

Spotting expressivity bottlenecks in neural networks and fixing them by optimal architecture growth

Manon Verbockhaven

Sylvain Chevallier, Guillaume Charpiat

TAU team, LISN, INRIA Saclay, Université Paris-Saclay

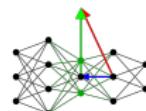
Rémi Gribonval

Hervé Luga

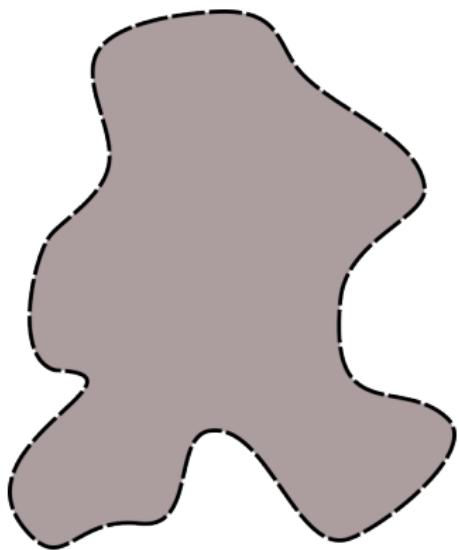
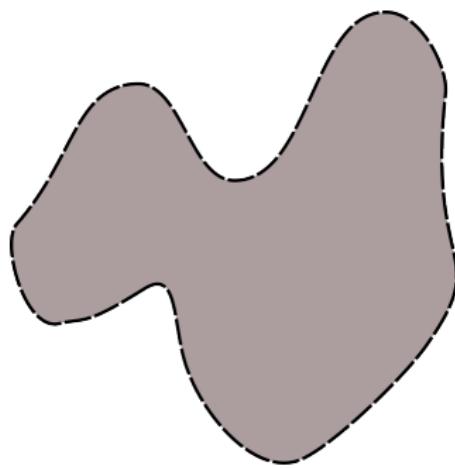
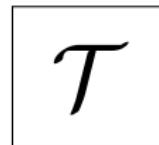
David Picard

Guillaume Lecué

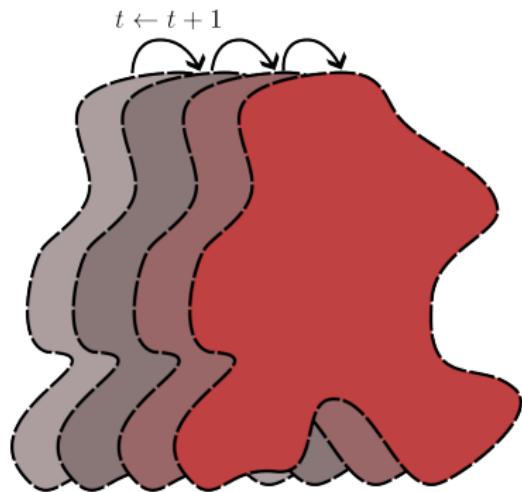
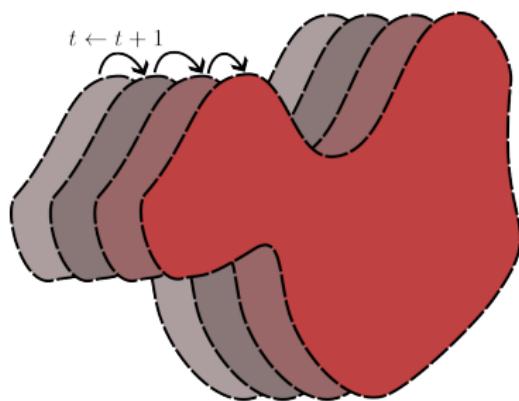
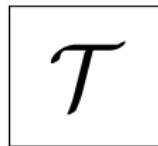
Aurélie Névéol



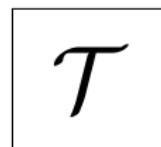
Introduction



Introduction



Introduction



eval ()

A red blob shape, which is irregular and roughly heart-shaped with a hole on the right side. It is enclosed in parentheses, suggesting it is the result of an evaluation.

eval ()

A red blob shape, which is more complex and elongated than the first one, with multiple lobes and a central indentation. It is enclosed in parentheses, suggesting it is the result of an evaluation.

Plan

- 1 Literature review on Neural Architecture Search (NAS)
 - The definitions
 - The One-Shot techniques
 - Some expressivity metrics
- 2 The expressivity bottleneck
 - The intuition and the definitions
 - The best parameter move and the best neurons to add
- 3 Reformulation and extension
 - Revisiting GradMax and NORTH
 - Move away from first order approximation

I Literature review on Neural Architecture Search (NAS)

- The definitions
- The One-Shot techniques
- Some expressivity metrics

NAS : the optimization problem

Consider a task \mathcal{T} and a corresponding dataset \mathcal{X} . Let \mathbb{A} a set of architecture and \mathcal{L} a loss function, we would like to solve :

$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(\mathcal{A}(\theta)) \quad (1)$$



$$\iff \arg \min_{\mathcal{A} \in \mathbb{A}} \mathcal{L}_{train}(\mathcal{A}(\theta^*)) \quad s.t. \quad \theta^* := \arg \min_{\theta} \mathcal{L}_{train}(\mathcal{A}(\theta))$$

NAS Optimization problem

$$\arg \min_{\mathcal{A} \in \mathbb{A}} \mathcal{L}_{val}(\mathcal{A}(\theta^*)) \quad s.t. \quad \theta^* := \arg \min_{\theta} \mathcal{L}_{train}(\mathcal{A}(\theta)) \quad (2)$$

NAS : a two-steps optimization procedure

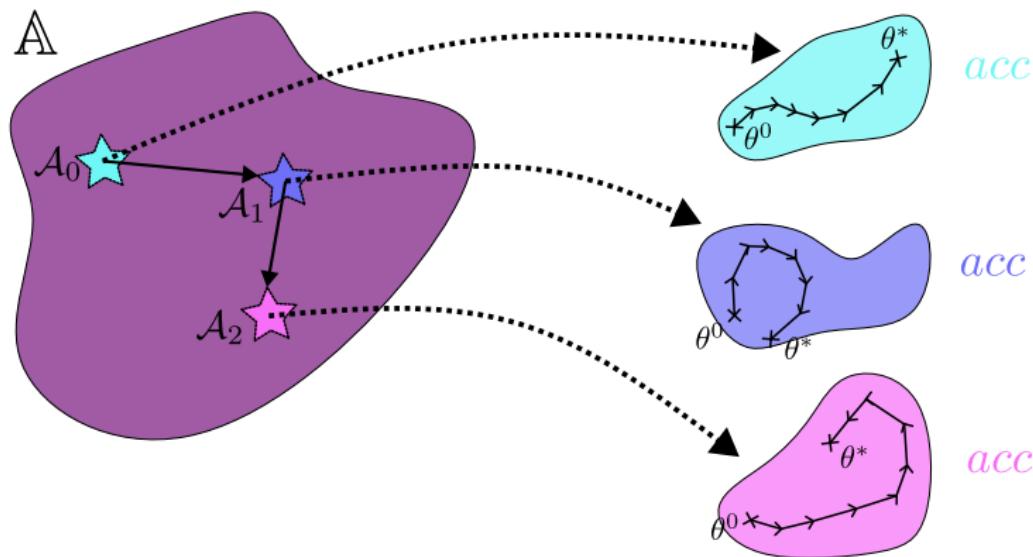
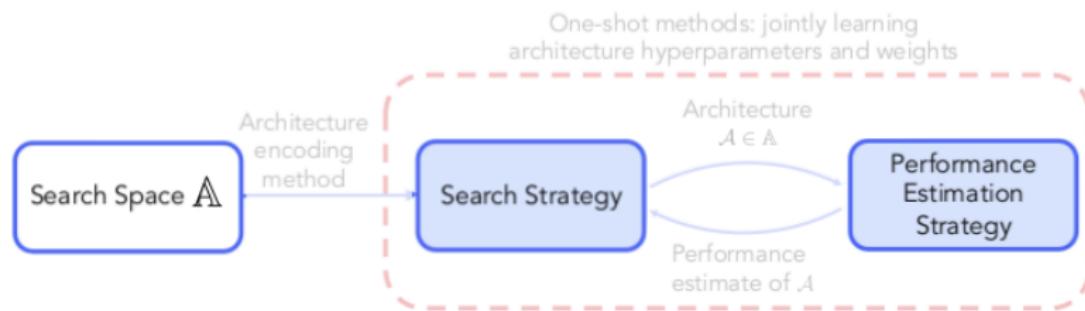
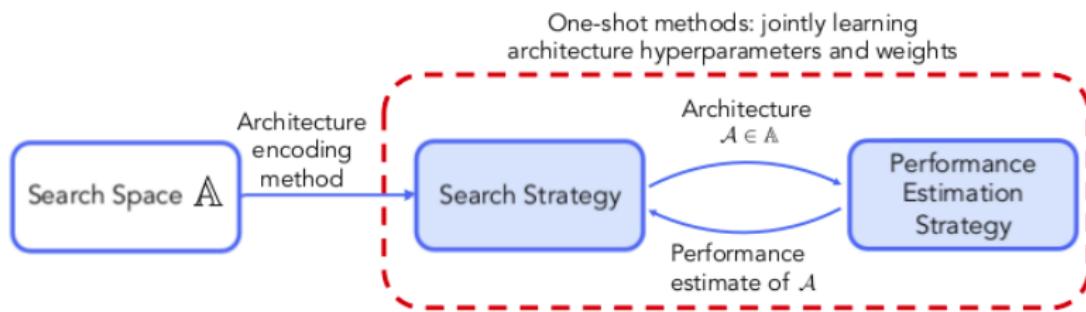


Figure: Optimization of a network architecture with the decoupling between the search into the space of architecture \mathbb{A} and the space of parameters $\mathbb{R}^{d(\mathcal{A})}$.

