
Growing Tiny Networks: Spotting Expressivity Bottlenecks and Fixing Them Optimally

Manon Verbockhaven, Théo Rudkiewicz, Sylvain Chevallier, Guillaume Charpiat

TAU team, LISN, Université Paris-Saclay, CNRS, Inria, 91405, Orsay, France

firstname.name@inria.fr

Abstract

Machine learning tasks are generally formulated as optimization problems, where one searches for an optimal function within a certain functional space. In practice, parameterized functional spaces are considered, in order to be able to perform gradient descent. Typically, a neural network architecture is chosen and fixed, and its parameters (connection weights) are optimized, yielding an architecture-dependent result. This way of proceeding however forces the evolution of the function during training to lie within the realm of what is expressible with the chosen architecture, and prevents any optimization across architectures. Costly architectural hyper-parameter optimization is often performed to compensate for this. Instead, we propose to adapt the architecture on the fly during training. We show that the information about desirable architectural changes, due to expressivity bottlenecks when attempting to follow the functional gradient, can be extracted from backpropagation. To do this, we propose a mathematical definition of expressivity bottlenecks, which enables us to detect, quantify and solve them while training, by adding suitable neurons. Thus, while the standard approach requires large networks, in terms of number of neurons per layer, for expressivity and optimization reasons, we are able to start with very small neural networks and let them grow appropriately. As a proof of concept, we show results on the CIFAR dataset, matching large neural network accuracy, with competitive training time, while removing the need for standard architectural hyper-parameter search.

1 Introduction

Issues with the fixed-architecture paradigm. Universal approximation theorems such as (Hornik et al., 1989; Cybenko, 1989) are historically among the first theoretical results obtained on neural networks, stating the family of neural networks with arbitrary width as a good candidate for a parameterized space of functions to be used in machine learning. However the current common practice in neural network training consists in choosing a fixed architecture, and training it, without any possible architecture modification meanwhile. This inconveniently prevents the direct application of these universal approximation theorems, as expressivity bottlenecks that might arise in a given layer during training will not be able to be fixed. There are two approaches to circumvent this in daily practice. Either one chooses a (very) large width, to be sure to avoid expressivity and optimization issues (Hanin & Rolnick, 2019b; Raghu et al., 2017), to the cost of extra computational power consumption for training and applying such big models; to mitigate this cost, model reduction techniques are often used afterwards, using pruning, tensor factorization, quantization (Louizos et al., 2017) or distillation (Hinton et al., 2015). Or one tries different architectures and keeps the most suitable one (in terms of performance-size compromise for instance), which multiplies the computational burden by the number of trials. This latter approach relates to the Auto-DeepLearning field (Liu et al., 2020), where different exploration strategies over the space of architecture hyper-parameters (among other ones) have been tested, including reinforcement learning (Baker et al., 2017; Zoph & Le, 2016), Bayesian optimization techniques (Mendoza et al., 2016), and evolutionary approaches (Miller et al., 1989; Stanley et al., 2009; Miikkulainen et al., 2017; Bennet et al., 2021), that all rely on random tries and consequently take time for exploration. Within that line, Net2Net (Chen et al., 2015), AdaptNet (Yang et al., 2018)

and MorphNet (Gordon et al., 2018) propose different strategies to explore possible variations of a given architecture, possibly guided by model size constraints. Instead, we aim at providing a way to locate precisely expressivity bottlenecks in a trained network, which might speed up neural architecture search significantly. Moreover, based on such observations, we aim at modifying the architecture *on the fly* during training, in a single run (no re-training), using first-order derivatives only, while avoiding neuron redundancy. Related work on architecture adaptation while training includes probabilistic edges (Liu et al., 2019) or sparsifying priors (Wolinski et al., 2020). Yet the training is done on the largest architecture allowed, which is resource-consuming. On the opposite we aim at starting from the simplest architecture possible.

Optimization properties. An important reason for common practice to choose wide architectures is the associated optimization properties: sufficiently larger networks are proved theoretically and shown empirically to be better optimized than small ones (Jacot et al., 2018). Typically, small networks exhibit issues with spurious local minima, while wide ones find good nearly-global minima. One of our goals is to train small networks without suffering from such optimization difficulties.

Neural architecture growth. A related line of work consists in growing networks neuron by neuron, by iteratively estimating the best possible neurons to add, according to a certain criterion. For instance, approaches such as (Wu et al., 2019) or Firefly (Wu et al., 2020) aim at escaping local minima by adding neurons that minimize the loss under neighborhood constraints. These neurons are found by gradient descent or by solving quadratic problems involving second-order derivatives. Other approaches (Causse et al., 2019; Bashtova et al., 2022), including GradMax (Evci et al., 2022), seek to minimize the loss as fast as possible and involve another quadratic problem. However the neurons added by these approaches are possibly redundant with existing neurons, especially if one does not wait for training convergence to a local minimum (which is time consuming) before adding neurons, therefore producing larger-than-needed architectures.

Redundancy. To our knowledge, the only approach tackling redundancy in neural architecture growth adds random neurons that are orthogonal in some sense to the ones already present (Maile et al., 2022). More precisely, the new neurons are picked within the *kernel* (preimage of $\{0\}$) of an application describing already existing neurons. Two such applications are proposed, respectively the matrix of fan-in weights and the pre-activation matrix, yielding two different notions of orthogonality. The latter formulation is close to the one of GradMax, in that both study first-order loss variations and use the same pre-activation matrix, with an important difference though: GradMax optimally decreases the loss without caring about redundancy, while the other one avoids redundancy but picks random directions instead of optimal ones. In this paper we bridge the gap between these two approaches, picking optimal directions that avoid redundancy in the pre-activation space.

Notions of expressivity. Several concepts of expressivity or complexity exist in the Machine Learning literature, ranging from Vapnik-Chervonenkis dimension (Vapnik & Chervonenkis, 1971) and Rademacher complexity (Koltchinskii, 2001) to the number of pieces in a piecewise affine function (as networks with ReLU activations are) (Serra et al., 2018; Hanin & Rolnick, 2019a). Bottlenecks have been also studied from the point of view of Information Theory, through mutual information between the activities of different layers (Tishby & Zaslavsky, 2015; Dai et al., 2018); this quantity is difficult to estimate though. Also relevant and from Information Theory, the Minimum Description Length paradigm and Kolmogorov complexity (Kolmogorov, 1965; Li et al., 2008) enable to define trade-offs between performance and model complexity.

In this article, we aim at measuring lacks of expressivity as the difference between what the backpropagation asks for and what can be done by a small parameter update (such as a gradient step), that is, between the desired variation for each activation in each layer (for each sample) and the best one that can be realized by a parameter update. Intuitively, differences arise when a layer does not have sufficient expressive power to realize the desired variation. Our main contributions are that we:

- adopt a functional analysis perspective on gradient descent in neural networks, advocating to follow the functional gradient. We not only optimize the weights of the current architecture but also dynamically adjust the architecture itself to progress towards a suitable parameterized functional

spaces. This approach mitigates optimization challenges like local minima that are due to thin architectures;

- properly define and quantify the concept of expressivity bottlenecks, both globally at the neural network output and locally at individual layers, in a computationally accessible manner. This methodology enables the localize expressivity bottlenecks within a neural network;
- formally define as a quadratic problem the best possible neurons to add to a given layer to decrease lacks of expressivity ; solve it and compute the associated expressivity gain;
- automatically adapt the architecture to the specific task by expanding it where necessary within a single run, maintaining competitive computational complexity compared to training a large model once. To remove the need for hyper-optimization of layer width, one could specify a target accuracy and stop neuron additions when reached.

2 Main concepts

2.1 Notations

Let \mathcal{F} be a functional space, e.g. $L_2(\mathbb{R}^p \rightarrow \mathbb{R}^d)$, and a loss function $\mathcal{L} : \mathcal{F} \rightarrow \mathbb{R}^+$ defined on it, of the form $\mathcal{L}(f) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\ell(f(\mathbf{x}), \mathbf{y})]$, where ℓ is the per-sample loss, assumed to be differentiable, and where \mathcal{D} is the sample distribution, from which the dataset $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ is sampled, with $\mathbf{x}_i \in \mathbb{R}^p$ and $\mathbf{y}_i \in \mathbb{R}^d$.

For the sake of simplicity we consider a feedforward neural network $f_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^d$ with L hidden layers, each of which consisting of an affine layer with weights \mathbf{W}_l followed by a differentiable activation function σ_l which satisfies $\sigma_l(0) = 0$. The network parameters are then $\theta := (\mathbf{W}_l)_{l=1 \dots L}$. The network iteratively computes:

$$\begin{aligned} \mathbf{b}_0(\mathbf{x}) &= \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \\ \forall l \in [1, L], \quad &\begin{cases} \mathbf{a}_l(\mathbf{x}) &= \mathbf{W}_l \mathbf{b}_{l-1}(\mathbf{x}) \\ \mathbf{b}_l(\mathbf{x}) &= \begin{pmatrix} \sigma_l(\mathbf{a}_l(\mathbf{x})) \\ 1 \end{pmatrix} \end{cases} \\ f_\theta(\mathbf{x}) &= \sigma_L(\mathbf{a}_L(\mathbf{x})) \end{aligned}$$

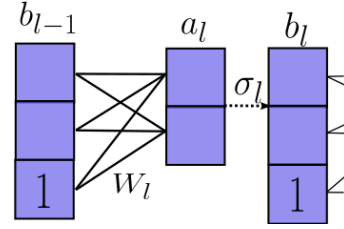


Figure 1: Notations

To any vector-valued function noted $\mathbf{t}(\mathbf{x})$ and any batch of inputs $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$, we associate the concatenated matrix $\mathbf{T}(\mathbf{X}) := (\mathbf{t}(\mathbf{x}_1) \dots \mathbf{t}(\mathbf{x}_n)) \in \mathbb{R}^{|\mathbf{t}(\cdot)| \times n}$. The matrices of pre-activation and post-activation activities at layer l over a minibatch \mathbf{X} are thus respectively: $\mathbf{A}_l(\mathbf{X}) = (\mathbf{a}_l(\mathbf{x}_1) \dots \mathbf{a}_l(\mathbf{x}_n))$ and $\mathbf{B}_l(\mathbf{X}) = (\mathbf{b}_l(\mathbf{x}_1) \dots \mathbf{b}_l(\mathbf{x}_n))$.

NB: convolutions can also be considered, with appropriate representations (cf matrix $\mathbf{b}_l^c(\mathbf{x})$ in 46).

2.2 Approach

Functional gradient descent. We take a functional perspective on the use of neural networks. Ideally in a machine learning task, one would search for a function $f : \mathbb{R}^p \rightarrow \mathbb{R}^d$ that minimizes the loss \mathcal{L} by gradient descent: $\frac{\partial f}{\partial t} = -\nabla_f \mathcal{L}(f)$ for some metric on the functional space \mathcal{F} (typically, $L_2(\mathbb{R}^p \rightarrow \mathbb{R}^d)$), where ∇_f denotes the functional gradient and t denotes the evolution time of the gradient descent. The descent direction $\mathbf{v}_{\text{goal}} := -\nabla_f \mathcal{L}(f)$ is a function of the same type as f and whose value at \mathbf{x} is easily computable as $\mathbf{v}_{\text{goal}}(\mathbf{x}) = -(\nabla_f \mathcal{L}(f))(\mathbf{x}) = -\nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y}(\mathbf{x}))|_{\mathbf{u}=f(\mathbf{x})}$ (see Appendix A.1 for more details). This direction \mathbf{v}_{goal} is the best infinitesimal variation in \mathcal{F} to add to f to decrease the loss \mathcal{L} .

Parametric gradient descent reminder. However in practice, to represent functions and to compute gradients, the infinite-dimensional functional space \mathcal{F} has to be replaced with a finite-dimensional parametric space of functions, which is usually done by choosing a particular neural network architecture \mathcal{A} with weights $\theta \in \Theta_{\mathcal{A}}$. The associated parametric search space $\mathcal{F}_{\mathcal{A}}$ then consists of all possible functions f_{θ} that can be represented with such a network for any parameter value θ . Under standard weak assumptions (see Appendix A.2), the gradient descent is of the form:

$$\frac{\partial \theta}{\partial t} = -\nabla_{\theta} \mathcal{L}(f_{\theta}) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}), \mathbf{y}) \right]. \quad (1)$$

Using the chain rule (on $\frac{\partial f_{\theta}}{\partial \theta}$ then on $\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}), \mathbf{y})$), these parameter updates yield a functional evolution:

$$\mathbf{v}_{\text{GD}} := \frac{\partial f_{\theta}}{\partial t} = \frac{\partial f_{\theta}}{\partial \theta} \frac{\partial \theta}{\partial t} = \frac{\partial f_{\theta}}{\partial \theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[\frac{\partial f_{\theta}}{\partial \theta}^T(\mathbf{x}) \mathbf{v}_{\text{goal}}(\mathbf{x}) \right] \quad (2)$$

which significantly differs from the original functional gradient descent. We will aim to augment the neural network architecture so that parametric gradient descents can get closer to the functional one.

Optimal move direction. We name $\mathcal{T}_{\mathcal{A}}^{f_{\theta}}$, or just $\mathcal{T}_{\mathcal{A}}$, the tangent space of $\mathcal{F}_{\mathcal{A}}$ at f_{θ} , that is, the set of all possible infinitesimal variations around f_{θ} under small parameter variations:

$$\mathcal{T}_{\mathcal{A}}^{f_{\theta}} := \left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta \theta \mid \text{s.t. } \delta \theta \in \Theta_{\mathcal{A}} \right\}$$

This linear space is a first-order approximation of the neighborhood of f_{θ} within $\mathcal{F}_{\mathcal{A}}$. The direction \mathbf{v}_{GD} obtained above by gradient descent is actually not the best one to consider within $\mathcal{T}_{\mathcal{A}}$. Indeed, the best move \mathbf{v}^* would be the orthogonal projection of the desired direction $\mathbf{v}_{\text{goal}} := -\nabla_{f_{\theta}} \mathcal{L}(f_{\theta})$ onto $\mathcal{T}_{\mathcal{A}}$. This projection is what a (generalization of the notion of) natural gradient would compute (Ollivier, 2017).

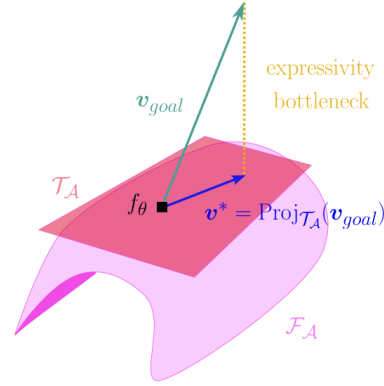


Figure 2: Expressivity bottleneck

Indeed, the parameter variation $\delta \theta^*$ associated to the functional variation $\mathbf{v}^* = \frac{\partial f_{\theta}}{\partial \theta} \delta \theta^*$ is the gradient $-\nabla_{\theta}^{\mathcal{T}_{\mathcal{A}}} \mathcal{L}(f_{\theta})$ of $\mathcal{L} \circ f_{\theta}$ w.r.t. parameters θ when considering the L_2 metric on *functional* variations $\|\frac{\partial f_{\theta}}{\partial \theta} \delta \theta\|_{L_2(\mathcal{T}_{\mathcal{A}})}$, not to be confused with the usual gradient $\nabla_{\theta} \mathcal{L}(f_{\theta})$, based on the L_2 metric on *parameter* variations $\|\delta \theta\|_{L_2(\mathbb{R}^{|\Theta_{\mathcal{A}}|})}$. This can be seen in a proximal formulation as:

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathcal{T}_{\mathcal{A}}} \|\mathbf{v} - \mathbf{v}_{\text{goal}}\|^2 = \arg \min_{\mathbf{v} \in \mathcal{T}_{\mathcal{A}}} \left\{ D_f \mathcal{L}(f)(\mathbf{v}) + \frac{1}{2} \|\mathbf{v}\|^2 \right\} \quad (3)$$

where D is the directional derivative (see details in Appendix A.3), or equivalently as:

$$\delta \theta^* = \arg \min_{\delta \theta \in \Theta_{\mathcal{A}}} \left\| \frac{\partial f_{\theta}}{\partial \theta} \delta \theta - \mathbf{v}_{\text{goal}} \right\|^2 = \arg \min_{\delta \theta \in \Theta_{\mathcal{A}}} \left\{ D_{\theta} \mathcal{L}(f_{\theta})(\delta \theta) + \frac{1}{2} \left\| \frac{\partial f_{\theta}}{\partial \theta} \delta \theta \right\|^2 \right\} =: -\nabla_{\theta}^{\mathcal{T}_{\mathcal{A}}} \mathcal{L}(f_{\theta}).$$

Lack of expressivity. When \mathbf{v}_{goal} does not belong to the reachable subspace $\mathcal{T}_{\mathcal{A}}$, there is a lack of expressivity, that is, the parametric space \mathcal{A} is not rich enough to follow the ideal functional gradient descent. This happens frequently with small neural networks (see Appendix A.4 for an example). The expressivity bottleneck is then quantified as the distance $\|\mathbf{v}^* - \mathbf{v}_{\text{goal}}\|$ between the functional gradient \mathbf{v}_{goal} and the optimal functional move \mathbf{v}^* given the architecture \mathcal{A} (in the sense of Eq. 3).

2.3 Generalizing to all layers

Ideal updates. The same reasoning can be applied to the pre-activations \mathbf{a}_l at each layer l , seen as functions $\mathbf{a}_l : \mathbf{x} \in \mathbb{R}^p \mapsto \mathbf{a}_l(\mathbf{x}) \in \mathbb{R}^{d_l}$ defined over the input space of the neural network. The optimal parameter update for a given layer l then follows the projection of the desired update $-\nabla_{\mathbf{a}_l} \mathcal{L}(f_\theta)$ of the pre-activation functions \mathbf{a}_l onto the linear subspace $\mathcal{T}_{\mathcal{A}}^{\mathbf{a}_l}$ of pre-activation variations that are possible with the architecture, as we will detail now.

Given a sample $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$, standard backpropagation already iteratively computes $\mathbf{v}_{\text{goal}}^l(\mathbf{x}) := -(\nabla_{\mathbf{a}_l} \mathcal{L}(f_\theta))(\mathbf{x}) = -\nabla_{\mathbf{u}} \ell(\sigma_L(\mathbf{W}_L \sigma_{L-1}(\mathbf{W}_{L-1} \dots \sigma_l(\mathbf{u}))), \mathbf{y})|_{\mathbf{u}=\mathbf{a}_l(\mathbf{x})}$, which is the derivative of the loss $\ell(f_\theta(\mathbf{x}), \mathbf{y})$ with respect to the pre-activations $\mathbf{u} = \mathbf{a}_l(\mathbf{x})$ of each layer. This is usually performed in order to compute the gradients w.r.t. model parameters \mathbf{W}_l , as $\nabla_{\mathbf{W}_l} \ell(f_\theta(\mathbf{x}), \mathbf{y}) = \frac{\partial \mathbf{a}_l(\mathbf{x})}{\partial \mathbf{W}_l} \nabla_{\mathbf{a}_l} \ell(f_\theta(\mathbf{x}), \mathbf{y})$.

$\mathbf{v}_{\text{goal}}^l(\mathbf{x}) := -(\nabla_{\mathbf{a}_l} \mathcal{L}(f_\theta))(\mathbf{x})$ indicates the direction in which one would like to change the layer pre-activations $\mathbf{a}_l(\mathbf{x})$ in order to decrease the loss at point \mathbf{x} . However, given a minibatch of points (\mathbf{x}_i) , most of the time no parameter move $\delta\theta$ is able to induce this progression for each \mathbf{x}_i simultaneously, because the θ -parameterized family of functions \mathbf{a}_l is not expressive enough.

Activity update resulting from a parameter change. Given a subset of parameters $\tilde{\theta}$ (such as the ones specific to a layer: $\tilde{\theta} = \mathbf{W}_l$), and an incremental direction $\delta\tilde{\theta}$ to update these parameters (e.g. the one resulting from a gradient descent: $\delta\tilde{\theta} = -\eta \sum_{(\mathbf{x}, \mathbf{y}) \in \text{minibatch}} \nabla_{\tilde{\theta}} \ell(f_\theta(\mathbf{x}), \mathbf{y})$ for some learning rate η), the impact of the parameter update $\delta\tilde{\theta}$ on the pre-activations \mathbf{a}_l at layer l at order 1 in $\delta\tilde{\theta}$ is $\mathbf{v}^l(\mathbf{x}, \delta\tilde{\theta}) := \frac{\partial \mathbf{a}_l(\mathbf{x})}{\partial \tilde{\theta}} \delta\tilde{\theta}$.

3 Expressivity bottlenecks

We now quantify expressivity bottlenecks at any layer l as the distance between the desired activity update $\mathbf{v}_{\text{goal}}^l(\cdot)$ and the best realizable one $\mathbf{v}^l(\cdot)$ (cf Figure 2):

Definition 3.1 (Lack of expressivity). *For a neural network f_θ and a minibatch of points $\mathbf{X} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, we define the lack of expressivity at layer l as how far the desired activity update $\mathbf{V}_{\text{goal}}^l = (\mathbf{v}_{\text{goal}}^l(\mathbf{x}_1), \mathbf{v}_{\text{goal}}^l(\mathbf{x}_2), \dots)$ is from the closest possible activity update $\mathbf{V}^l = (\mathbf{v}^l(\mathbf{x}_1), \mathbf{v}^l(\mathbf{x}_2), \dots)$ realizable by a parameter change $\delta\theta$:*

$$\Psi^l := \min_{\mathbf{v}^l \in \mathcal{T}_{\mathcal{A}}^{\mathbf{a}_l}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{v}^l(\mathbf{x}_i) - \mathbf{v}_{\text{goal}}^l(\mathbf{x}_i)\|^2 = \min_{\delta\theta} \frac{1}{n} \|\mathbf{V}^l(\mathbf{X}, \delta\theta) - \mathbf{V}_{\text{goal}}^l(\mathbf{X})\|_{\text{Tr}}^2 \quad (4)$$

where $\|\cdot\|$ stands for the L_2 norm, $\|\cdot\|_{\text{Tr}}$ for the Frobenius norm, and $\mathbf{V}^l(\mathbf{X}, \delta\theta)$ is the activity update resulting from parameter change $\delta\theta$ as defined in previous section. In the two following parts we fix the minibatch \mathbf{X} and simplify notations accordingly by removing the dependency on \mathbf{X} .

3.1 Best move without modifying the architecture of the network

Let $\delta\mathbf{W}_l^*$ be the solution of 4 when the parameter variation $\delta\theta$ is restricted to involve only layer l parameters, i.e. \mathbf{W}_l . This move is sub-optimal in that it does not result from an update of all architecture parameters but only of the current layer ones:

$$\delta\mathbf{W}_l^* = \arg \min_{\delta\mathbf{W}} \frac{1}{n} \|\mathbf{V}^l(\delta\mathbf{W}) - \mathbf{V}_{\text{goal}}^l\|_{\text{Tr}}^2 \quad (5)$$

Proposition 3.1. *The solution of Problem (5) is: $\delta\mathbf{W}_l^* = \frac{1}{n} \mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T (\frac{1}{n} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T)^+$ where P^+ denotes the generalized inverse of matrix P .*

¹Note: given a learning rate η , in the sequel we will rather consider $\mathbf{v}_{\text{goal}}^l(\mathbf{x}) := -\eta \nabla_{\mathbf{a}_l} \mathcal{L}(f_\theta)(\mathbf{x})$

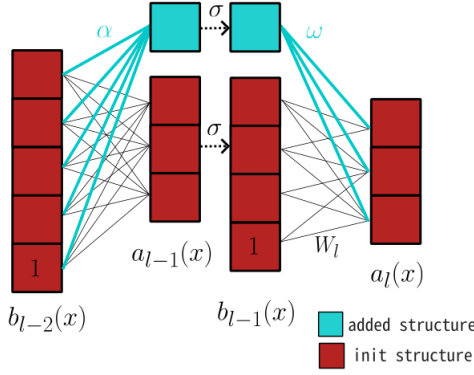


Figure 3: Adding one neuron to layer l in cyan ($K = 1$), with connections in cyan. Here, $\alpha \in \mathbb{R}^5$ and $\omega \in \mathbb{R}^3$.

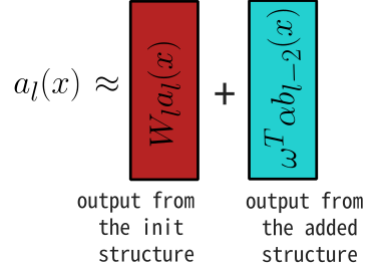


Figure 4: Sum of functional moves

This update $\delta \mathbf{W}_l^*$ is not equivalent to the usual gradient descent update, whose form is $\delta \mathbf{W}_l^{\text{GD}} \propto \mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T$. In fact the associated activity variation, $\delta \mathbf{W}_l^* \mathbf{B}_{l-1}$, is the projection of $\mathbf{V}_{\text{goal}}^l$ on the post-activation matrix of layer $l-1$, that is to say onto the span of all possible post-activation directions, through the projector $\frac{1}{n} \mathbf{B}_{l-1}^T (\frac{1}{n} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T)^+ \mathbf{B}_{l-1}$. To increase expressivity if needed, we will aim at increasing this span with the most useful directions to close the gap between this best update and the desired one. Note that the update $\delta \mathbf{W}_l^*$ consists of a standard gradient ($\mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T$) and of a (kind of) natural gradient only for the last part (projector), as we consider metrics in the pre-activation space.

