# A round-trip journey in pruned artificial neural networks

*Andrea Bragagnolo[+], Enzo Tartaglione[*], Gianluca Dalmasso[+,x], Marco Grangetto[x]*

[+] **Synesthesia s.r.l. Torino**
[*]**TELECOM Paris, France**
[x]**Università di Torino**

# Outline

**EIDOSLAB**

**A round-trip journey in pruned artificial neural networks**
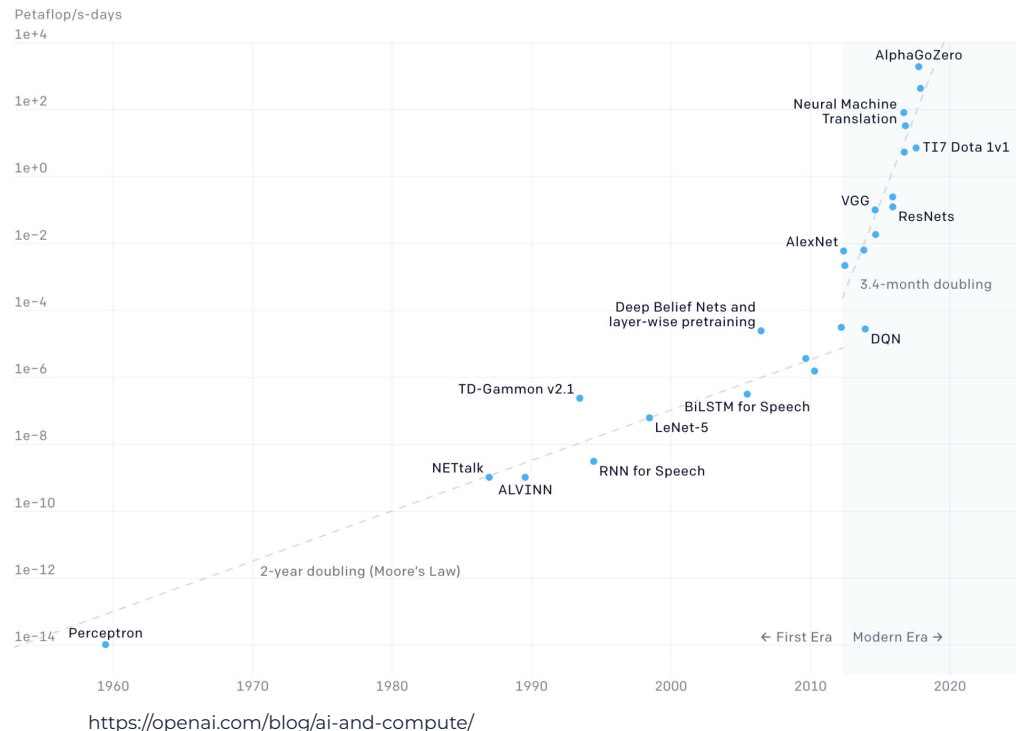
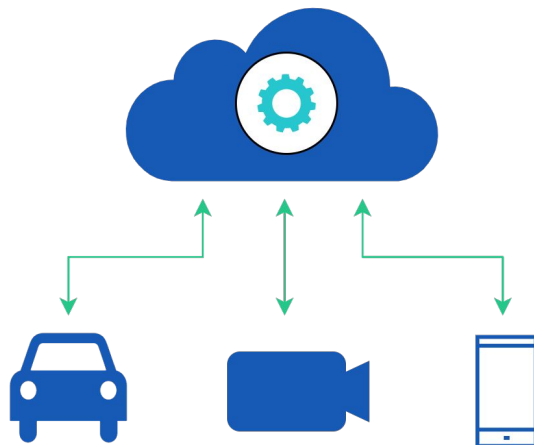# Deep models scale and **training cost**

State-of-the-art performance on complex tasks.

Very high number of parameters (hundreds of millions and even more…).

