

# Extracting Effective Subnetworks with Gumbel-Softmax

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

Sorbonne Université & Netatmo  
Netatmo  
Sorbonne Université  
Netatmo

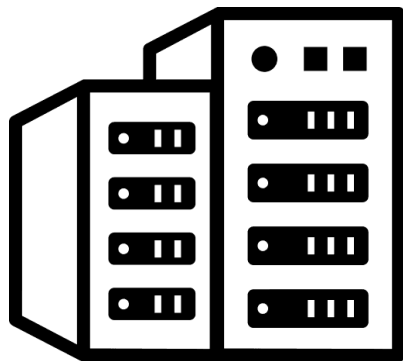
# Table of Content

- Why Lightweight Neural Networks ?
- Lightweight Networks Design via Pruning
- Our Method
- Results
- Sum Up

**Why Lightweight Neural Networks ?**

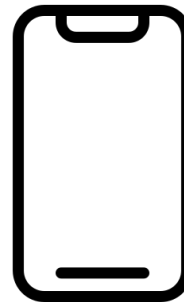
# Why Lightweight Neural Networks ?

Server



- Powerful 💪
- Handle full size models 🧠

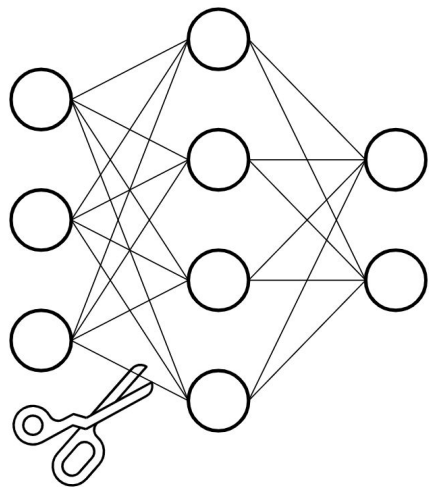
Embedded device



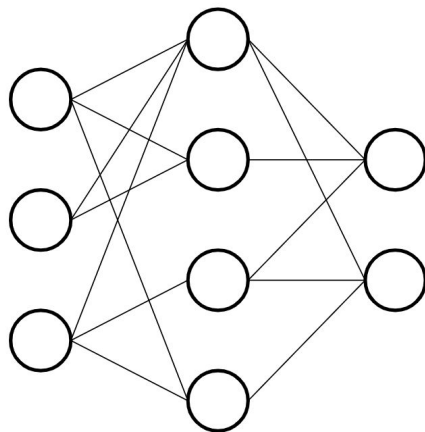
- Limited resources 🤖
- Require lightweight models 🪶

# **Lightweight Networks Design via Pruning**

# Lightweight Networks Design via Pruning



Before pruning

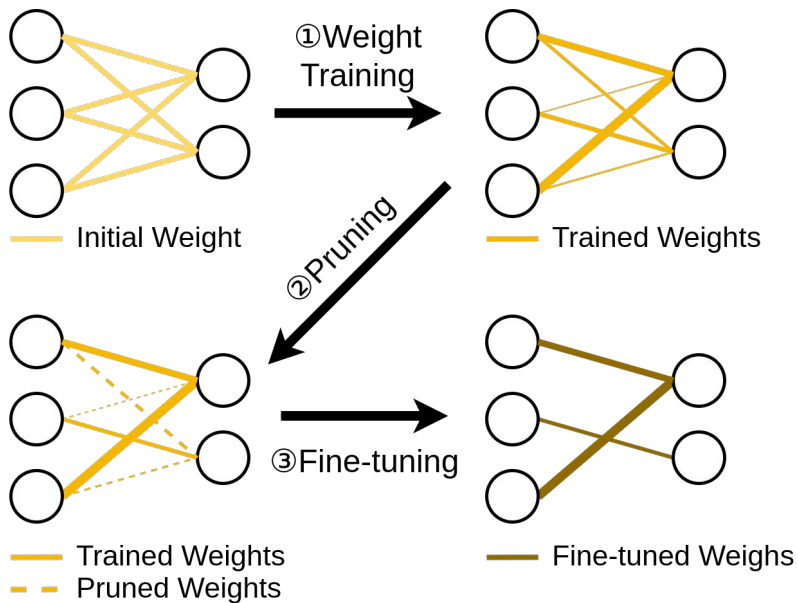


After pruning

- **Pruning** removes **weights**
- **Unstructured** pruning
- Yields **sparse** and **lightweight** models