3.2 Reducing expressivity bottleneck by modifying the architecture

To get as close as possible to $\mathbf{V}_{\text{goal}}^l$ and to increase the expressive power of the current neural network, we modify each layer of its structure. At layer $l-1$, we add K neurons n_1, \dots, n_K with input weights $\alpha_1, \dots, \alpha_K$ and output weights $\omega_1, \dots, \omega_K$ (cf Figure 3). We have the following expansions by concatenation: $\mathbf{W}_{l-1}^T \leftarrow (\mathbf{W}_{l-1}^T \quad \alpha_1 \quad \dots \quad \alpha_K)$ and $\mathbf{W}_l \leftarrow (\mathbf{W}_l \quad \omega_1 \quad \dots \quad \omega_K)$. We note this architecture modification $\theta \leftarrow \theta \oplus \theta_{\leftrightarrow}^K$ where \oplus is the concatenation sign and $\theta_{\leftrightarrow}^K := (\alpha_k, \omega_k)_{k=1}^K$ are the K added neurons.

The added neurons could be chosen randomly, as in usual neural network initialization, but this would not yield any guarantee regarding the impact on the system loss. Another possibility would be to set either input weights $(\alpha_k)_{k=1}^K$ or output weights $(\omega_k)_{k=1}^K$ to 0, so that the function $f_\theta(\cdot)$ would not be modified, while its gradient w.r.t. θ would be enriched from the new parameters. Another option is to solve a optimization problem as in the previous section with the modified structure $\theta \leftarrow \theta \oplus \theta_{\leftrightarrow}^K$ and jointly search for both the optimal new parameters $\theta_{\leftrightarrow}^K$ and the optimal variation $\delta \mathbf{W}$ of the old ones:

$$\arg \min_{\theta_{\leftrightarrow}^K, \delta \mathbf{W}} \left\| \mathbf{V}^l(\delta \mathbf{W} \oplus \theta_{\leftrightarrow}^K) - \mathbf{V}_{\text{goal}}^l \right\|_{\text{Tr}}^2 \quad (6)$$

As shown in figure 4, the displacement \mathbf{V}^l at layer l is actually a sum of the moves induced by the neurons already present ($\delta \mathbf{W}$) and by the added neurons ($\theta_{\leftrightarrow}^K$), our problem rewrites as :

$$\arg \min_{\theta_{\leftrightarrow}^K, \delta \mathbf{W}} \left\| \mathbf{V}^l(\theta_{\leftrightarrow}^K) + \mathbf{V}^l(\delta \mathbf{W}) - \mathbf{V}_{\text{goal}}^l \right\|_{\text{Tr}}^2 \quad (7)$$

with $\mathbf{v}^l(\mathbf{x}, \theta_{\leftrightarrow}^K) := \sum_{k=1}^K \omega_k (b_{l-2}(\mathbf{x}))^T \alpha_k$ (See A.5). We choose $\delta \mathbf{W}$ as the best move of already-existing parameters as defined in Proposition 3.1 and we note $\mathbf{V}_{\text{goal}_{proj}}^l := \mathbf{V}_{\text{goal}}^l - \mathbf{V}^l(\delta \mathbf{W}^*)$. We are looking for the solution $(K^*, \theta_{\leftrightarrow}^{K*})$ of the optimization problem :

$$\arg \min_{K, \theta_{\leftrightarrow}^K} \left\| \mathbf{V}^l(\theta_{\leftrightarrow}^K) - \mathbf{V}_{\text{goal}_{proj}}^l \right\|_{\text{Tr}}^2. \quad (8)$$

This quadratic optimization problem can be solved thanks to the low-rank matrix approximation theorem (Eckart & Young, 1936), using matrices $\mathbf{N} := \frac{1}{n} \mathbf{B}_{l-2} (\mathbf{V}_{\text{goalproj}}^l)^T$ and $\mathbf{S} := \frac{1}{n} \mathbf{B}_{l-2} \mathbf{B}_{l-2}^T$. As \mathbf{S} is semi-positive definite, let its truncated SVD be $\mathbf{S} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^T$, and define $\mathbf{S}^{-\frac{1}{2}} := \mathbf{U} \sqrt{\mathbf{\Sigma}}^{-1} \mathbf{U}^T$, with the convention that the inverse of 0 eigenvalues is 0. Finally, consider the truncated SVD of matrix $\mathbf{S}^{-\frac{1}{2}} \mathbf{N} = \sum_{k=1}^R \lambda_k \mathbf{u}_k \mathbf{v}_k^T$, where R is the rank of the matrix $\mathbf{S}^{-\frac{1}{2}} \mathbf{N}$. Then:

Proposition 3.2. *The solution of Problem (8) is:*

- optimal number of neurons: $K^* = R$
- their optimal weights: $\theta_{\leftrightarrow}^{K^*} = (\boldsymbol{\alpha}_k^*, \boldsymbol{\omega}_k^*)_k^{K^*} = \left(\sqrt{\lambda_k} \mathbf{S}^{-\frac{1}{2}} \mathbf{u}_k, \sqrt{\lambda_k} \mathbf{v}_k \right)_k^{K^*}$

Moreover for any number of neurons $K \leq R$, and associated weights $\theta_{\leftrightarrow}^{K,*}$, the expressivity gain can be quantified very simply as a function of the singular values λ_k :

$$\Psi_{\theta \oplus \theta_{\leftrightarrow}^{K,*}}^l \leq \Psi_{\theta}^l - \sum_{k=1}^K \lambda_k^2 \quad (9)$$

Proposition 3.3. *If \mathbf{S} is positive definite, then solving 8 is equivalent to taking $\boldsymbol{\omega}_k = \mathbf{N} \boldsymbol{\alpha}_k$ and finding the K first eigenvectors $\boldsymbol{\alpha}_k$ associated to the K largest eigenvalues λ of the generalized eigenvalue problem :*

$$\mathbf{N} \mathbf{N}^T \boldsymbol{\alpha}_k = \lambda \mathbf{S} \boldsymbol{\alpha}_k$$

Corollary 1. *For all integers m, m' such that $m + m' \leq R$, at order one in \mathbf{V} , adding $m + m'$ neurons simultaneously according to the previous method is equivalent to adding m neurons then m' neurons by applying successively the previous method twice.*

Note: Problems (7) and (8) are generally not equivalent, though similar (cf C.4).

Note 2: Solving 8 is equivalent to minimizing the loss \mathcal{L} at order one in \mathbf{V}^l . Furthermore performing an update of architecture with $\delta \mathbf{W}^*$ (5) and a neuron addition with $\theta_{\leftrightarrow}^{K,*}$ (3.2), has an impact on the loss at first order in $\|\mathbf{V}^l(\delta \mathbf{W}^*) + \mathbf{V}^l(\theta_{\leftrightarrow}^{K,*})\|$ as :

$$\mathcal{L}(f_{\theta \oplus \theta_{\leftrightarrow}^{K,*}}) \approx \mathcal{L}(f_{\theta}) - \frac{1}{\eta n} \left(\sigma'_{l-1}(0) \Delta_{\theta_{\leftrightarrow}^{K,*}} + \Delta_{\delta \mathbf{W}^*} \right) \quad (10)$$

With

$$\Delta_{\theta_{\leftrightarrow}^{K,*}} := \left\langle \mathbf{V}_{\text{goalproj}}^l, \mathbf{V}^l(\theta_{\leftrightarrow}^{K,*}) \right\rangle_{\text{Tr}} = \sum_{k=1}^K \lambda_k^2 \quad (11)$$

$$\Delta_{\delta \mathbf{W}^*} := \left\langle \mathbf{V}_{\text{goal}}^l, \mathbf{V}^l(\delta \mathbf{W}^*) \right\rangle_{\text{Tr}} \geq 0 \quad (12)$$

The family $\{\mathbf{V}^{l+1}((\boldsymbol{\alpha}_k, \boldsymbol{\omega}_k))\}_{k=1}^K$ of pre-activity variations induced by adding the neurons $\theta_{\leftrightarrow}^{K,*}$ is orthogonal for the trace scalar product. We could say that the added neurons are orthogonal to each other (and to the already-present ones) in that sense. Interestingly, the GradMax method (Evci et al., 2022) also aims at minimizing the loss 10, but without avoiding redundancy (see Appendix B.1 for more details).

Addition of new neurons. In practice before adding new neurons (α, ω) , we multiply them by an amplitude factor γ found by a simple line search (see Appendix E.3), i.e. we add $(\sqrt{\gamma} \alpha, \sqrt{\gamma} \omega)$. The addition of each neuron k has an impact on the loss of the order of $\gamma \lambda_k^2$ provided γ is small. This performance gain could be used in a selection criterion realizing a trade-off with computational complexity. A selection based on statistical significance of singular values can also be performed. The full algorithm and its complexity are detailed in Appendices E.4 and E.5.

4 About greedy growth sufficiency and TINY convergence

One might wonder whether a greedy approach on layer growth might get stuck in a non-optimal state. By *greedy* we mean that every neuron added has to decrease the loss. Since in this work we add neurons layer per layer independently, we study here the case of a single hidden layer network, to spot potential layer growth issues. For the sake of simplicity, we consider the task of least square regression towards an explicit continuous target f^* , defined on a compact set. That is, we aim at minimizing the loss:

$$\inf_f \sum_{\mathbf{x} \in \mathcal{D}} \|f(\mathbf{x}) - f^*(\mathbf{x})\|^2 \quad (13)$$

where $f(\mathbf{x})$ is the output of the neural network and \mathcal{D} is the training set. Proofs and supplementary propositions are deferred to Appendix D, in particular D.4 and D.7.

First, if one allows only adding neurons but no modification of already existing ones:

Proposition 4.1 (Exponential convergence to 0 training error by ReLU neuron additions). *It is possible to decrease the loss exponentially fast with the number t of added neurons, i.e. as $\gamma^t \mathcal{L}(f)$, towards 0 training loss, and this in a greedy way, that is, such that each added neuron decreases the loss. The factor γ is $\gamma = 1 - \frac{1}{n^3 d'} \left(\frac{d_m}{d_M}\right)^2$, where d_m and d_M are quantities solely dependent on the dataset geometry, d' is the output dimension of the network, and n is the dataset size.*

In particular, there exists no situation where one would need to add many neurons simultaneously to decrease the loss: it is always feasible with a single neuron.

TINY might get stuck when no correlation between inputs \mathbf{x}_i and desired output variations $f^*(\mathbf{x}_i) - f(\mathbf{x}_i)$ can be found anymore. To prevent this, one can choose an auxiliary method to add neurons in such cases, for instance random neurons (with a line search over their amplitude, cf. Appendix D.3), or locally-optimal neurons found by gradient descent, or solutions of higher-order expressivity bottleneck formulations using further developments of the activation function. We will name *completed-TINY* the completion of TINY by any such auxiliary method.

Now, if we also update already existing weights when adding new neurons, we get a stronger result:

Proposition 4.2 (Completed-TINY reaches 0 training error in at most n neuron additions). *Under certain assumptions (full batch optimization, updating already existing parameters, and, more technically: polynomial activation function of order $\geq n^2$), completed-TINY reaches 0 training error in at most n neuron additions almost surely.*

Hence we see the importance of updating existing parameters on the convergence speed. This optimization protocol is actually the one we follow in practice when training neural networks with TINY (except when comparing with other methods using their protocol).

Note that our approach shares similarity with gradient boosting Friedman (2001) somehow, as we grow the architecture based on the gradient of the loss. Note also that finding the optimal neuron to add is actually NP-hard (Bach, 2017), but that we do not need new neuron optimality to converge to 0 training error.

5 Results

5.1 Comparison with GradMax on CIFAR-100

The closest growing method to TINY is GradMax (Evci et al. (2022)), as it solves a quadratic problem similar to (8). By construction, the objective of GradMax is to decrease the loss as fast as possible considering an infinitesimal increment of new neurons. The main difference is that GradMax does not take into account the expressivity of the current architecture as TINY does in (8) by projecting v_{goal} . In-depth details about the difference between the GradMax and TINY are provided in Appendix B.1.

In this section, we show on the CIFAR-100 dataset that solving (8) instead of B.1 (defined by GradMax) to grow a network allows better final performance and almost full expressivity power. To do so, we have

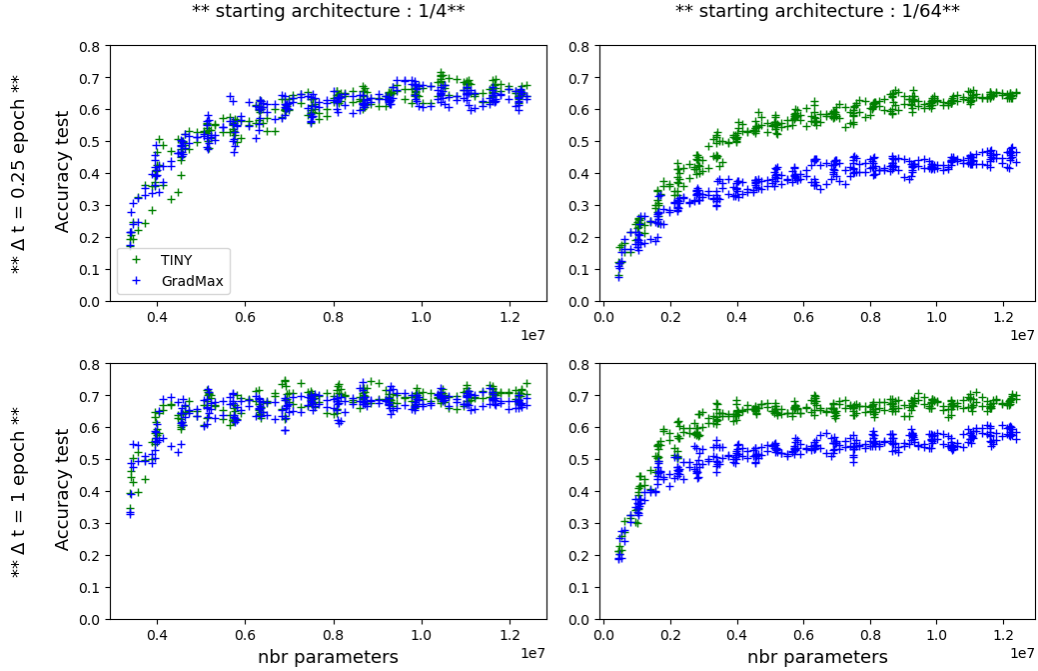


Figure 5: Test accuracy as a function of the number of parameters during architecture growth from ResNet_s to ResNet_{18} . The left (resp. right) column is for the starting architecture $\text{ResNet}_{1/4}$ (resp. $\text{ResNet}_{1/64}$). The upper (resp. lower) row is for Δt equal to 0.25 (resp. 1) epoch.

re-implemented the GradMax method and mimicked its growing process which consists in increasing the architecture of a thin ResNet_{18} until it reaches the architecture of the usual ResNet_{18} . This process is described in the pseudo code 1, where two parameters can be chosen : the relative thinness s of the starting architecture, w.r.t. the usual ResNet_{18} architecture 3 ($s = 1/4$ or $s = 1/64$), and the amount of training time between consecutive neuron additions ($\Delta t = 1$ or $\Delta t = 0.25$ epochs). Then the number of parameters and the performance of the growing network are evaluated at regular intervals to plot Figure 5.

Once the models have reached the final architecture ResNet_{18} , they are trained for 250 epochs (or 500 epochs if they have not converged on the training set). We have summarized the final performance in Table 1. We also added the column Reference, which gives the performance of a ResNet_{18} trained from scratch by usual gradient descent with all its neurons. We do not expect TINY or GradMax to achieve the performance of the reference as its architecture and optimisation process have been optimised for years.

The details of the protocol can be found in the annexes F.1, as well as other technical details such as the dynamic of the learning batch size E.2, the number of examples used to solve the expressivity bottleneck 8 and the complexity of the algorithms E.5. For both methods, all the latter apply so that the main differences between GradMax and TINY in this experiment is the mathematical definition of the new neurons.

For $s = 1/64$, we observe a significant difference in performance between TINY and GradMax methods. While TINY models almost achieve the reference’s performance, GradMax remain stuck 10 points below. This suggests that the framework proposed by GradMax is not sufficient to be able to start with an architecture far from full expressivity, i.e. $\text{ResNet}_{1/64}$, while TINY is able to handle it. As for the setting $s = 1/4$, both methods seem equivalent in terms of final performance and achieve full expressivity.

The curves on Figure 6, which are extracted from Figure 13 in the appendix, show that TINY models have converged at the end of the growing process, while GradMax ones have not. This latter effect contrasts with GradMax formulation which is to accelerate the gradient descent as fast as possible by adding neurons. Furthermore GradMax needs extra training to achieve full expressivity: for the particular setting $s = 1/64, \Delta t = 1$, the extra training time required by GradMax is twice as high as TINY’s, as shown in Figure

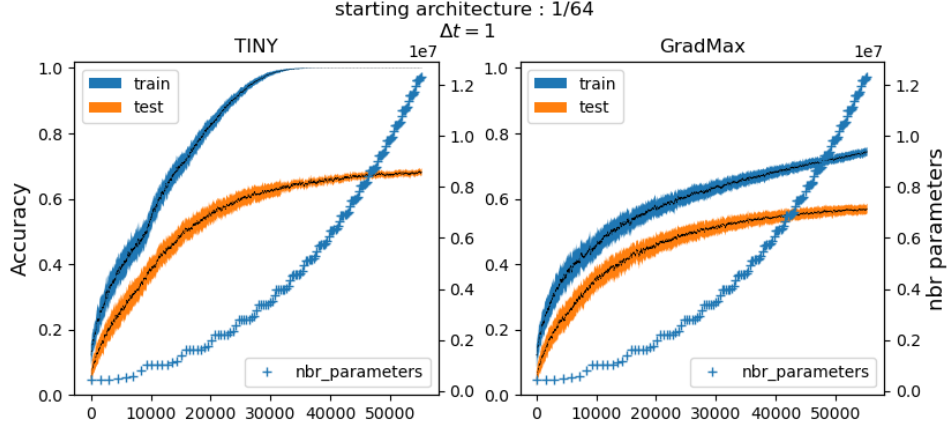


Figure 6: Evolution of accuracy and number of parameters as a function of gradient step for the setting $\Delta t = 1$, $s = 1/64$ for TINY and GradMax, mean and standard deviation over two runs. Other settings in the annexes 13

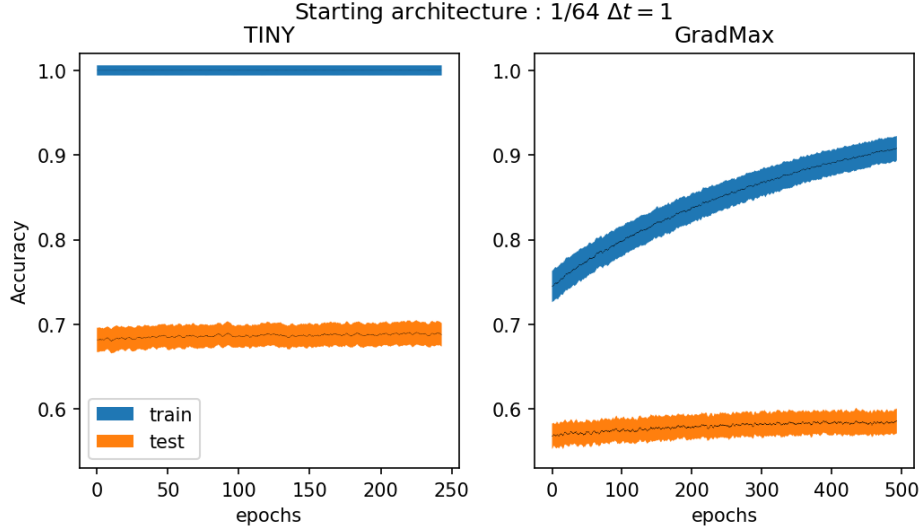


Figure 7: Evolution of accuracy and number of parameters as a function of gradient step for the setting $\Delta t = 1$, $s = 1/64$ during extra training for TINY and GradMax, mean and standard deviation over two runs. Other settings in the annexes 15 and 14.

$s \backslash \Delta t$	TINY		GradMax		Reference
	0.25	1	0.25	1	
1/4	67.0 \pm 0.1	71.0 \pm 0.1	65.0 \pm 0.1	69.0 \pm 0.1	72.9 \pm 0.1 ^{5*}
1/4	70.0 \pm 0.2 ^{5*}	71.0 \pm 0.2^{5*}	67.0 \pm 0.2 ^{5*}	69.0 \pm 0.1 ^{5*}	
1/64	66.0 \pm 0.1	68.0 \pm 0.4	45.0 \pm 0.2	57.0 \pm 0.2	
1/64	69.0 \pm 0.1 ^{5*}	69.0 \pm 0.6^{5*}	57.0 \pm 0.3 ^{10*}	59.0 \pm 0.1 ^{10*}	

Table 1: Final accuracy on test of ResNet18 after the architecture growth (*grey*) and after convergence (*black*). The number of stars indicates the multiple of 50 epochs needed to achieve convergence. With the starting architecture ResNet_{1/64} and $\Delta t = 0.25$, the method TINY achieves 66.0 \pm 0.1 on test after its growth and it reaches 69.0 \pm 0.1^{5*} after 5* := 5 \times 50 epochs (examples of training curves for the extra training in Figure 14). Mean and standard deviation are performed on 2 runs for each setting.

7. This need for extra training also appears for all settings in Table 1. In particular for $s = 1/64, \Delta t = 0.25$, the difference in performance after and before extra training goes up to 20 % above the initial performance while it is only of 6% for TINY.

5.2 Comparison with Random on CIFAR-100 : initialisation impact

In this section, we focus on the impact of the new neurons’ initialization. To do so, we consider as a baseline the Random method, which initializes the new neurons according to a Gaussian distribution: $(\alpha_k^*, \omega_k^*)_{k=1}^K \sim \mathcal{N}(0, Id)$. Also, when adding new neurons, instead of normalizing them as previously, we search for the best scaling using a line-search on the loss. Thus, we perform the operation $\theta_{\leftrightarrow}^K \leftarrow \gamma^* \theta_{\leftrightarrow}^K$, with the amplitude factor $\gamma^* \in \mathbb{R}$ defined as :

$$\gamma^* := \arg \min_{\gamma \in [-L, L]} \sum_i \mathcal{L}(f_{\theta \oplus \gamma \theta_{\leftrightarrow}^K}(\mathbf{x}_i), \mathbf{y}_i) \quad \text{with } \gamma \theta_{\leftrightarrow}^{K*} = (\gamma \alpha_k^*, \gamma \omega_k^*)_k^K \quad (14)$$

with L a positive constant. More details can be found in Appendix E.3.2 and in Algorithm 1. With such an amplitude factor, one can measure the quality of the directions generated by TINY and Random by quantifying the maximal decrease of loss in these directions.

To better measure the impact of the initialisation method, and to distinguish it from the optimization process, we do not perform any gradient descent. This contrasts with the previous section where long training time after architecture growth was modifying the direction of the added neurons, dampening initialization impact with training time, especially as they were added with a small amplitude factor (cf Section E.3.1).

With these two modifications to the protocol of previous section, we obtain Figure 8. We see the crucial impact of TINY initialization compared to the Random one. Indeed, TINY reaches more than 17% accuracy just by adding neurons (without any further update), which accounts for about one quarter of the total accuracy with the full method (69% in Table 1 using in plus gradient descent). On the opposite, the Random initialization does not contribute to the accuracy through the growing process (just about 1%); this can be explained and quantified as follows. In the random setting, we model $\mathbf{v}(X)$ and $\mathbf{v}_{\text{goal}}(X)$ as independent Gaussian variables following $\mathcal{N}(0_d, I_d \frac{1}{d})$ where d is the dimension of \mathbf{v}_{goal} and \mathbf{v} . From Equation 10, the scalar product $\langle \mathbf{V}(X), \mathbf{V}_{\text{goal}}(X) \rangle := \frac{1}{n} \sum_i \mathbf{v}_{\text{goal}}(\mathbf{x}_i)^T \mathbf{v}(\mathbf{x}_i)$ is a proxy of the expected decrease of loss after each architecture growth. This quantity can be approximated by its standard deviation, ie $\frac{1}{\sqrt{nd}}$, which makes the expected relative gain of loss (for a gradient step) of the order of magnitude of $\frac{1}{\sqrt{64}}$ for the first layer and $\frac{1}{\sqrt{512}}$ for the last layer when compared to the true gradient, and consequently when compared to TINY. Furthermore, one can take into account the effect of a line search over the random direction: in that case the expected relative loss gain is quadratic in the angle between the directions and therefore of the order of magnitude of $\frac{1}{64}$ or $\frac{1}{512}$ respectively (see Appendix D.3).

Note that the search interval of equation 14 for can be shrunk to $[0, L]$ with TINY initializations, as the first order development of the loss in Equation 10 is positive. This property is the direct consequence of the

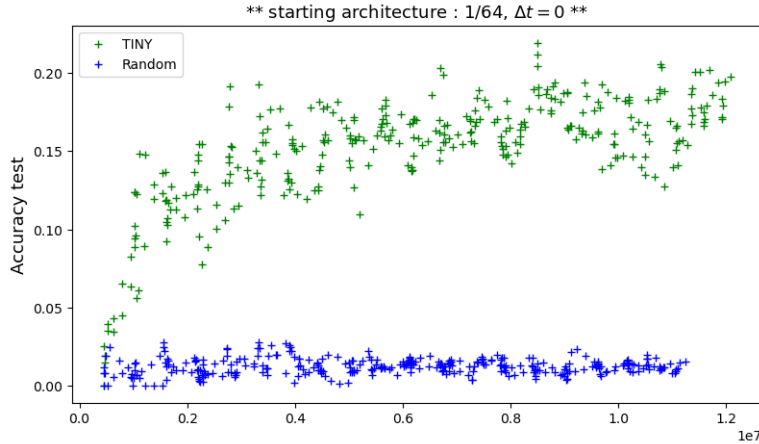


Figure 8: Test accuracy as a function of the number of parameters during architecture growth from ResNet_{1/64} to ResNet₁₈.

definition of \mathbf{V}^* as the minimizer of the expressivity bottleneck equation 8. One can also note that we do not include GradMax in Figure 8, because its protocol initializes the on-going weights to zero ($\alpha_k \leftarrow 0$) and imposes a small norm on its out-going weights ($\|\omega_k\| = \varepsilon$). Those two aspects make the amplitude factor γ^* meaningless and the impact of the new neurons initialization invisible without gradient descent.

