



WEIGHT REPARAMETRIZATION FOR BUDGET-AWARE NETWORK PRUNING

ICIP 2021

Robin Dupont - *Sorbonne Université & Netatmo*

Hichem Sahbi - *Sorbonne Université*

Guillaume Michel - *Netatmo*



1 PRUNING AND SPARSITY

2 OUR METHOD

3 RESULTS

4 SUM UP

PRUNING AND SPARSITY

Pruning - Overview

Pruning

- Group of methods that aim to design **lightweight architectures** (*small memory footprint or fast inference time*)

Pruning - Overview

Pruning

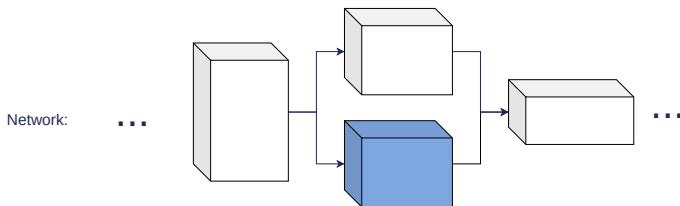
- Group of methods that aim to design **lightweight architectures** (*small memory footprint or fast inference time*)
- Introduce **sparsity** by removing **redundant** or **unnecessary weights** of the network

Pruning - Overview

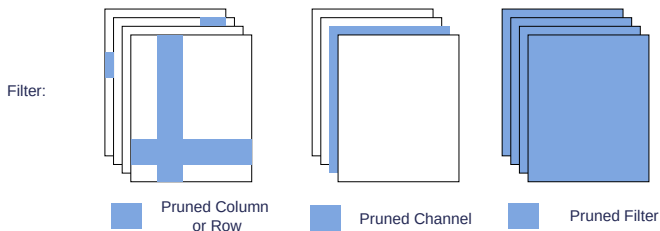
Pruning

- Group of methods that aim to design **lightweight architectures** (*small memory footprint or fast inference time*)
- Introduce **sparsity** by removing **redundant** or **unnecessary weights** of the network
- Can be applied at **different granularity** (*fine-grained vs coarse-grained*)

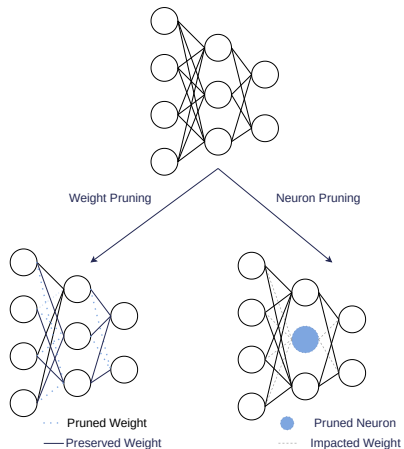
Pruning - Coarse-grained



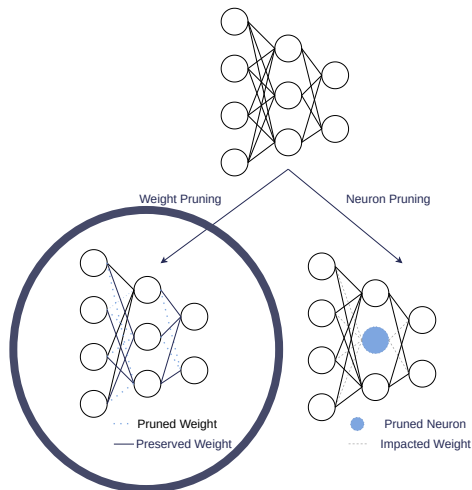
■ Pruned Subnetwork



Pruning - Fine-grained



Pruning - Fine-grained

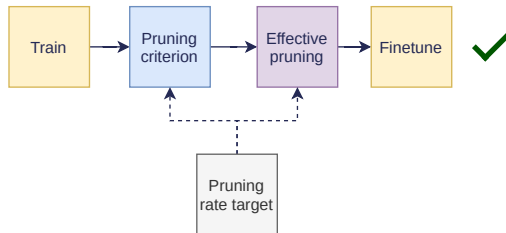


Our method belongs to the the **fine-grained weight pruning** category.