NAS : the definitions



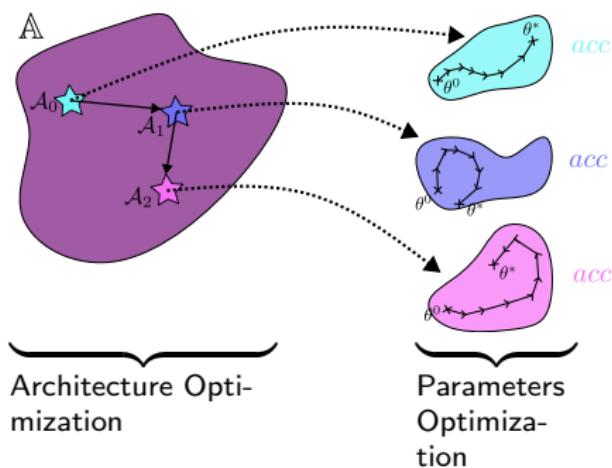
NAS : the definitions



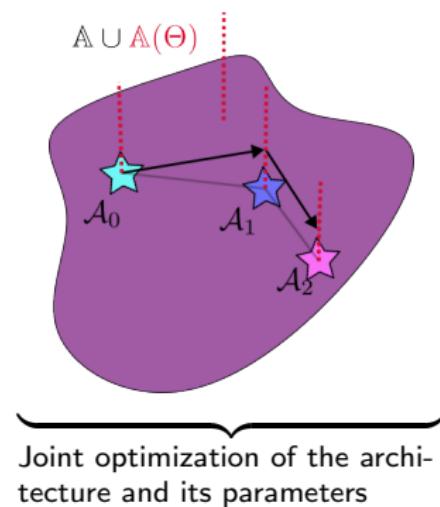
NAS : the search strategy

Probabilistic techniques

Search \rightleftarrows Evaluation



One-shot techniques



First formulation

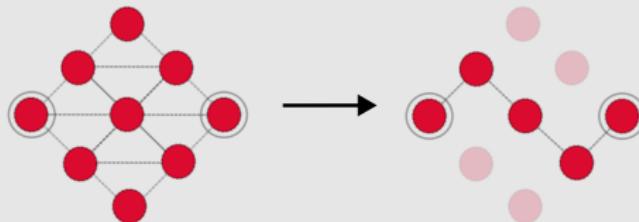
$$\arg \min_{\mathcal{A} \in \mathbb{A}} \mathcal{L}_{train}(\mathcal{A}(\theta^*)) \quad s.t. \quad \theta^* := \arg \min_{\theta} \mathcal{L}_{train}(\mathcal{A}(\theta))$$



$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(\mathcal{A}(\theta))$$

NAS : categories of one-shot techniques

Supernet



Gradient-based

A gradient-based method searches in \mathbb{A} using the back-products of the standard backpropagation (gradient estimation method to update network parameters) (Evci et al. [2022], Dai et al. [2019], Wu et al. [2020], Verbockhaven et al. [2024]).

Examples of gradient-based methods

GradMax

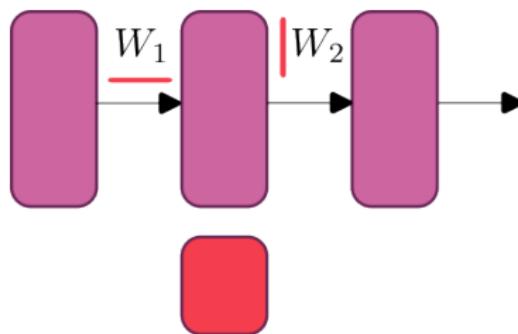


Figure: Adding one neuron (red) to a layer (pink) as in paper Evcı et al. [2022]

Firefly

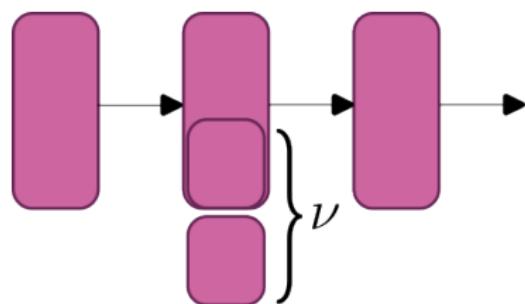
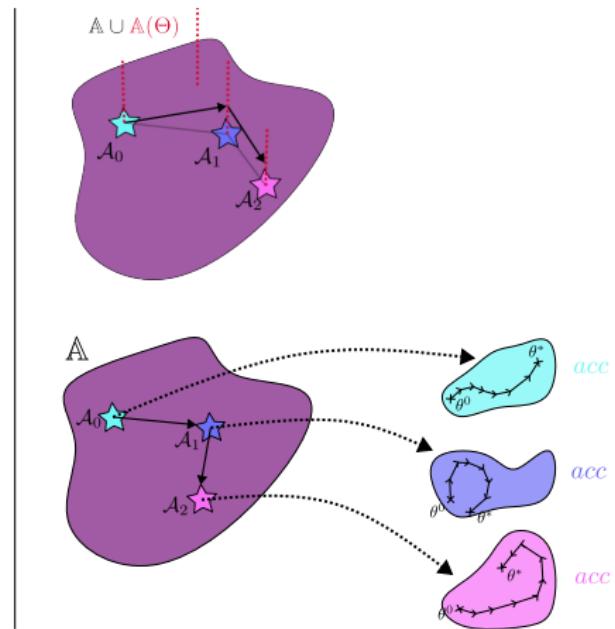


Figure: Dividing one neuron into two neurons (ν) as in paper Wu et al. [2020]

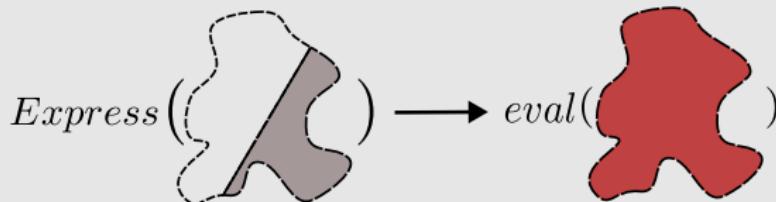
Orient the search

- Probabilistic techniques & One-Shot methods
↪ optimization in \mathbb{A} and θ



The expressivity metrics

Expressivity



$$\arg \min_{\mathcal{A} \in \mathbb{A}} \mathcal{L}_{val}(\mathcal{A}(\theta^*)) \quad s.t. \quad \theta^* := \arg \min_{\theta} \mathcal{L}_{train}(\mathcal{A}(\theta))$$

The expressivity metrics

Expressivity

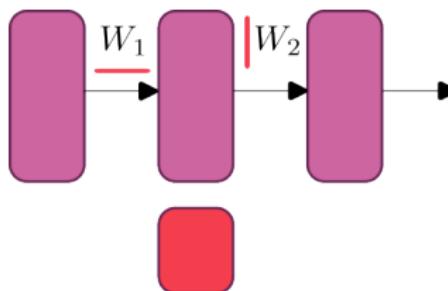
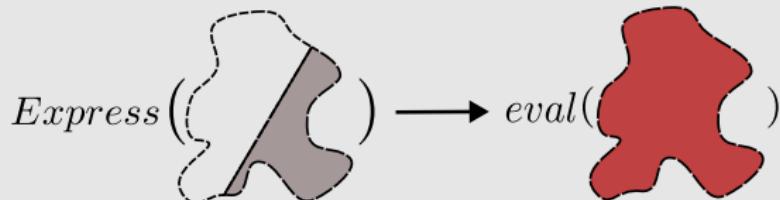
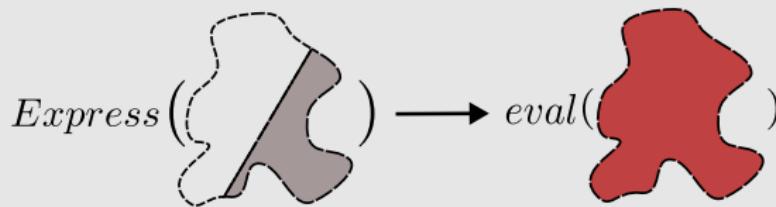


Figure: Adding one neuron (red) to a layer (pink) as in papers of Maile et al. [2022]

The expressivity metrics

Expressivity



	easy to compute	specific to the task \mathcal{T}
VC-dimension	✗	✗
Rademacher Complexity	✗	✓

Part II

I Literature review on Neural Architecture Search (NAS)

The definitions

The One-Shot techniques

Some expressivity metrics

II The expressivity bottleneck

- The intuition and the definitions
- The best parameter move and the best neurons to add

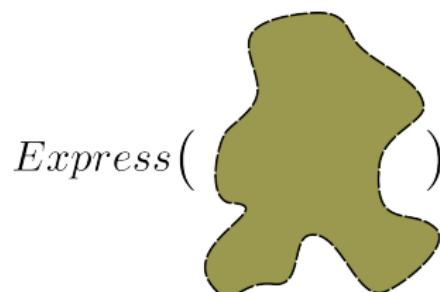
the problem statement

We would like to solve :

$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(\mathcal{A}(\theta))$$

- Jointly optimize the architecture and its parametrization.

