

Extracting Effective Subnetworks with Gumbel-Softmax

Robin **Dupont**

Mohammed Amine **Alaoui**

Hichem **Sahbi**

Alice **Lebois**

Sorbonne Université & Netatmo

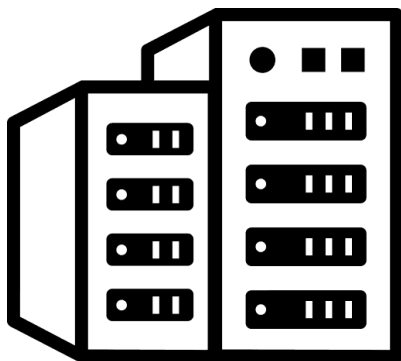
Netatmo

Sorbonne Université

Netatmo

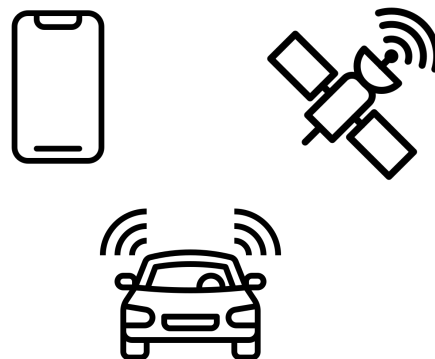
Embedded devices Need Lightweight Neural Networks

Servers



- Powerful 💪
- Handle **full size** models 🧠

Embedded devices



- Limited resources 🤖
- Require **lightweight** models 🪶

How to design **lightweight networks ?**

Existing network **compression** methods

- Neural network **distillation**
- Weight **quantization**
- Neural **Architecture Search**
- **Pruning**
 - Structured
 - **Unstructured**

Existing network **compression** methods

- Neural network **distillation**
- Weight **quantization**
- Neural **Architecture Search**
- **Pruning**
 - Structured
 - **Unstructured**

Our method **relies** on **unstructured** pruning

Our method is **not a typical pruning method**

 **Automatic** pruning rate

 **No** weight training

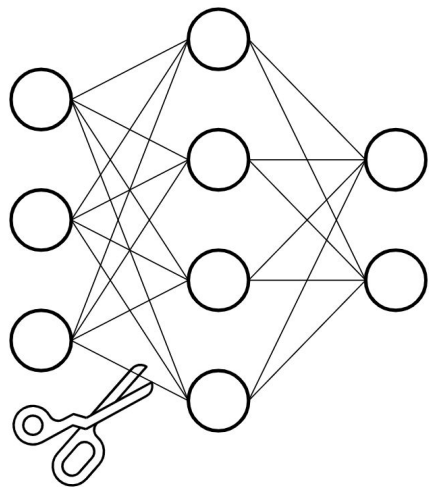
Our method is **not a typical pruning** method

 **Automatic** pruning rate

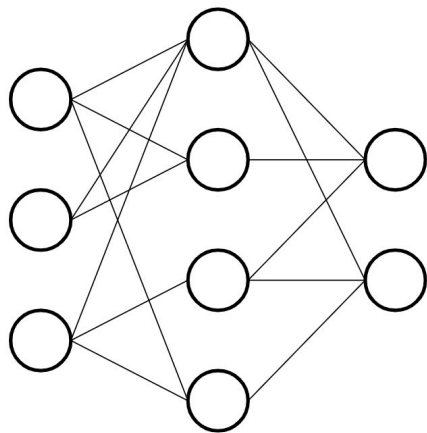
 **No weight training**  performs **topology selection**

How unstructured weight pruning works ?