The code will be available on git; for now the code is available as supplementary materials (zip file TINY_code.zip).

6 Conclusion

We provided the theoretical principles of TINY, a method to grow the architecture of a neural net while training it; starting from a very thin architecture, TINY adds the neurons where needed and yields a fully trained architecture at the end. Our method relies on the functional gradient to find new directions that tackle the expressivity bottleneck, even for small networks, by expanding their space of parameters. This way, we combine in the same framework gradient descent and neural architecture search, that is optimizing the network parameters and its architectures at the same time, and this, in a way that guarantees convergence to 0 training error, thus escaping expressivity bottlenecks indeed.

The method is generic for all architectures and is instantiated for linear and convolutional layers. Extension to self-attention mechanism (transformers) is part of future works. Although the common architectures consist of a succession of layers, a research direction is to develop tools handling general computational graphs (such as U-net, Inception, Dense-Net), which offers the possibility to let the architecture graph grow and bypass manual design.

Another possible development would be to study the statistical reliability of the TINY method, for instance using tools borrowed from random matrix theory. Indeed statistical tests can be applied on intermediate computations to obtain the new neurons. An interesting byproduct of this approach would be to define a threshold to select neurons found by 3.2, based on statistical significance.

References

- Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017.
- Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=S1c2cvqee>.

-
- Kateryna Bashtova, Mathieu Causse, Cameron James, Florent Masmoudi, Mohamed Masmoudi, Houcine Turki, and Joshua Wolff. Application of the topological gradient to parsimonious neural networks. In *Computational Sciences and Artificial Intelligence in Industry*, pp. 47–61. Springer, 2022.
- Pauline Bennet, Carola Doerr, Antoine Moreau, Jeremy Rapin, Fabien Teytaud, and Olivier Teytaud. Nevgrad: black-box optimization platform. *ACM SIGEVOlution*, 14(1):8–15, 2021.
- Mathieu Causse, Cameron James, Mohamed Slim Masmoudi, and Houcine Turki. Parsimonious neural networks. In *Cesar Conference*, 2019. URL https://www.cesar-conference.org/wp-content/uploads/2019/10/s5_p1_21_1330.pdf.
- Tianqi Chen, Ian Goodfellow, and Jonathon Shlens. Net2net: Accelerating learning via knowledge transfer. *arXiv preprint arXiv:1511.05641*, 2015.
- G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Bin Dai, Chen Zhu, Baining Guo, and David Wipf. Compressing neural networks using the variational information bottleneck. In *International Conference on Machine Learning*, pp. 1135–1144. PMLR, 2018.
- Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, September 1936. ISSN 1860-0980. doi: 10.1007/BF02288367. URL <https://doi.org/10.1007/BF02288367>.
- Utku Evci, Bart van Merriënboer, Thomas Unterthiner, Fabian Pedregosa, and Max Vladymyrov. Gradmax: Growing neural networks using gradient information. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=qjN4h_wwU0.
- Harley Flanders. Differentiation under the integral sign. *The American Mathematical Monthly*, 80(6):615–627, 1973. ISSN 00029890, 19300972. URL <http://www.jstor.org/stable/2319163>.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Training pruned neural networks. *CoRR*, abs/1803.03635, 2018. URL <http://arxiv.org/abs/1803.03635>.
- Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pp. 1189–1232, 2001.
- Ariel Gordon, Elad Eban, Ofir Nachum, Bo Chen, Hao Wu, Tien-Ju Yang, and Edward Choi. Morphnet: Fast & simple resource-constrained structure learning of deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1586–1595, 2018.
- Boris Hanin and David Rolnick. Complexity of linear regions in deep networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2596–2604. PMLR, 09–15 Jun 2019a. URL <https://proceedings.mlr.press/v97/hanin19a.html>.
- Boris Hanin and David Rolnick. Deep relu networks have surprisingly few activation patterns. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019b. URL <https://proceedings.neurips.cc/paper/2019/file/9766527f2b5d3e95d4a733fcfb77bd7e-Paper.pdf>.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.

-
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf>.
- Andrei N Kolmogorov. Three approaches to the quantitative definition of information'. *Problems of information transmission*, 1(1):1–7, 1965.
- Vladimir Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*, 47(5):1902–1914, 2001.
- Ming Li, Paul Vitányi, et al. *An introduction to Kolmogorov complexity and its applications*, volume 3. Springer, 2008.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable architecture search. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=S1eYHoC5FX>.
- Zhengying Liu, Zhen Xu, Shangeth Rajaa, Meysam Madadi, Julio CS Jacques Junior, Sergio Escalera, Adrien Pavao, Sebastien Treguer, Wei-Wei Tu, and Isabelle Guyon. Towards automated deep learning: Analysis of the autodl challenge series 2019. In *NeurIPS 2019 Competition and Demonstration Track*, pp. 242–252. PMLR, 2020.
- Christos Louizos, Karen Ullrich, and Max Welling. Bayesian compression for deep learning. *Advances in neural information processing systems*, 30, 2017.
- Kaitlin Maile, Emmanuel Rachelson, Hervé Luga, and Dennis George Wilson. When, where, and how to add new neurons to ANNs. In *First Conference on Automated Machine Learning (Main Track)*, 2022. URL <https://openreview.net/forum?id=SW0g-arIg9>.
- Hector Mendoza, Aaron Klein, Matthias Feurer, Jost Tobias Springenberg, and Frank Hutter. Towards automatically-tuned neural networks. In Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren (eds.), *Proceedings of the Workshop on Automatic Machine Learning*, volume 64 of *Proceedings of Machine Learning Research*, pp. 58–65, New York, New York, USA, 24 Jun 2016. PMLR. URL https://proceedings.mlr.press/v64/mendoza_towards_2016.html.
- Risto Miiikkulainen, Jason Zhi Liang, Elliot Meyerson, Aditya Rawal, Daniel Fink, Olivier Francon, Bala Raju, Hormoz Shahrzad, Arshak Navruzian, Nigel Duffy, and Babak Hodjat. Evolving deep neural networks. *CoRR*, abs/1703.00548, 2017. URL <http://arxiv.org/abs/1703.00548>.
- Geoffrey F. Miller, Peter M. Todd, and Shailesh U. Hegde. Designing neural networks using genetic algorithms. In *ICGA*, 1989.
- Yann Ollivier. True asymptotic natural gradient optimization, 2017.
- Allan Pinkus. Approximation theory of the mlp model in neural networks. *ACTA NUMERICA*, 8:143–195, 1999.
- 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. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/raghu17a.html>.
- Thiago Serra, Christian Tjandraatmadja, and Srikumar Ramalingam. Bounding and counting linear regions of deep 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. 4558–4566. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/serra18b.html>.

-
- Kenneth O Stanley, David B D'Ambrosio, and Jason Gauci. A hypercube-based encoding for evolving large-scale neural networks. *Artificial life*, 15(2):185–212, 2009.
- Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. *CoRR*, abs/1503.02406, 2015. URL <http://dblp.uni-trier.de/db/journals/corr/corr1503.html#TishbyZ15>.
- V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971. doi: 10.1137/1116025. URL <http://link.aip.org/link/?TPR/16/264/1>.
- P. Wolinski, G. Charpiat, and O. Ollivier. Asymmetrical scaling layers for stable network pruning. *OpenReview Archive*, 2020.
- Lemeng Wu, Dilin Wang, and Qiang Liu. Splitting steepest descent for growing neural architectures. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/3a01fc0853eb94fde4d1cc6fb842a-Paper.pdf>.
- Lemeng Wu, Bo Liu, Peter Stone, and Qiang Liu. Firefly neural architecture descent: a general approach for growing neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 22373–22383. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/fdbe012e2e11314b96402b32c0df26b7-Paper.pdf>.
- Tien-Ju Yang, Andrew Howard, Bo Chen, Xiao Zhang, Alec Go, Mark Sandler, Vivienne Sze, and Hartwig Adam. Netadapt: Platform-aware neural network adaptation for mobile applications. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 285–300, 2018.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning, 2016. URL <http://arxiv.org/abs/1611.01578>. cite arxiv:1611.01578.

Appendix outline

- Appendix A details the theoretical approach of TINY.
- Appendix B compares the theoretical approach of TINY with other approaches.
- Appendix C proves the propositions of the paper.
- Appendix E provides the hyper parameters for learning.
- Appendix F gives additional graphics associated to the result part.

For part B and C we use the trace scalar product and its associated norm. We note $\langle \cdot, \cdot \rangle := \langle \cdot, \cdot \rangle_{\text{Tr}}$ and $\| \cdot \| := \| \cdot \|_{\text{Tr}}$. One should remark that $\| \cdot \| = \| \cdot \|_{\text{Tr}} = \| \cdot \|_2$.

A Theoretical details for part 2

A.1 Functional gradient

The functional loss \mathcal{L} is a functional that takes as input a function $f \in \mathcal{F}$ and outputs a real score:

$$\mathcal{L} : f \in \mathcal{F} \mapsto \mathcal{L}(f) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\ell(f(\mathbf{x}), \mathbf{y})] \in \mathbb{R}.$$

The function space \mathcal{F} can typically be chosen to be $L_2(\mathbb{R}^p \rightarrow \mathbb{R}^d)$, which is a Hilbert space. The directional derivative (or Gateaux derivative, or Fréchet derivative) of functional \mathcal{L} at function f in direction v is defined as:

$$D\mathcal{L}(f)(v) = \lim_{\varepsilon \rightarrow 0} \frac{\mathcal{L}(f + \varepsilon v) - \mathcal{L}(f)}{\varepsilon}$$

if it exists. Here v denotes any function in the Hilbert space \mathcal{F} and stands for the direction in which we would like to update f , following an infinitesimal step (of size ε), resulting in a function $f + \varepsilon v$.

If this directional derivative exists in all possible directions $v \in \mathcal{F}$ and moreover is continuous in v , then the Riesz representation theorem implies that there exists a unique direction $v^* \in \mathcal{F}$ such that:

$$\forall v \in \mathcal{F}, \quad D\mathcal{L}(f)(v) = \langle v^*, v \rangle.$$

This direction v^* is named the *gradient* of the functional \mathcal{L} at function f and is denoted by $\nabla_f \mathcal{L}(f)$.

Note that while the inner product $\langle \cdot, \cdot \rangle$ considered is usually the L_2 one, it is possible to consider other ones, such as Sobolev ones (e.g., H^1). The gradient $\nabla_f \mathcal{L}(f)$ depends on the chosen inner product and should consequently rather be denoted by $\nabla_f^{L_2} \mathcal{L}(f)$ for instance.

Note that continuous functions from \mathbb{R}^p to \mathbb{R}^d , as well as C^∞ functions, are dense in $L_2(\mathbb{R}^p \rightarrow \mathbb{R}^d)$.

Let us now study properties specific to our loss design: $\mathcal{L}(f) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\ell(f(\mathbf{x}), \mathbf{y})]$. Assuming sufficient ℓ -loss differentiability and integrability, we get, for any function update direction $v \in \mathcal{F}$ and infinitesimal step size $\varepsilon \in \mathbb{R}$:

$$\begin{aligned} \mathcal{L}(f + \varepsilon v) - \mathcal{L}(f) &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\ell(f(\mathbf{x}) + \varepsilon v(\mathbf{x}), \mathbf{y}) - \ell(f(\mathbf{x}), \mathbf{y})] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[\nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y})|_{\mathbf{u}=f(\mathbf{x})} \cdot \varepsilon v(\mathbf{x}) + O(\varepsilon^2 \|v(\mathbf{x})\|^2) \right] \end{aligned}$$

using the usual gradient of function ℓ at point $(\mathbf{u} = f(\mathbf{x}), \mathbf{y})$ w.r.t. its first argument \mathbf{u} , with the standard Euclidean dot product \cdot in \mathbb{R}^p . Then the directional derivative is:

$$D\mathcal{L}(f)(v) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[\nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y})|_{\mathbf{u}=f(\mathbf{x})} \cdot v(\mathbf{x}) \right] = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\mathbb{E}_{\mathbf{y} \sim \mathcal{D}|\mathbf{x}} \left[\nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y})|_{\mathbf{u}=f(\mathbf{x})} \right] \cdot v(\mathbf{x}) \right]$$

and thus the functional gradient for the inner product $\langle v, v' \rangle_{\mathbb{E}} := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [v(\mathbf{x}) \cdot v'(\mathbf{x})]$ is the function:

$$\nabla_f^{\mathbb{E}} \mathcal{L}(f) : \mathbf{x} \mapsto \mathbb{E}_{\mathbf{y} \sim \mathcal{D}|\mathbf{x}} \left[\nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y})|_{\mathbf{u}=f(\mathbf{x})} \right]$$

which simplifies into:

$$\nabla_f^{\mathbb{E}} \mathcal{L}(f) : \mathbf{x} \mapsto \nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y}(\mathbf{x}))|_{\mathbf{u}=f(\mathbf{x})}$$

if there is no ambiguity in the dataset, i.e. if for each \mathbf{x} there is a unique $\mathbf{y}(\mathbf{x})$.

Note that by considering the $L_2(\mathbb{R}^p \rightarrow \mathbb{R}^d)$ inner product $\int v \cdot v'$ instead, one would respectively get:

$$\nabla_f^{L_2} \mathcal{L}(f) : \mathbf{x} \mapsto p_{\mathcal{D}}(\mathbf{x}) \mathbb{E}_{\mathbf{y} \sim \mathcal{D}|\mathbf{x}} \left[\nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y})|_{\mathbf{u}=f(\mathbf{x})} \right]$$

and

$$\nabla_f^{L_2} \mathcal{L}(f) : \mathbf{x} \mapsto p_{\mathcal{D}}(\mathbf{x}) \nabla_{\mathbf{u}} \ell(\mathbf{u}, \mathbf{y}(\mathbf{x}))|_{\mathbf{u}=f(\mathbf{x})}$$

instead, where $p_{\mathcal{D}}(\mathbf{x})$ is the density of the dataset distribution at point \mathbf{x} . In practice one estimates such gradients using a minibatch of samples (\mathbf{x}, \mathbf{y}) , obtained by picking uniformly at random within a finite dataset, and thus the formulas for the two inner products coincide (up to a constant factor).

A.2 Differentiation under the integral sign

Let X be an open subset of \mathbb{R} , and Ω be a measure space. Suppose $f : X \times \Omega \rightarrow \mathbb{R}$ satisfies the following conditions:

- $f(x, \omega)$ is a Lebesgue-integrable function of ω for each $x \in X$.
- For almost all $\omega \in \Omega$, the partial derivative $\frac{\partial}{\partial x} f$ of f according to x exists for all $x \in X$.
- There is an integrable function $\theta : \Omega \rightarrow \mathbb{R}$ such that $|\frac{\partial}{\partial x} f(x, \omega)| \leq \theta(\omega)$ for all $x \in X$ and almost every $\omega \in \Omega$.

Then, for all $x \in X$,

$$\frac{\partial}{\partial x} \int_{\Omega} f(x, \omega) d\omega = \int_{\Omega} \frac{\partial}{\partial x} f(x, \omega) d\omega \quad (15)$$

See proof and details :Flanders (1973).

A.3 Gradients and proximal point of view

Gradients with respect to standard variables such as vectors are defined the same way as functional gradients above: given a sufficiently smooth loss $\tilde{\mathcal{L}} : \theta \in \Theta_{\mathcal{A}} \mapsto \tilde{\mathcal{L}}(\theta) = \mathcal{L}(f_{\theta}) \in \mathbb{R}$, and an inner product \cdot in the space $\Theta_{\mathcal{A}}$ of parameters θ , the gradient $\nabla_{\theta} \tilde{\mathcal{L}}(\theta)$ is the unique vector $\boldsymbol{\tau} \in \Theta_{\mathcal{A}}$ such that:

$$\forall \delta\theta \in \Theta_{\mathcal{A}}, \quad \boldsymbol{\tau} \cdot \delta\theta = D_{\theta} \tilde{\mathcal{L}}(\theta)(\delta\theta)$$

where $D_{\theta} \tilde{\mathcal{L}}(\theta)(\delta\theta)$ is the directional derivative of $\tilde{\mathcal{L}}$ at point θ in the direction $\delta\theta$, defined as in the previous section. This gradient depends on the inner product chosen, which can be highlighted by the following property. The opposite $-\nabla_{\theta} \tilde{\mathcal{L}}(\theta)$ of the gradient is the unique solution of the problem:

$$\arg \min_{\delta\theta \in \Theta_{\mathcal{A}}} \left\{ D_{\theta} \tilde{\mathcal{L}}(\theta)(\delta\theta) + \frac{1}{2} \|\delta\theta\|_P^2 \right\}$$

where $\|\cdot\|_P$ is the norm associated to the chosen inner product. Changing the inner product obviously changes the way candidate directions $\delta\theta$ are penalized, leading to different gradients. This proximal formulation can be obtained as follows. For any $\delta\theta$, its distance to the gradient descent direction is:

$$\begin{aligned} \left\| \delta\theta - \left(-\nabla_{\theta} \tilde{\mathcal{L}}(\theta) \right) \right\|^2 &= \|\delta\theta\|^2 + 2 \delta\theta \cdot \nabla_{\theta} \tilde{\mathcal{L}}(\theta) + \left\| \nabla_{\theta} \tilde{\mathcal{L}}(\theta) \right\|^2 \\ &= 2 \left(\frac{1}{2} \|\delta\theta\|^2 + D_{\theta} \tilde{\mathcal{L}}(\theta)(\delta\theta) \right) + K \end{aligned}$$

where K does not depend on $\delta\theta$. For the above to hold, the inner product used has to be the one from which the norm is derived. By minimizing this expression with respect to $\delta\theta$, one obtains the desired property.

In our case of study, for the norm over the space $\Theta_{\mathcal{A}}$ of parameter variations, we consider a norm in the space of associated functional variations, i.e.:

$$\|\delta\theta\|_P := \left\| \frac{\partial f_{\theta}}{\partial \theta} \delta\theta \right\|$$

which makes more sense from a physical point of view, as it is more intrinsic to the task to solve and depends as little as possible on the parameterization (i.e. on the architecture chosen). This results in a functional move that is the projection of the functional one to the set of possible moves given the architecture. On the opposite, the standard gradient (using Euclidean parameter norm $\|\delta\theta\|$ in parameter space) yields a functional move obtained not only by projecting the functional gradient but also by multiplying it by a matrix $\frac{\partial f_\theta}{\partial \theta} \frac{\partial f_\theta}{\partial \theta}^T$ which can be seen as a strong architecture bias over optimization directions.

We consider here that the loss \mathcal{L} to be minimized is the real loss that the user wants to optimize, possibly including regularizers to avoid overfitting, and since the architecture is evolving during training, possibly to architectures far from usual manual design and never tested before, one cannot assume architecture bias to be desirable. We aim at getting rid of it in order to follow the functional gradient descent as closely as possible.

Searching for

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathcal{T}_A} \|\mathbf{v} - \mathbf{v}_{\text{goal}}\|^2 = \arg \min_{\mathbf{v} \in \mathcal{T}_A} \left\{ D\mathcal{L}(f)(\mathbf{v}) + \frac{1}{2} \|\mathbf{v}\|^2 \right\} \quad (16)$$

or equivalently for:

$$\delta\theta^* = \arg \min_{\delta\theta \in \Theta_A} \left\| \frac{\partial f_\theta}{\partial \theta} \delta\theta - \mathbf{v}_{\text{goal}} \right\|^2 = \arg \min_{\delta\theta \in \Theta_A} \left\{ D_\theta \mathcal{L}(f_\theta)(\delta\theta) + \frac{1}{2} \left\| \frac{\partial f_\theta}{\partial \theta} \delta\theta \right\|^2 \right\} =: -\nabla_\theta^{\mathcal{T}_A} \mathcal{L}(f_\theta) \quad (17)$$

then appears as a natural goal.

A.4 Example of expressivity bottleneck

Example. Suppose one tries to estimate the function $y = f_{\text{true}}(x) = 2\sin(x) + x$ with a linear model $f_{\text{predict}}(x) = ax + b$. Consider $(a, b) = (1, 0)$ and the square loss \mathcal{L} . For the dataset of inputs $(x_0, x_1, x_2, x_3) = (0, \frac{\pi}{2}, \pi, \frac{3\pi}{2})$, there exists no parameter update $(\delta a, \delta b)$ that would improve prediction at x_0, x_1, x_2 and x_3 simultaneously, as the space of linear functions $\{f : x \rightarrow ax + b \mid a, b \in \mathbb{R}\}$ is not expressive enough. To improve the prediction at x_0, x_1, x_2 **and** x_3 , one should look for another, more expressive functional space such that for $i = 0, 1, 2, 3$ the functional update $\Delta f(x_i) := f^{t+1}(x_i) - f^t(x_i)$ goes into the same direction as the functional gradient $\mathbf{v}_{\text{goal}}(x_i) := -\nabla_{f(x_i)} \mathcal{L}(f(x_i), y_i) = -2(f(x_i) - y_i)$ where $y_i = f_{\text{true}}(x_i)$.

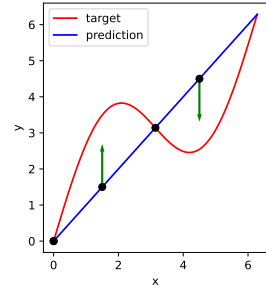


Figure 9: Linear interpolation

A.5 Problem formulation and choice of pre-activities

There are several ways to design the problem of adding neurons, which we discuss now, in order to explain our choice of the pre-activities to express expressivity bottlenecks.

Suppose one wishes to add K neurons $\theta_{\leftrightarrow}^K := (\alpha_k, \omega_k)_{k=1}^K$ to layer $l-1$, which impacts the activities \mathbf{a}_l at the next layer, in order to improve its expressivity. These neurons could be chosen to have only nul weights, or nul input weights α_k and non-nul output weights ω_k , or the opposite, or both non-nul weights. Searching for the best neurons to add for each of these cases will produce different optimization problems.

Let us remind first that adding such K neurons with weights $\theta_{\leftrightarrow}^K := (\alpha_k, \omega_k)_{k=1}^K$ changes the activities \mathbf{a}_l of the (next) layer by

$$\delta \mathbf{a}_l = \sum_{k=1}^K \omega_k \sigma(\alpha_k^T \mathbf{b}_{l-2}(x)) \quad (18)$$

Small weights approximation Under the hypothesis of small input weights α_k , the activity variation 18 can be approximated by:

$$\sigma'(0) \sum_{k=1}^K \omega_k \alpha_k^T \mathbf{b}_{l-2}(\mathbf{x}) \quad (19)$$

at first order in $\|\alpha_k\|$. We will drop the constant $\sigma'(0)$ in the sequel.

This quantity is linear both in α_k and ω_k , therefore the first-order parameter-induced activity variations are easy to compute:

$$\begin{aligned} \mathbf{v}^l(\mathbf{x}, (\alpha_k)_{k=1}^K) &= \frac{\partial \mathbf{a}_l(\mathbf{x})}{\partial ((\alpha_k)_{k=1}^K)} \Big|_{(\alpha_k)_{k=1}^K=0} (\alpha_k)_{k=1}^K = \sum_{k=1}^K \omega_k \mathbf{b}_{l-2}(\mathbf{x})^T \alpha_k \\ \mathbf{v}^l(\mathbf{x}, (\omega_k)_{k=1}^K) &= \frac{\partial \mathbf{a}_l(\mathbf{x})}{\partial ((\omega_k)_{k=1}^K)} \Big|_{(\omega_k)_{k=1}^K=0} (\omega_k)_{k=1}^K = \sum_{k=1}^K \omega_k \mathbf{b}_{l-2}(\mathbf{x})^T \alpha_k \end{aligned}$$

so with a slight abuse of notation we have:

$$\mathbf{v}^l(\mathbf{x}, \theta_{\leftrightarrow}^K) = \sum_{k=1}^K \omega_k \alpha_k^T \mathbf{b}_{l-2}(\mathbf{x})$$

Note also that technically the quantity above is first-order in α_k and in ω_k but second-order in the joint variable $\theta_{\leftrightarrow}^K = (\alpha_k, \omega_k)$.

Adding neurons with 0 weights (both input and output weights). In that case, one increases the number of neurons in the layer, but without changing the function (since only nul quantities are added) and also without changing the gradient with respect to the parameters, thus not improving expressivity. Indeed, the added quantity (Eq. 18) involves 0×0 multiplications, and consequently the derivative $\frac{\partial \mathbf{a}_l(\mathbf{x})}{\partial \theta_{\leftrightarrow}^K} \Big|_{\theta_{\leftrightarrow}^K=0}$ w.r.t. these new parameters, that is, $\mathbf{b}_{l-2}(\mathbf{x})^T \alpha_k$ w.r.t. ω_k and $\omega_k \mathbf{b}_{l-2}(\mathbf{x})^T$ w.r.t. α_k is 0, as both α_k and ω_k are 0.

Adding neurons with non-0 input weights and 0 output weights or the opposite. In these cases, the addition of neurons will not change the function (because of multiplications by 0), but just the gradient. One of the 2 gradients (w.r.t. α_k or w.r.t. ω_k) will be non-0, as the variable that is 0 has non-0 derivatives.

The question is then how to pick the best non-0 variable, (α_k or ω_k) such that the added gradient will be the most useful. The problem can then be formulated similarly as what is done in the paper.