With the use of :



Express(

- Easy to compute.
- Provides guidance in \mathbb{A} without redundancy.

Intuition

Notations

- Dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ iid
- Neural Network $f_\theta := \mathcal{A}(\theta)$
- Loss function \mathcal{L}

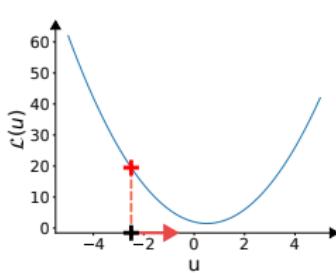
$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(f_\theta)$$

Intuition

Notations

- Dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ iid
- Neural Network $f_\theta := \mathcal{A}(\theta)$
- Loss function \mathcal{L}

$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(f_\theta)$$



Minimize $\mathcal{L} : U \rightarrow \mathbb{R}^+$ on $(U, \langle \cdot, \cdot \rangle_\dagger)$ with

$$u^{t+1} = u^t - \eta \nabla_u^\dagger \mathcal{L}(u)|_{u=u^t}$$

Where $\nabla_u^\dagger \mathcal{L}(u)$ verifies :

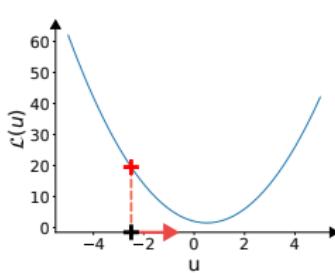
$$\mathcal{L}(u + \delta u) = \mathcal{L}(u) + \langle \nabla_u^\dagger \mathcal{L}(u), \delta u \rangle_\dagger + o(\|\delta u\|_\dagger)$$

Intuition

Notations

- Dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ iid
- Neural Network $f_\theta := \mathcal{A}(\theta)$
- Loss function \mathcal{L}

$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(f_\theta)$$



Minimize $\mathcal{L} : \Theta \rightarrow \mathbb{R}^+$ on $(\Theta, \langle \cdot, \cdot \rangle_{\mathbb{R}^p})$ with

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta}^{\mathbb{R}^p} \mathcal{L}(\theta)|_{\theta=\theta^t}$$

Where $\nabla_{\theta}^{\mathbb{R}^p} \mathcal{L}(\theta)$ verifies :

$$\mathcal{L}(\theta + \delta\theta) = \mathcal{L}(\theta) + \langle \nabla_{\theta}^{\mathbb{R}^p} \mathcal{L}(\theta), \delta\theta \rangle_{\mathbb{R}^p} + o(\|\delta\theta\|_{\mathbb{R}^p})$$

Intuition

Notations

- Dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ iid
- Neural Network $f_\theta := \mathcal{A}(\theta)$
- Loss function \mathcal{L}

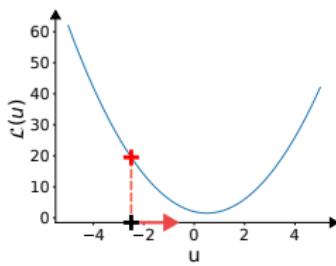
$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(f_\theta)$$

Minimize $\mathcal{L} : \mathcal{F} \rightarrow \mathbb{R}^+$ on $(\mathcal{F}, \langle \cdot, \cdot \rangle_{\textcolor{red}{L_2}})$ with

$$f^{t+1} = f^t - \eta \nabla_f^{\textcolor{red}{L_2}} \mathcal{L}(f)_{|f=f^t}$$

Where $\nabla_f^{\textcolor{red}{L_2}} \mathcal{L}(f)$ verifies :

$$\mathcal{L}(f + \delta f) = \mathcal{L}(f) + \langle \nabla_f^{\textcolor{red}{L_2}} \mathcal{L}(f), \delta f \rangle_{\textcolor{red}{L_2}} + o(\|\delta f\|_{\textcolor{red}{L_2}})$$

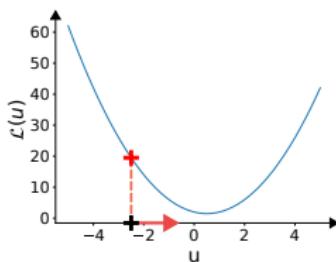


Notations

- Dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ iid
 - Neural Network $f_\theta := \mathcal{A}(\theta)$
 - Loss function \mathcal{L}

$$\arg \min_{\mathcal{A} \in \mathbb{A}, \theta \in \mathbb{R}^{d(\mathcal{A})}} \mathcal{L}_{train}(f_\theta)$$

Minimize $\mathcal{L} : \mathcal{F} \rightarrow \mathbb{R}^+$ on $(\mathcal{F}, \langle \cdot, \cdot \rangle_{\textcolor{red}{L_2}})$ with



$$f^{t+1} = f^t \underbrace{-\eta \nabla_f \mathcal{L}(f)}_{v\nabla}_{|f=f^t}$$

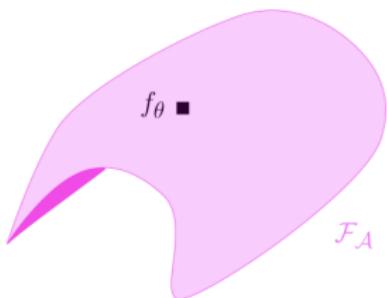
Where $\nabla_f^{\text{L}_2} \mathcal{L}(f)$ verifies :

$$\mathcal{L}(f + \delta f) = \mathcal{L}(f) + \left\langle \nabla_f^{\textcolor{red}{L_2}} \mathcal{L}(f), \delta f \right\rangle_{\textcolor{red}{L_2}} + o(\|\delta f\|_{\textcolor{red}{L_2}})$$

The expressivity bottleneck

Objects

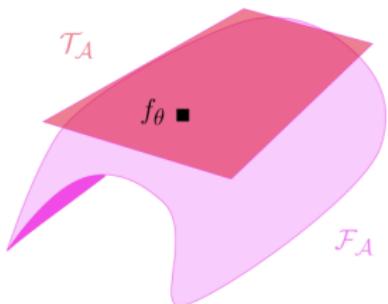
- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$



The expressivity bottleneck

Objects

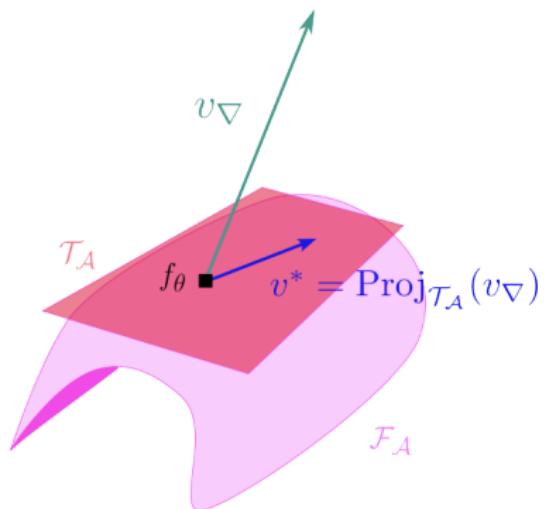
- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
- $\mathcal{T}_{\mathcal{A}}^{f_{\theta}}$ Tangent space in f_{θ} , ie
 $\left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta \theta \mid \text{s.t. } \delta \theta \in \Theta \right\}$



The expressivity bottleneck

Objects

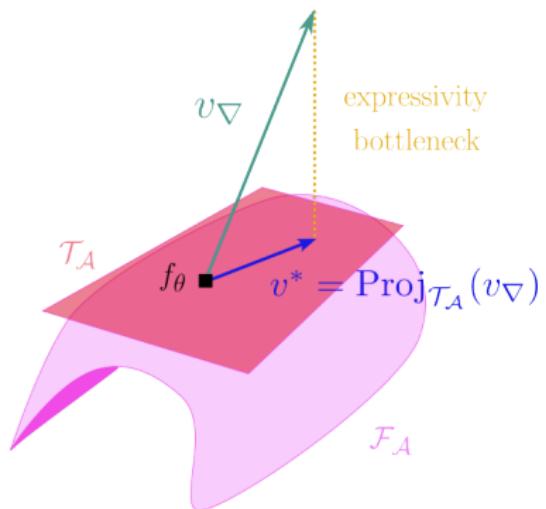
- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
- $\mathcal{T}_{\mathcal{A}}^{f_{\theta}}$ Tangent space in f_{θ} , ie
 $\left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta \theta \mid \text{s.t. } \delta \theta \in \Theta \right\}$



The expressivity bottleneck

Expressivity Bottleneck

$$\min_{\mathbf{v} \in \mathcal{T}_{\mathcal{A}}} \mathbb{E}_{(\mathbf{x}, \mathbf{y})} [\|\mathbf{v}_{\nabla}(\mathbf{x}) - \mathbf{v}(\mathbf{x})\|^2]$$



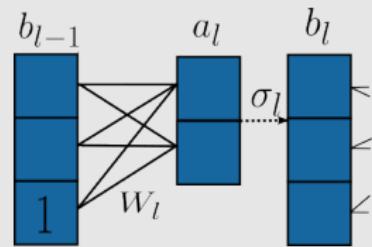
Definition (Feed-forward network)