Unstructured weight pruning yields lightweights networks



Before pruning

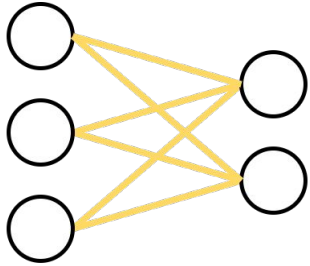


After pruning

- weights removed **individually**
- **Flexible** ✓
- **High sparsity rate** ✓

Typical pruning pipelines rely on weight training

Typical pruning pipelines are composed of 3 steps



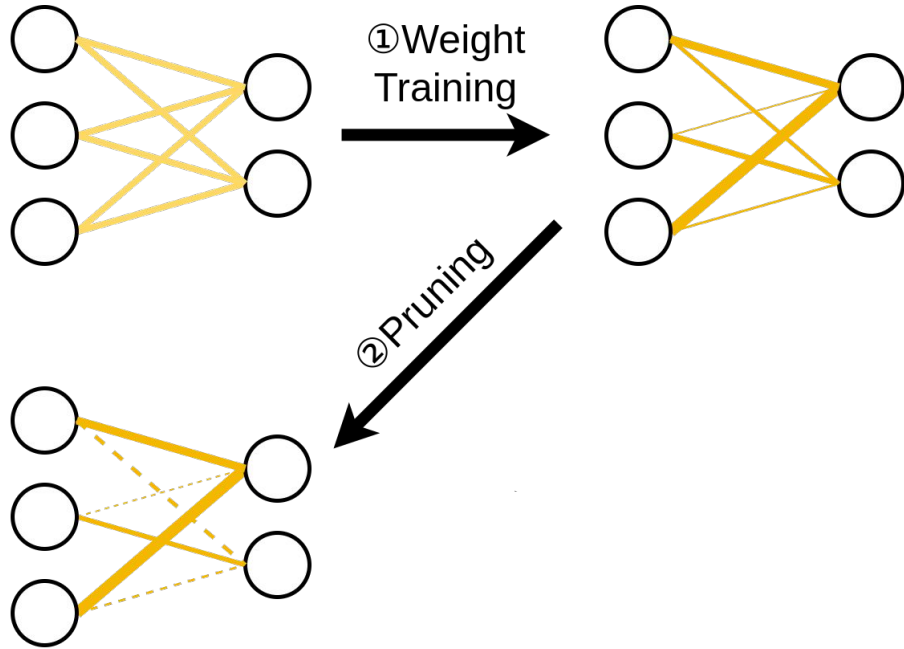
- 3 steps procedure
train - prune - **fine-tune**
- Pruning **criterion** depends on the **method**
- **Fine-tuning** needed

Typical pruning pipelines are composed of 3 steps



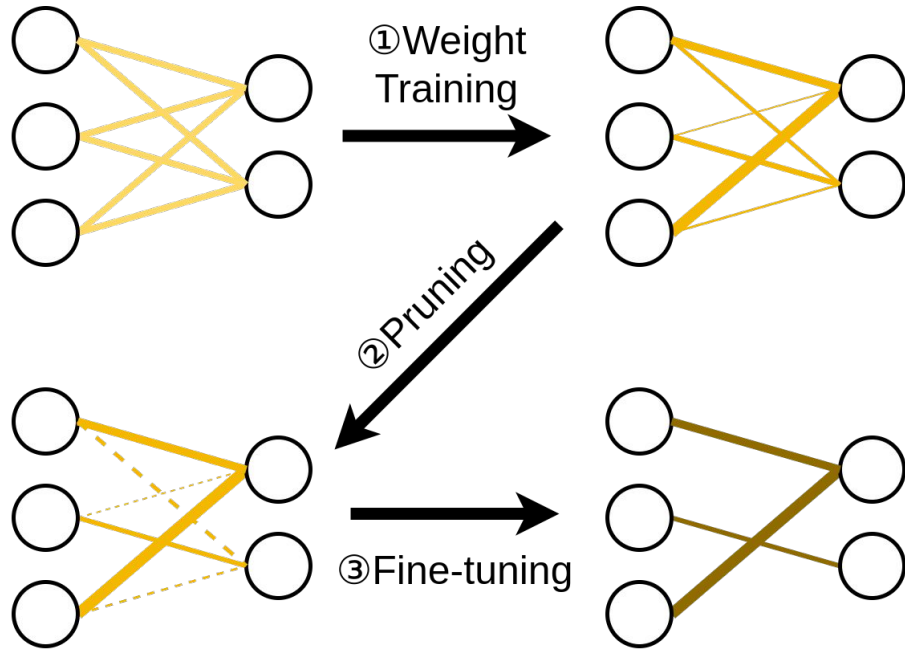
- 3 steps procedure
train - prune - **fine-tune**
- Pruning **criterion** depends on the **method**
- **Fine-tuning** needed

Typical pruning pipelines are composed of 3 steps



- 3 steps procedure
train - prune - **fine-tune**
- Pruning **criterion** depends on the **method**
- **Fine-tuning** needed

Typical pruning pipelines are composed of 3 steps

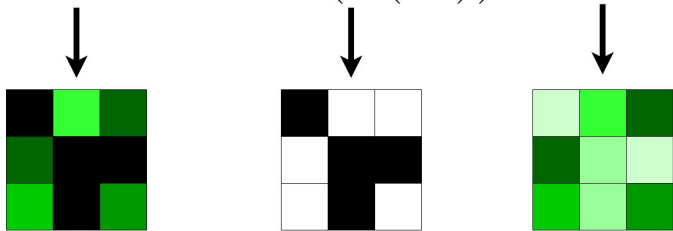


- 3 steps procedure
train - prune - **fine-tune**
- Pruning **criterion** depends on the **method**
- **Fine-tuning** needed

What if we do not train the weights ?

