}(DW)^{\otimes 2}))^{-1} = \sum_{n=0}^{\infty} (m^{-1}(DW)^{\otimes 2})^n. \tag{A.6}$$

This is again only justified if

$$\|(m^{-1}(DW)^{\otimes 2})^n\|_F \rightarrow 0. \tag{A.7}$$

By submultiplicativity it holds that

$$\|(m^{-1}(DW)^{\otimes 2})^n\|_F \leq \|m^{-1}(DW)^{\otimes 2}\|_F^n. \tag{A.8}$$

For the Kronecker product it holds that $\text{tr}(A \otimes B) = \text{tr}(A)\text{tr}(B)$ and thus $\|(DW)^{\otimes 2}\|_F = \text{tr}((DW)^*(DW)) = \|DW\|_F^2$ by definition of the Frobenius norm. Therefore, by our assumption of $\|D\|_F \|W\|_F < \sqrt{m}$, condition (A.7) holds and (A.6) is valid. To simplify the notation we let $\Sigma_{\text{sum}} = \Sigma_{\text{spl}} = 0$, leading to the representation of X_\star as

$$\text{vec}(X_\star) = \sum_{n=0}^{\infty} m^{-n} ((DW)^{\otimes 2})^n \text{vec}(m^{-1}D\Sigma_w D^\top + \Sigma_a).$$

Returning to the matrix notation we have

$$X_\star = \sum_{n=0}^{\infty} m^{-n} (DW)^n (m^{-1}D\Sigma_w D^\top + \Sigma_a) ((DW)^n)^\top.$$

That is it. \square