We note $f_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^d$ a feed-forward network with L hidden layers. Formally $\theta := (\mathbf{W}_l)_{l=1\dots L}$, and the function f_θ iteratively computes :

$$\mathbf{b}_0(\mathbf{x}) = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}$$

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

$$f_\theta(\mathbf{x}) = \sigma_L(\mathbf{a}_L(\mathbf{x}))$$



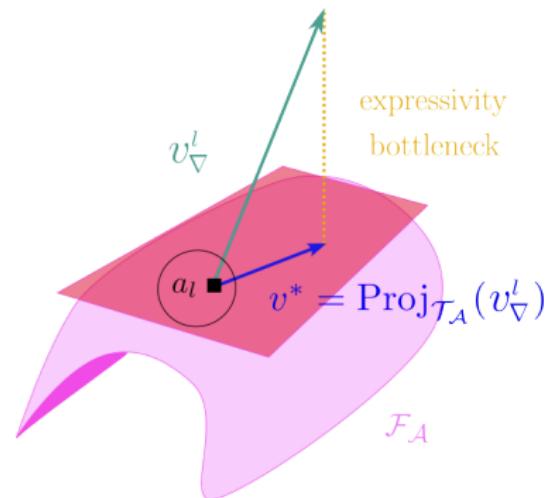
The expressivity bottleneck at layer l

Expressivity Bottleneck at l

$$\min_{\mathbf{v} \in \mathcal{T}_{\mathcal{A}}} \mathbb{E}_{(\mathbf{x}, \mathbf{y})} [\|\mathbf{v}_{\nabla}^l(\mathbf{x}) - \mathbf{v}^l(\mathbf{x})\|^2]$$

$$\Psi^l := \min_{\mathbf{v}^l \in \mathcal{T}_{\mathcal{A}}} \frac{1}{n} \|\mathbf{V}_{\nabla}^l - \mathbf{V}^l\|_{\text{Tr}}^2$$

$$\mathbf{V}_{\nabla} = (\mathbf{v}_{\nabla}(\mathbf{x}_1) \quad \dots \quad \mathbf{v}_{\nabla}(\mathbf{x}_n)) \text{ and } \mathbf{V} = (\mathbf{v}(\mathbf{x}_1) \quad \dots \quad \mathbf{v}(\mathbf{x}_n))$$



Updating the current architecture

$$\delta \mathbf{W}^* := \arg \min_{\delta \mathbf{W}^*} \frac{1}{n} \| \mathbf{V}_\nabla^l - \mathbf{V}^l(\delta \mathbf{W}) \|^2$$

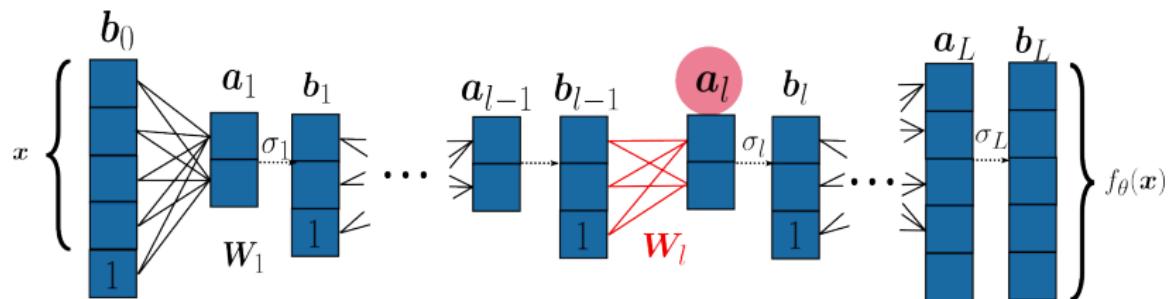


Figure: Approximating the expressivity bottleneck by modifying the parameters at layer l .

Adding neurons

$$\mathbf{A}^*, \Omega^*, K^* := \arg \min_{\mathbf{A}, \Omega, K} \frac{1}{n} \left\| \mathbf{V}_{\perp}^l - \mathbf{V}^l(\delta \mathbf{W}^*) \right\|^2$$

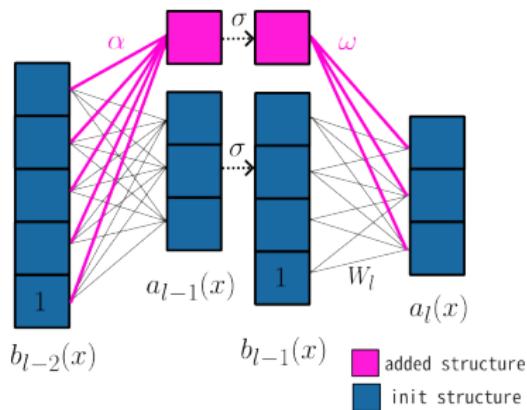
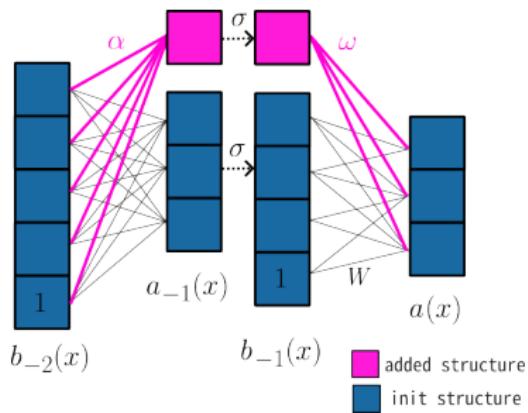


Figure: Adding one neuron to layer l in cyan, with connections in cyan. Here, $\alpha \in \mathbb{R}^5$ and $\omega \in \mathbb{R}^3$.

The optimization problems

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

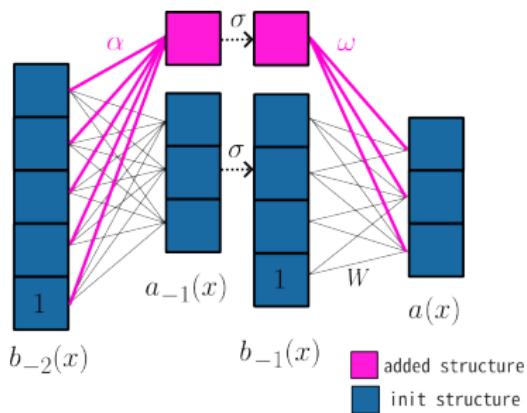
$$\mathbf{A}^*, \Omega^*, K^* = \arg \min_{K, \text{rk}(\Omega \mathbf{A}^T) = K} \frac{1}{n} \| \mathbf{V}_\perp - \mathbf{V}((\mathbf{A}, \Omega)) \|_{\text{Tr}}^2$$



The optimization problems

$$\delta \mathbf{W}^* = \arg \min_{\delta \mathbf{W}} \frac{1}{n} \| \mathbf{V}_{\nabla} - \delta \mathbf{W} \mathbf{B}_{-1} \|_{\text{Tr}}^2$$

$$\mathbf{A}^*, \boldsymbol{\Omega}^*, K^* = \arg \min_{K, \text{rk}(\boldsymbol{\Omega} \mathbf{A}^T) = K} \frac{1}{n} \| \mathbf{V}_{\perp} - \boldsymbol{\Omega} \mathbf{A}^T \mathbf{B}_{-2} \|_{\text{Tr}}^2$$



The optimization problems

$$\delta \mathbf{W}^* = \arg \min_{\delta \mathbf{W}} \frac{1}{n} \| \mathbf{V}_{\nabla} - \delta \mathbf{W} \mathbf{B}_{-1} \|_{\text{Tr}}^2$$

$$\mathbf{A}^*, \Omega^*, K^* = \arg \min_{K, \text{rk}(\Omega \mathbf{A}^T) = K} \frac{1}{n} \| \mathbf{V}_{\perp} - \Omega \mathbf{A}^T \mathbf{B}_{-2} \|_{\text{Tr}}^2$$

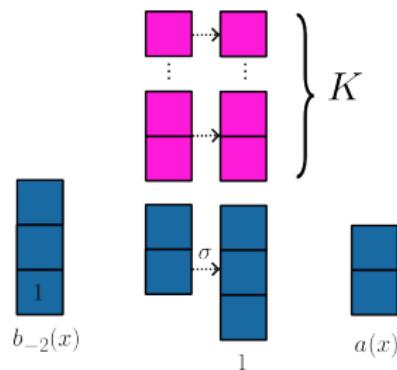
Best Update and new neurons

The optimal solution is $K^* = \text{rank}(\mathbf{V}_{\perp} \mathbf{B}_{-2}^T)$ and :

$$\delta \mathbf{W}^* = \mathbf{V}_{\nabla} \mathbf{B}_{-1}^T (\mathbf{B}_{-1} \mathbf{B}_{-1}^T)^+$$

$$\Omega^* \mathbf{A}^{*T} = \mathbf{V}_{\perp} \mathbf{B}_{-2}^T (\mathbf{B}_{-2} \mathbf{B}_{-2}^T)^+$$

The new neurons

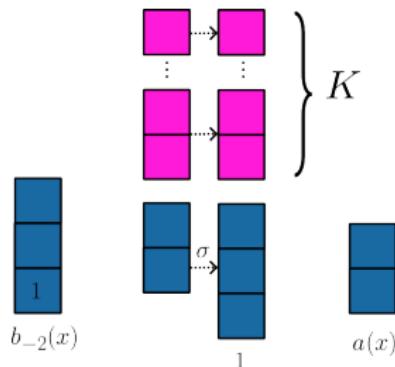


New neurons

The optimal solution is $K^* = \text{rank}(\mathbf{V}_\perp \mathbf{B}_{-2}^T)$ and :

$$\Omega^* \mathbf{A}^{*T} = \mathbf{V}_\perp \mathbf{B}_{-2}^T (\mathbf{B}_{-2} \mathbf{B}_{-2}^T)^+$$