Pruning Pipelines

Pruning - Standard Pipeline

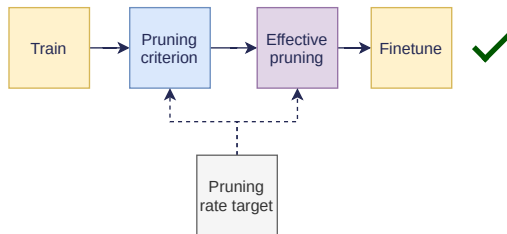
Standard Pruning Pipeline



Effective Pruning: Setting pruned weights to 0 and freezing them

Pruning - Standard Pipeline

Standard Pruning Pipeline

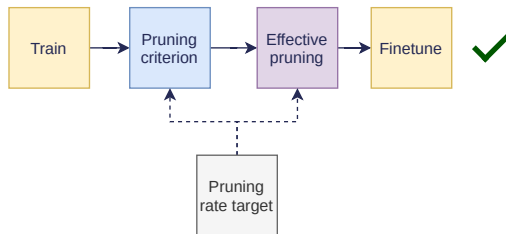


Standard pipelines requires the application of a **pruning criterion** to determine which weights will be **pruned**. This is done **after training**.

Effective Pruning: Setting pruned weights to 0 and freezing them

Pruning - Standard Pipeline

Standard Pruning Pipeline



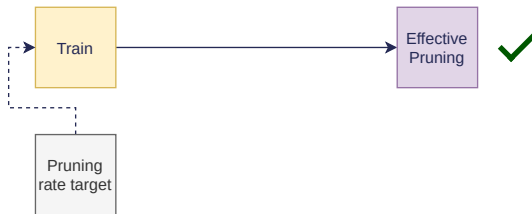
Standard pipelines requires the application of a **pruning criterion** to determine which weights will be **pruned**. This is done **after training**.

Initial training does **not take into account final weight budget**.

Effective Pruning: Setting pruned weights to 0 and freezing them

Pruning - Our Pipeline

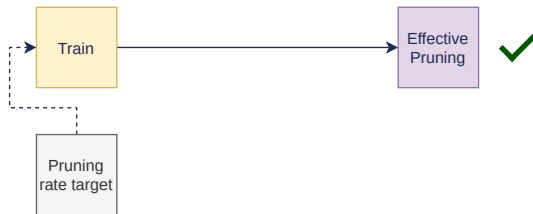
Our Pipeline



Effective Pruning: Setting pruned weights to 0 and freezing them

Pruning - Our Pipeline

Our Pipeline

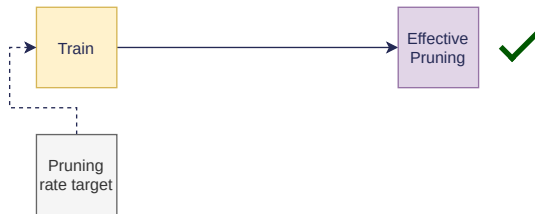


Our pipeline takes the **final pruning rate** as an input **during initial training**.

Effective Pruning: Setting pruned weights to 0 and freezing them

Pruning - Our Pipeline

Our Pipeline



Our pipeline takes the **final pruning rate** as an input **during initial training**.

Topology is optimized with the budget constraint **from the start**. This helps **preventing disconnections** in the network.

Effective Pruning: Setting pruned weights to 0 and freezing them

OUR METHOD

Weight Reparametrization

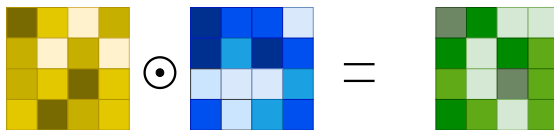
Weight Reparametrization

Weight Reparametrization

 Weight
 w

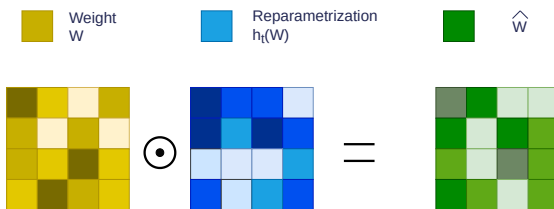
 Reparametrization
 $h_t(w)$

 \hat{w}



Weight Reparametrization

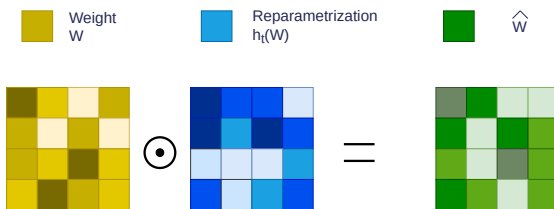
Weight Reparametrization



- New weights \hat{w} are defined by $\hat{w} = w \odot h_t(w)$.