# Lightweight Networks Design via Pruning

## Standard pruning pipelines



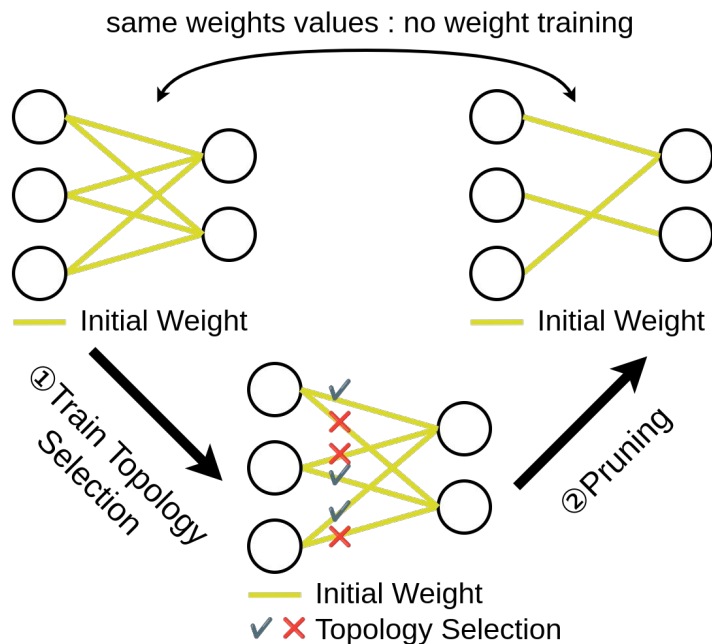
- 3 steps procedure  
train - prune - fine-tune
- Pruning **criterion** depends  
on the **method**
- **Fine-tuning** needed

# **Our Method**



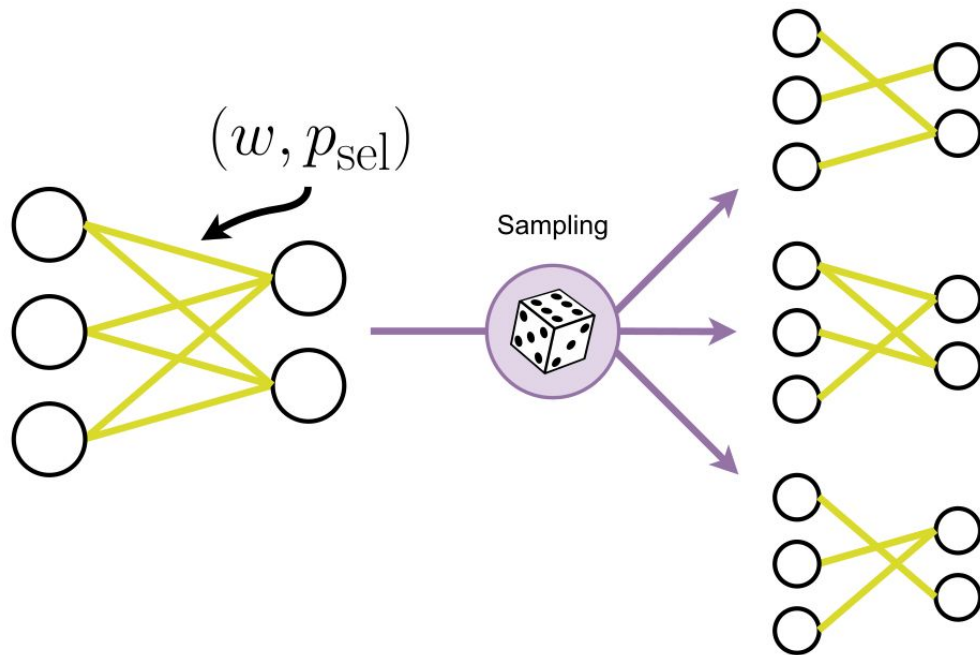
# Our Method

## Our pruning pipeline



- **No weight training** ⚠
- **Topology selection only**



# Our Method



# Our Method

Layer equation :

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

binary masks tensor  weights tensor 

coefficient follows bernoulli distribution :

$$m \sim \mathcal{B}(p_{\text{sel}})$$

Sampling is **not differentiable** ⚠

 Probability for a **weight** of  
**being selected**

# Our Method

Probability reparametrization :

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

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

Learnt variable

# Our Method

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

Naive Gumbel-Softmax 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**

# Our Method - ASLP

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

## Our Method - ASLP

Our formulation :