Training requires substantial computing power.
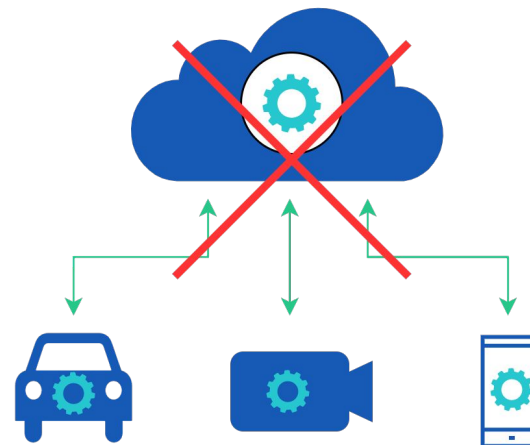
Petaflop/s-days



https://openai.com/blog/ai-and-compute/

# Deploying models

**Cloud AI**

**Edge AI**



Inferences are performed remotely.

Inferences are performed locally.

**A round-trip journey in pruned artificial neural networks**

# Advantages of Edge AI

Costs & Energy

Low latency

Privacy

Connectivity

**A round-trip journey in pruned artificial neural networks**

# Outline

**A round-trip journey in pruned artificial neural networks**

# Solution for efficient deployment
## Neural Network Pruning
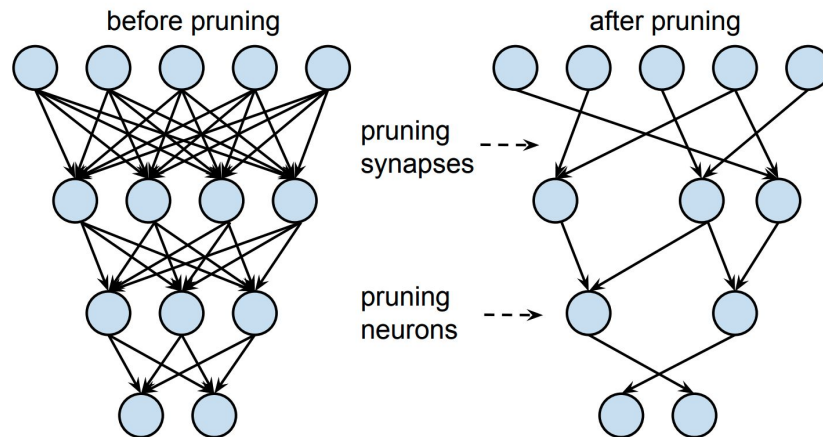


Neural Network Pruning

# Neural Network Pruning

## Intuition

Removes less influential elements while preserving the generalization capabilities.

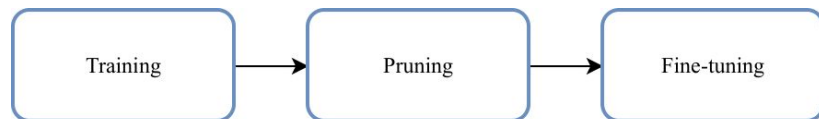Reduces the resources required to use the model.

Studied since the late '80s has seen a resurgence in 2015.

before pruning

after pruning

pruning synapses - - →

pruning neurons - - →

**A round-trip journey in pruned artificial neural networks**
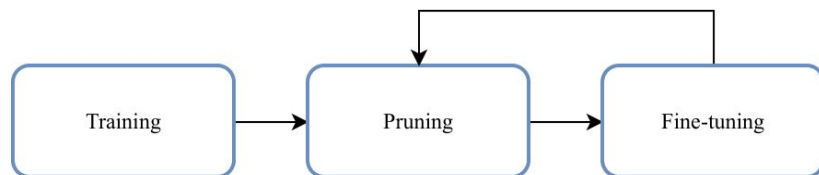
# Categorization of Pruning Procedures
## One-shot vs. Iterative



One-Shot

Performs a single pruning step.

Fine-tuning to recover performance.