Adding neurons with small yet non-0 weights. In this case, both the function and its gradient will change when adding the neurons. Fortunately, Proposition 3.2 states that the best neurons to add in terms of expressivity (to get the gradient closer to the variation desired by the backpropagation) are also the best neurons to add to decrease the loss, i.e. the function change they will imply goes into the right direction.

For each family $(\omega_k)_{k=1}^K$, the tangent space in \mathbf{a}_l restricted to the family $(\alpha_k)_{k=1}^K$, ie $\mathcal{T}_{\mathcal{A}}^{\mathbf{a}_l} := \left\{ \frac{\partial \mathbf{a}_l}{\partial (\alpha_k)_{k=1}^K} (\cdot) (\alpha_k)_{k=1}^K \Big| (\alpha_k)_{k=1}^K \in (\mathbb{R}^{|\mathbf{b}_{l-2}(\mathbf{x})|})^K \right\}$ varies with the family $(\omega_k)_{k=1}^K$, ie $\mathcal{T}_{\mathcal{A}}^{\mathbf{a}_l} := \mathcal{T}_{\mathcal{A}}^{\mathbf{a}_l}((\omega_k)_{k=1}^K)$. Optimizing w.r.t. the ω_k is equivalent to search for the best tangent space for the α_k , while symmetrically optimizing w.r.t. the α_k is equivalent to find the best projection on the tangent space defined by the ω_k .

Pre-activities vs. post-activities. The space of pre-activities \mathbf{a}_l is a natural space for this framework, as they are formed with linear operations and we compute first-order variation quantities. Considering the space of post-activities $\mathbf{b}_l = \sigma(\mathbf{a}_l)$ is also possible, though computing variations will be more complex. Indeed, without first-order approximation, the obtained problem is not manageable, because of the non-linear activation function σ added in front of all quantities (while in the case of pre-activations, quantity 18

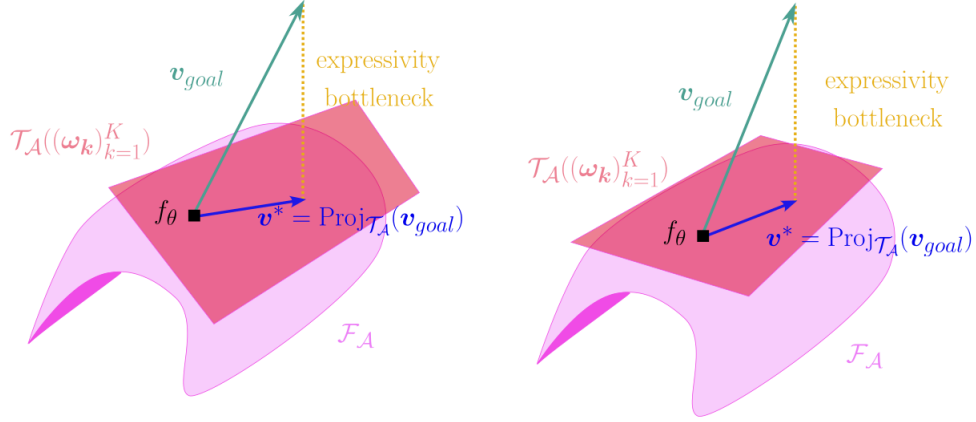


Figure 10: Changing the tangent space with different values of $(\omega_k)_{k=1}^K$.

is linear in ω_k and thus does not require approximation in ω_k , which allow considering large ω_k), and, with first-order approximation, it would add the derivative of the activation function, taken at various locations $\sigma'(\mathbf{a}_l)$ (while in the previous case the derivatives of the activation function were always taken at 0).

A.6 Adding convolutional neurons

To add a convolutional neuron at layer $l-1$, one should add a kernel at layer $l-1$ and expand one dimension to all the kernels in layer l to match the new dimension of the post-activity.

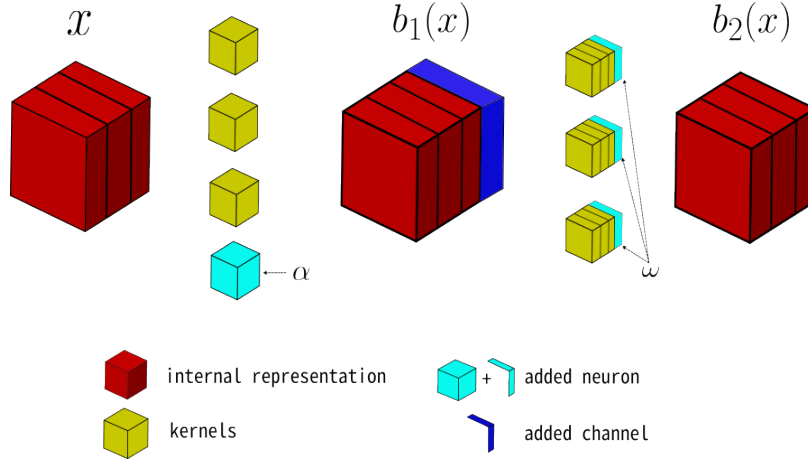


Figure 11: Adding one convolutional neuron at layer one for a input with tree channels.

B Theoretical comparison with other approaches

B.1 GradMax method

To facilitate reading we remove the layer index of each quantity, ie $\mathbf{b} := \mathbf{b}_{l-2}$, $\mathbf{B} := \mathbf{B}_{l-2}$, $\mathbf{V}_{\text{goal}} := \mathbf{V}_{\text{goal}}^l$ and $\mathbf{V}_{\text{goal}_{proj}} := \mathbf{V}_{\text{goal}_{proj}}^l$.

The theoretical approach of GradMax is to add neurons with zero fan-in and choose the fan-out which would decrease the loss as much as possible after one gradient step. We note $\mathbf{\Omega}$ the fan-out of such neurons and perform the addition at time t at layer l . After one gradient step, ie $t \rightarrow t + 1$, the decrease of loss is :

$$\mathcal{L}^{t+1} \approx \mathcal{L}^t - \|\nabla_{\theta} \mathcal{L}\|^2 - \|\nabla_{\Omega} \mathcal{L}\|^2$$

The solution of GradMax as formulated in the paper Evci et al. (2022) is :

$$(\omega_1^*, \dots, \omega_K^*) := \mathbf{\Omega}^* = \arg \max_{\mathbf{\Omega}} \|\nabla_{\Omega} \mathcal{L}\|^2 \quad s.t. \quad \|\mathbf{\Omega}\|^2 \leq c \quad (20)$$

we remark that :

$$\|\nabla_{\Omega} \mathcal{L}\|^2 = \left\| \sum_i b(\mathbf{x}_i) \mathbf{v}_{\text{goal}}^T(\mathbf{x}_i) \mathbf{\Omega} \right\|^2 \quad (21)$$

$$= \left\| \mathbf{B} \mathbf{V}_{\text{goal}}^T \mathbf{\Omega} \right\|^2 \quad (22)$$

$$= \left\| \tilde{\mathbf{N}} \mathbf{\Omega} \right\|^2 \quad \tilde{\mathbf{N}} := \mathbf{B} \mathbf{V}_{\text{goal}}^T \quad (23)$$

It follows that the fan-out of the neurons computed by GradMax are the solution of the problem that we shall prove later:

$$\mathbf{\Omega}^* := \arg \max_{\mathbf{\Omega}} \left\| \tilde{\mathbf{N}} \mathbf{\Omega} \right\| \quad s.t. \quad \|\mathbf{\Omega}\|^2 \leq c \quad (24)$$

To compare this optimization problem with TINY, we use the following proposition:

Proposition B.1.

$$\forall \mathbf{a} \in \mathbb{R}^d, \mathbf{M} \in \mathbb{R}^{j,d}, \exists c \in \mathbb{R} \text{ s.t. } \arg \min_{\mathbf{b}} \|\mathbf{a} - \mathbf{M}\mathbf{b}\|^2 = \arg \max_{\mathbf{b}, \|\mathbf{M}\mathbf{b}\|^2 \leq c} \langle \mathbf{a}, \mathbf{M}\mathbf{b} \rangle$$

Taking \mathbf{V}_{goal} as \mathbf{a} and \mathbf{V} as $\mathbf{M}\mathbf{b}$, we can reformulate TINY optimization problem 8 as :

$$\mathbf{A}^*, \mathbf{\Omega}^* = \arg \max_{\mathbf{A}, \mathbf{\Omega}} \langle \mathbf{V}(\mathbf{A}, \mathbf{\Omega}), \mathbf{V}_{\text{goal}_{proj}} \rangle \quad s.t. \quad \|\mathbf{V}(\mathbf{A}, \mathbf{\Omega})\|^2 \leq c \quad (25)$$

We remark that

$$\langle \mathbf{V}(\mathbf{A}, \mathbf{\Omega}), \mathbf{V}_{\text{goal}_{proj}} \rangle = \langle \mathbf{\Omega} \mathbf{A}^T \mathbf{B}, \mathbf{V}_{\text{goal}_{proj}} \rangle \quad (26)$$

$$= \text{Tr}(\mathbf{B}^T \mathbf{A} \mathbf{\Omega}^T \mathbf{V}_{\text{goal}_{proj}}) \quad (27)$$

$$= \text{Tr}(\mathbf{A} \mathbf{\Omega}^T \mathbf{V}_{\text{goal}_{proj}} \mathbf{B}^T) \quad (28)$$

$$= \langle \mathbf{\Omega} \mathbf{A}^T, \mathbf{V}_{\text{goal}_{proj}} \mathbf{B}^T \rangle \quad (29)$$

We perform the change of variable $\tilde{\mathbf{A}} := \mathbf{S}^{\frac{1}{2}} \mathbf{A}$ and re-write the constrain as :

$$\|\mathbf{V}(\mathbf{A}, \mathbf{\Omega})\|^2 = \|\mathbf{\Omega} \mathbf{A}^T \mathbf{B}\|^2 \quad (30)$$

$$= \text{Tr}(\mathbf{A} \mathbf{\Omega}^T \mathbf{\Omega} \mathbf{A}^T \mathbf{S}) \quad \mathbf{S} = \mathbf{B} \mathbf{B}^T \quad (31)$$

$$= \left\| \mathbf{\Omega} (\mathbf{S}^{\frac{1}{2}} \mathbf{A})^T \right\|^2 \quad (32)$$

$$= \left\| \mathbf{\Omega} \tilde{\mathbf{A}}^T \right\|^2 \quad (33)$$

Then we define $\mathbf{N} := \mathbf{B} \mathbf{V}_{\text{goal}_{proj}}^T$. The initial scalar product is then :

$$\langle \mathbf{V}(\mathbf{A}, \mathbf{\Omega}), \mathbf{V}_{\text{goal}_{proj}} \rangle = \langle \mathbf{\Omega} \tilde{\mathbf{A}}^T \mathbf{S}^{-\frac{1}{2}}, \mathbf{N}^T \rangle = \langle \mathbf{\Omega} \tilde{\mathbf{A}}^T, \mathbf{N}^T \mathbf{S}^{-\frac{1}{2}} \rangle = \langle \tilde{\mathbf{A}} \mathbf{\Omega}^T, \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \rangle \quad (34)$$

To maximise the scalar product, we choose $\tilde{\mathbf{A}}\mathbf{\Omega}^T = \mathbf{S}^{-\frac{1}{2}}\mathbf{N}$. A solution for $(\tilde{\mathbf{A}}, \mathbf{\Omega})$ is the (left, right) eigenvectors of the matrix $\mathbf{S}^{-\frac{1}{2}}\mathbf{N}$. It implies that :

$$\mathbf{\Omega}^* := \arg \max_{\mathbf{\Omega}} \left\| \mathbf{S}^{-\frac{1}{2}}\mathbf{N}\mathbf{\Omega} \right\|^2 \quad s.t. \|\mathbf{\Omega}\| \leq \tilde{c} \quad (35)$$

One can note three differences between those optimization problems:

- First, the matrix $\tilde{\mathbf{N}}$ is not defined using the projection of the desired update $\mathbf{V}_{\text{goal}_{proj}}^{l+1}$. As a consequence, GradMax does not take into account redundancy, and on the opposite will actually try to add new neurons that are as redundant as possible with the part of the goal update that is already feasible with already-existing neurons.
- Second, the constraint lies in the weight space for GradMax method while it lies in the pre-activation space in our case. The difference is that GradMax relies on the Euclidean metric in the space of parameters, which arguably offers less meaning than the Euclidean metric in the space of activities. Essentially this is the same difference as between the standard L2 gradient w.r.t. parameters and the natural gradient, which takes care of parameter redundancy and measures all quantities in the output space in order to be independent from the parameterization. In practice we do observe that the "natural" gradient direction improves the loss better than the usual L2 gradient.
- Third, our fan-in weights are not set to 0 but directly to their optimal values (at first order).

We now prove the proposition B.1

Proof. Indeed,

$$\arg \min_{\mathbf{b}} \|\mathbf{a} - \mathbf{M}\mathbf{b}\|^2 = \arg \min_{\mathbf{b}} \|\mathbf{M}\mathbf{b}\|^2 - \langle \mathbf{a}, \mathbf{M}\mathbf{b} \rangle \quad (36)$$

$$= \arg \min_{\mathbf{b}} \|\mathbf{M}\mathbf{b}\| \left(\|\mathbf{M}\mathbf{b}\| - \left\langle \mathbf{a}, \frac{\mathbf{M}\mathbf{b}}{\|\mathbf{M}\mathbf{b}\|} \right\rangle \right) \quad (37)$$

$$= \arg \min_c \arg \min_{\|\mathbf{M}\mathbf{b}\|=c, \mathbf{u}=\frac{\mathbf{M}\mathbf{b}}{\|\mathbf{M}\mathbf{b}\|}} c - \langle \mathbf{a}, \mathbf{u} \rangle \quad (38)$$

With the convention that $\frac{0}{\|0\|} = 0$. □

B.2 NORTH Preactivation

In paper Maile et al. (2022), fan-out weights are initialized to 0 while fan-in weights are initialized as $\alpha_i = \mathbf{S}^{-1}\mathbf{B}_{l-1}\mathbf{V}_{\mathbf{Z}_{l-1}}^T r_i$ where r_i is a random vector and $\mathbf{V}_{\mathbf{Z}_{l-1}} \in \mathcal{M}(|\text{Ker}(\mathbf{B}_{l-1}^T)|, |\mathbf{b}_{l-2}(\mathbf{x})|)$ is a matrix consisting of the orthogonal vectors of the kernel of pre-activations \mathbf{B}_{l-1} , i.e $\{z \mid \mathbf{B}_{l-1}^T z = 0\}$. In our paper fan-in weights are initialized as $\alpha_i = \mathbf{S}^{-1}\mathbf{B}_{l-2}\mathbf{V}_{\text{goal}_{proj}}^T \mathbf{v}_i = \mathbf{S}^{-1}\mathbf{B}_{l-2}\mathbf{V}_{\text{goal}}^T \mathbf{V}_{\mathbf{Z}_{l-1}} \mathbf{V}_{\mathbf{Z}_{l-1}}^T \mathbf{v}_i$, where the \mathbf{v}_i are right eigenvectors of the matrix $\mathbf{S}^{-\frac{1}{2}}\mathbf{N}$.

The main difference is thus that we use the backpropagation to find the best \mathbf{v}_i or r_i directly, while the NORTH approach tries random directions r_i to explore the space of possible neuron additions.

C Proofs of Part 3

C.1 Proposition 3.1

Denoting by \mathbf{M}^+ the generalized (pseudo-)inverse of \mathbf{M} , we have:

$$\delta \mathbf{W}_l^* = \frac{1}{n} \mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T \left(\frac{1}{n} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T \right)^+ \text{ and } \mathbf{V}_0^l = \frac{1}{n} \mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T \left(\frac{1}{n} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T \right)^+ \mathbf{B}_{l-1}$$

Proof

Consider the function

$$g(\delta \mathbf{W}) = \|\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1}\|^2 \quad (39)$$

then:

$$g(\delta \mathbf{W} + \mathbf{H}) = \|\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1} - \mathbf{H} \mathbf{B}_{l-1}\|^2 \quad (40)$$

$$= g(\delta \mathbf{W}) - 2\langle \mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1}, \mathbf{H} \mathbf{B}_{l-1} \rangle + o(\|\mathbf{H}\|) \quad (41)$$

$$= g(\delta \mathbf{W}) - 2 \text{Tr}((\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1})^T \mathbf{H} \mathbf{B}_{l-1}) + o(\|\mathbf{H}\|) \quad (42)$$

$$= g(\delta \mathbf{W}) - 2 \text{Tr}(\mathbf{B}_{l-1} (\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1})^T \mathbf{H}) + o(\|\mathbf{H}\|) \quad (43)$$

$$= g(\delta \mathbf{W}) - 2\langle (\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1}) \mathbf{B}_{l-1}^T, \mathbf{H} \rangle + o(\|\mathbf{H}\|) \quad (44)$$

By identification $\nabla_{\delta \mathbf{W}} g(\delta \mathbf{W}) = -2(\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W} \mathbf{B}_{l-1}) \mathbf{B}_{l-1}^T$, and thus:

$$\nabla_{\delta \mathbf{W}} g(\delta \mathbf{W}) = 0 \implies \mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T = \delta \mathbf{W} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T$$

Using that g is convex and the definition of the generalized inverse, we get:

$$\delta \mathbf{W}_l^* = \frac{1}{n} \mathbf{V}_{\text{goal}}^l \mathbf{B}_{l-1}^T \left(\frac{1}{n} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T \right)^+$$

as one solution. **For convolutional layer**, we defined as b_i^c the input associated to the activation $a_l(X_i) \in \mathbb{R}^{k,p,p}$, such that for a convolution layer with one output channel, noted with parameter \mathbf{W} , we have :

$$\text{Conv}(a_l(X_i)) = \mathbf{B}_i^c \text{vect}(\mathbf{W}) \quad (45)$$

Example : considering the kernel of Conv to be $(2, 2)$, then :

$$\mathbf{b}_i^k = \begin{pmatrix} b_i^{1,k} & b_i^{2,k} & \cdot & b_i^{p,k} \\ b_i^{p+1,k} & b_i^{p+2,k} & \cdot & b_i^{2p,k} \\ \cdot & \cdot & \cdot & \cdot \\ b_i^{p(p-1)+1,k} & \cdot & \cdot & b_i^{p^2,k} \end{pmatrix} \quad (46)$$

$$\mathbf{B}_i^c = \begin{pmatrix} b_i^{1,1} & b_i^{2,1} & b_i^{p+1,1} & b_i^{p+2,1} & b_i^{1,2} & b_i^{2,2} & b_i^{p+1,2} & b_i^{p+2,2} & b_i^{1,3} & \cdot \\ b_i^{2,1} & b_i^{3,1} & b_i^{p+2,1} & b_i^{p+3,1} & b_i^{2,2} & b_i^{3,2} & b_i^{p+2,2} & b_i^{p+3,2} & b_i^{2,3} & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}$$

Then the function to minimize is

$$g(\delta \mathbf{W}) = \|\mathbf{V}_{\text{goal}}^l - \mathbf{B}^c \delta \mathbf{W}\|^2$$

where $\mathbf{B}^c := (\mathbf{B}_1^c \quad \dots \quad \mathbf{B}_n^c)$ and $\delta \mathbf{W}$ is the concatenation over the output channels. \square

C.2 Proposition 3.2

We define the matrices $\mathbf{N} := \frac{1}{n} \mathbf{B}_{l-2} (\mathbf{V}_{\text{goal}}^l \text{proj})^T$ and $\mathbf{S} := \frac{1}{n} \mathbf{B}_{l-2} \mathbf{B}_{l-2}^T$. Let us denote its SVD by $\mathbf{S} = \mathbf{U} \Sigma \mathbf{U}^T$, and note $\mathbf{S}^{-\frac{1}{2}} := \mathbf{U} \sqrt{\Sigma}^{-1} \mathbf{U}^T$ and consider the SVD of the matrix $\mathbf{S}^{-\frac{1}{2}} \mathbf{N} = \sum_{k=1}^R \lambda_k \mathbf{u}_k \mathbf{v}_k^T$ with $\lambda_1 \geq \dots \geq \lambda_R \geq 0$, where R is the rank of the matrix \mathbf{N} . Then:

Proposition C.1 (3.2). *The solution of (7) can be written as:*

- *optimal number of neurons: $K^* = R$*
- *their optimal weights: $(\boldsymbol{\alpha}_k^*, \boldsymbol{\omega}_k^*) = (\sqrt{\lambda_k} \mathbf{S}^{-\frac{1}{2}} \mathbf{u}_k, \sqrt{\lambda_k} \mathbf{v}_k)$ for $k = 1, \dots, R$.*

Moreover for any number of neurons $K \leq R$, and associated scaled weights $\theta_{\leftrightarrow}^{K,*}$, the expressivity gain and the first order in η of the loss improvement due to the addition of these K neurons are equal and can be quantified very simply as a function of the eigenvalues λ_k :

$$\Psi_{\theta \oplus \theta_{\leftrightarrow}^{K,*}}^l \leq \Psi_{\theta}^l - \sum_{k=1}^K \lambda_k^2$$

Proof

To facilitate reading we remove the layer index of each quantity, ie $\mathbf{B} := \mathbf{B}_{l-2}$, $\mathbf{V}^l(\mathbf{A}, \boldsymbol{\Omega}) := \mathbf{V}(\mathbf{A}, \boldsymbol{\Omega})$ and $\mathbf{V}_{\text{goal}_{proj}}^l := \mathbf{V}_{\text{goal}_{proj}}$.

To solve this problem, we consider the input of the incoming connexions \mathbf{B} and the desired change in the output of the outgoing connexions $\mathbf{V}_{\text{goal}_{proj}}$. Hence if we note $L(\mathbf{A})$ and $L(\boldsymbol{\Omega})$ the additional connexions of the expanded representation and σ the non linearity, we optimize the following proxy problem:

$$\arg \min_{\mathbf{A}, \boldsymbol{\Omega}} \left\| (L(\boldsymbol{\Omega}) \circ \sigma \circ L(\mathbf{A}))(\mathbf{B}) - \mathbf{V}_{\text{goal}_{proj}} \right\|_{\text{Tr}} \quad (47)$$

We solve this problem at first order by linearizing the non linearity σ . We denote $\mathbf{Lin}_{(a,b)}(W)$ the fully connected layer with input size a , output size b and weight matrix \mathbf{W} . We also note $C[+1]$ and $C[-1]$ the layer width at layer $l+1$ and $l-1$ with the convention that $C[0]$ is the dimension of the input X . With those notations, for fully connected layers, we have for the additions of K neurons:

$$\arg \min_{\mathbf{A}, \boldsymbol{\Omega}} \left\| \mathbf{Lin}_{(C[+1], K)}(\boldsymbol{\Omega})(\mathbf{Lin}_{(K, C[-1])}(\mathbf{A})(\mathbf{B})) - \mathbf{V}_{\text{goal}_{proj}} \right\|_2 \quad (48)$$

With the same notations, for convolutional layers, we have for the additions of K intermediate channels:

$$\arg \min_{\mathbf{A}, \boldsymbol{\Omega}} \left\| \mathbf{Conv}_{(C[+1], K)}(\boldsymbol{\Omega})(\mathbf{Conv}_{(K, C[-1])}(\mathbf{A})(\mathbf{B})) - \mathbf{V}_{\text{goal}_{proj}} \right\|_2 \quad (49)$$

If we note $\mathbf{V}(\mathbf{A}, \boldsymbol{\Omega})$ the result of \mathbf{B} after applying the layers parametrized by \mathbf{A} and $\boldsymbol{\Omega}$, in both cases we aim to optimize:

$$\arg \min_{\mathbf{A}, \boldsymbol{\Omega}} \left\| \mathbf{V}(\mathbf{A}, \boldsymbol{\Omega}) - \mathbf{V}_{\text{goal}_{proj}} \right\|_{\text{Tr}} \quad (50)$$

First we will transform the resolution of the problem in solving the following optimization problem:

$$\arg \min_{\mathbf{A}, \boldsymbol{\Omega}} \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{A} \boldsymbol{\Omega} - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|_2 \quad (51)$$

where \mathbf{S} depends of \mathbf{B} and \mathbf{N} of \mathbf{B} and $\mathbf{V}_{\text{goal}_{proj}}$.

If we note $\mathbf{S} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$ the SVD of \mathbf{S} , we define the square root of \mathbf{S} as $\mathbf{S}^{\frac{1}{2}} := \mathbf{U} \sqrt{\boldsymbol{\Lambda}} \mathbf{U}^\top$ and $\mathbf{S}^{-\frac{1}{2}} := \mathbf{U} \sqrt{\boldsymbol{\Lambda}^{-1}} \mathbf{U}^\top$ with the convention $0^{-1} = 0$.

Lemma C.1. For $S \in \mathbb{R}(s, s)$, $N \in \mathbb{R}(s, t)$, $A \in \mathbb{R}(s, k)$, $B \in \mathbb{R}(k, t)$.

If $N = S^{\frac{1}{2}} S^{-\frac{1}{2}} N$, we have:

$$\langle AB, SAB \rangle - 2 \langle N, AB \rangle = \left\| S^{\frac{1}{2}} AB - S^{-\frac{1}{2}} N \right\|^2 - \left\| S^{-\frac{1}{2}} N \right\|^2 \quad (52)$$

Proof.