Weight Reparametrization

Weight Reparametrization



- New weights \hat{w} are defined by $\hat{w} = w \odot h_t(w)$.
- w are the standard neural network weights.

Reparametrization Function

Reparametrization function

C_1 and C_2 ensure
 $0 \leq h_t(x) \leq 1$

$$h_t(x) = C_1 \left(\exp \left\{ - \frac{1}{(tx)^n + 1} \right\} - C_2 \right)$$

t controls the bandwidth of the pit.
It is a learnt parameter

n controls the sharpness of the falling and raising edges

$+1$ to numerically stabilize $h_t(x)$

Reparametrization function

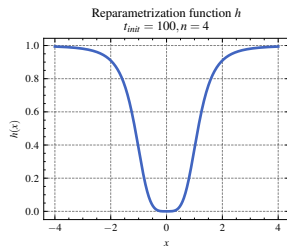
\mathbf{C}_1 and \mathbf{C}_2 ensure
 $0 \leq h_t(x) \leq 1$

$$h_t(x) = C_1 \left(\exp \left\{ - \frac{1}{(tx)^n + 1} \right\} - C_2 \right)$$

\mathbf{t} controls the bandwidth of the pit.
It is a learnt parameter

\mathbf{n} controls the sharpness of the falling and raising edges

$+1$ to numerically stabilize $h_t(x)$



Reparametrization function

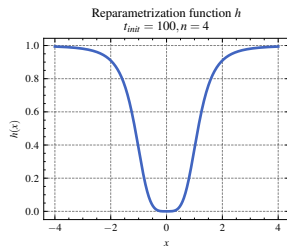
C_1 and C_2 ensure
 $0 \leq h_t(x) \leq 1$

$$h_t(x) = C_1 \left(\exp \left\{ - \frac{1}{(tx)^n + 1} \right\} - C_2 \right)$$

t controls the bandwidth of the pit. It is a learnt parameter

n controls the sharpness of the falling and raising edges

$+1$ to numerically stabilize $h_t(x)$



- t is a **learnt parameter**, optimized with SGD. t is initialized to 100.

Reparametrization function

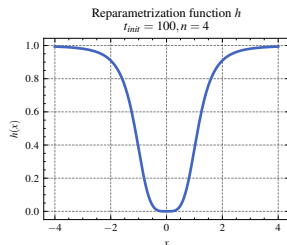
C_1 and C_2 ensure
 $0 \leq h_t(x) \leq 1$

$$h_t(x) = C_1 \left(\exp \left\{ - \frac{1}{(tx)^n + 1} \right\} - C_2 \right)$$

t controls the bandwidth of the pit. It is a learnt parameter

n controls the sharpness of the falling and raising edges

+1 to numerically stabilize $h_t(x)$

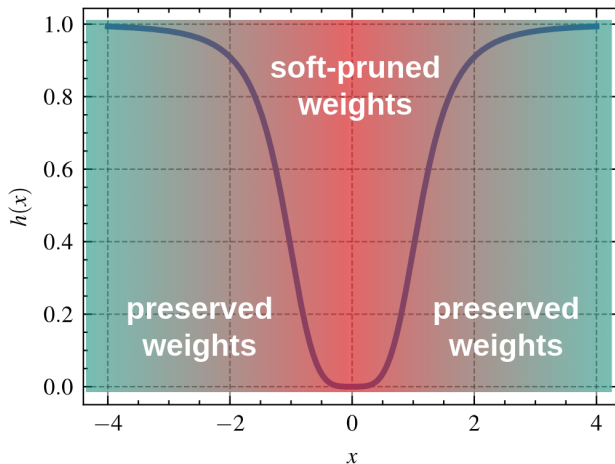


- t is a **learnt parameter**, optimized with SGD. t is initialized to 100.
- n is fixed to 4.