The new neurons

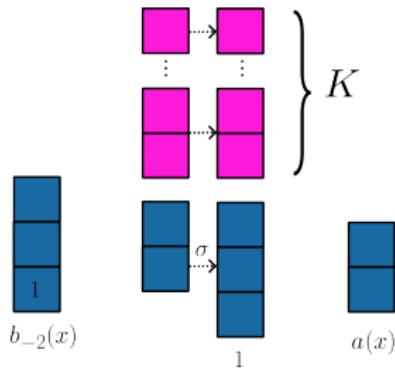


New neurons

Noting $\mathbf{S} := \frac{1}{n} \mathbf{B}_{-2} \mathbf{B}_{-2}^T$, $\mathbf{N}_\perp := \frac{1}{n} \mathbf{B}_{-2} \mathbf{V}_\perp^T$, the optimal solution is $K^* = \text{rank}(\mathbf{N}_\perp)$ and :

$$\Omega^* \mathbf{A}^{*T} = \mathbf{N}_\perp^T \mathbf{S}^+$$

The new neurons

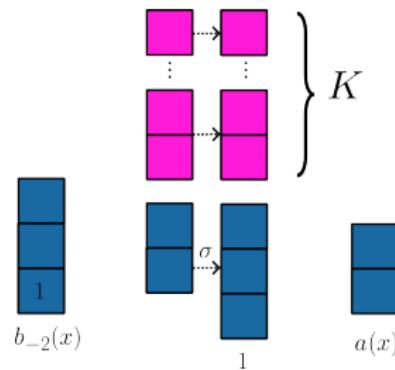


$$\mathbf{A}^*, \Omega^* = \underset{rk(\Omega \mathbf{A}^T) = K}{\arg \min} \frac{1}{n} \left\| \mathbf{V}_\perp - \Omega \mathbf{A}^T \mathbf{B}_{-2} \right\|_{\text{Tr}}^2$$

Theorem

Imposing $K \leq \text{rank}(\mathbf{N}_\perp)$, then the optimal solution is :

$$(\Omega_K^* \mathbf{A}_K^{*T}) = \left(\mathbf{N}_\perp^T \mathbf{S}^{-\frac{1}{2}} \right)_{|\text{Approx rk}(K)} \mathbf{S}^{-\frac{1}{2}}$$



Theorem : the expressivity bottleneck decrease

If the new neurons are initialised as previously, then noting
 $\lambda_1 \geq \dots \geq \lambda_K$ the K first eigenvalues of $N_\perp^T S^{-\frac{1}{2}}$, then :

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

Part III

I Literature review on Neural Architecture Search (NAS)

The definitions

The One-Shot techniques

Some expressivity metrics

II The expressivity bottleneck

The intuition and the definitions

The best parameter move and the best neurons to add

- Reformulation and extension

- Revisiting GradMax and NORTH
- Move away from first order approximation

Similar new neurons initialization methods

GradMax (Evci et al. [2022]) & NORTH (Maile et al. [2022])

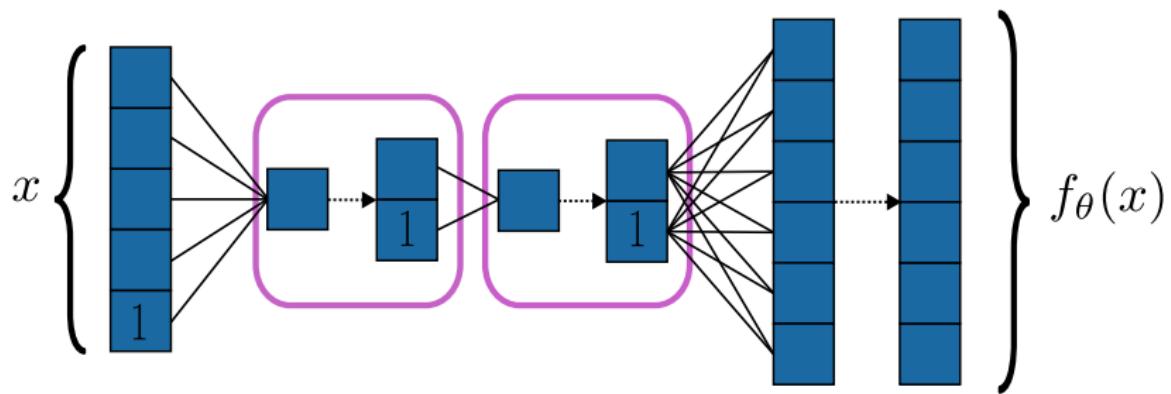
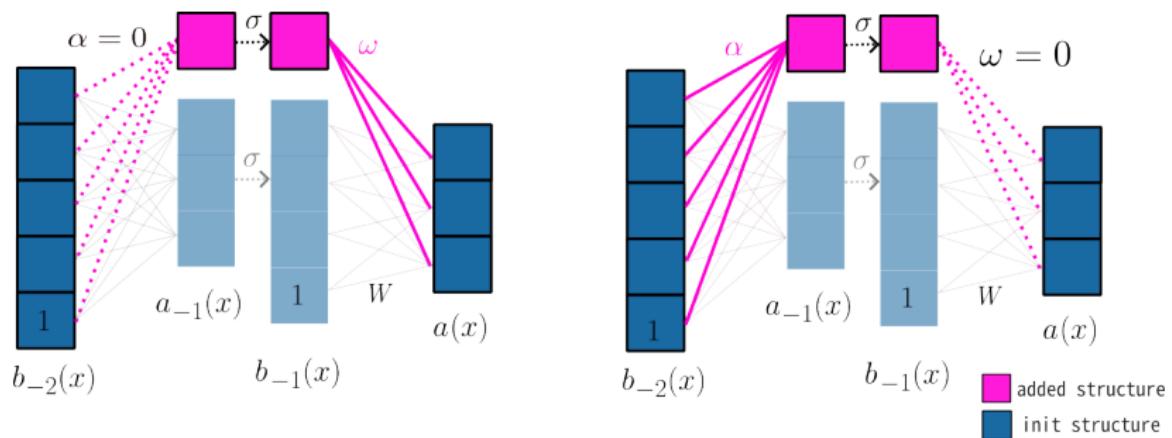


Figure: Thin feed forward network with two hidden layers.

GradMax & NORTH



GradMax initialization philosophy

objective : Accelerate the training of a network.

Updating the network parameters with $\delta\theta$:

$$\mathcal{L}^{t+1} = \mathcal{L}^t + \langle \nabla_{\theta} \mathcal{L}, \delta\theta \rangle + o(\|\delta\theta\|)$$

Choosing the gradient direction, i.e. $\delta\theta = -\eta \nabla_{\theta} \mathcal{L}$:

$$\mathcal{L}^{t+1} = \mathcal{L}^t - \eta \|\nabla_{\theta} \mathcal{L}\|^2 + o(\eta)$$

Idea : Maximizing the gradient according to the new parameters.

GradMax trick

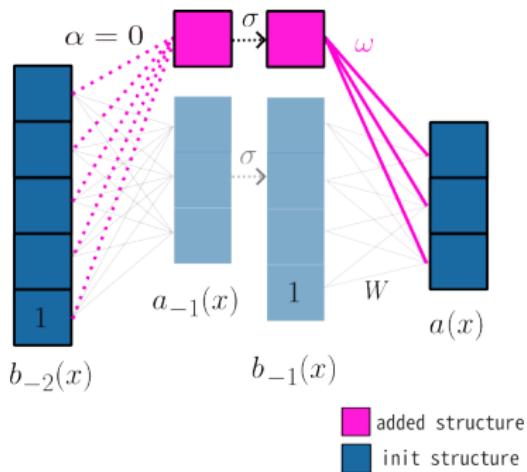
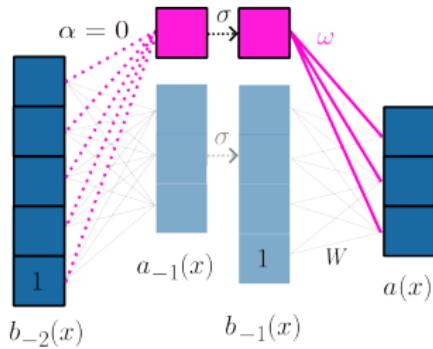


Figure: GradMax new neurons initialization.

$$\mathcal{L}^{t+1} \approx \mathcal{L}^t - \eta \|\nabla_{\theta}\mathcal{L}\|^2 - \eta \left(\underbrace{\|\nabla_{\Omega}\mathcal{L}\|^2}_{=0} + \underbrace{\|\nabla_{\mathbf{A}}\mathcal{L}\|^2}_{\text{quadratic in } \Omega} \right)$$