- For the first term we have:

$$\langle AB, SAB \rangle = \left\langle AB, S^{\frac{1}{2}} S^{\frac{1}{2}} AB \right\rangle \quad (53)$$

$$= \left\langle S^{\frac{1}{2}} AB, S^{\frac{1}{2}} AB \right\rangle \quad (54)$$

$$= \left\| S^{\frac{1}{2}} AB \right\|^2 \quad (55)$$

- For the second term we have:

$$\langle N, AB \rangle = \left\langle S^{\frac{1}{2}} S^{-\frac{1}{2}} N, AB \right\rangle \quad (56)$$

$$= \left\langle S^{-\frac{1}{2}} N, S^{\frac{1}{2}} AB \right\rangle \quad (57)$$

Hence we have that:

$$\langle AB, SAB \rangle - 2 \langle N, AB \rangle = \left\| S^{\frac{1}{2}} AB \right\|^2 - 2 \langle N, AB \rangle + \left\| S^{-\frac{1}{2}} N \right\|^2 - \left\| S^{-\frac{1}{2}} N \right\|^2 \quad (58)$$

$$= \left\| S^{\frac{1}{2}} AB - S^{-\frac{1}{2}} N \right\|^2 - \left\| S^{-\frac{1}{2}} N \right\|^2 \quad (59)$$

□

C.2.1 Fully connected layers

For a fully connected layer, we have

$$V(A, \Omega) = \text{Lin}_{(C[+1], K)}(\Omega)(\text{Lin}_{(K, C[-1])}(A)(B)) = BA\Omega \quad (60)$$

We will use the following lemma to get the desired result:

Lemma C.2. Let $Y \in \mathbb{R}(n, t)$, $X \in \mathbb{R}(n, s)$ and $S := X^T X \in \mathbb{R}(s, s)$.

$$S^{\frac{1}{2}} S^{-\frac{1}{2}} X^T Y = X^T Y \quad (61)$$

Proof. Let decompose Y on $\text{Im}(X) \oplus_{\perp} \ker(X^T)$: $Y = XI + K$.

$$X^T Y = X^T XI + X^T K = X^T XI = SI$$

Hence $X^T Y \in \text{Im}(S)$, hence as $S|_{\text{Im}(S)}$ is invertible, we have: $S^{-\frac{1}{2}} S^{\frac{1}{2}} X^T Y = X^T Y$. □

Lemma C.3. Let $Y \in \mathbb{R}(n, t)$, $X \in \mathbb{R}(n, s)$, $k \leq \min(s, t)$, $A \in \mathbb{R}(s, k)$, $B \in \mathbb{R}(k, t)$.

We define:

$$S := X^T X \in \mathbb{R}(s, s) \quad (62)$$

$$N := X^T Y \in \mathbb{R}(s, t) \quad (63)$$

$$\|Y - XAB\|^2 = \left\| S^{\frac{1}{2}} AB - S^{-\frac{1}{2}} X^T Y \right\|^2 - \left\| S^{-\frac{1}{2}} X^T Y \right\|^2 + \|Y\|^2 \quad (64)$$

Proof. By developing the scalar product we get:

$$\|Y - XAB\|^2 = \|Y\|^2 - 2\langle Y, XAB \rangle + \|XAB\|^2 \quad (65)$$

$$= \|Y\|^2 - 2\langle Y, XAB \rangle + \langle XAB, XAB \rangle \quad (66)$$

$$= \|Y\|^2 - 2\langle Y, XAB \rangle + \langle AB, X^\top XAB \rangle \quad (67)$$

$$(68)$$

Using theorem C.2 we can apply theorem C.1 and immediately get the result. \square

Hence using theorem C.3 we have:

$$\begin{aligned} & \left\| \text{Lin}_{(C[+1], K)}(\Omega)(\text{Lin}_{(K, C[-1])}(A)(B)) - V_{\text{goal}_{proj}} \right\|^2 \\ &= \\ & \left\| S^{\frac{1}{2}} A \Omega - S^{-\frac{1}{2}} N \right\|^2 - \left\| S^{-\frac{1}{2}} N \right\|^2 + \left\| V_{\text{goal}_{proj}} \right\|^2 \end{aligned} \quad (69)$$

With:

$$\begin{aligned} S &:= B^\top B \in \mathbb{R}(C[-1], C[-1]) \\ N &:= B^\top V_{\text{goal}_{proj}} \in \mathbb{R}(C[-1], C[+1]) \end{aligned}$$

C.2.2 Convolutional connected layers

We have $A \in \mathbb{R}(K, C[-1], d, d)$ and $\Omega \in \mathbb{R}(C[+1], K, d[+1], d[+1])$ where $d, d[+1]$ is the kernel size at l and $l + 1$. We note $A_F \in \mathbb{R}(K, C[-1]dd)$ the flatten version of A and $A_k := A_F[k, :]$. We also note $\Omega_{m,k} \in \mathbb{R}(d[+1]d[+1])$ the flatten version of $\Omega[m, k, :, :]$. Using this we define $\Omega_k := (\Omega_{k,1} \ \cdots \ \Omega_{k,C[+1]})$

and $\Omega_F := \begin{pmatrix} \Omega_1 \\ \vdots \\ \Omega_K \end{pmatrix}.$

We define the tensor T such that for a pixel j of the output of the convolutional layer, T_j is a linear application that select the pixels of the input of the convolutional layer that are used to compute the pixel j of the output in a flatten version image (flatten only on the space not on the channels). $T \in \mathbb{R}(H[+1]W[+1], d[+1]d[+1], HW)$ where H and W are the height and width of the intermediate image and $H[+1]$ and $W[+1]$ are the height and width of the output image.

As previously, we have B^c the unfolded version of B such that $B^c \in \mathbb{R}(n, C[-1]dd, HW)$ satisfying $\text{Conv}(B_i)$ is equal with the correct reshape to AB_i^c .

In addition, we use j as an index on the space of pixel instead of having a couple h, w for height and width. With those notations we have:

$$V(A, \Omega)[i, m, j] = \text{Conv}_{(C[+1], K)}(\Omega)(\text{Conv}_{(K, C[-1])}(A)(B_i))[m, j] \quad (70)$$

$$= \sum_k^K \Omega_{m,k}^T T_j B_i^c A_k \quad (71)$$

In the following for simplicity we note $B_{i,j}^t := T_j B_i^c$. To find the best neurons to add we solve the expressivity bottleneck as :

$$\arg \min_{A, \Omega} \sum_i \sum_j \sum_m \|V_{\text{goal}_{proj}_i}^{(j,m)} - \sum_{k=1}^K \omega_{m,k}^T B_{i,j}^t \alpha_k\|^2 \quad (72)$$

$$(73)$$

Using the properties of the trace, it follows that :

$$\sum_{i,j,m} \|\mathbf{V}_{\text{goalproj}_i}^{(j,m)} - \sum_{k=1}^K \boldsymbol{\omega}_{m,k}^T \mathbf{B}_{i,j}^t \boldsymbol{\alpha}_k\|^2 = \sum_{i,j,m} \|\mathbf{V}_{\text{goalproj}_i}^{(j,m)} - \sum_k \text{Tr}(\mathbf{B}_{i,j}^t \boldsymbol{\alpha}_k \boldsymbol{\omega}_{m,k}^T)\|^2 \quad (74)$$

$$= \sum_{i,j,m} \|\mathbf{V}_{\text{goalproj}_i}^{(j,m)} - \text{Tr}(\mathbf{B}_{i,j}^t \overbrace{\sum_k \boldsymbol{\alpha}_k \boldsymbol{\omega}_{m,k}^T}^{\mathbf{F}})\|^2 \quad (75)$$

$$= \sum_{i,j,m} \|\mathbf{V}_{\text{goalproj}_i}^{(j,m)} - \text{flat}(\mathbf{B}_{i,j}^t)^T \text{flat}(\mathbf{F}_m)\|^2 \quad (76)$$

$$= \sum_{i,j} \|\mathbf{V}_{\text{goalproj}_i}^{(j)} - \text{flat}(\mathbf{B}_{i,j}^t)^T \mathbf{F}\|^2 \quad (77)$$

$$(78)$$

With $\mathbf{F} := (\text{flat}(\mathbf{F}_1) \ \dots \ \text{flat}(\mathbf{F}_{C[+1]}))$.

We remark that $\mathbf{V}(\mathbf{A}, \boldsymbol{\Omega})$ is a linear function of the matrix \mathbf{F} which implies that the solution of 72 is the same as for linear layer. Replacing $\boldsymbol{\Omega}\mathbf{A}$ by \mathbf{F} in 60 and following the same reasoning as for linear layer, it follows that 72 is equivalent to :

$$\arg \min_{\mathbf{F}} \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{F} - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|_2 \quad (79)$$

with $\mathbf{S} := \sum_{i,j} \text{flat}(\mathbf{B}_{i,j}^t) \text{flat}(\mathbf{B}_{i,j}^t)^T$ and $\mathbf{N} := \sum_{i,j} \mathbf{V}_{\text{goalproj}_i}^{(j)} \text{flat}(\mathbf{B}_{i,j}^t)^T$.

However, we remark that the dimension of $\mathbf{N} \in \mathbb{R}(C[-1]d[+1]^2d^2, d[+1])$ is quite large and that computing the SVD of such matrix is costly. To avoid expensive computation, we approximate 72 by defining the matrix \mathbf{S} and \mathbf{N} as 80 and 82. We now prove that 3.2, 3.3 and equation (10) still hold with such new definitions of \mathbf{S} and \mathbf{N} .

Lemma C.4. *Let $H = \max(\mathbf{B}_{i,j}^t.\text{shape})$, we define:*

$$\mathbf{S} := H \sum_{i=1}^n \sum_{j=1}^{H[+1]W[+1]} (\mathbf{B}_{i,j}^t)^\top (\mathbf{B}_{i,j}^t) \in (C[-1]dd, C[-1]dd) \quad (80)$$

$$\mathbf{N}_m := \sum_{i,j}^{n, H[+1]W[+1]} \mathbf{V}_{\text{goalproj}_{i,j,m}} (\mathbf{B}_{i,j}^t)^\top \in (C[-1]dd, d[+1]d[+1]) \quad (81)$$

$$\mathbf{N} := (\mathbf{N}_1 \cdots \mathbf{N}_{C[+1]}) \in (C[-1]dd, C[+1]d[+1]d[+1]) \quad (82)$$

We have:

$$\left\| \mathbf{V}(\mathbf{A}, \boldsymbol{\Omega}) - \mathbf{V}_{\text{goalproj}} \right\|^2 \leq \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{A}_F \boldsymbol{\Omega}_F - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|^2 - \left\| \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|^2 + \left\| \mathbf{V}_{\text{goalproj}} \right\|^2 \quad (83)$$

Proof.

$$\left\| \mathbf{V}(\mathbf{A}, \mathbf{\Omega}) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 = \sum_m^{C[+1]} \sum_{i,j}^{n, H[+1]W[+1]} \left(\sum_k^K \mathbf{\Omega}_{k,m}^T \mathbf{B}_{i,j}^t \boldsymbol{\alpha}_k - \mathbf{V}_{\text{goal}_{proj}_{i,j,m}} \right)^2 \quad (84)$$

$$(Q :=) = \sum_m^{C[+1]} \sum_{i,j}^{n, H[+1]W[+1]} \left(\sum_k^K \mathbf{\Omega}_{k,m}^T \mathbf{B}_{i,j}^t \boldsymbol{\alpha}_k \right)^2 \quad (85)$$

$$(-2L :=) - 2 \sum_m^{C[+1]} \sum_{i,j}^{n, H[+1]W[+1]} \sum_k^K \mathbf{\Omega}_{k,m}^T \mathbf{B}_{i,j}^t \boldsymbol{\alpha}_k \mathbf{V}_{\text{goal}_{proj}_{i,j,m}} \quad (86)$$

$$+ \sum_m^{C[+1]} \sum_{i,j}^{n, H[+1]W[+1]} \mathbf{V}_{\text{goal}_{proj}_{i,j,m}}^2 \quad (87)$$

$$= Q - 2L + \|\mathbf{V}_{\text{goal}_{proj}}\|^2 \quad (88)$$

We will use the following lemma to simplify this expression.

Lemma C.5. *For any square matrix $\mathbf{A} \in \mathbb{R}^{(n,n)}$, $\text{Tr}(\mathbf{A})^2 \leq \text{rank}(\mathbf{A}) \|\mathbf{A}\|^2$.*

Proof. Using the truncated SVD we have $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}$ with $\mathbf{\Sigma}$ a diagonal and $\mathbf{U} \in \mathbb{R}^{(n, \text{rank}(\mathbf{A}))}$, $\mathbf{V} \in \mathbb{R}^{(\text{rank}(\mathbf{A}), n)}$ truncated orthonormal matrices.

We have:

$$\text{Tr}(\mathbf{A})^2 = \text{Tr}(\mathbf{U}\mathbf{\Sigma}\mathbf{V})^2 \quad (89)$$

$$= \text{Tr}(\mathbf{V}\mathbf{U}\mathbf{\Sigma})^2 \quad (90)$$

$$= \langle \mathbf{V}\mathbf{U}, \mathbf{\Sigma} \rangle^2 \quad (91)$$

$$\text{(Cauchy-Swarz)} \leq \|\mathbf{V}\mathbf{U}\|^2 \|\mathbf{\Sigma}\|^2 \quad (92)$$

As \mathbf{U}, \mathbf{V} are truncated orthonormal matrices, we have:

$$\|\mathbf{V}\mathbf{U}\|^2 = \text{Tr}(\mathbf{U}^T \mathbf{V}^T \mathbf{V} \mathbf{U}) = \text{Tr}(\mathbf{U}^T \mathbf{U}) = \text{Tr}(\mathbf{I}_{\text{rank}(\mathbf{A})}) = \text{rank}(\mathbf{A})$$

Hence:

$$\text{Tr}(\mathbf{A})^2 \leq \text{rank}(\mathbf{A}) \|\mathbf{\Sigma}\|^2 \quad (93)$$

As \mathbf{U}, \mathbf{V} are truncated orthonormal matrices, we have:

$$\|\mathbf{\Sigma}\|^2 = \text{Tr}(\mathbf{\Sigma}^T \mathbf{\Sigma}) = \text{Tr}((\mathbf{V}\mathbf{\Sigma}\mathbf{U})\mathbf{U}\mathbf{\Sigma}\mathbf{V}) = \text{Tr}(\mathbf{A}^T \mathbf{A}) = \|\mathbf{A}\|^2$$

We conclude that:

$$\text{Tr}(\mathbf{A})^2 \leq \text{rank}(\mathbf{A}) \|\mathbf{A}\|^2 \quad (94)$$

□

Lemma C.6. *For $(\mathbf{M}_i)_{i \in I} \in (m, n)^I$, $(\mathbf{u}_k)_{k \in \llbracket K \rrbracket} \in (m)^K$, $(\mathbf{v}_k)_{k \in \llbracket K \rrbracket} \in (n)^K$ and with $\mathbf{W} := \sum_{k \in \llbracket K \rrbracket} \mathbf{v}_k \mathbf{u}_k^T \in (n, m)$ we have:*

$$\sum_{i \in I} \left\| \sum_{k \in \llbracket K \rrbracket} \mathbf{u}_k^T \mathbf{M}_i \mathbf{v}_k \right\|^2 \leq \left\langle \mathbf{W}, \sum_{i \in I} \mathbf{M}_i^T \mathbf{M}_i \mathbf{W} \right\rangle \quad (95)$$

Proof. Let $i \in I$:

$$\| \sum_{k \in \llbracket K \rrbracket} \mathbf{v} \mathbf{u}_k^T \mathbf{M}_i \mathbf{v}_k \|^2 = \left(\sum_{k \in \llbracket K \rrbracket} \mathbf{u}_k^T \mathbf{M}_i \mathbf{v}_k \right)^2 \quad (96)$$

$$= \text{Tr} \left(\sum_{k \in \llbracket K \rrbracket} \mathbf{u}_k^T \mathbf{M}_i \mathbf{v}_k \right)^2 \quad (97)$$

$$= \text{Tr}(\mathbf{M}_i \sum_{k \in \llbracket K \rrbracket} \mathbf{v}_k \mathbf{u}_k^T)^2 \quad (98)$$

$$= \text{Tr}(\mathbf{M}_i \mathbf{W})^2 \quad (99)$$

$$\text{(theorem C.5)} \leq \text{rank}(\mathbf{M}_i \mathbf{W}) \|\mathbf{M}_i \mathbf{W}\|^2 \quad (100)$$

$$\leq H \|\mathbf{M}_i \mathbf{W}\|^2 \quad H := \min(\mathbf{M}_i.\text{shape}) \quad (101)$$

Hence we have:

$$\sum_{i \in I} \left\| \sum_{k \in \llbracket K \rrbracket} \mathbf{u}_k^T \mathbf{M}_i \mathbf{v}_k \right\|^2 \leq H \sum_{i \in I} \|\mathbf{M}_i \mathbf{W}\|^2 \quad (102)$$

$$= H \sum_{i \in I} \langle \mathbf{M}_i \mathbf{W}, \mathbf{M}_i \mathbf{W} \rangle \quad (103)$$

$$= H \sum_{i \in I} \langle \mathbf{W}, \mathbf{M}_i^\top \mathbf{M}_i \mathbf{W} \rangle \quad (104)$$

$$= H \left\langle \mathbf{W}, \sum_{i \in I} \mathbf{M}_i^\top \mathbf{M}_i \mathbf{W} \right\rangle \quad (105)$$

□

Using theorem C.6 we have:

$$Q = \sum_m^{C[+1]} \sum_{i,j}^{n,H[+1]W[+1]} \left\| \sum_k^K \boldsymbol{\Omega}_{k,m}^T (\mathbf{B}_{i,j}^t) \boldsymbol{\alpha}_k \right\|^2 \quad (106)$$

$$\leq \sum_m^{C[+1]} \left\langle \sum_k^K \boldsymbol{\alpha}_k \boldsymbol{\Omega}_{k,m}^T, H \sum_{i,j}^{n,H[+1]W[+1]} (\mathbf{B}_{i,j}^t)^\top (\mathbf{B}_{i,j}^t) \sum_k^K \boldsymbol{\alpha}_k \boldsymbol{\Omega}_{k,m}^T \right\rangle \quad (107)$$

$$= \sum_m^{C[+1]} \langle \mathbf{A}_F \boldsymbol{\Omega}_m, \mathbf{S} \mathbf{A}_F \boldsymbol{\Omega}_m \rangle \quad (108)$$

$$= \langle \mathbf{A}_F \boldsymbol{\Omega}_F, \mathbf{S} \mathbf{A}_F \boldsymbol{\Omega}_F \rangle \quad (109)$$

Lemma C.7. For $\mathbf{M} \in (m, n)$, $\mathbf{u} \in (m)$, $\mathbf{v} \in (n)$ we have:

$$\mathbf{u}^T \mathbf{M} \mathbf{v} = \langle \mathbf{v} \mathbf{u}^T, \mathbf{M}^\top \rangle \quad (110)$$

Proof.

$$\mathbf{u}^T \mathbf{M} \mathbf{v}_k = (\mathbf{u}^T \mathbf{M} \mathbf{v})^\top \quad (111)$$

$$= \mathbf{v}^\top \mathbf{M}^\top \mathbf{u} \quad (112)$$

$$= \langle \mathbf{v}, \mathbf{M}^\top \mathbf{u} \rangle \quad (113)$$

$$= \langle \mathbf{v} \mathbf{u}^T, \mathbf{M}^\top \rangle \quad (114)$$

□

We have:

$$L = \sum_m^{C[+1]} \sum_{i,j}^{n,H[+1]W[+1]} \sum_k^K \Omega_{k,m}^T \mathbf{B}_{i,j}^t \alpha_k \mathbf{V}_{\text{goal}_{proj_{i,j,m}}} \quad (115)$$

$$(\mathbf{V}_{\text{goal}_{proj_{i,j,m}}} \in (1)) = \sum_m^{C[+1]} \sum_k^K \Omega_{k,m}^T \sum_{i,j}^{n,H[+1]W[+1]} (\mathbf{B}_{i,j}^t \mathbf{V}_{\text{goal}_{proj_{i,j,m}}}) \alpha_k \quad (116)$$

$$(\text{theorem C.7}) = \sum_m^{C[+1]} \left\langle \sum_k^K \alpha_k \Omega_{k,m}^T, \sum_{i,j}^{n,H[+1]W[+1]} \mathbf{V}_{\text{goal}_{proj_{i,j,m}}} (\mathbf{B}_{i,j}^t)^\top \right\rangle \quad (117)$$

$$= \sum_m^{C[+1]} \langle \mathbf{A}_F \Omega_m, \mathbf{N}_m \rangle \quad (118)$$

$$= \langle \mathbf{A}_F \Omega_F, \mathbf{N} \rangle \quad (119)$$

In total, we get:

$$\left\| \mathbf{V}(\mathbf{A}, \Omega) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \leq \langle \mathbf{A}_F \Omega_F, \mathbf{S} \mathbf{A}_F \Omega_F \rangle - 2 \langle \mathbf{A}_F \Omega_F, \mathbf{N} \rangle + \left\| \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \quad (120)$$

If we suppose that \mathbf{S} is invertible, we can apply theorem C.1 and get the result.

□

Lemma C.8. For $\mathbf{S} \in \mathbb{R}(s, s)$, $\mathbf{N} \in \mathbb{R}(s, t)$, $\mathbf{A} \in \mathbb{R}(s, K)$, $\mathbf{B} \in \mathbb{R}(K, t)$.

We note $\mathbf{U} \Lambda \mathbf{V}$ the singular value decomposition of $\mathbf{S}^{-\frac{1}{2}} \mathbf{N}$ and \mathbf{U}_K the first K columns of \mathbf{U} , \mathbf{V}_K the first K lines of \mathbf{V} , Λ_K the first K singular values of Λ and Λ_{K+1} the other singular values of Λ .

We define:

$$\mathbf{A}^* := \mathbf{S}^{-\frac{1}{2}} \mathbf{U}_K \sqrt{\Lambda_K} \quad (121)$$

$$\Omega^* := \sqrt{\Lambda_K} \mathbf{V}_K \quad (122)$$

$$\min_{\mathbf{A}, \Omega} \left\| \mathbf{V}(\mathbf{A}, \Omega) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \leq \left\| \mathbf{V}(\mathbf{A}^*, \Omega^*) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 = -\|\Lambda_K\|^2 + \left\| \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \quad (123)$$

with equality for the linear case.

$$\Psi_{\theta \oplus \theta_{\leftrightarrow}^{K*}}^l \leq \Psi_{\theta}^l - \sum_{k=1}^K \lambda_k^2 \quad (124)$$

Proof. Using theorem C.3 and theorem C.4:

$$\left\| \mathbf{V}(\mathbf{A}, \Omega) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \leq \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{A} \Omega - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|^2 - \left\| \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|^2 + \overbrace{\left\| \mathbf{V}_{\text{goal}_{proj}} \right\|^2}^{\Psi_{\theta}^l} \quad (125)$$

Hence we have:

$$\arg \min_{\mathbf{A}, \Omega} \left\| \mathbf{V}(\mathbf{A}, \Omega) - \mathbf{V}_{\text{goal}_{proj}} \right\| = \arg \min_{\mathbf{A}, \Omega} \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{A} \Omega - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\| \quad (126)$$

As we suppose that S is invertible, we can use the change of variable $\tilde{\mathbf{A}} = \mathbf{S}^{\frac{1}{2}} \mathbf{A}$ thus we have:

$$\min_{\mathbf{A}, \mathbf{\Omega}} \left\| \mathbf{S}^{\frac{1}{2}} \mathbf{A} \mathbf{\Omega} - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\| = \min_{\tilde{\mathbf{A}}, \mathbf{\Omega}} \left\| \tilde{\mathbf{A}} \mathbf{\Omega} - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\| \quad (127)$$

The solution of such problems is given by the paper Eckart & Young (1936) and is:

$$\tilde{\mathbf{A}}^* = \mathbf{U}_K \sqrt{\Lambda_K} \quad (128)$$

$$\mathbf{\Omega}^* = \sqrt{\Lambda_K} \mathbf{V}_K \quad (129)$$

To recover \mathbf{A}^* we simply have to multiply by $\mathbf{S}^{-\frac{1}{2}}$ on the left side of $\tilde{\mathbf{A}}^*$. By definition of the SVD and the construction of $(\mathbf{A}^*, \mathbf{\Omega}^*)$ we have:

$$\left\| \mathbf{S}^{\frac{1}{2}} \mathbf{A}^* \mathbf{\Omega}^* - \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|^2 - \left\| \mathbf{S}^{-\frac{1}{2}} \mathbf{N} \right\|^2 = \|\Lambda_{K+1}\|^2 - \|\Lambda\|^2 = -\|\Lambda_K\|^2 \quad (130)$$

Using this and equation (125) we immediately get the desired equation (123). To conclude, we can also rewrite this with the bottleneck expression:

$$\Psi_{\theta \oplus \theta_{\leftrightarrow}^K} := \min_{\mathbf{A}, \mathbf{\Omega}} \left\| \mathbf{V}(\mathbf{A}, \mathbf{\Omega}) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \leq \Psi_{\theta}^l - \sum_{k=1}^K \lambda_k^2 \quad (131)$$

□

Note 2: Considering one update of architecture $\delta \mathbf{W}^*$ at layer l and adding neurons $\theta_{\leftrightarrow}^{K,*}$ at layer $l-1$ the loss at first order in $\left\| \mathbf{V}^l(\delta \mathbf{W}^*) + \mathbf{V}^l(\theta_{\leftrightarrow}^{K,*}) \right\|$ is :

$$\mathcal{L}(f_{\theta \oplus \theta_{\leftrightarrow}^{K,*}}) \approx \mathcal{L}(f_{\theta}) - \frac{\sigma'_l(0)}{\eta} \left(\sum_{k=1}^K \lambda_k^2 + \langle \mathbf{V}_{\text{goal}}, \mathbf{V}^l(\delta \mathbf{W}^l) \rangle \right) \quad (132)$$

To prove note 2 we use the following lemma :