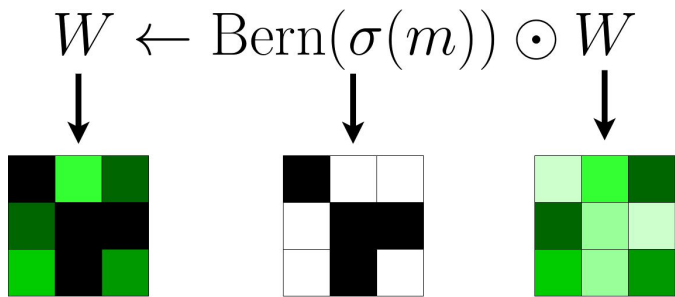
Zhou et al. train the Supermask

$$W \leftarrow \text{Bern}(\sigma(m)) \odot W$$



- mask **sampled** from Bernoulli distribution
- optimize **only** m

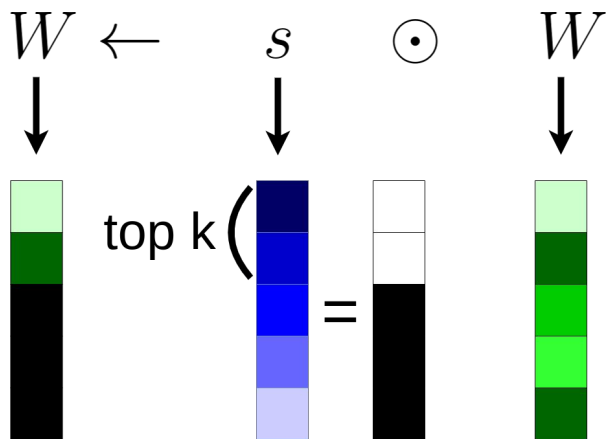
Zhou et al. train the Supermask



- mask **sampled** from Bernoulli distribution
- optimize **only** m

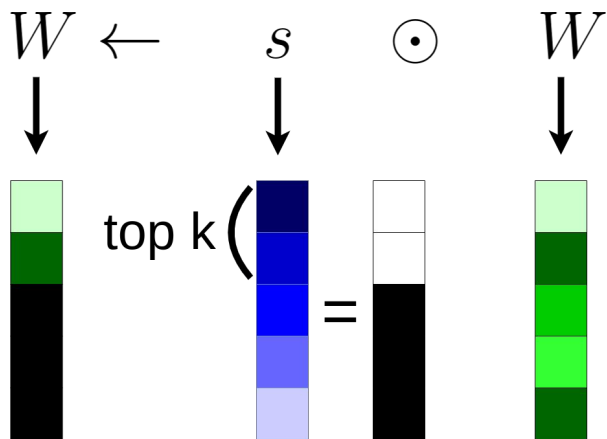
✗ Basic mask parametrization with Straight Through
✗ Cumbersome weight rescaling

Ramanujan et al. take the **top-k** weights in each layers



- **top-k** element of s are selected
- top-k elements are chosen **per layer**
- k depends on the **pruning rate**

Ramanujan et al. take the **top-k** weights in each layers



- **top-k** element of s are selected
- top-k elements are chosen **per layer**
- k depends on the **pruning rate**

- ✗ Pruning rate **has** to be **known in advance**
- ✗ **Same** pruning rate for **all layer**
- ✗ Best pruning rate if found by **grid search** → computationally **expensive** 💰
- ✗ Best pruning rate **depends** on the network **architecture**

Our Method

Our method characteristics

- We use **Gumbel-Softmax** for differentiable sampling (Jang et al. 2016)
 better performances than Straight-Through





Our method characteristics

- We use **Gumbel-Softmax** for differentiable sampling (Jang et al. 2016)
 - ✓ better performances than Straight-Through
- Pruning rate is **not needed**
 - ✓ No computationally intensive grid search