Reparametrization function

Reparametrization function h

$t_{init} = 100, n = 4$



Budget Loss

Budget Loss

Budget Loss

Budget: ℓ_0 -norm is not differentiable. h is used as a surrogate.

Budget Loss

Budget: ℓ_0 -norm is not differentiable. h is used as a surrogate.

$$C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) = \sum_{i=1}^L h(\mathbf{w}_i)$$

Current cost
(sum of weight reparametrizations)

$$\mathcal{L}_{\text{budget}} = \left(\frac{C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) - C_{\text{target}}}{C_{\text{initial}}} \right)^2$$

Target cost

Initial cost

Budget Loss

Budget: ℓ_0 -norm is not differentiable. h is used as a surrogate.

$$C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) = \sum_{i=1}^L h(\mathbf{w}_i)$$

Current cost
(sum of weight reparametrizations)

$$\mathcal{L}_{\text{budget}} = \left(\frac{C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) - C_{\text{target}}}{C_{\text{initial}}} \right)^2$$

Target cost

Initial cost

- **Current cost** is computed at **each step**.

Budget Loss

Budget: ℓ_0 -norm is not differentiable. h is used as a surrogate.

$$C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) = \sum_{i=1}^L h(\mathbf{w}_i)$$

Current cost
(sum of weight reparametrizations)

$$\mathcal{L}_{\text{budget}} = \left(\frac{C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) - C_{\text{target}}}{C_{\text{initial}}} \right)^2$$

Target cost

Initial cost

- **Current cost** is computed at **each step**.
- **Target cost** is the **number of parameters** that will be **kept**.

Budget Loss

Budget: ℓ_0 -norm is not differentiable. h is used as a surrogate.

$$C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) = \sum_{i=1}^L h(\mathbf{w}_i)$$

Current cost
(sum of weight reparametrizations)

$$\mathcal{L}_{\text{budget}} = \left(\frac{C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) - C_{\text{target}}}{C_{\text{initial}}} \right)^2$$

Target cost

Initial cost

- **Current cost** is computed at **each step**.
- **Target cost** is the **number of parameters** that will be **kept**.
- **Initial cost** is the sum of weight reparametrizations **before** the first step.

Budget Loss

Budget: ℓ_0 -norm is not differentiable. h is used as a surrogate.

$$C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) = \sum_{i=1}^L h(\mathbf{w}_i)$$

Current cost
(sum of weight reparametrizations)

$$\mathcal{L}_{\text{budget}} = \left(\frac{C(\{\mathbf{w}_1, \dots, \mathbf{w}_L\}) - C_{\text{target}}}{C_{\text{initial}}} \right)^2$$

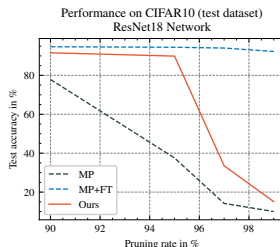
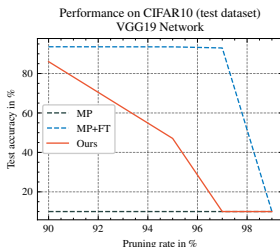
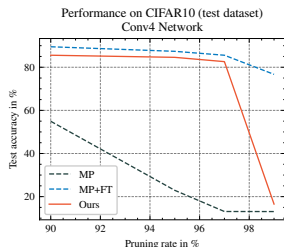
Target cost

Initial cost

- **Current cost** is computed at **each step**.
- **Target cost** is the **number of parameters** that will be **kept**.
- **Initial cost** is the sum of weight reparametrizations **before** the first step.
- Total loss is $\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda \mathcal{L}_{\text{budget}}$, with $\lambda > 0$