Lemma C.9. *We note $\mathbf{V}(\mathbf{A}, \mathbf{\Omega})$ the result of \mathbf{B} after applying the layers parameterized by \mathbf{A} and $\mathbf{\Omega}$. We note $\mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*)$ where \mathbf{A}^* and $\mathbf{\Omega}^*$ as define in 3.2*

$$\langle \mathbf{V}_{\text{goal}_{proj}}, \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*) \rangle = \|\Lambda_K\|^2 \quad (133)$$

Proof. Starting from theorem C.8

$$\left\| \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*) - \mathbf{V}_{\text{goal}_{proj}} \right\|^2 = -\|\Lambda_K\|^2 + \left\| \mathbf{V}_{\text{goal}_{proj}} \right\|^2 \quad (134)$$

Hence by developing the norm, we have:

$$\left\| \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*) \right\|^2 - 2 \langle \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*), \mathbf{V}_{\text{goal}_{proj}} \rangle = -\|\Lambda_K\|^2 \quad (135)$$

Moreover by construction we have $\left\| \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*) \right\| = \|\Lambda_K\|$ and therefore we get:

$$-2 \langle \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*), \mathbf{V}_{\text{goal}_{proj}} \rangle = -2 \|\Lambda_K\|^2 \quad (136)$$

which conclude the proof. □

We now prove the note 2. Using the first order approximation of the loss function at \mathbf{a}^l , we have the following:

$$\mathcal{L}(\mathbf{a}^l + \delta \mathbf{a}^l) = \mathcal{L}(\mathbf{a}^l) + \langle \nabla_{\mathbf{a}^l} \mathcal{L}, \delta \mathbf{a}^l \rangle + o(\|\delta \mathbf{a}^l\|) \quad (137)$$

On one hand, performing an update of architecture, ie $\mathbf{W}^* \leftarrow \mathbf{W} + \delta \mathbf{W}^*$ change the activation function \mathbf{a}^l as $\delta \mathbf{a}_1^l := \mathbf{V}(\delta \mathbf{W}^*)$. Then, as explained in appendix A.5, adding neurons at layer $l - 1$ change the activation function \mathbf{a}^l as :

$$\delta \mathbf{a}_2^l = \sigma'(0) \mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*) + o(\|\mathbf{V}(\mathbf{A}^*, \mathbf{\Omega}^*)\|) \quad (138)$$

Which combined give us that:

$$\mathcal{L}(\mathbf{A}^*, \mathbf{\Omega}^*) = \mathcal{L} + \langle \nabla_{\mathbf{a}^l} \mathcal{L}, \delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l \rangle + o(\|\delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l\|) \quad (139)$$

Using that $\mathbf{V}_{\text{goal}} := -\eta \nabla_{\mathbf{a}^l} \mathcal{L}$ we have that :

$$\mathcal{L}(\mathbf{A}^*, \mathbf{\Omega}^*) = \mathcal{L} - \frac{1}{\eta} \langle \mathbf{V}_{\text{goal}}, \delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l \rangle + o(\|\delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l\|) \quad (140)$$

$$= \mathcal{L} - \frac{1}{\eta} (\langle \mathbf{V}_{\text{goal}} - \delta \mathbf{a}_1^l, \delta \mathbf{a}_2^l \rangle + \langle \delta \mathbf{a}_1^l, \delta \mathbf{a}_2^l \rangle + \langle \mathbf{V}_{\text{goal}}, \delta \mathbf{a}_1^l \rangle) + o(\|\delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l\|) \quad (141)$$

Remarking that $\langle \delta \mathbf{a}_1^l, \delta \mathbf{a}_2^l \rangle = o(\|\delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l\|)$ and using C.9 we have :

$$\mathcal{L}(\mathbf{A}^*, \mathbf{\Omega}^*) = \mathcal{L} - \frac{1}{\eta} \left(\sigma'(0) \sum_{k=1}^K \lambda_k^2 + \langle \mathbf{V}_{\text{goal}}, \delta \mathbf{a}_1^l \rangle \right) + o(\|\delta \mathbf{a}_1^l + \delta \mathbf{a}_2^l\|) \quad (142)$$

Note on the approximation for convolutional layer. By developing the expression $\|\mathbf{V} - \mathbf{V}_{\text{goal}_{proj}}\|^2$, we remark that minimising $\|\mathbf{V} - \mathbf{V}_{\text{goal}_{proj}}\|^2$ over \mathbf{V} is equivalent to maximising $\langle \mathbf{V}, \mathbf{V}_{\text{goal}_{proj}} \rangle$ with a constrain on the norm of \mathbf{V} . This constrain lies in the functional space of the activities and can be reformulated in the parameter space with the matrix \mathbf{S} as $\|\mathbf{A}\mathbf{\Omega}^T\|_{\mathbf{S}} = \|\mathbf{V}\|$. By changing the matrix \mathbf{S} , we modify the metric on \mathbf{V} and obtain a pseudo solution $\mathbf{S}_{pseudo}^{-1} \mathbf{N}$ which is still positively correlated with $\mathbf{S}^{-1} \mathbf{N}$ as $\langle \mathbf{S}^{-1} \mathbf{N}, \mathbf{S}_{pseudo}^{-1} \mathbf{N} \rangle \geq 0$.

C.3 Proposition and remark 3.3 and 3.2

Proposition C.2. 3.3 Suppose \mathbf{S} is semi definite, we note $\mathbf{S} = \mathbf{S}^{\frac{1}{2}} \mathbf{S}^{\frac{1}{2}}$. Solving (7) is equivalent to find the K first eigenvectors α_k associated to the K largest eigenvalues λ of the generalized eigenvalue problem :

$$\mathbf{N} \mathbf{N}^T \alpha_k = \lambda \mathbf{S} \alpha_k \quad (143)$$

Corollary 2. (3.2) For all integers m, m' such that $m + m' \leq R$, at order one in η , adding $m + m'$ neurons simultaneously according to the previous method is equivalent to adding m neurons then m' neurons by applying successively the previous method twice.

Proof

To prove 3.3, we show that the solution of 143 and the formula of 3.2 are collinear.

Solving 143 is equivalent to maximizing the following generalized Rayleigh quotient (which is solvable by the LOBPCG technique):

$$\alpha^* = \arg \max_{\alpha} \frac{\alpha^T N N^T \alpha}{\alpha^T S \alpha} \quad (144)$$

$$p^* = \arg \max_{p=S^{1/2}\alpha} \frac{p^T S^{-\frac{1}{2}} N N^T S^{-\frac{1}{2}} p}{p^T p} \quad (145)$$

$$p^* = \arg \max_{\|p\|=1} \|N^T S^{-\frac{1}{2}} p\| \quad (146)$$

$$\alpha^* = S^{-\frac{1}{2}} p^* \quad (147)$$

Considering the SVD of $N^T S^{-\frac{1}{2}} = \sum_{r=1}^R \lambda_r e_r f_r^T$, then $S^{-\frac{1}{2}} N N^T S^{-\frac{1}{2}} = \sum_{r=1}^R \lambda_r^2 f_r f_r^T$, because $j \neq i \implies e_i^T e_j = 0$ and $f_i^T f_j = 0$. Hence maximizing the first quantity is equivalent to $p_k^* = f_k$, then $\alpha_k = S^{-\frac{1}{2}} f_k$, which match the formula of proposition 3.2. The same reasoning can be applied on ω_k .

We prove second corollary 3.2 by induction. Note that $v(\theta_{\leftrightarrow}^{K,*}, x) = o(\eta)$, then for $m = m' = 1$:

$$a_l(x)^{t+1} = a_l(x)^t + v(\theta_{\leftrightarrow}^{1,*}, x) + o(\eta) \quad (148)$$

Remark that $v_{\text{goal}}(x)$ is a function of $a_l(x)$, ie $v_{\text{goal}}(x) := g(a_l(x))$. Then suppose that $\mathcal{L}(f(x), y)$ is twice differentiable in $a_l(x)$. It follows that $g(a_l(x))$ is differentiable and :

$$v_{\text{goal}}^{t+1}(x) = g(a_l^t(x) + v(\theta_{\leftrightarrow}^{1,*}, x)) \quad (149)$$

$$= g(a_l^t(x)) + \nabla_{a_l^t(x)} g(a_l^t(x))^T v(\theta_{\leftrightarrow}^{1,*}, x) + o(\eta^2) \quad (150)$$

$$= v_{\text{goal}}^t(x) + \eta \frac{\partial^2 \mathcal{L}(f_{\theta}(x), y)}{\partial a_l^t(x)^2} v(\theta_{\leftrightarrow}^{1,*}, x) + o(\eta^2) \quad (151)$$

$$= v_{\text{goal}}^t(x) + o(\eta) \quad (152)$$

Adding the second neuron we obtain the minimization problem:

$$\arg \min_{\alpha_2, \omega_2} \|V_{\text{goal}}^t - V(\alpha_2, \omega_2)\| + o(\eta) \quad (153)$$

□

C.4 About equivalence of quadratic problems

Problems 8 and 7 are generally not equivalent, but might be very close, depending on layer sizes and number of samples. The difference between the two problems is that in one case one minimizes the quadratic quantity:

$$\left\| V^l(\theta_{\leftrightarrow}^K) + V^l(M) - V_{\text{goal}}^l \right\|^2$$

w.r.t. M and $\theta_{\leftrightarrow}^K$ **jointly**, while in the other case the problem is first minimized w.r.t. M and then w.r.t. $\theta_{\leftrightarrow}^K$. The latter process, being greedy, might thus provide a solution that is not as optimal as the joint optimization.

We chose this two-step process as it intuitively relates to the spirit of improving upon a standard gradient descent: we aim at adding neurons that complement what the other ones have already done. This choice is debatable and one could solve the joint problem instead, with the same techniques.

The topic of this section is to check how close the two problems are. To study this further, note that $V^l(M) = \delta W_l B_{l-1}$ while $V^l(\theta_{\leftrightarrow}^K) = \sum_{k=1}^K \omega_k B_{l-2}^T \alpha_k$. The rank of B_{l-1} is $\min(n_S, n_{l-1})$ where n_S is the number of samples and n_{l-1} the number of neurons (post-activities) in layer $l-1$, while the rank of B_{l-2} is $\min(n_S, n_{l-2})$ where n_{l-2} is the number of neurons (post-activities) in layer $l-2$. Note also that the number of degrees of freedom in the optimization variables δW_l and $\theta_{\leftrightarrow}^K = (\omega_k, \alpha_k)$ is much larger than these ranks.

Small sample case. If the number n_S of samples is lower than the number of neurons n_{l-1} and n_{l-2} (which is potentially problematic, see Section E.1), then it is possible to find suitable variables $\delta \mathbf{W}_l$ and $\theta_{\leftrightarrow}^K$ to form any desired $\mathbf{V}^l(\mathbf{M})$ and $\mathbf{V}^l(\theta_{\leftrightarrow}^K)$. In particular, if $n_S \leq n_{l-1} \leq n_{l-2}$, one can choose $\mathbf{V}^l(\theta_{\leftrightarrow}^K)$ to be $\mathbf{V}_{\text{goal}}^l - \mathbf{V}^l(\mathbf{M})$ and thus cancel any effect due to the greedy process in two steps. The two problems are then equivalent.

Large sample case. On the opposite, if the number of samples is very large (compared to the number of neurons n_{l-1} and n_{l-2}), then the lines of matrices \mathbf{B}_{l-1} and \mathbf{B}_{l-2} become asymptotically uncorrelated, under the assumption of their independence (which is debatable, depending on the type of layers and activation functions). Thus the optimization directions available to $\mathbf{V}^l(\mathbf{M})$ and $\mathbf{V}^l(\theta_{\leftrightarrow}^K)$ become orthogonal, and proceeding greedily does not affect the result, the two problems are asymptotically equivalent.

In the general case, matrices \mathbf{B}_{l-1} and \mathbf{B}_{l-2} are not independent, though not fully correlated, and the number of samples (in the minibatch) is typically larger than the number of neurons; the problems are then different.

Note that technically the ranks could be lower, in the improbable case where some neurons are perfectly redundant, or, e.g., if some samples yield exactly the same activities.

D Section About greedy growth sufficiency and TINY convergence with more details and proofs

One might wonder whether a greedy approach on layer growth might get stuck in a non-optimal state. By *greedy* we mean that every neuron added has to decrease the loss. We derive the following series of propositions in this regard. Since in this work we add neurons layer per layer independently, we study here the case of a single hidden layer network, to spot potential layer growth issues. For the sake of simplicity, we consider the task of least square regression towards an explicit continuous target f^* , defined on a compact set. That is, we aim at minimizing the loss:

$$\inf_{\mathbf{x} \in \mathcal{D}} \|f(\mathbf{x}) - f^*(\mathbf{x})\|^2 \quad (154)$$

where $f(\mathbf{x})$ is the output of the neural network and \mathcal{D} is the training set.

We start with an optional introductive section D.1 about greedy growth possibilities, then prepare lemmas in Sections D.2 and D.3 that will be used in Section D.4 to show that one can keep on adding neurons to a network (without modifying already existing weights) to make it converge exponentially fast towards the optimal function. Then in Section D.6 we present a growth method that explicitly overfits each dataset sample one by one, thus requiring only n neurons, thanks to existing weights modification. Finally, more importantly, in Section D.7, we show that actually any reasonable growth method that follows a certain optimization protocol (this includes TINY completed by random neuron additions if necessary) will reach 0 training error in at most n neuron additions.

D.1 Possibility of greedy growth

Proposition D.1 (Greedy completion of an existing network). *If f is not f^* yet, there exists a set of neurons to add to the hidden layer such that the new function f' will have a lower loss than f .*

One can even choose the added neurons such that the loss is arbitrarily well minimized.

Proof. The classic universal approximation theorem about neural networks with one hidden layer Pinkus (1999) states that for any continuous function g^* defined on a compact set ω , for any desired precision γ , and for any activation function σ provided it is not a polynomial, then there exists a neural network g with one hidden layer (possibly quite large when γ is small) and with this activation function σ , such that

$$(.*) \quad (155)$$

We apply this theorem to the case where $g^* = f^* - f$, which is continuous as f^* is continuous, and f is a shallow neural network and as such is a composition of linear functions and of the function σ , that we will suppose to be continuous for the sake of simplicity. We will suppose that f is real-valued for the sake of simplicity as well, but the result is trivially extendable to vector-valued functions (just concatenate the networks obtained for each output independently). We choose $\gamma = \frac{1}{10}\|f^* - f\|_{L2}$, where $\langle a | b \rangle_{L2} = \frac{1}{|\omega|} \int_{\mathbf{x} \in \omega} a(\mathbf{x}) b(\mathbf{x}) d\mathbf{x}$. This way we obtain a one-hidden-layer neural network g with activation function σ such that:

$$\forall \mathbf{x} \in \omega, \quad -\gamma \leq g(\mathbf{x}) - g^*(\mathbf{x}) \leq \gamma \quad (156)$$

$$\forall \mathbf{x} \in \omega, \quad g(\mathbf{x}) = f^*(\mathbf{x}) - f(\mathbf{x}) + a(\mathbf{x}) \quad (157)$$

with $\forall \mathbf{x} \in \omega, |a(\mathbf{x})| \leq \gamma$.

Then:

$$\forall \mathbf{x} \in \omega, \quad f^*(\mathbf{x}) - (f(\mathbf{x}) + g(\mathbf{x})) = -a(\mathbf{x}) \quad (158)$$

$$\forall \mathbf{x} \in \omega, \quad (f^*(\mathbf{x}) - h(\mathbf{x}))^2 = a^2(\mathbf{x}) \quad (159)$$

with h being the function corresponding to a neural network consisting in concatenating the hidden layer neurons of f and g , and consequently summing their outputs.

$$\|f^* - h\|_{L2}^2 = \|a\|_{L2}^2 \quad (160)$$

$$\|f^* - h\|_{L2}^2 \leq \gamma^2 = \frac{1}{100} \|f^* - f\|_{L2}^2 \quad (161)$$

and consequently the loss is reduced indeed (by a factor of 100 in this construction).

The same holds in expectation or sum over a training set, by choosing $\gamma = \frac{1}{10} \sqrt{\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \|f(\mathbf{x}) - f^*(\mathbf{x})\|^2}$, as Equation (159) then yields:

$$\sum_{\mathbf{x} \in \mathcal{D}} (f^*(\mathbf{x}) - h(\mathbf{x}))^2 = \sum_{\mathbf{x} \in \mathcal{D}} a^2(\mathbf{x}) \leq \frac{1}{100} \sum_{\mathbf{x} \in \mathcal{D}} (f^*(\mathbf{x}) - f(\mathbf{x}))^2 \quad (162)$$

which proves the proposition as stated.

For more general losses, one can consider order-1 (linear) developpment of the loss and ask for a network g that is close to (the opposite of) the gradient of the loss.

□

Proof of the additional remark. The proof in Pinkus (1999) relies on the existence of real values c_n such that the n -th order derivatives $\sigma^{(n)}(c_n)$ are not 0. Then, by considering appropriate values arbitrarily close to c_n , one can approximate the n -th derivative of σ at c_n and consequently the polynomial c^n of order n . This standard proof then concludes by density of polynomials in continuous functions.

Provided the activation function σ is not a polynomial, these values c_n can actually be chosen arbitrarily, in particular arbitrarily close to 0. This corresponds to choosing neuron input weights arbitrarily close to 0. □

Proposition D.2 (Greedy completion by one single neuron). *If f is not f^* yet, there exists a neuron to add to the hidden layer such that the new function f' will have a lower loss than f .*

Proof. From the previous proposition, there exists a finite set of neurons to add such that the loss will be decreased. In this particular setting of $L2$ regression, or for more general losses if considering small function moves, this means that the function represented by this set of neurons has a strictly negative component over the gradient g of the loss ($g = -2(f^* - f)$ in the case of the $L2$ regression). That is, denoting by $a_i \sigma(\mathbf{W}_i \cdot \mathbf{x})$ these N neurons:

$$\left\langle \sum_{i=1}^N a_i \sigma(\mathbf{w}_i \cdot \mathbf{x}) \mid g \right\rangle_{L2} = K < 0 \quad (163)$$

i.e.

$$\sum_{i=1}^N \langle a_i \sigma(\mathbf{w}_i \cdot \mathbf{x}) | g \rangle_{L_2} = K < 0 \quad (164)$$

We have:

$$0 > \frac{1}{N} K = \frac{1}{N} \sum_{i=1}^N \langle a_i \sigma(\mathbf{w}_i \cdot \mathbf{x}) | g \rangle_{L_2} \geq \min_{i=1}^N \langle a_i \sigma(\mathbf{w}_i \cdot \mathbf{x}) | g \rangle_{L_2} \quad (165)$$

Then necessarily at least one of the N neurons satisfies

$$\langle a_i \sigma(\mathbf{w}_i \cdot \mathbf{x}) | g \rangle_{L_2} \leq \frac{1}{N} K < 0 \quad (166)$$

and thus decreases the loss when added to the hidden layer of the neural network representing f . Moreover this decrease is at least $\frac{1}{N}$ of the loss decrease resulting from the addition of all neurons. \square

As a consequence, there exists no situation where one would need to add many neurons simultaneously to decrease the loss: it is always feasible with a single neuron. Note that finding the optimal neuron to add is actually NP-hard (Bach, 2017), so we will not necessarily search for the optimal one. A constructive lower bound on how much the loss can be improved will be given later in this section.

Proposition D.3 (Greedy completion by one infinitesimal neuron). *The neuron in the previous proposition can be chosen to have arbitrarily small input weights.*

Proof. This is straightforward, as, following a previous remark, the neurons found to collectively decrease the loss can be supposed to all have arbitrarily small input weights. \square

This detail is important in that our approach is based on the tangent space of the function f and thus manipulates infinitesimal quantities. Our optimization problem indeed relies on the linearization of the activation function by requiring the added neuron to have infinitely small input weights, to make the problem easier to solve. This proposition confirms that such neuron exists indeed.

Correlations and higher orders. Note that, as a matter of fact, our approach exploits linear correlations between inputs of a layer and desired output variations. It might happen that the loss is not minimized yet but there is no such correlation to exploit anymore. In that case the optimization problem (8) will not find neurons to add. Yet following Prop. D.3 there does exist a neuron with arbitrarily small input weights that can reduce the loss. This paradox can be explained by pushing further the Taylor expansion of that neuron output in terms of weight amplitude (single factor ε on all of its input weights), for instance $\sigma(\varepsilon \boldsymbol{\alpha} \cdot \mathbf{x}) \simeq \sigma(0) + \sigma'(0) \varepsilon \boldsymbol{\alpha} \cdot \mathbf{x} + \frac{1}{2} \sigma''(0) \varepsilon^2 (\boldsymbol{\alpha} \cdot \mathbf{x})^2 + O(\varepsilon^3)$. Though the linear term $\boldsymbol{\alpha} \cdot \mathbf{x}$ might be uncorrelated over the dataset with desired output variation $v(\mathbf{x})$, i.e. $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{x} v(\mathbf{x})] = 0$, the quadratic term $(\boldsymbol{\alpha} \cdot \mathbf{x})^2$, or higher-order ones otherwise, might be correlated with $v(\mathbf{x})$. Finding neurons with such higher-order correlations can be done by increasing accordingly the power of $(\boldsymbol{\alpha} \cdot \mathbf{x})$ in the optimization problem (7). Note that one could consider other function bases than the polynomials from Taylor expansion, such as Hermite or Legendre polynomials, for their orthogonality properties. In all cases, one does not need to solve such problems exactly but just to find an approximate solution, i.e. a neuron improving the loss.

Adding random neurons. Another possibility to suggest additional neurons, when expressivity bottlenecks are detected but no correlation (up to order p) can be exploited anymore, is to add random neurons. The first p order Taylor expansions will show 0 correlation with desired output variation, hence no loss improvement nor worsening, but the correlation of the $p + 1$ -th order will be non-0, with probability 1, in the spirit of random projections. Furthermore, in the spirit of common neural network training practice, one could consider brute force combinatorics by adding many random neurons and hoping some will be close enough to the desired direction (Frankle & Carbin, 2018). The difference with usual training is that we would perform such computationally-costly searches only when and where relevant, exploiting all simple information first (linear correlations in each layer).

D.2 Loss decrease with a line search on a quadratic energy

Let \mathcal{L} be a quadratic loss over \mathbb{R}^d and g be a vector in \mathbb{R}^d . The loss \mathcal{L} can be written as:

$$\mathcal{L}(g) = g^T Q g + v^T g + K \quad (167)$$

where Q is a matrix that we will suppose to be symmetric positive definite. This is to ensure that all eigenvalues of Q are positive, hence modelling a local minimum without saddle point. v is a vector in \mathbb{R}^d and K is a real constant.

For instance, the mean square loss $\mathbb{E}_{x \in \mathcal{D}} [\|f(x) - f^*(x)\|_S^2]$, where \mathcal{D} is a finite dataset of N samples, f^* a target function, and S is a symmetric positive definite matrix used as a metric, fits these hypotheses, considering $g = (f(x_1), f(x_2), \dots)$ as a vector. Indeed this loss rewrites as

$$\sum_{i=1}^N f(x_i)^T S f(x_i) - 2 \sum_i f^{*T}(x_i) S f(x_i) + K = g^T Q g + v^T g + K \quad (168)$$

by flattening and concatenating the vectors $f(x_i)$ and considering $Q = S \otimes S \otimes S \otimes \dots$ the tensor product of N times the same matrix S , i.e. a diagonal-block matrix with N identical blocks S . Note that for the standard regression with the L^2 metric, this matrix Q is just the Identity.

Starting from point g , and given a direction $h \in \mathbb{R}^d$, the question is to perform a line search in that direction, i.e. to optimize the factor $\lambda \in \mathbb{R}$ in order to minimise $\mathcal{L}(g + \lambda h)$.

Developing that expression, we get:

$$\mathcal{L}(g + \lambda h) = (g + \lambda h)^T Q (g + \lambda h) + v^T (g + \lambda h) + K = \lambda^2 h^T Q h + \lambda (2h^T Q g + v^T h) + \mathcal{L}(g) \quad (169)$$

which is a second-order polynomial in λ with positive quadratic coefficient. Note that the linear coefficient is $h^T \nabla_g \mathcal{L}(g)$, where $\nabla_g \mathcal{L}(g) = 2Qg + v$ is the gradient of \mathcal{L} at point g . The unique minimum of the polynomial in λ is then:

$$\lambda^* = -\frac{1}{2} \frac{h^T \nabla_g \mathcal{L}(g)}{h^T Q h} \quad (170)$$

which leads to

$$\min_{\lambda} \mathcal{L}(g + \lambda h) = \lambda^{*2} h^T Q h + \lambda^* h^T \nabla_g \mathcal{L}(g) + \mathcal{L}(g) \quad (171)$$

$$= \mathcal{L}(g) - \frac{1}{4} \frac{(h^T \nabla_g \mathcal{L}(g))^2}{h^T Q h} \quad (172)$$

$$= \mathcal{L}(g) - \frac{1}{4} \left\langle \frac{h}{\|h\|_Q} \middle| \nabla_g^Q \mathcal{L}(g) \right\rangle_Q^2. \quad (173)$$

Thus the loss gain obtained by a line search in a direction h is quadratic in the angle between that direction and the gradient of the loss, in the sense of the Q norm (and it is also quadratic in the norm of the gradient). Note that inner products with the gradient do not depend on the metric, in the sense that $\langle h | \nabla_g \mathcal{L}(g) \rangle_{L^2} = \langle h | \nabla_g^S \mathcal{L}(g) \rangle_S \quad \forall h$ for any metric S , i.e. any symmetric definite positive matrix S , associated to the norm $\|h\|_S^2 = h^T S h$ and to the gradient $\nabla_g^S \mathcal{L}(g) = S^{-1} \nabla_g^{L^2} \mathcal{L}(g)$.