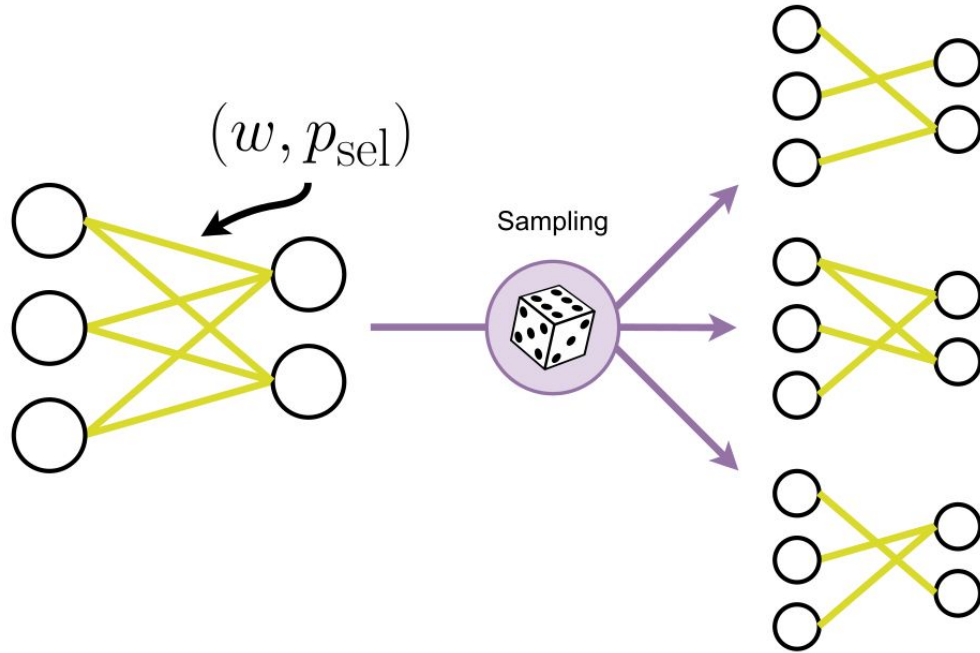
Our method characteristics

- We use **Gumbel-Softmax** for differentiable sampling (Jang et al. 2016)
 - ✓ better performances than Straight-Through
- Pruning rate is **not needed**
 - ✓ No computationally intensive grid search
- **Learnt** weight rescaling factor
 - ✓ Faster inference time and lower overhead

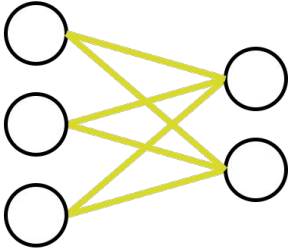
Our method characteristics

- We use **Gumbel-Softmax** for differentiable sampling (Jang et al. 2016)
 better performances than Straight-Through
- Pruning rate is **not needed**
 No computationally intensive grid search
- **Learnt** weight rescaling factor
 Faster inference time and lower overhead
- **No** weight training
 Topology selection only