Faster procedure.
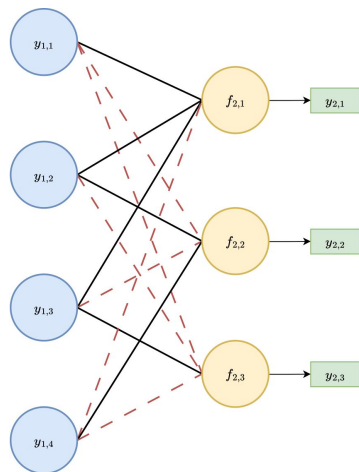


Iterative

Performs multiple pruning step.

Successive training and pruning steps.

Prunes more parameters.

# Categorization of Pruning Procedures

Unstructured vs. Structured



### Unstructured

Removes many parameters from the network.

Highly reduces the compressed model size.

### Structured

Removes entire neurons in the network.

Reduces the number of operations.

# Categorization of Pruning Procedures
## Unstructured vs. Structured



## Unstructured

No guarantee in removing neurons.

We still consider the entire matrix to define the output

## Structured

Removes entire rows from the matrix.

The rank of the final matrix is lower.

# Learning Both Weights and Connections

## A Template for Modern Techniques



Han et al., *Learning both weights and connections for efficient neural network* *(2015)*.

Kick-started modern pruning research.

Unstructured, iterative.

Acts as the foundation for the proposed procedures.

# Sensitivity Regularization

Standard Feed Forward Neural Network.

$w_{n,i,j}$

$y_N$

# Sensitivity Regularization



Standard Feed Forward Neural Network.

Our goal: assess to which extent changes in the value of the weight $w_{n,i,j}$ would affect the output $y_N$.

# Sensitivity Regularization



Standard Feed Forward Neural Network.

Our goal: assess to which extent changes in the value of the weight $w_{n,i,j}$ would affect the output $y_N$.

Intuitively:

- If the change in $y_N$ is big, $w_{n,i,j}$ has a high Sensitivity.

# Sensitivity Regularization



Standard Feed Forward Neural Network.

Our goal: assess to which extent changes in the value of the weight $w_{n,i,j}$ would affect the output $y_N$.

Intuitively:

- If the change in $y_N$ is big, $w_{n,i,j}$ has a high Sensitivity.
- If the change in $y_N$ is small, $w_{n,i,j}$ has a low Sensitivity.

# Sensitivity Regularization



Standard Feed Forward Neural Network.

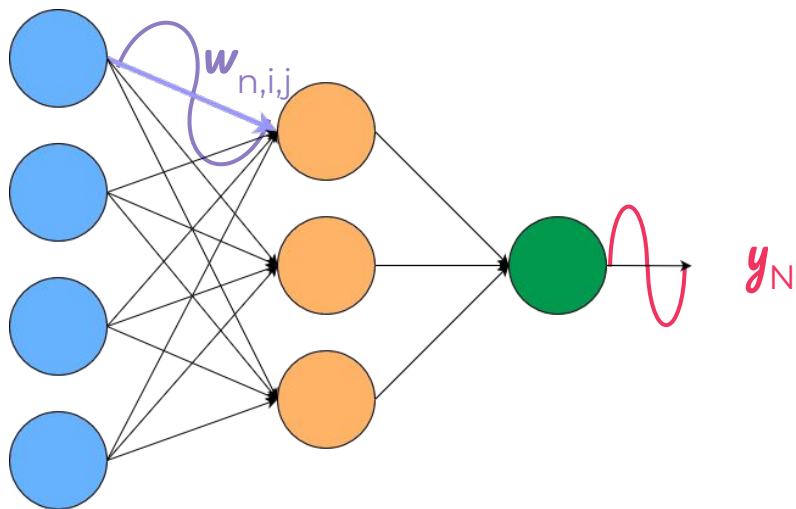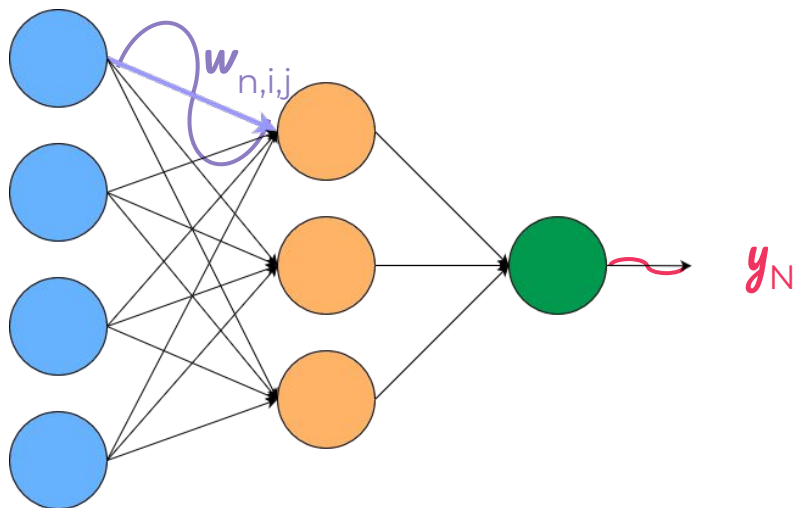Our goal: assess to which extent changes in the value of the weight $w_{n,i,j}$ would affect the output $y_N$.