In the case of a standard L^2 regression this boils down to:

$$\min_{\lambda} \|g + \lambda h\|_{L^2}^2 = \|g\|^2 - \left\langle \frac{h}{\|h\|} \middle| g \right\rangle_{L^2}^2 \quad (174)$$

i.e. considering $\mathcal{L}(f) := \mathbb{E}_{x \in \mathcal{D}} [\|f(x) - f^*(x)\|^2]$:

$$\min_{\lambda} \mathcal{L}(f + \lambda h) = \mathcal{L}(f) - \left\langle \frac{h}{\|h\|} \middle| f^* - f \right\rangle_{L^2}^2 = \mathcal{L}(f) - \frac{\mathbb{E}_{x \in \mathcal{D}} [(f^* - f) h]^2}{\mathbb{E}_{x \in \mathcal{D}} [\|h\|^2]}. \quad (175)$$

D.3 Expected loss gain with a line search in a random direction

Using Appendix D.2 above, the loss gain when performing a line search on a quadratic loss is quadratic in the angle $\alpha = \left\langle \frac{\mathbf{V}(X)}{\|\mathbf{V}(X)\|} \mid \frac{\mathbf{V}_{\text{goal}}(X)}{\|\mathbf{V}_{\text{goal}}(X)\|} \right\rangle_{L2}$ between the random search direction $\mathbf{V}(X)$ and the gradient $\mathbf{V}_{\text{goal}}(X)$.

This angle has average 0 and is of standard deviation $\frac{1}{nd}$, as described in Section 5.2. The loss gain is thus of the order of magnitude of $\frac{1}{d}$ in the best case (single-sample minibatch).

D.4 Exponential convergence to 0 training error

Considering a regression to a target f^* with the quadratic loss, the function f represented by the current neural network (fully-connected, one hidden layer, with ReLU activation function) can be improved to reach 0 loss by an addition of n neurons $(h_i)_{1 \leq i \leq n}$, with n is the dataset size, using Zhang et al. (2017). Unfortunately there is no guarantee that if one adds each of these neurons one by one, the loss decreases each time. We will prove that one of these neurons does decrease the loss, and we will quantify by how much, relying on the explicit construction in Zhang et al. (2017). This decrease will actually be a constant factor of the loss, thus leading to exponential convergence towards the target f^* on the training set.

As in the proof of Proposition D.2 in Appendix D, at least one of the added neurons satisfies that its inner product with the gradient direction is at least $1/n$. While one could consequently hope for a loss gain in $O(\frac{1}{n})$, one has to see that this decrease would be the one of a gradient step, which is multiplied by a step size η , and asks for multiple steps to be done. Instead in our approach we actually perform a line search over the direction of the new neuron. In both cases (line search or multiple small gradient steps), one has to take into account at least order-2 changes of the loss to compute the line search or estimate suitable η and/or its associated number of steps. Luckily in our case of least square regression, the loss is exactly equal to its second order Taylor development, and all following computations are exact.

We consider the mean square regression loss $\mathcal{L}(f) = \mathbb{E}_{x \in \mathcal{D}} [\|f(x) - f^*(x)\|_S^2]$, where \mathcal{D} is a finite training dataset of N samples. Its functional gradient $\nabla \mathcal{L}(f)$ at point f is $2(f - f^*)$, which is proportional to the optimal change to add to f , that is, $f^* - f$. The n neurons $(h_i)_{1 \leq i \leq n}$ to be added to f following Zhang et al. (2017) satisfy $\sum_i h_i = f^* - f = -\frac{1}{2} \nabla \mathcal{L}(f)$. Thus

$$\left\langle \sum_i h_i \mid f^* - f \right\rangle_{L2} = \|f^* - f\|_{L2}^2 = \mathcal{L}(f). \quad (176)$$

Then like in the proof of D.2 we use that the maximum is greater or equal to the mean to get that there exists a neuron h_i that satisfies:

$$\langle h_i \mid f^* - f \rangle_{L2} \geq \mathcal{L}(f)/n. \quad (177)$$

By applying Appendix D.2 one obtains that the new loss after line search into the direction of h_i yields:

$$\min_{\lambda} \mathcal{L}(f + \lambda h_i) = \mathcal{L}(f) - \frac{\langle h_i \mid f^* - f \rangle_{L2}^2}{\|h_i\|^2} \leq \mathcal{L}(f) \times \left(1 - \frac{\mathcal{L}(f)}{n^2 \|h_i\|^2} \right). \quad (178)$$

From the particular construction in Zhang et al. (2017) it is possible to bound the square norm of the neuron $\|h_i\|^2$ by $n d' \left(\frac{d_M}{d_m} \right)^2 \mathcal{L}(f)$, where d_M is related to the maximum distance between 2 points in the dataset, d_m is another geometric quantity related to the minimum distance, and d' is the network output dimension. To ease the reading of this proof, we defer the construction of this bound to next section, Appendix D.5.

Then the loss at each neuron addition decreases by a factor which is at least $\gamma = 1 - \frac{1}{n^3 d'} \left(\frac{d_m}{d_M} \right)^2 < 1$. This factor is a constant, as it is a bound that depends only on the geometry of the dataset (not on f).

Thus it is possible to decrease the loss exponentially fast with the number t of added neurons, i.e. $\mathcal{L}(f_t) \leq \gamma^t \mathcal{L}(f)$, towards 0 training loss, and this in a greedy way, that is, by adding neuron one by one, with the property that each neuron addition decreases the loss.

Note that, in the proof of Zhang et al. (2017), the added neurons could be chosen to have arbitrarily small input weights. This corresponds to choosing a with small norm instead of unit norm in Equation 179.

The number of neuron additions expected to reach good performance according to this bound is in the order of magnitude of n^3 , which is to be compared to n (number of neurons needed to overfit the dataset, without the constraint that each addition decreases the loss). This bound might be improved using other constructions than Zhang et al. (2017), though with this proof the bound cannot be better than n^2 (supposing $\|h_i\|$ can be made not to depend on n).

Note also that with ReLU activation functions, all points that are on the convex hull of the dataset (which is necessarily the case of all points if the input dimension is higher than the number of points) can easily in turn be perfectly predicted (0 loss) by just one neuron addition each (without changing the outputs for the other points), by choosing an hyperplane that separates the current convex hull point from the rest of the dataset, and setting a ReLU neuron in that direction.

D.5 Bound on the norm of the neurons

Here we prove that the neurons obtained by Zhang et al. (2017) can be chosen so as to bound the square norm of any neuron $\|h_i\|^2$ by $n d' \left(\frac{d_M}{d_m}\right)^2 \mathcal{L}(f)$, where d_M is related to the maximum distance between 2 points in the dataset, and d_m is another geometric quantity related to the minimum distance. For the sake of simplicity, we first consider the case where the output dimension is $d' = 1$.

In Zhang et al. (2017), the n neurons are obtained by solving $y = Aw$, where $y = (y_1, y_2, \dots)$ is the target function (here $(f^* - f)$ at each x_j), A is the matrix given by $A_{jk} = \text{ReLU}(a \cdot x_j - b_k)$, representing neuron activations, and a is any vector that separates the dataset points, i.e. $a \cdot x_j \neq a \cdot x_{j'} \forall j \neq j'$, that is, a could be almost any vector in \mathbb{R}^d (in the sense of random projections, that is, the set of vectors that do not satisfy this is of measure 0).

Here we will pick a particular unit direction a , one that maximizes the distance between any two samples after projection:

$$a \in \arg \max_{\|a\|=1} \min_{j, j'} |a \cdot (x_j - x_{j'})| \quad (179)$$

and let us denote d'_m the associated value: $d'_m = \min_{j, j'} |a \cdot (x_j - x_{j'})|$ for that a . Note that $d'_m \leq \min_{j, j'} \|x_j - x_{j'}\|$ and that it depends only on the training set. The quantity d'_m is likely to be also lower-bounded (over all possible datasets) by $\min_{j, j'} \|x_j - x_{j'}\|$ times a factor depending on the embedding dimension d and the number of points n .

Now, let us sort the samples according to increasing $a \cdot x_j$, that is, let us re-index the samples such that $(a \cdot x_j)$ now grows with j . By definition of a , the difference between any two consecutive $a \cdot x_j$ is at least d'_m .

We now choose biases $b_j = a \cdot x_j - d'_m + \varepsilon$ for some very small ε . The neurons are then defined as $h_k(x) = w_k \text{ReLU}(a \cdot x - b_k)$. The induced activation matrix $A_{jk} = \text{ReLU}(a \cdot x_j - b_k)$ then satisfies $\forall j < k; A_{jk} = 0$ and $\forall j \geq k; A_{jk} \geq d'_m - \varepsilon$. The matrix A is lower triangular with diagonal elements above $d'_m := d'_m - \varepsilon$, hence invertible. Recall that $y = Aw$.

Consequently, $w = A^{-1}y$, and hence $\|w\|^2 \leq \|A^{-1}\|^2 \|y\|^2$, that is,

$$\|w\|^2 \leq \frac{1}{d_m^2} \mathcal{L}(f) \quad (180)$$

as the target y is the vector $f^* - f$ in our case. Consequently, for any neuron h_i , one has:

$$w_i^2 \leq \frac{1}{d_m^2} \mathcal{L}(f). \quad (181)$$

As the norm of the neuron is $\|h_i\|^2 = w_i^2 \sum_j A_{ji}^2$, one still has to bound the activities $A_{ji} = \text{ReLU}(a \cdot x_j - b_i)$. As a was chosen a unit direction, the values $a \cdot x_j$ span a domain smaller than the diameter of the dataset \mathcal{D} : $|a \cdot (x_j - x_{j'})| \leq \|x_j - x_{j'}\| \leq \text{diam}(\mathcal{D}) \forall j, j'$. Hence all values $\forall i, j, |A_{ij}| = |a \cdot x_i - b_j| = |a \cdot x_i - a \cdot x_j + d_m| < d_M := \text{diam}(\mathcal{D}) + d_m$. Note that d_M depends only the dataset geometry, as for d_m .

We now have:

$$\|h_i\|^2 = w_i^2 \sum_j A_{ji}^2 \leq n \frac{d_M^2}{d_m^2} \mathcal{L}(f) \quad (182)$$

which ends the proof.

For higher output dimensions d' , one vector w of output weights is estimated per dimension, independently, leading to the same bound for each dimension. The square norms of neurons are summed over all dimensions and thus multiplied by at most d' .

D.6 Reaching 0 training error in n neuron additions by overfitting each dataset sample in turn

If one allows updating already existing output weights at the same time as one adds new neurons, then it is possible to reach 0 training error in only n steps (where n is the size of the dataset) while decreasing the loss at each addition.

This scenario is closer to the one we consider with TINY, as we compute the optimal update of existing weights inside the layer, as a byproduct of new neuron estimation, and apply them.

However the existence proof here follows a very different way to create new neurons, tailored to obtain a constructive proof, and inspired by the previous section. See Appendix D.7 for another, more generic proof, applicable to a wide range of growth methods.

Here we consider the same approach as in Appendix D.5 above, but introducing neurons one by one instead of n neurons at once. After computing a and the biases b_j , thus forming the activity matrix A , we add only the last neuron h_n . The activity of this neuron is 0 for all input samples x_j except for the last one, for which it is $A_{nn} > 0$. Thus, the neuron h_n separates the sample x_n from the rest of the dataset, and it is easy to find w_n so that the loss gets to 0 on that training sample, without changing the outputs for other samples.

Similarly, one can then add neuron h_{n-1} , which is active only for samples x_{n-1} and x_n . However designing w_{n-1} so that the loss becomes 0 at point x_{n-1} disturbs the output for point x_n (and for that point only). Luckily if one allows to update w_n then there exists a (unique) solution (w_{n-1}, w_n) to achieve 0 loss at both points. This is done exactly as previously, by solving $y = Aw$, but considering only the 2 last lines and rows of A , leading to a smaller 2×2 system which is also lower-triangular with positive diagonal.

Proceeding iteratively this way adds neuron one by one in a way that sends each time one more sample to 0 loss. Thus adding n neurons is sufficient to achieve 0 loss on the full training set, and this in a way that each time decreases the loss.

Note that updating existing output weights w_i while adding a new neuron, to decrease optimally the loss, is actually what TINY does. However, the construction in this Appendix completely overfits each sample in turn, by design, without being to generalize to new test points. On the opposite, TINY exploits correlations over the whole dataset to extract the main tendencies.

D.7 TINY reaches 0 training error in n neuron additions

We will now show that the TINY approach, as well as any other suitable greedy growth method, implemented within the right optimization procedure, reaches 0 training error in at most n steps (where n is the size of the dataset), almost surely.

Before stating it formally, we need to introduce the optimization protocol, growth completion and a probability measure over activation functions.

Optimization protocol. For this we consider the following optimization protocol conditions, that has to be applied at least during the last, n -th addition step:

- a **full batch** approach,
- when adding new neurons, also compute and add the **optimal moves of already existing parameters** (i.e. of output weights w).

The first point is to ensure that all dataset samples will be taken into account in the loss during the n -th update. Otherwise, for instance if using minibatches instead, the optimization of output weights w will not be able to overfit the training loss.

The second point is to make sure that, after update, the output weights w will be optimal for the training loss. Note that in the mean square regression case, this is easy to do, as the loss is quadratic in w : the optimal move (leading to the global optimum f^*) can be obtained by line search over the natural gradient (which is obtained for free as a by-product of TINY’s projection of \mathbf{V}_{goal} , and is proportional to $f^* - f$). This is precisely what we do in practice with TINY when training networks (except when comparing with other methods and using their own protocol).

Growth completion. For this proof to make sense, we will need the growth method to actually be able to perform n neuron additions, if it has not reached 0 training loss before. A counter-example would be a growth method that gets stuck at a place where the training loss is not 0 while being unable to propose new neuron to add. In the case of TINY, this can happen when no correlation between inputs x_i and desired output variations $f^*(x_i) - f(x_i)$ can be found anymore. To prevent this, one can choose any auxiliary method to add neurons in such cases, for instance random directions, solutions of higher-order expressivity bottleneck formulations using further developments of the activation function, or locally-optimal neurons found by gradient descent. Some auxiliary methods are guaranteed to further decrease the loss by a neuron addition (cf. Appendices D.2, D.3, D.4), while any other one will be guaranteed not to increase the loss if combined with a line search along that neuron direction.

We will name *completed-TINY* the completion of TINY by any such auxiliary method.

Activation function. For technical reasons, the result will stand *almost surely* only, depending on the invertibility of a certain matrix, namely, the activation matrix A , defined as $A_{ij} = \sigma(\mathbf{v}_j \cdot \mathbf{x}_i + b_j)$, indexed by samples i and neurons j .

Generally speaking, kernels induced by neurons $k_j : \mathbf{x} \mapsto \sigma(\mathbf{v}_j \cdot \mathbf{x} + b_j)$ form free families, in the sense that they are linearly independent (to the notable exception of the linear kernel). This linear independency means that a linear combination of kernels cannot be equal, as a function, to another kernel with different parameters. Equality is to be understood as *for all possible points \mathbf{x} ever*. However here we will evaluate the functions only at a finite number n of points (the dataset samples), therefore linear independence will be considered between the rows of the activation matrix A . This notion of linear dependence is much weaker: kernels might form a free family as functions but be linearly dependent once restricted to the dataset samples, by mere chance. While this is not likely (over dataset samples), this is not impossible in general (though of measure 0), and it is difficult to express an explicit, simple condition on the activation function to be sure that the activation matrix A is *always* invertible (up to slight changes of parameters). Thus instead we will express results *almost surely* over activation functions and neuron parameters.

For most activation functions in the space of smooth functions, the activation matrix A will be invertible almost surely over all possible datasets. In the unlucky case where the matrix is not invertible, an infinitesimal move of the neurons’ parameters will be sufficient to make it invertible. For some activation functions, however, such as linear or piecewise-linear ones (e.g., ReLU), the matrix might remain non-invertible over a wide range of parameter variations (unless further assumptions are made on the neurons added by the growth process). Yet, in such cases, slight perturbations of the activation function (i.e., choosing another, smooth, activation function, arbitrarily close to the original one) will yield invertibility.

To properly define "almost surely" regarding activation functions, let us restrict the activation function σ to belong to the space \mathfrak{P} of polynomials of order at least n^2 , that is:

$$\sigma(x) = \sum_{k=0}^K \gamma_k x^k \quad (183)$$

with $n^2 \leq K < +\infty$, and non-0 highest-order amplitude $\gamma_K \neq 0$. This set \mathfrak{P} is dense in the set of all continuous functions over the set $\Omega = [-r_M, r_M]^d$ which is an hypercube of sufficient radius r_M to cover all samples from the given dataset. One can define probability distributions over \mathfrak{P} , for instance consider the density $p(\sigma) = \frac{\alpha}{K^2} \prod_{k=0}^K \frac{e^{-\gamma_k^2}}{\sqrt{2\pi}}$ with a factor $\alpha = \left(\frac{\pi^2}{6} - \sum_{k < n^2} \frac{1}{k^2}\right)^{-1}$ to normalize the distribution, and where K is the order of the polynomial and thus depends on σ . This density is continuous in the space of parameters γ_k (though not continuous in usual functional metric spaces). Note that the decomposition of any $\sigma \in \mathfrak{P}$ as a finite-order polynomial is unique, as monomials of different orders are linearly independent.

We can now state the following lemma (that we will prove later):

Lemma D.1 (Invertibility of the activation matrix). *Let $\mathcal{D} = \{\mathbf{x}_i, 1 \leq i \leq n\}$ be a dataset of n distinct points, and let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a function in \mathfrak{P} , that is, a polynomial of order at least n^2 . Then with probability 1 over function and neuron parameters (γ_k) , (\mathbf{v}_j) and (b_j) , the activity matrix A defined by $A_{ij} = \sigma(\mathbf{v}_j \cdot \mathbf{x}_i + b_j)$ is full rank.*

and the following proposition:

Proposition D.4 (Reaching 0 training error in at most n neuron additions). *Under the assumptions above (polynomial activation function of order $\geq n^2$, full batch optimization and computation of the optimal moves of already existing parameters), completed-TINY reaches 0 training error in at most n neuron additions almost surely.*

Proof. If the growth method reaches 0 training error before n neuron additions, the proof is done. Otherwise, let us consider the n -th neuron addition. We will show in Lemma D.1 that the activity matrix A , defined by $A_{ij} = \sigma(\mathbf{v}_j \cdot \mathbf{x}_i + b_j)$, indexed by samples i and neurons j , is invertible. Then there exists a unique $\mathbf{w} \in \mathbb{R}^n$ such that $A\mathbf{w} = f^*$, i.e. $\sum_j w_j \sigma(\mathbf{v}_j \cdot \mathbf{x}_i + b_j) = f^*(\mathbf{x}_i)$ for each point \mathbf{x}_i of the dataset. This vector of output parameters \mathbf{w} realises the global minimum of the loss over already existing weights: $\inf_{\mathbf{w}} \mathcal{L}(f_{\mathbf{v}}, \mathbf{w}) = \inf_{\mathbf{w}} \|A\mathbf{w} - f^*\|^2$. They are also the ones found by a natural gradient step over the loss (up to a factor 2, that can easily be found by line search as the loss is convex). Then after that update the training loss is exactly 0. \square

Note: piecewise-linear activation functions such as ReLU are not covered by this proposition. However the result might still hold with further assumptions over the growth process. For instance with the method in Zhang et al. (2017), the ReLU neurons are chosen in such a way that the matrix A is full rank by construction.

Proof of Lemma D.1. Let us first show that if, unluckily, for a given activation function σ and given parameters (\mathbf{v}_j, b_j) , the matrix A is not full rank, then upon infinitesimal variation of the parameters, the matrix A becomes full rank.

Indeed, if all pre-activities $a_{i,j} := \mathbf{v}_j \cdot \mathbf{x}_i + b_j$ are not distinct for all i, j , then an infinitesimal variation of the vectors \mathbf{v}_j can make them distinct. For this, one can see that the set of directions \mathbf{v}_j on which any two dataset points \mathbf{x}_i and $\mathbf{x}_{i'}$ have the same projection is finite (since it has to be the direction of $\mathbf{x}_i - \mathbf{x}_{i'}$, for a given pair of dataset samples (i, i')) and thus of measure 0. As a consequence with probability 1 over neuron parameters \mathbf{v}_j and b_j , all pre-activities are distinct.

Now, if the matrix A is not invertible, as invertible matrices are dense in the space of matrices, one can easily find an infinitesimal change δA to apply to A to make it invertible. This corresponds to changing the activation function σ accordingly at each of the n^2 distinct pre-activity values. Since σ has more than n^2 parameters, this is doable. For instance one can select the n^2 first parameters and search for a suitable variation $\mathbf{g} := (\delta\gamma_k)_{0 \leq k < n^2}$ of them by solving the linear system $S\mathbf{g} = \delta A$ where the $n^2 \times n^2$ matrix S is

defined by $S_{ij,k} = a_{i,j}^k = (\mathbf{v}_j \cdot \mathbf{x}_i + b_j)^k$. This matrix S is invertible because any \mathbf{g} such that $S\mathbf{g} = 0$ would induce:

$$\forall i, j, \sum_{k=0}^{n^2-1} \delta \gamma_k a_{i,j}^k = 0 \quad (184)$$

and thus the polynomial $P(x) = \sum_{k=0}^{n^2-1} \delta \gamma_k x^k$ has at least n^2 roots while being of order at most $n^2 - 1$. Thus $S\mathbf{g} = 0 \implies \mathbf{g} = 0$ and S is invertible. Note that as δA is infinitesimal, $\mathbf{g} = S^{-1} \delta A$ will be infinitesimal as well, and so is the change brought to the activation function σ .

Consequently we have that the set of activation functions σ and neuron parameters (\mathbf{v}_j, b_j) for which the matrix A is full rank is dense in the set of polynomials \mathfrak{P} of order at least n^2 and of neuron parameters \mathfrak{N} .

Now, the function $\det : \mathfrak{P} \times \mathfrak{N} \rightarrow \mathbb{R}$, $((\gamma_k)_k, (\mathbf{v}_j, b_j)_j) \mapsto \det A = \det(\sigma_\gamma(\mathbf{v}_j \cdot \mathbf{x}_i + b_j))$ is smooth as a function of its input parameters (the determinant being a polynomial function of the matrix coefficients). As this continuous function is non-0 on a dense set of its inputs, the pre-image $\det^{-1}\{0\}$ is closed and contains no open subset. This is not yet sufficient to prove that this pre-image is of measure 0 (e.g., fat Cantor set).

For a fixed order K , one can see this function as a polynomial of its inputs γ_k and \mathbf{v}_j, b_j , and conclude² that the set of its roots is of measure 0. As a consequence, the probability, over coefficients γ_k or equivalently over polynomials σ of order K , that $\det A$ is non-0, is 1. As this stands for all K , we have that the probability that the matrix A is invertible is at least the mass of polynomials of all orders K , i.e. $\sum_{k \geq n^2} \frac{\alpha}{k^2} = 1$. Thus A is invertible with probability 1. \square

E Technical details

E.1 Batch size to estimate the new neuron and the best update

In this section we study the variance of the matrices \mathbf{M}^* and $\mathbf{S}^{-1/2}\mathbf{N}$ computed using a minibatch of n samples, seeing the samples as random variables, and the matrices computed as estimators of the true matrices one would obtain by considering the full distribution of samples. Those two matrices are the solutions of the multiple linear regression problems defined in (39) and in (51), as we are trying to regress the desired update noted Y onto the span of the activities noted X . We suppose we have the following setting :

$$Y \sim \mathbf{A}X + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2), \quad \mathbb{E}[\varepsilon|X] = 0$$

where the (X_i, Y_i) are *i.i.d.* and A is the oracle for \mathbf{M}^* or matrix $\mathbf{S}^{-1/2}\mathbf{N}$. If Y is multidimensional, the the total variance of our estimator can be seen as the sum of the variances of the estimator on each dimension of Y .

We now suppose that $Y \in \mathbb{R}$. The estimator $\hat{A} := (\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}Y^T$ has variance $\text{var}(\hat{A}) = \sigma^2(\mathbf{X}\mathbf{X}^T)^{-1}$. If n is large, and if matrix $\frac{1}{n}\mathbf{X}\mathbf{X}^T \rightarrow \mathbf{Q}$, with \mathbf{Q} non singular, then, asymptotically, we have $\hat{A} \sim \mathcal{N}(\mathbf{A}, \sigma^2 \frac{\mathbf{Q}^{-1}}{n})$, which is equivalent to $(\hat{A} - \mathbf{A}) \frac{\sqrt{n}}{\sigma} \mathbf{Q}^{1/2} \sim \mathcal{N}(0, I)$. Then $\|(\hat{A} - \mathbf{A}) \frac{\sqrt{n}}{\sigma} \mathbf{Q}^{1/2}\|^2 \sim \chi^2(k)$ where k is the dimension of X . It follows that $\mathbb{E}[\|(\hat{A} - \mathbf{A}) \mathbf{Q}^{1/2}\|^2] = \frac{k\sigma^2}{n}$ and as $\mathbf{Q}^{1/2}\mathbf{Q}^{1/2^T}$ is positive definite, we conclude that $\text{var}(\hat{A}) \leq \frac{k\sigma^2}{n\lambda_{\min}(\mathbf{Q})}$.

In practice and to keep the variance of our estimators stable during architecture growth, for the estimation of the best neuron to add we use batch size

$$n \propto \frac{(SW)^2}{P},$$

with the notations defined in Figure 12, since the matrices we estimate have side size SW and that each input sample contains P values, i.e. P quantities that each play the role of X here.

²See for instance a proof by recurrence that roots of a polynomial are always of measure 0: <https://math.stackexchange.com/questions/1920302/the-lebesgue-measure-of-zero-set-of-a-polynomial-function-is-zero>.