Each weight as a probability of being selected

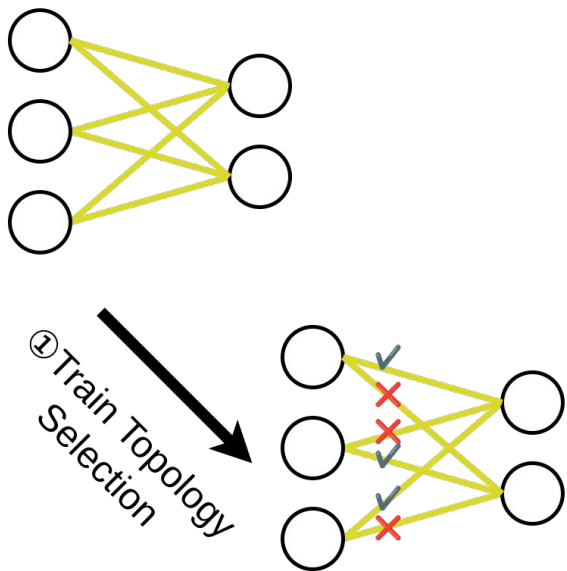


Our method performs **topology selection only**



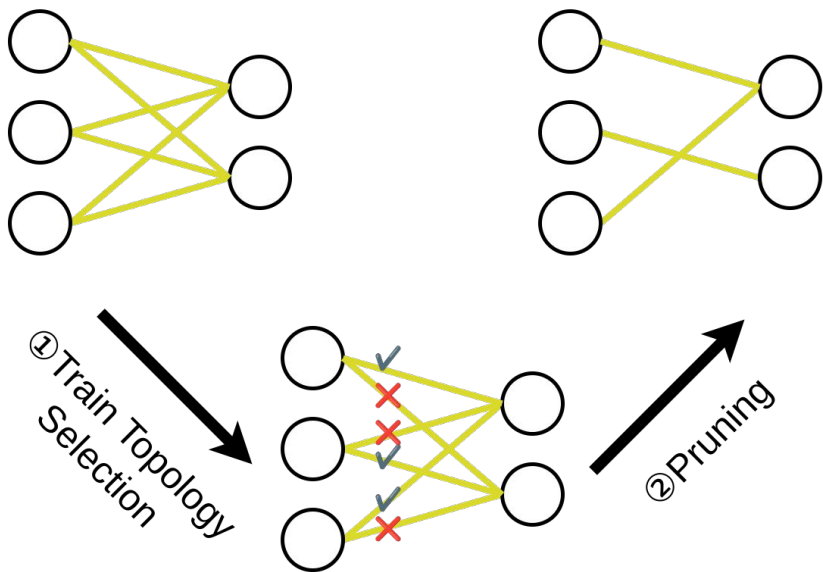
- No weight **training** or **fine-tuning** ⚠
- **Topology selection only**

Our method performs **topology selection** only



- No weight **training** or **fine-tuning** ⚠
- **Topology selection** only

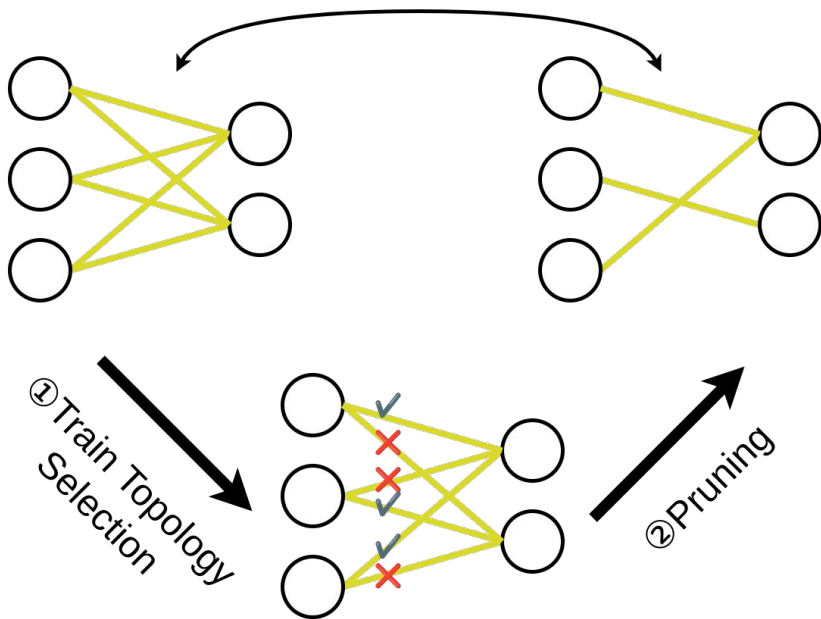
Our method performs **topology selection** only



- No weight **training** or **fine-tuning** ⚠
- **Topology selection** only

Our method performs **topology selection only**

same weights values : no weight training



- No weight **training** or **fine-tuning** ⚠
- **Topology selection only**

Topology are selected with a **stochastic** mask

Layer equation :

$$\mathbf{z}_\ell = g_\ell((\mathbf{m}_\ell \odot \mathbf{w}_\ell) \otimes \mathbf{z}_{\ell-1})$$



binary masks tensor

weights tensor

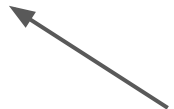
Topology are selected with a **stochastic** mask

Layer equation :

$$\mathbf{z}_\ell = g_\ell((\mathbf{m}_\ell \odot \mathbf{w}_\ell) \otimes \mathbf{z}_{\ell-1})$$

binary masks tensor  weights tensor 



each coefficient follows Bernoulli distribution : $m \sim \mathcal{B}(p_{\text{sel}})$

 Probability for a **weight** to be selected

Topology are selected with a **stochastic** mask

Layer equation :

$$\mathbf{z}_\ell = g_\ell((\mathbf{m}_\ell \odot \mathbf{w}_\ell) \otimes \mathbf{z}_{\ell-1})$$

binary masks tensor  weights tensor 

each coefficient follows Bernoulli distribution : $m \sim \mathcal{B}(p_{\text{sel}})$