Intuitively:

- If the change in $y_N$ is big, $w_{n,i,j}$ has a high Sensitivity.
- If the change in $y_N$ is small, $w_{n,i,j}$ has a low Sensitivity.

$$S(y, w_j) = \sum_{k=1}^{C} \alpha_k \left| \frac{\partial y_k}{\partial w_j} \right|$$

Tartaglione et al., Learning sparse neural networks via sensitivity-driven regularization (2018).

**A round-trip journey in pruned artificial neural networks**

# Our pruning techniques

## LOBSTER

Contribution of the parameters to the loss of the network.

$$S(\mathcal{L}, w_{n,i,j}) = \left| \frac{\partial \mathcal{L}}{\partial w_{n,i,j}} \right|$$

**Unstructured.**

"Loss-based sensitivity regularization: towards deep sparse neural networks." **Neural Networks** (2022).



(a)          (b)

## SeReNe

Contribution of the neuron to the output of the network.

$$S_{n,i}(\mathbf{y}_N, p_{n,i}) = \frac{1}{C} \sum_{k=1}^{C} \left| \frac{\partial y_{N,k}}{\partial p_{n,i}} \right|$$

**Structured.**

"SeReNe: Sensitivity-based regularization of neurons for structured sparsity in neural networks." **IEEE Transactions on Neural Networks and Learning Systems** (2021).

**A round-trip journey in pruned artificial neural networks**

# Pruning potential



ResNet-18 on ImageNet

Legend:
- LOBSTER
- FSKD
- ABCPruner
- Liebenwein et al. (2020)
- $\ell_2$ + pruning
- SCOP
- PFP

Y-axis: Top-1 Error (%)
X-axis: Pruned Parameters (%)

# Pruning in the MPEG-7 Part 17 pipeline

## Experimental Setup



Pruning → Simplification → Entropy Coding → Bit stream → Decompression

Evaluate the benefits of structured pruning approaches
within the MPEG-7 Part 17 neural network compression pipeline.

"On the role of structured pruning for neural network compression." 2021 IEEE International Conference on Image Processing (2021).