E.2 Batch size for learning

We adjust the batch size for gradient descent as follow : the batch size is set to $b_{t=0} = 32$ at the beginning of each experiment, and it is scheduled to increase as the square root of the complexity of the model (ie number of parameters). If at time t the network has complexity C_t parameters, then at time $t+1$ the training batch size is equal to $b_{t+1} = b_t \times \sqrt{\frac{C_{t+1}}{C_t}}$.

E.3 Normalization

E.3.1 Figures 5 and 16 : Usual normalization

For the GradMax method of figure 5 and 16, before adding the new neurons to the architecture, we normalize the out-going weight of the new neurons according to Evci et al. (2022), ie :

$$\alpha_k^* \leftarrow 0 \quad (185)$$

$$\text{for } 5 \quad \omega_k^* \leftarrow \omega_k^* \times \frac{1e-3}{\sqrt{\|(\omega_j^*)_{j=1}^{n_d}\|_2^2/n_d}} \quad (186)$$

$$\text{for } 16 \quad \omega_k^* \leftarrow \omega_k^* \times \sqrt{\frac{1e-3}{\|(\omega_j^*)_{j=1}^{n_d}\|_2^2/n_d}} \quad (187)$$

For TINY method of both figures, the previous normalization process is mimicked by normalizing the in and out going weights by theirs norms and multiplying them by $\sqrt{1e-3}$, ie :

$$\alpha_k \leftarrow \alpha_k^* \times \sqrt{\frac{1e-3}{\|(\alpha_j^*)_{j=1}^{n_d}\|_2^2/n_d}} \quad (188)$$

$$\omega_k \leftarrow \omega_k^* \times \sqrt{\frac{1e-3}{\|(\omega_j^*)_{j=1}^{n_d}\|_2^2/n_d}} \quad (189)$$

E.3.2 Figure 8 : Amplitude Factor

For the Random and the TINY methods of figure 8, we first normalize the parameters as :

$$\begin{aligned} &\text{For the new neurons} \\ \alpha_k^* &\leftarrow \alpha_k^* \times \frac{1}{\sqrt{\|(\alpha_j^*)_{j=1}^{n_d}\|_2^2/n_d}} \\ \omega_k^* &\leftarrow \omega_k^* \times \frac{1}{\sqrt{\|(\omega_j^*)_{j=1}^{n_d}\|_2^2/n_d}} \end{aligned}$$

$$\begin{aligned} &\text{For the best update} \\ \mathbf{W}^* &\leftarrow \mathbf{W}^* \times \frac{1}{\sqrt{\|\mathbf{W}^*\|_2^2/n_d}} \end{aligned}$$

Then, we multiply them by the amplitude factor γ^* :

$$\begin{aligned} &\text{For the new neurons :} \\ \alpha_k^*, \omega_k^* &\leftarrow \alpha_k^* \gamma^*, \omega_k^* \gamma^* \\ \gamma^* &:= \arg \min_{\gamma \in [-L, L]} \sum_i \mathcal{L}(f_{\theta \oplus \gamma \theta_{\leftrightarrow}^K}(\mathbf{x}_i), \mathbf{y}_i) \end{aligned}$$

$$\begin{aligned} &\text{For the best update :} \\ \mathbf{W}_l^* &\leftarrow \gamma^* \mathbf{W}_l \\ \gamma^* &:= \arg \min_{\gamma \in [-L, L]} \sum_i \mathcal{L}(f_{\theta + \gamma \mathbf{W}^*}(\mathbf{x}_i), \mathbf{y}_i) \end{aligned}$$

Where the operation $\gamma \theta_{\leftrightarrow}^{K*} = (\gamma \alpha_k^*, \gamma \omega_k^*)^K$ is the concatenation of the neural network with the new neurons and $\theta + \gamma \mathbf{W}^*$ is the update of one layer with its best update. The batch on which γ^* is computed is different from the one used to estimate the new parameters and its size is fixed to 1000 for all experiments.

E.4 Full algorithm

In this section we describe in detail the pseudo code to plot 5 and 8. The function $\text{NewNeurons}(l)$, in Algorithm 2, computes the new neurons defined at Proposition 3.2 for layer l sorted by decreasing eigenvalues. The function $\text{BestUpdate}(l)$, in Algorithm 4 computes the best update at Proposition 3.1 for layer l .

Algorithm 1: Algorithm to plot Figure 5 and 8.

```

1 for each method [TINY, MethodToCompareWith] do
2   Start from neural network  $N$  with initial structure  $s \in \{1/4, 1/64\}$ ;
3   while  $N$  architecture doesn't match ResNet18 do
4     for  $d$  in {depths to growth} do
5        $\theta_{\leftrightarrow}^{K^*} = \text{NewNeurons}(d, \text{method} = \text{method})$  ;
6       Normalize  $\theta_{\leftrightarrow}^{K^*}$  according to E.3;
7       Add the neurons at  $d$  ;
8       Train  $N$  for  $\Delta t$  epochs ;
9       Save model  $N$  and its performance ;
10    end
11  end
12 end

```

Algorithm 2: NewNeurons

Data: $l, \text{method} = \text{TINY}$

Result: Best neurons at l

```

1 if  $\text{method} = \text{TINY}$  then
2    $M = \text{BestUpdate}(l + 1)$ ;
3    $S, N = \text{MatrixSN}(l - 1, l + 1, M = M)$ ;
4   Compute the SVD of  $S := U\Sigma U^T$ ;
5   Compute the SVD of
       $U\sqrt{\Sigma}^{-1}UN := A\Lambda\Omega$ ;
6   Use the columns of  $A$ , the lines of  $\Omega$  and
      the diagonal of  $\Lambda$  to construct the new
      neurons of Prop. 3.2;
7 else if  $\text{method} = \text{GradMax}$  then
8    $M = \text{None}$  ;
9    $\_, N = \text{MatrixSN}(l - 1, l + 1, M = M)$  ;
10  Compute the SVD of  $N^T N$  ;
11  Use the eigenvectors to define the new
      out-going weights ;
12  Set the new in-going weight to 0;
13 else if  $\text{method} = \text{Random}$  then
14    $(\alpha_k, \omega_k)_{k=1}^{n_d} \sim \mathcal{N}(0, Id)$ ;
15 end

```

Algorithm 3: MatrixSN

Data: p_1, p_2 (layer indexes), $M = \text{None}$

Result: Construct matrices S and N

```

1 Take a minibatch  $X$  of size  $\propto \frac{(SW)^2}{P}$ ;
2 Propagate and backpropagate  $X$ ;
3 Compute  $V_{goal}$  at  $p_2$ , ie  $-\frac{\partial \mathcal{L}^{tot}}{\partial A_{p_2}}$ ;
4 if  $M \neq \text{None}$  then
5    $V_{goal-} = MB_{p_1}$ 
6 end
7  $S, N = B_{p_1} B_{p_1}^T, B_{p_1} V_{goal}^T$ ;

```

Algorithm 4: BestUpdate

Data: l , index of a layer

Result: Best update at l

```

1 Take a minibatch  $X$  of size  $\propto \frac{(SW)^2}{P}$ ;
2 Compute  $(S, N)$  with  $\text{MatrixSN}(l, l)$ ;
3  $M = N^T S^{-1}$ ;

```

E.5 Computational complexity

We estimate here the computational complexity of the above algorithm for architecture growth.

Theoretical estimate. We use the following notations:

- number of layers: L
- layer width, or number of kernels if convolutions: W (assuming for simplicity that all layers have same width or kernels)

- number of pixels in the image: P ($P = 1$ for fully-connected)
- kernel filter size: S ($S = 1$ if fully-connected)
- minibatch size used for standard gradient descent: M
- minibatch size used for new neuron estimation: M'
- minibatch size used in the line-search to estimate amplitude factor: M''
- number of classical gradients steps performed between 2 addition tentatives: T

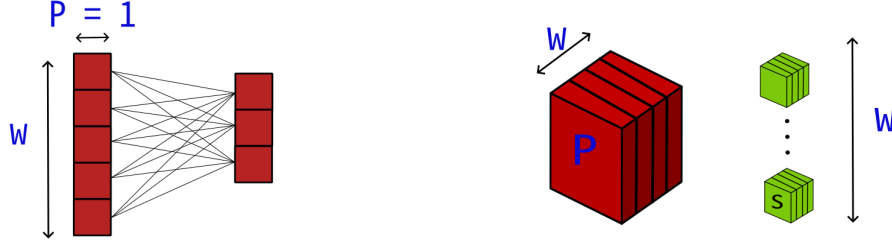


Figure 12: Notation and size for convolutional and linear layers

Complexity, estimated as the number of basic operations, cumulated over all calls of the functions:

- of the standard training part: $TMLW^2SP$
- of the computation of matrices of interest (function MatrixSN): $LM'(SW)^2P$
- of SVD computations (function NewNeurons): $L(SW)^3$
- of line-searches (function AmplitudeFactor): $L^2M''W^2SP$
- of weight updates (function BestUpdate): LSW

The relative added complexity w.r.t. the standard training part is thus:

$$M'S/TM + S^2W/TMP + M''L/TM + 1/WTMP.$$

SVD cost is negligible. The relative cost of the SVD w.r.t. the standard training part is S^2W/TMP . In the fully-connected network case, $S = 1$, $P = 1$, and the relative cost of the SVD is then W/TM . It is then negligible, as layer width W is usually much smaller than TM , which is typically 10×100 for instance. In the convolutional case, $S = 9$ for 3×3 kernels, and $P \approx 1000$ for CIFAR, $P \approx 100000$ for ImageNet, so the SVD cost is negligible as long as layer width $W \ll 10000$ or $1\,000\,000$ respectively. So one needs no worrying about SVD cost.

Likewise, the update of existing weights using the “optimal move” (already computed as a by-product) is computationally negligible, and the relative cost of the line searches is limited as long as the network is not extremely deep ($L < TM/M''$).

On the opposite, the estimation of the matrices (to which SVD is applied) can be more resource demanding. The factor $M'S/TM$ can be large if the minibatch size M' needs to be large for statistical significance reasons. One can show that an upper bound to the value required for M' to ensure estimator precision (see Appendix E.1) is $(SW)^2/P$. In that case, if $W > \sqrt{TMP/S^3}$, these matrix estimations will get costly. In the fully-connected network case, this means $W > \sqrt{TM} \approx 30$ for $T = 10$ and $M = 100$. In the convolutional case, this means $W > \sqrt{TMP/S^3} \approx 30$ for CIFAR and ≈ 300 for ImageNet. We are working on finer variance estimation and on other types of estimators to decrease M' and consequently this cost. Actually $(SW)^2/P$ is just an upper bound on the value required for M' , which might be much lower, depending on the rank of computed matrices.

In practice. In practice the cost of a full training with our architecture growth approach is similar (sometimes a bit faster, sometimes a bit slower) than a standard gradient descent training using the final architecture from scratch. This is great as the right comparison should take into account the number of different architectures to try in the classical neural architecture search approach. Therefore we get layer width hyper-optimization for free.

F Additional experimental results and remarks

F.1 ResNet18 on CIFAR-100

Figures. In all plots the black line represents the average performance over two independent runs, and the colored regions indicate the confidence interval.

technical details of figure 5 and 8 The experiment were performed on 1 GPU. The optimizer is SGD($lr = 1e-2$) with the starting batch size 32 E.2. At each depth l we set the number n_l of neurons to be added at this depth 2. These numbers do not depend on the starting architecture and have been chosen such that each depth will reach its final width with the same number of layer extensions. For the initial structure $s = 1/4$, resp. $1/64$, we set the number of layer extensions to 16, resp. 21, such that at depth 2 (named Conv2 in Table 3), $n_2 = (\text{Size}_2^{final} - \text{Size}_2^{start})/\text{nb of layer extensions} = (64 - 16)/16 = (64 - 1)/21 = 3$. The initial architecture is described in Table 3.

depth l	Conv2	Conv3	Conv5	Conv6	Conv8	Conv9	Conv11	Conv12
n_l	3	3	6	6	12	12	24	24

Table 2: Number of neurons to add per layer. The depth is identified by its name on tab 3

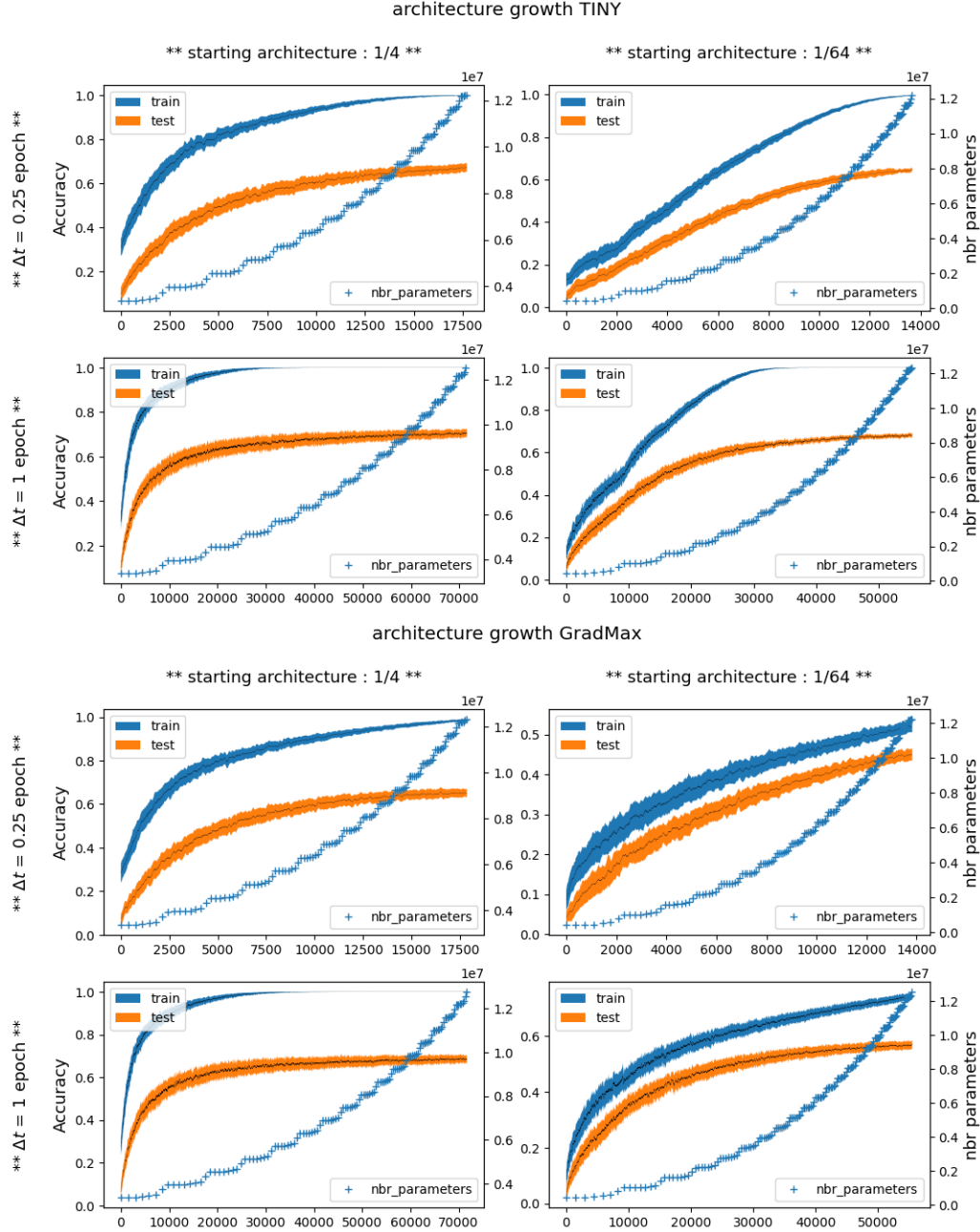


Figure 13: Accuracy and number of parameters during architecture growth for methods TINY and GradMax as a function of gradient step.

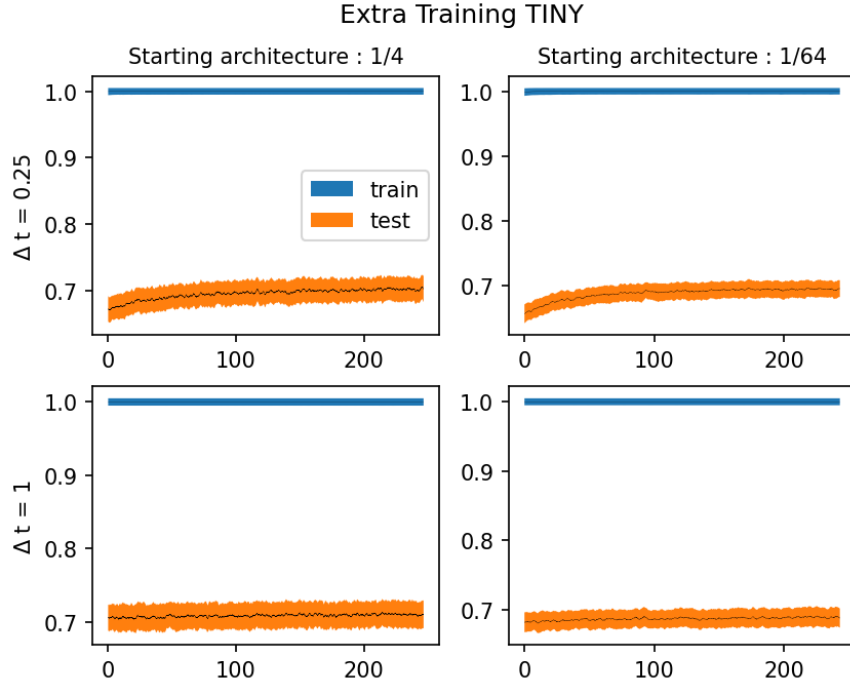


Figure 14: Accuracy as a function of the number of epochs during extra training for TINY.

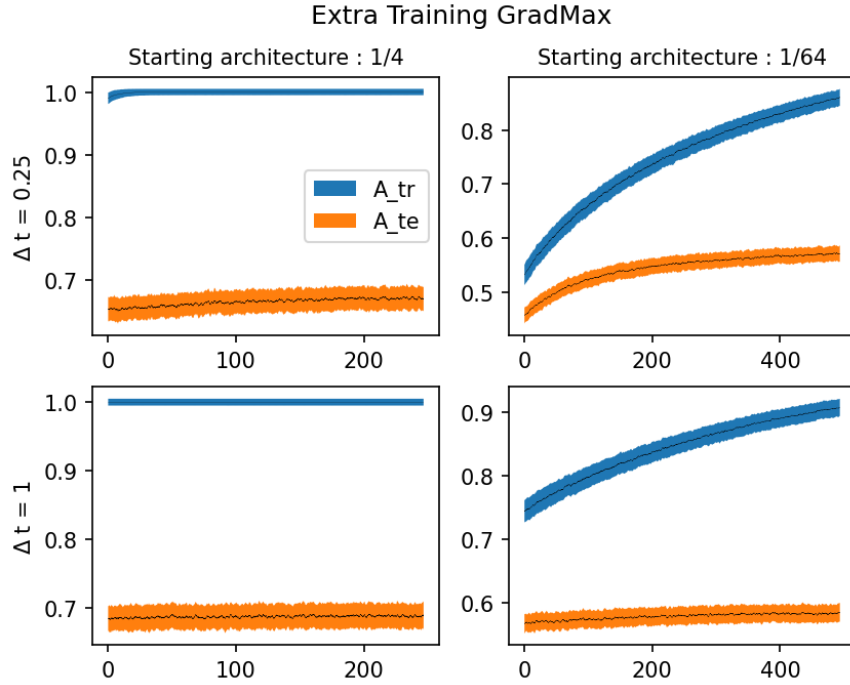


Figure 15: Accuracy curves as a function of the number of epochs during extra training for GradMax.

Table 3: Initial and final architecture for the models of Figure 5. Numbers in color indicate where the methods were allowed to add neurons (middle of ResNet blocks). In **blue** the initial structure for the model 1/64 and in **green** the initial structure for the model 1/4, ie **1/16** indicates that the model 1/64 started with 1 neuron at this layer while the model 1/4 started with 16 neurons at the same layer.

ResNet18					
Layer name	Output size	Initial layers (kernel=(3,3), padd.=1)		Final layers (end of Fig 5)	
Conv 1	$32 \times 32 \times 64$	$3 \times 3,$		$3 \times 3, 64$	
Conv 2	$32 \times 32 \times 64$	$3 \times 3, 64$ $3 \times 3, \text{1/16}$	$3 \times 3, \text{1/16}$ $3 \times 3, 64$	$3 \times 3, 64$ $3 \times 3, \text{64}$	$3 \times 3, \text{64}$ $3 \times 3, 64$
Conv 3	$32 \times 32 \times 64$	$3 \times 3, 64$ $3 \times 3, \text{1/16}$	$3 \times 3, \text{1/16}$ $3 \times 3, 64$	$3 \times 3, 64$ $3 \times 3, \text{64}$	$3 \times 3, \text{64}$ $3 \times 3, 64$
Conv 4	$16 \times 16 \times 64$	$3 \times 3, 128$		$3 \times 3, 128$	
Conv 5	$16 \times 16 \times 128$	$3 \times 3, 128$ $3 \times 3, \text{2/32}$	$3 \times 3, \text{2/32}$ $3 \times 3, 128$	$3 \times 3, 128$ $3 \times 3, \text{128}$	$3 \times 3, \text{128}$ $3 \times 3, 128$
Conv 6	$16 \times 16 \times 128$	$3 \times 3, 128$ $3 \times 3, \text{2/32}$	$3 \times 3, \text{2/32}$ $3 \times 3, 128$	$3 \times 3, 128$ $3 \times 3, \text{128}$	$3 \times 3, \text{128}$ $3 \times 3, 128$
Conv 7	$8 \times 8 \times 256$	$3 \times 3, 256$		$3 \times 3, 256$	
Conv 8	$8 \times 8 \times 256$	$3 \times 3, 256$ $3 \times 3, \text{4/64}$	$3 \times 3, \text{4/64}$ $3 \times 3, 256$	$3 \times 3, 256$ $3 \times 3, \text{256}$	$3 \times 3, \text{256}$ $3 \times 3, 256$
Conv 9	$8 \times 8 \times 256$	$3 \times 3, 256$ $3 \times 3, \text{4/64}$	$3 \times 3, \text{4/64}$ $3 \times 3, 256$	$3 \times 3, 256$ $3 \times 3, \text{256}$	$3 \times 3, \text{256}$ $3 \times 3, 256$
Conv 10	$4 \times 4 \times 512$	$3 \times 3, 512$		$3 \times 3, 512$	
Conv 11	$4 \times 4 \times 512$	$3 \times 3, 512$ $3 \times 3, \text{8/128}$	$3 \times 3, \text{8/128}$ $3 \times 3, 512$	$3 \times 3, 512$ $3 \times 3, \text{512}$	$3 \times 3, \text{512}$ $3 \times 3, 512$
Conv 12	$4 \times 4 \times 512$	$3 \times 3, 512$ $3 \times 3, \text{8/128}$	$3 \times 3, \text{8/128}$ $3 \times 3, 512$	$3 \times 3, 512$ $3 \times 3, \text{512}$	$3 \times 3, \text{512}$ $3 \times 3, 512$
AvgPool2d	$1 \times 1 \times 512$				
FC 1	100	512×100		256×100	
SoftMax	100				

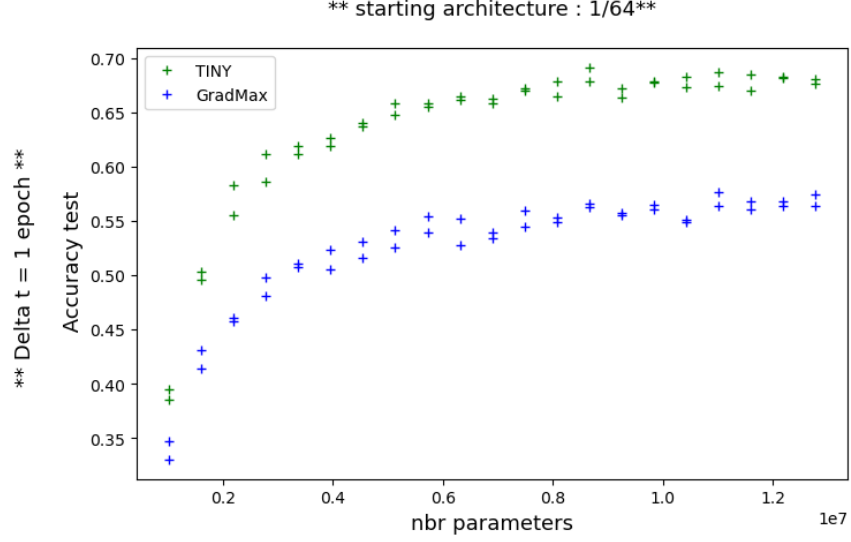


Figure 16: Accuracy on test of as a function of the number of parameters during architecture growth from ResNet_{1/64} to ResNet18. **The normalization for GradMax is $\sqrt{10^{-3}}$**

Δt s	TINY	GradMax	Baseline
	1	1	
1/64	68.0 ± 0.4	57.2 ± 0.3	72.9 ± 0.1 ^{5*}
1/64	69.0 ± 0.6 ^{5*}	57.7 ± 0.3 ^{3*}	

Table 4: Final accuracy on test of ResNet18 of 16 after the architecture growth (*grey*) and after convergence (*black*). The number of start indicated the multiple of 50 epochs needed to achieve convergence. With the starting architecture ResNet_{1/64} and $\Delta t = 1$ the method TINY achieves 68.0 ± 0.4 on test after its growth and it reaches **69.0 ± 0.6** ^{5*} after $* := 5 \times 50$ epochs.