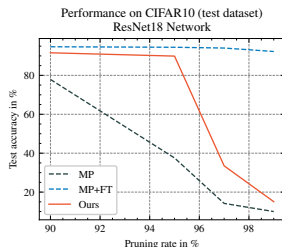
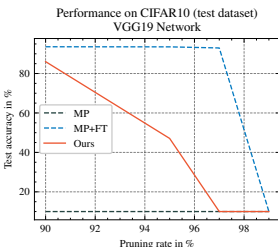
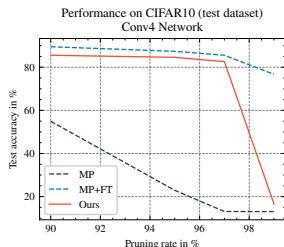
RESULTS

CIFAR 10

Results on CIFAR 10 for **Conv4**, **VGG19** and **ResNet18** networks

MP: Magnitude Pruning, MP+FT: Finetuned Magnitude Pruning

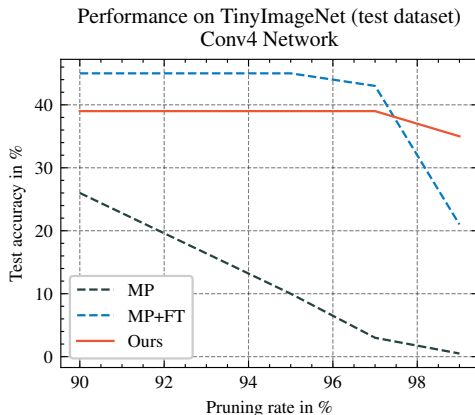
CIFAR 10

Results on CIFAR 10 for **Conv4**, **VGG19** and **ResNet18** networks

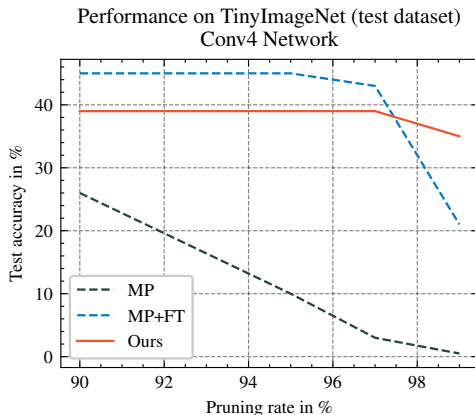
MP: Magnitude Pruning, MP+FT: Finetuned Magnitude Pruning



Our method **does not need finetuning.**

Results on TINYIMAGENET for **Conv4** network

MP: Magnitude Pruning, MP+FT: Finetuned Magnitude Pruning

Results on TINYIMAGENET for **Conv4** network

MP: Magnitude Pruning, MP+FT: Finetuned Magnitude Pruning



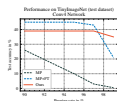
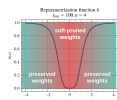
Our method **does not need finetuning**.

SUM UP

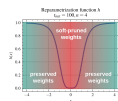
Key points

Key points

- Our reparametrization acts as a **regularizer** and a **saliency indicator**, which **induce sparsity** by **soft-pruning** the smallest weights.



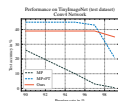
Key points



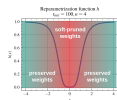
- Our reparametrization acts as a **regularizer** and a **saliency indicator**, which **induce sparsity** by **soft-pruning** the smallest weights.



- It allows to optimize **both topology and weights** under **budget constraints**.



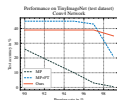
Key points



- Our reparametrization acts as a **regularizer** and a **saliency indicator**, which **induce sparsity** by **soft-pruning** the smallest weights.



- It allows to optimize **both topology and weights** under **budget constraints**.



- Our method significantly **overperforms magnitude pruning without finetuning**, and performs better than finetuned magnitude pruning for **very high pruning rates** on more complex datasets.

Perspectives

- Evaluate our method on **larger** and **more complex datasets**.

- Evaluate our method on **larger** and **more complex datasets**.
- Improve performances to outperform **finetuned** magnitude pruning more consistently.

- Evaluate our method on **larger** and **more complex datasets**.
- Improve performances to outperform **finetuned** magnitude pruning more consistently.
- Assess the **impact of the reparametrization** function and test other functions.

Thank you!

Robin DUPONT

*Sorbonne Université
& Netatmo*

Hichem SAHBI

Sorbonne Université

Guillaume MICHEL

Netatmo

 `robin.dupont@netatmo.com`