 Probability for a **weight** to be selected

Sampling is **not differentiable** ⚠

Probability is reparametrized

Probability reparametrization :

$$p_{\text{sel}} = \sigma(\hat{m})$$

Sigmoid ensures $0 \leq p_{\text{sel}} \leq 1$

Learnt variable

Naive Straight Through Gumbel-Softmax is **flawed**

Probability reparametrization : $p_{\text{sel}} = \sigma(\hat{m})$

Naive **S**traight **T**hrough **G**umbel-**S**oftmax
formulation :

$$m = \text{STGS} \left(\begin{bmatrix} \log(\sigma(\hat{m})) \\ \log(1 - \sigma(\hat{m})) \end{bmatrix} \right)$$

Naive Straight Through Gumbel-Softmax is **flawed**

Probability reparametrization : $p_{\text{sel}} = \sigma(\hat{m})$

Naive **S**traight **T**hrough **G**umbel-**S**oftmax
formulation :

$$m = \text{STGS} \left(\begin{bmatrix} \log(\sigma(\hat{m})) \\ \log(1 - \sigma(\hat{m})) \end{bmatrix} \right)$$

Combination of log and exponential functions :

- ✗ Numerical **instabilities**
- ✗ Computationally **intensive**

ASLP is simpler and resolves issues

Probability reparametrization : $p_{\text{sel}} = \sigma(\hat{m})$

Our formulation **Arbitrarily Shifted Log Parameterization** (ASLP)

$$m = \text{STGS} \left(\begin{bmatrix} \hat{m} \\ 0 \end{bmatrix} \right)$$

ASLP is simpler and resolves issues

Probability reparametrization : $p_{\text{sel}} = \sigma(\hat{m})$

Our formulation **Arbitrarily Shifted Log Parameterization** (ASLP)

$$m = \text{STGS} \left(\begin{bmatrix} \hat{m} \\ 0 \end{bmatrix} \right)$$

- ✓ Numerically **stable**
- ✓ **Less** computationally intensive

ASLP formulation implies the same parametrization for P_{sel}

Our formulation : $m = \text{STGS} \left(\begin{bmatrix} \hat{m} \\ 0 \end{bmatrix} \right)$

Arbitrary unknown constant that shifts log probabilities

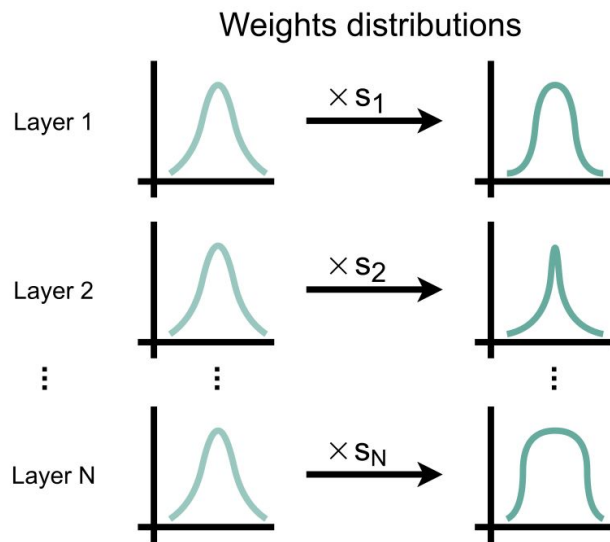
$$\begin{bmatrix} \hat{m} \\ 0 \end{bmatrix} = \begin{bmatrix} \log(p_{sel}) + c \\ \log(1 - p_{sel}) + c \end{bmatrix} \Rightarrow p_{sel} = \sigma(\hat{m})$$

💡 Adding a constant **does not change** the result of STGS

Same reparametrization

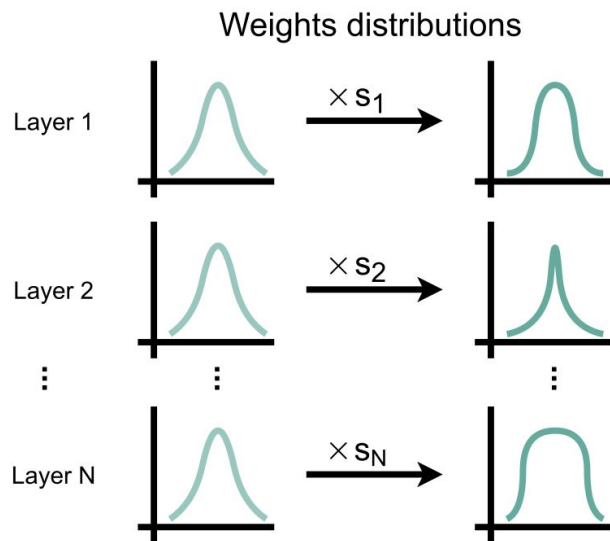
Pruning weights affects the signal propagation dynamic