GradMax vs TINY



$$\mathbf{S} := \frac{1}{n} \mathbf{B}_{-2} \mathbf{B}_{-2}^T$$

$$N_{\perp} := \frac{1}{n} B_{-2} V_{\perp}^T$$

$$N_{\nabla} := \frac{1}{n} B_{-2} V_{\nabla}^T$$

GradMax	TINY
$\arg \max_{\Omega} \ \nabla_A \mathcal{L}(f_\theta)\ ^2 \text{ s.t. } \ \Omega\ \leq c$ <p style="text-align: center;">\Downarrow</p> $\arg \max_{\Omega} \ \textcolor{red}{N}_\nabla \Omega\ ^2 \text{ s.t. } \ \Omega\ ^2 \leq c$	$\arg \max_{A, \Omega} \langle A \Omega^T, N_\perp \rangle \text{ s.t. } \ \Omega A^T B\ ^2 \leq c$ <p style="text-align: center;">\Downarrow</p> $\arg \max_{\Omega} \ \textcolor{red}{N}_\perp \Omega\ _{\textcolor{blue}{S^{-1}}}^2 \text{ s.t. } \ \Omega\ \leq \tilde{c}$

GradMax & NORTH

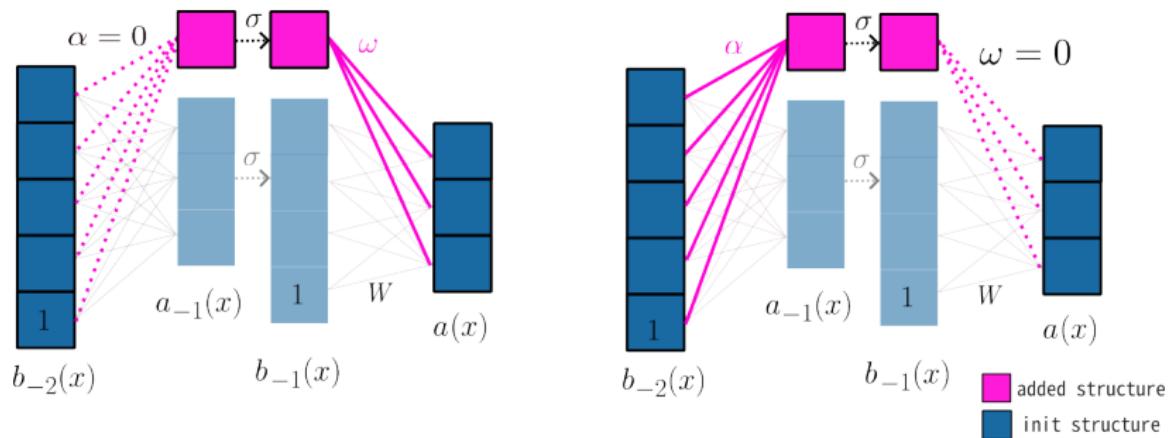
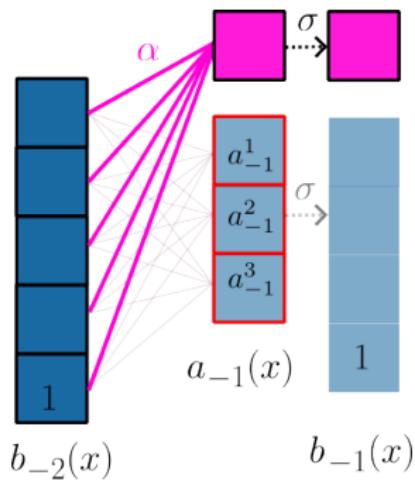


Figure: left : GradMax initialization; right : NORTH initialization

NORTH : activation-based



$$\mathbf{v}(\cdot, \alpha) \in \{\mathbf{a}_{-1}^1(\cdot), \mathbf{a}_{-1}^2(\cdot), \mathbf{a}_{-1}^3(\cdot)\}^\perp$$

NORTH vs TINY

$$S := \frac{1}{n} B_{-2} B_{-2}^T$$

NORTH	TINY
$S^+ B_{-2} V_{\text{Ker}(A_{-1})} r, \quad r \sim \mathcal{N}(0, I)$	$S^+ B_{-2} \overbrace{\mathcal{P}_{\text{Ker}(B_{-1})}}^{V_\perp^T} \overbrace{V_\nabla^T}^{V_i^T} v_i$

Conclusion : Initialization differences

Input weights A :

	NORTH	TINY
Init.	$S^+ B_{-2} V_{\text{Ker}(A_{-1})} r, \quad r \sim \mathcal{N}(0, I)$	$S^+ B_{-2} P_{\text{Ker}(B_{-1})} V_\nabla^T v_i$
Ortho. Grad.	✓ ✗	✓ ✓

Output weights Ω :

	GradMax	TINY
Init.	$\arg \max_{\Omega, \ \Omega\ ^2 \leq c} \ N_\nabla \Omega\ ^2$	$\arg \max_{\Omega, \ \Omega\ \leq \tilde{c}} \ N_\perp \Omega\ _{S^{-1}}^2$
Ortho. Grad.	✗ ✓	✓ ✓

Comparaison on CIFAR-100

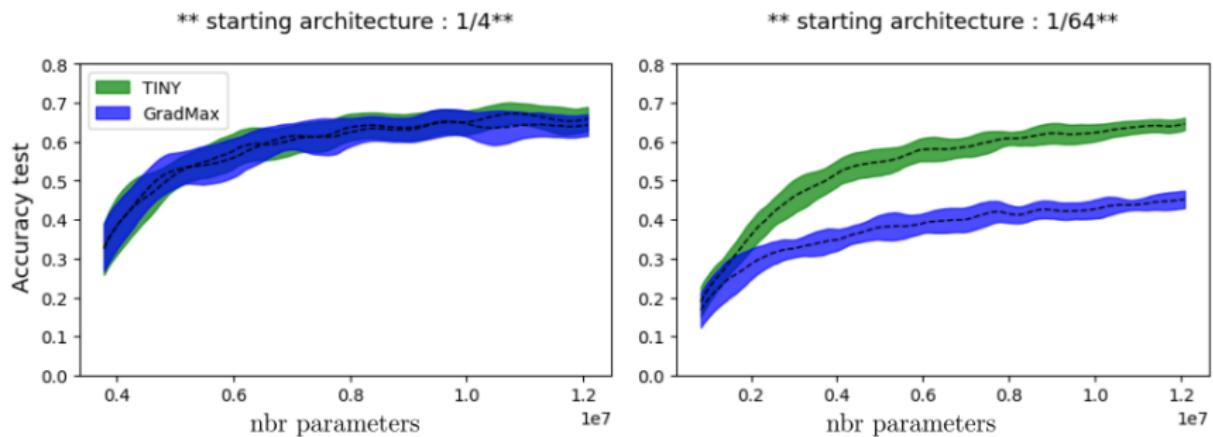


Figure: Test accuracy as a function of the number of parameters during architecture growth from ResNet_s to ResNet18; 4 runs are performed for each setting.

The amplitude factor (go faster !)

The amplitude factor γ

A new neuron with input and output weights α, ω is multiplied by the amplitude factor γ with the operation $\gamma(\alpha, \omega)$ as :

$$(\alpha, \omega) \leftarrow (\sqrt{\gamma}\alpha, \sqrt{\gamma}\omega).$$

For an input x , then :

$$v(x, (\alpha, \omega)) = \sqrt{\gamma}\omega\sigma(\sqrt{\gamma}\alpha^T b_{-2}(x)).$$

And for an input x , if $\|\gamma\| \ll 1$, then :

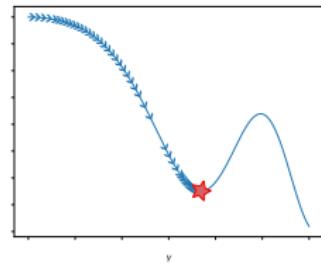
$$v(x, (\alpha, \omega)) = O(\|\gamma\|).$$

The best amplitude factor (go faster !)

The best amplitude factor γ^*

For a batch \mathbf{X}_γ , we can compute the best amplitude factor γ^* as :

$$\gamma^* := \arg \min_{\gamma \in [0, u]} \sum_{\mathbf{x}_i \in \mathbf{X}_\gamma} \mathcal{L}(f_{\theta \oplus \gamma(\mathbf{A}, \Omega)})(\mathbf{x}_i), \mathbf{y}_i)$$



For large γ^* :

$$\mathcal{L}(f_{\theta \oplus \gamma^*(\mathbf{A}, \Omega)}) = \mathcal{L}(f_\theta) - \gamma^* \frac{1}{n} \langle \mathbf{V}_\nabla, \mathbf{V}(\mathbf{A}, \Omega) \rangle_{\text{Tr}} + o(\gamma^*)$$

GradMax & NORTH : no impact

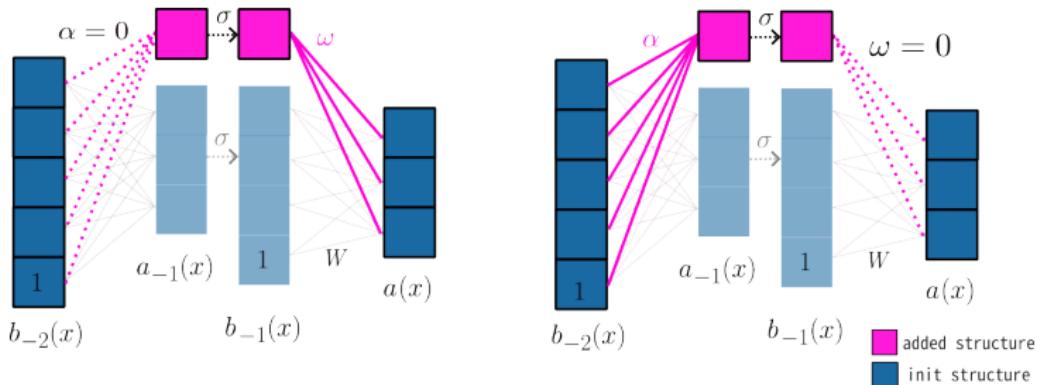


Figure: right : GradMax initialization; left : NORTH initialization

The amplitude factor for the strategy

Where to first add neurons ?