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

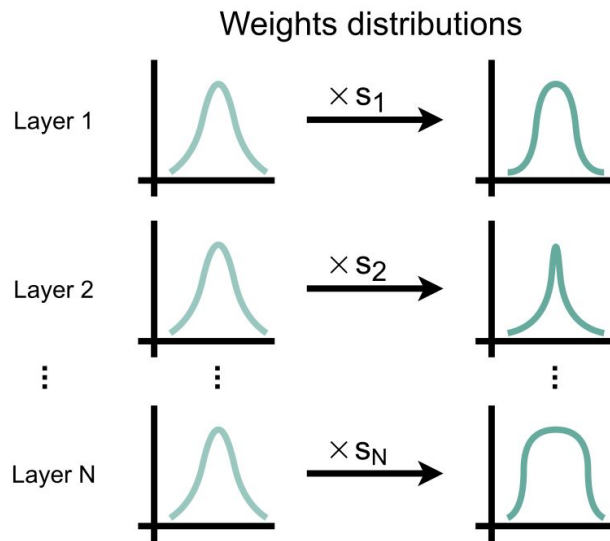
Arbitrary unknown constant that shifts log probabilities

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

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

Same reparametrization

# Our Method - Smart rescale



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

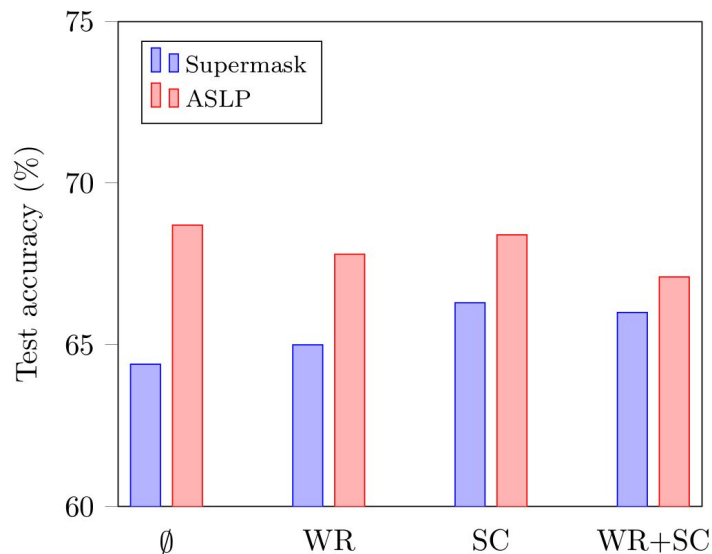


# Results

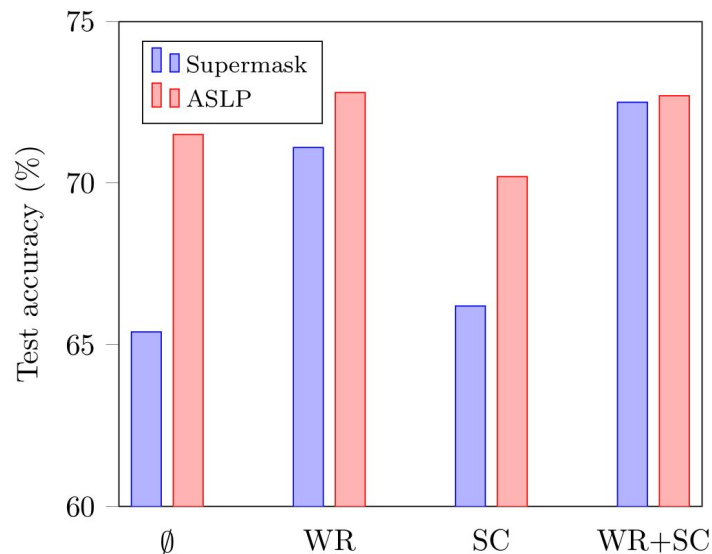
# Results

## CIFAR 10

Comparison of supermask and ASLP  
Conv2 – w/o data augmentation



Comparison of supermask and ASLP  
Conv4 – w/o data augmentation



WR = Weight Rescale, SC = Signed Constant

# Results

## CIFAR 100

|      | Conv2       | Conv4       | Conv6       |
|------|-------------|-------------|-------------|
| EP   | 40.9        | 51.1        | <b>53.2</b> |
| ASLP | <b>43.4</b> | <b>51.7</b> | 52.8        |

Table 1: Edge Popup and ASLP on CIFAR100

## Results for WR+SC

**Sum Up**

# Sum Up

- Lightweight networks are **useful** for **embedded devices**
- Our method prunes **untrained** networks - **topology selection** only
- We use **Gumbel Softmax** for 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**.

# Thanks you!

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

Sorbonne Université & Netatmo  
Netatmo  
Sorbonne Université  
Netatmo