Smart rescale is a simpler and useful weight rescaling



- Scaling learnt per layer
- **Mitigates** the change of **variance** due to **pruning**
- Improves **performances**
- **Reduces** number of epochs needed for **convergence**

Smart rescale is a simpler and useful weight rescaling



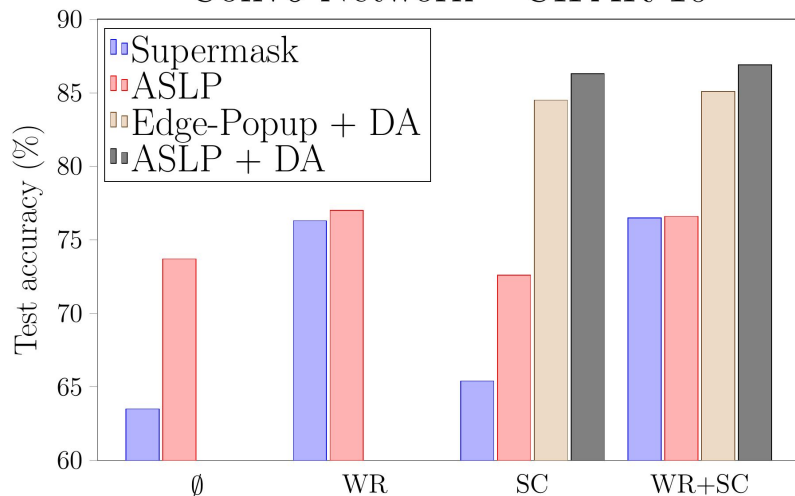
- Scaling learnt per layer
- **Mitigates** the change of **variance** due to **pruning**
- Improves **performances**
- **Reduces** number of epochs needed for **convergence**