Using first order information, i.e. $\gamma \ll 1$:

$$l^* := \arg \max_l \lambda_1^l$$

Using non-linear information :

$$\Delta_l := \mathcal{L}(f_\theta) - \mathcal{L}(f_{\theta \oplus \gamma^*(A, \Omega)})$$

$$l := \arg \max_l \Delta_l$$

Experiment : Algorithm (2)

Algorithm 1: Grow a network with the strategy s

Data: $AddStrat = \text{depth selection strategy}$

Result: Neural Network N

Initialize a *thin* network with *few* neurons per layer;

for t in $[1, \dots, T]$ **do**

depth = $AddStrat(t)$;

Compute the first neuron α_1, ω_1 of at *depth*;

Compute the amplitude factor γ^* for this new neuron;

Multiply the neuron with its amplitude factor;

Add it to the architecture;

Evaluate the performances of the network;

for d in $[1, \dots, \text{deep}]$ **do**

Compute the best update at layer d and its amplitude factor;

Multiply the best update with its amplitude factor;

Update the architecture with this best update;

Evaluate the performances of the network;

end

end

} Grow
} Train

Experiment : Algorithm (2)

Algorithm 2: Select depth to grow the network

Data: s = Strategy, t = iteration

Result: $depth$ where to add a neuron

```
if  $s = \text{Naive}$  then
     $depth = np.mod(t, deep) + 1;$ 
else if  $s = \text{MaxEigenValue}$  then
    for  $d$  in  $[1, \dots, deep]$  do
        Compute the first neuron  $\alpha_1, \omega_1$  at  $d$ ;
        Compute the first order information  $\lambda_1^d = \langle \mathbf{V}, \mathbf{V}_\perp \rangle$  at  $d$ ;
    end
     $depth = \arg \max_{d=1, \dots, deep} \{\lambda_1^d\};$ 
else if  $s = \text{AmplitudeFactor}$  then
    for  $d$  in  $[1, \dots, depth]$  do
        Compute the first neuron  $\alpha_1, \omega_1$  at  $d$ ;
        Compute the amplitude factor  $\gamma^*$  for this new neuron;
        Evaluate the decrease of loss  $\Delta_d$  if the architecture is expanded with
         $(\sqrt{\gamma^*} \alpha_1, \sqrt{\gamma^*} \omega_1)$ 
    end
     $depth = \arg \max_{d=1, \dots, deep} \{\Delta_d\};$ 
return  $depth$ ;
```

Experiment (2)

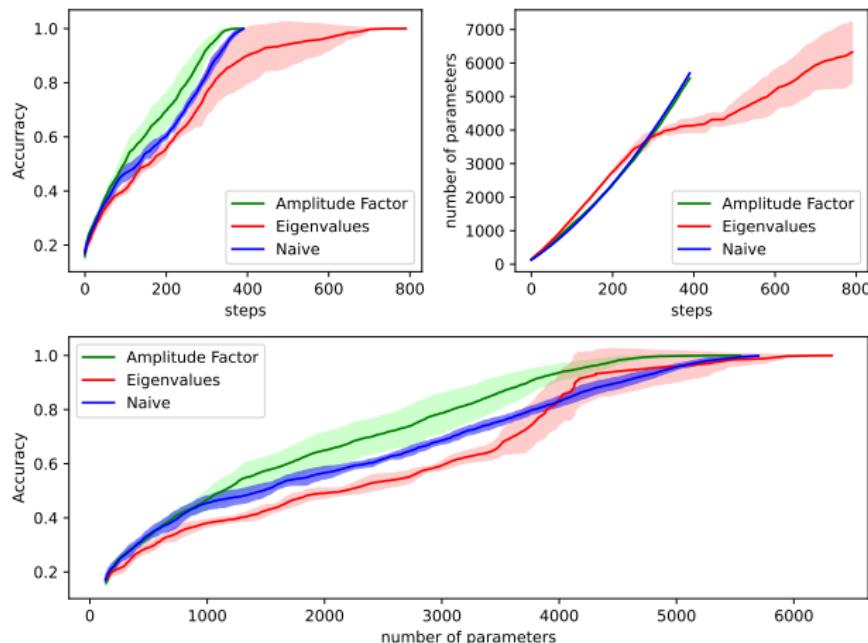


Figure: Accuracy and number of parameters of growing networks with different addition strategies for $\mathbf{x} \sim \mathcal{N}(0, I_{50})$ and $\mathbf{y} = \arg \max_{k=1, \dots, 10} \left\{ \sin(k \sum_{j=1}^{50} \mathbf{x}[j]) \right\}$.

Conclusion

- The expressivity bottleneck to update the parameters and initialize the new neurons
- First order information to select the new neurons at a given layer
- Formal comparison between TINY, GradMax and NORTH
- The amplitude factor to select the depth at which neurons should be added first and build a strategy across layers.

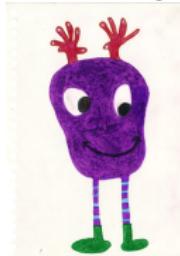
Open-source TINYpub, statistical significance of the estimators, MNIST, CIFAR10

Growing Tiny Networks:
Spotting Expressivity Bottlenecks and Fixing Them Optimally Verbockhoven et al. [2024]

Growth strategies for arbitrary DAG neural architectures Douka et al. [2025]

Neural Architecture Search

Suppose that God's apprentices would like to have a personal chef and decide to create it. To do so, they need to *assemble parts of the human body* into a human puppet, which will be *trained* to cook. Considering all possible layouts of the body sections, they try a first idea and place the hands on top of the head. After several hours of training, they observe, quite embarrassed, that their puppet is incapable of cooking the simplest recipes. One apprentice suggests that this failure might be linked to the puppet's inability to see its actions while using its hands. Considering this remark, they decide to revise the joining of the puppet and restart the training...



Replace the action *assemble the parts of the human body* by *construct the architecture of a neural network*, and the apprentices become a computer program of some researchers in the NAS field.

À Á Â Ã Ä Å

Ł Ł Ł Ł Ł Ł Ł
≤ 1 1 ł ł | l 1

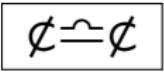
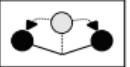
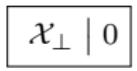
Method	Indicators				Dataset
	C_0	C_∞	T	P	
DARTS 	×	$3.3e6$	$1.5 + 2.5$	97.0 ± 0.1	CIFAR-10
	×	$3.3e6$	$4 + 2.5$	97.8 ± 0.1	
	×	$4.7e6$	4	73.3	ImageNet
Firefly 	$1.4e6$	$1.4e7$	$\approx 1 ?$	~ 94	CIFAR-10
	×	$5.8e5$	×	~ 71	CIFAR-100
GradMax 	$1e6$	$9e6$	×	84.4 ± 0.4	CIFAR-10
	$1e5$	$4e5$	×	91.1 ± 0.1	
	$1e5$	$4e5$	×	66.8 ± 0.2	CIFAR-100 ImageNet
	$1e6$	$4.2e6$	×	68.6 ± 0.2	
NORTH 	$1e6$	$3e7$	3.3	~ 76.2	CIFAR-10
	$3e4$	$3e5$	0.4	~ 82.5	
	$1e6$	$3e7$	1.6	~ 55.2	CIFAR-100

Figure: T : GPU days, P : accuracy on test (%) by-hand mean on all tested strategy, see figure 4 of paper Maile et al. [2022] for more precise results.

NAS : the search space

Search Space \mathbb{A}

$$\mathbb{A} := \{\text{macro}(\text{micro}) \mid \text{micro} \in S, \text{ macro} \in S'\} \quad (3)$$

A layer, noted , is of the form
 $x \rightarrow \sigma(Wx)$.

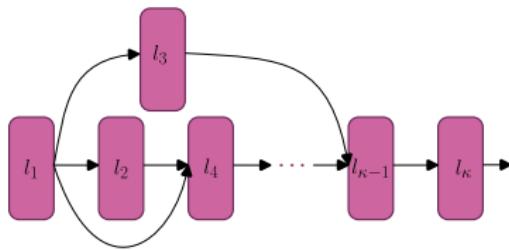


Figure: Example of a chain-structured architecture.

NAS : the search space

Search Space \mathbb{A}

$$\mathbb{A} := \{\text{macro}(\text{micro}) \mid \text{micro} \in S, \text{ macro} \in S'\} \quad (3)$$

A layer, noted , is of the form $x \rightarrow \sigma(Wx)$.

A cell, noted  ..., is a composition of layers.

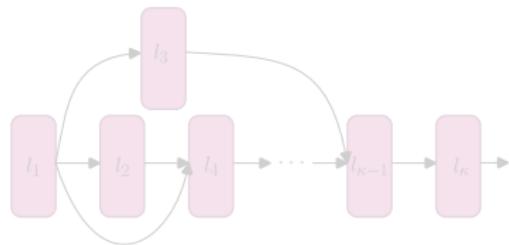


Figure: Example of a chain-structured architecture.