✓ **Less computationally intensive** than Zhou et al. weight rescaling

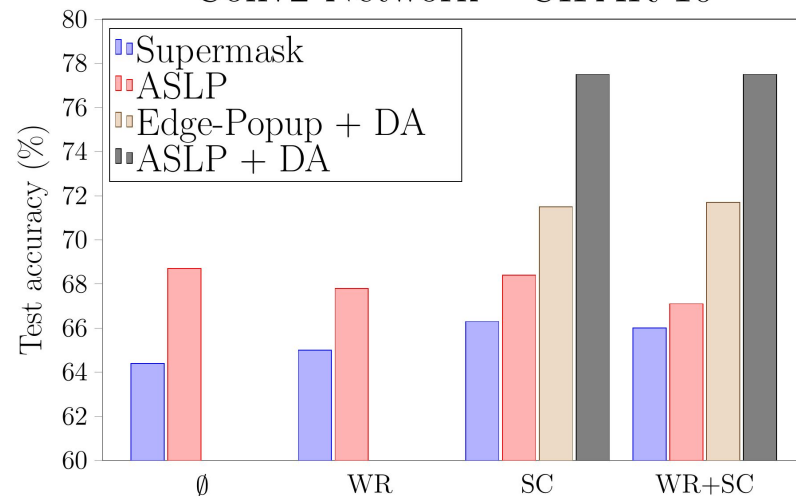
Results

Results : Our method performs better on various cases

Comparison of Supermask, EP and ASLP
Conv6 Network – CIFAR 10



Comparison of Supermask, EP and ASLP
Conv2 Network – CIFAR 10



WR = Weight Rescale, SC = Signed Constant, DA = Data Augmentation

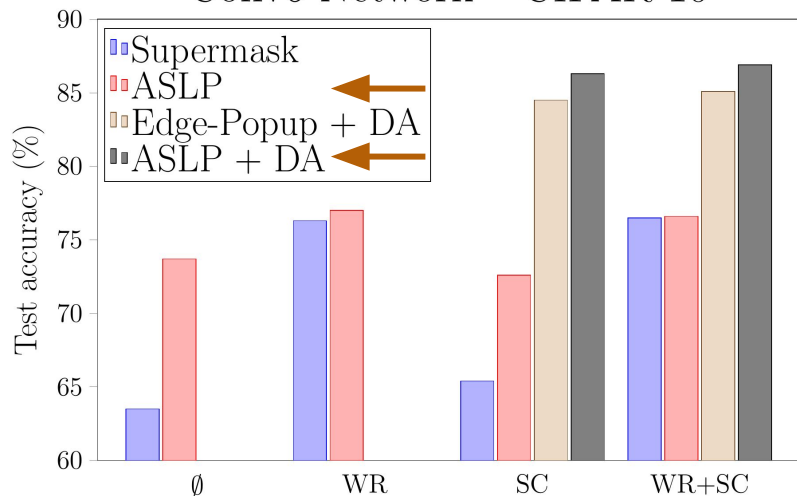
EP (Edge-Popup)¹, Supermask²

[1] H. Zhou, et al. "Deconstructing lottery tickets: Zeros, signs, and the supermask," in NeurIPS, 2019

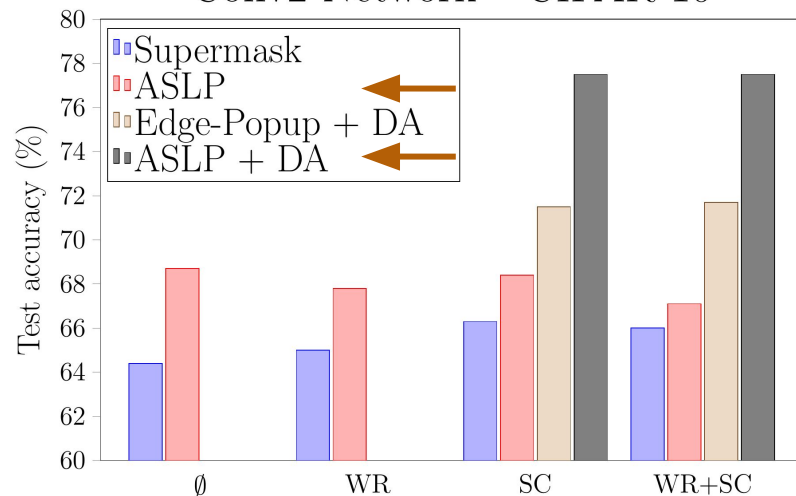
[2] V. Ramanujan, et al. "What's hidden in a randomly weighted neural network?," in CVPR, 2020

Results : Our method performs better on various cases

Comparison of Supermask, EP and ASLP
Conv6 Network – CIFAR 10



Comparison of Supermask, EP and ASLP
Conv2 Network – CIFAR 10



WR = Weight Rescale, SC = Signed Constant, DA = Data Augmentation

EP (Edge-Popup)¹, Supermask²

[1] H. Zhou, et al. "Deconstructing lottery tickets: Zeros, signs, and the supermask," in NeurIPS, 2019

[2] V. Ramanujan, et al. "What's hidden in a randomly weighted neural network?," in CVPR, 2020

Smart Rescale is faster than DWR¹ and accelerate convergence

- Smart rescale overhead **0.13s**
- DWR¹ overhead **0.2s**

Smart Rescale is faster than DWR¹ and accelerate convergence

- Smart rescale overhead **0.13s**

✓ SR overhead **35% faster**

- DWR¹ overhead **0.2s**

Smart Rescale is faster than DWR¹ and accelerate convergence

- Smart rescale overhead **0.13s**

 SR overhead **35% faster**

- DWR¹ overhead **0.2s**

- # Epochs reduction :
 - **8.2%** (Conv2)
 - **19.7%** (Conv4)
 - **14.0%** (Conv6)

Results : On CIFAR100 ASLP performs better on most cases








	Conv2	Conv4	Conv6
EP	40.9	51.1	53.2
ASLP	43.4	51.7	52.8

Table 1: Edge Popup and ASLP on CIFAR100

Results with Weight Rescale (WR) and Signed Constant (SC)

Sum Up

Sum Up

- Prune **untrained** networks  **topology selection** only
- **Gumbel Softmax**  differentiable sampling
- **ASLP** :
 - simpler formulation
 -  **less** computationally intensive
 -  numerically **stable**
- **Smart Rescale** :
 -  **improves** performances
 -  **reduces** number of **epochs**
- Our method yields **lightweight** networks , **without weight training.**

A few perspectives

- Test ASLP on other **network architectures** and **datasets**
- **Reduce** training time
- Test ASLP on **other context and applications**

Thank you!



Code available at:

github.com/n0ciple/aslp

Robin **Dupont**
Mohammed Amine **Alaoui**
Hichem **Sahbi**
Alice **Lebois**



Sorbonne Université & Netatmo
Netatmo
Sorbonne Université
Netatmo