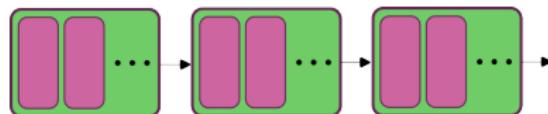


Figure: Example of a cell-based architecture.

Definition

- The dataset $\mathcal{D} := \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$
- the neural network or predictor f
- the class of functions $\mathcal{C} := \{f \text{ with architecture } \mathcal{A}\}$

Expr. & Compl. : noise fitting

$$\mathcal{C} := \{f \text{ for } f \in \mathcal{F}\} \quad \text{and} \quad f(\mathbf{x}) \in \{a, b\}$$

The Vapnik–Chervonenkis (VC) dimension

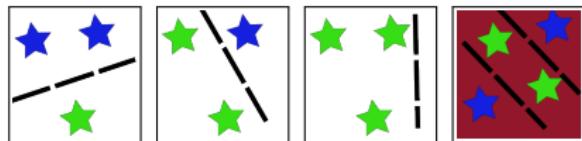


Figure: Example of VC-dimension when \mathcal{F} is the class of linear separator in a two-dimensional plane.

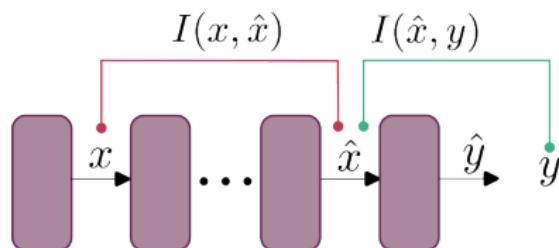
Rademacher complexity

Let $\sigma \sim \mathcal{B}(\frac{1}{2})$, then :

$$\mathbb{E}_{\mathcal{D}} \left[\frac{1}{n} \widehat{\mathbb{E}}_{\sigma} \left[\sup_{f \in \mathcal{C}} \sigma^T f(\mathbf{X}) \mid \mathcal{D} \right] \right]$$

Expr. & Compl. : minimization criterion

the Information Bottleneck



$$IB(f) := I(\mathbf{x}, \hat{\mathbf{x}}) - \beta I(\hat{\mathbf{x}}, \mathbf{y}).$$

the Kolmogorov Complexity

0 1 1 0 0 0 ... 1 1


shortest binary program generating f_θ

$H(\mathbf{x})$ the Shannon entropy, then $I(\mathbf{x}; \hat{\mathbf{x}}) := H(\mathbf{x}) - H(\mathbf{x}|\hat{\mathbf{x}})$.

From TINY to GradMax

First order development of the loss :

$$\begin{aligned}\mathcal{L}(f_{\theta \oplus (\mathbf{A}, \boldsymbol{\Omega})}) &:= \frac{1}{n} \sum_i \ell(f_{\theta}(\mathbf{x}_i)) \\ &= \mathcal{L}(f_{\theta}) - \frac{1}{n} \langle \mathbf{V}, \mathbf{V}_{\nabla} \rangle_{\text{Tr}} + o\left(\frac{1}{n} \|\mathbf{V}\|\right)\end{aligned}$$

The expressivity bottleneck :

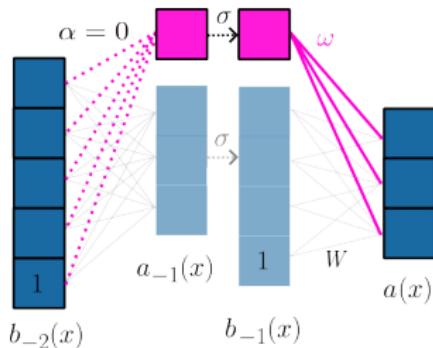
$$\Psi = \min_{\delta \theta} \frac{1}{n} \|\mathbf{V}((\mathbf{A}, \boldsymbol{\Omega})) - \mathbf{V}_{\perp}\|_{\text{Tr}}^2$$

Proposition

$\forall \mathbf{D} \in \mathbb{R}^{p,q}, \mathbf{B} \in \mathbb{R}^{k,q},$

$$\exists \tilde{c} \in \mathbb{R} \quad \text{s.t.} \quad \arg \min_{\mathbf{H}} \|\mathbf{D} - \mathbf{H}\mathbf{B}\|^2 = \arg \max_{\mathbf{H}, \|\mathbf{H}\mathbf{B}\|^2 \leq \tilde{c}} \langle \mathbf{D}, \mathbf{H}\mathbf{B} \rangle$$

GradMax initialization



$$\begin{aligned} S &:= \frac{1}{n} B_{-2} B_{-2}^T \\ N_{\perp} &:= \frac{1}{n} B_{-2} V_{\perp}^T \\ N_{\nabla} &:= \frac{1}{n} B_{-2} V_{\nabla}^T \end{aligned}$$

GradMax initialization :

$$\Omega^* := \arg \max_{\Omega} \|N_{\nabla} \Omega\|^2 \quad s.t. \quad \|\Omega\|^2 \leq c \quad (4)$$

TINY initialization :

$$\exists \tilde{c} \quad \Omega^* := \arg \max_{\Omega} \|N_{\perp} \Omega\|_{S^{-1}}^2 \quad s.t. \quad \|\Omega\| \leq \tilde{c} \quad (5)$$

Experiment : Algorithm (1)

Algorithm 3: Grow a network with *InitStrat*.

Data: *InitStrat* = Initialization strategy

Result: Neural Network *N*

Initialize a *thin* network with *few* neurons per layer;

for *t* in $[1, \dots, T]$ **do**

depth = *np.mod(t, deep)* + 1;

Using *InitStrat*, compute the first neuron α_1, ω_1 of at *depth*;

Compute the amplitude factor γ^* for this new neuron;

Multiply the neuron with its amplitude factor;

Add it to the architecture;

Evaluate the performances of the network;

end

}

Grow

Experiment (1) : Random vs TINY

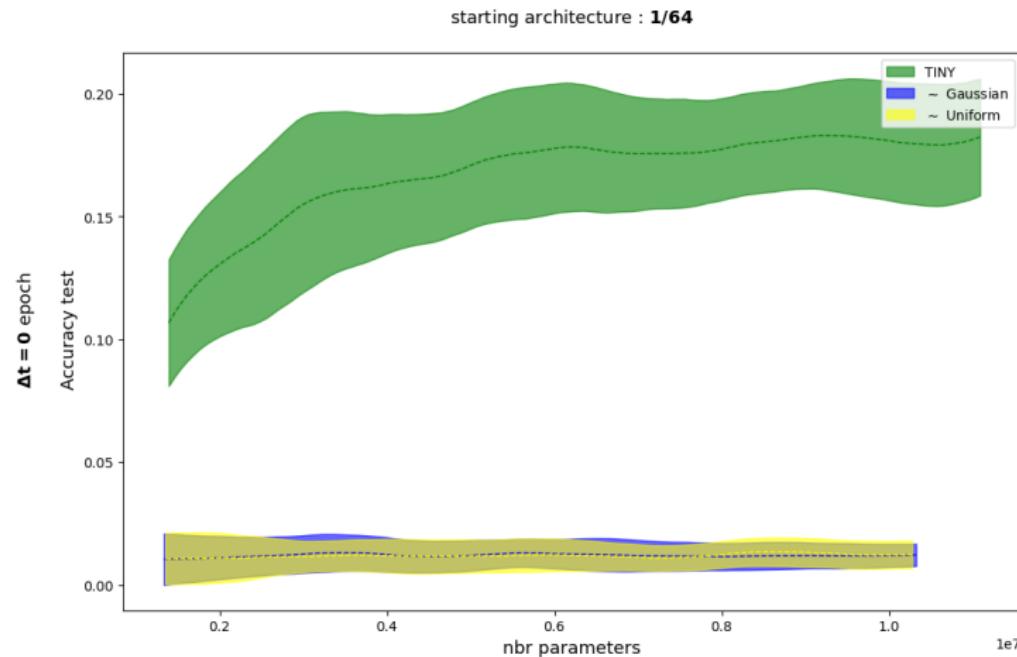


Figure: Test accuracy as a function of the number of parameters during pure architecture growth (no gradient steps performed, only neuron addition with a given initialization) from ResNet_{1/64} to ResNet₁₈, averaged over four independent runs.

Xiaoliang Dai, Hongxu Yin, and Niraj K. Jha. Nest: A neural network synthesis tool based on a grow-and-prune paradigm. *IEEE Transactions on Computers*, 68(10):1487–1497, 2019. doi: 10.1109/TC.2019.2914438.

Stella Douka, Manon Verbockhaven, Théo Rudkiewicz, Stéphane Rivaud, François P Landes, Sylvain Chevallier, and Guillaume Charpiat. Growth strategies for arbitrary dag neural architectures, 2025. URL <https://arxiv.org/abs/2501.12690>.

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.

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>.



Manon Verbockhaven, Théo Rudkiewicz, Guillaume Charpiat, and Sylvain Chevallier. Growing tiny networks: Spotting expressivity bottlenecks and fixing them optimally. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=hbtG6s6e7r>.

Lemeng Wu, Bo Liu, Peter Stone, and Qiang Liu. Firefly Neural Architecture Descent: a General Approach for Growing Neural Networks. In *Advances in Neural Information Processing Systems*, volume 33, pages 22373–22383. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/fdbe012e2e11314b96402b32c0df26b7-Abstract.html>.