# Pruning in the MPEG-7 Part 17 pipeline
## Experimental Setup



Pruning → Simplification → Entropy Coding → Bit stream → Decompression

MPEG-7 Part 17 standard

**A round-trip journey in pruned artificial neural networks**

# Pruning in the MPEG-7 Part 17 pipeline

## Experimental Setup



Pruning — Simplification — Entropy Coding — Bit stream — Decompression

Evaluation of pruned
models on mobile devices

# Pruning in the MPEG-7 Part 17 pipeline

## Experimental Setup



Pruning → Simplification → Entropy Coding → Bit stream → Decompression

Parameter pruning:
SeReNe vs LOBSTER

# Pruning in the MPEG-7 Part 17 pipeline

## Experimental Setup



Removal of pruned neurons: Simplify

**A round-trip journey in pruned artificial neural networks**

# Pruning in the MPEG-7 Part 17 pipeline

Evaluation Results

| Dataset | Architecture | Pruning | Pruning ratio [%] | Simplified topology [MB] | Compressed bitstream [MB] | Inference time [ms] | | | |
|---------|--------------|---------|-------------------|--------------------------|----------------------------|---------------------|------|------|------|
| | | | | | | RPi 3B | P20 | MI9 | S6L |
| CIFAR-10 | VGG-16 | No pruning | - | 60.0 | 3.6 | 647 | 204 | 153 | 251 |
| | | LOBSTER | **92.44** | 58.61 | 1.61 | 610 | 191 | 146 | 242 |
| | | SeReNe | 47.16 | **31.02** | **0.34** | **594** | **99** | **85** | **106** |
| | ResNet-32 | No pruning | - | 2.0 | 0.30 | 580 | 32 | 30 | 31 |
| | | LOBSTER | **81.19** | 1.96 | 0.12 | 545 | 32 | 26 | 30 |
| | | SeReNe | 52.80 | **1.0** | **0.09** | **536** | **25** | **17** | **25** |
| CIFAR-100 | AlexNet | No pruning | - | 94.6 | 10.1 | 246 | 131 | 84 | 168 |
| | | LOBSTER | **98.90** | 48.84 | 0.40 | 224 | 95 | 67 | 120 |
| | | SeReNe | 59.87 | **37.07** | **0.20** | **186** | **75** | **53** | **96** |
| ImageNet | ResNet-101 | No pruning | - | 178.4 | 26.24 | 11919 | 958 | 416 | 1008 |
| | | LOBSTER | **87.39** | 173.87 | 9.24 | 11879 | 956 | 403 | 985 |
| | | SeReNe | 1.09 | **172.53** | **7.51** | **11699** | **929** | **371** | **974** |

Even removing less parameters, SeReNe produces smaller and faster models.

**A round-trip journey in pruned artificial neural networks**

# Outline

**A round-trip journey in pruned artificial neural networks**
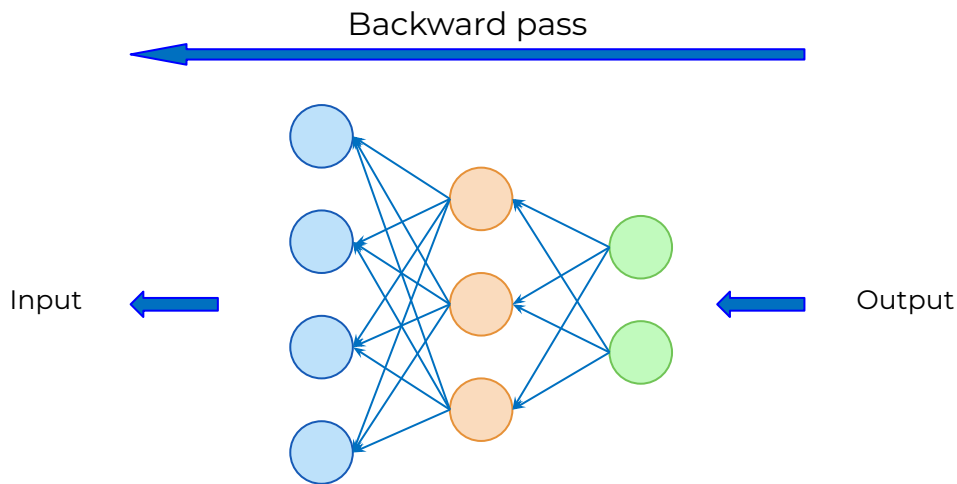
# Reducing the Training Cost

We have seen how pruning can reduce the resources required by a deployed model, but what about the training process?

It is true that a model will perform thousands of inferences but the training is still very expensive.
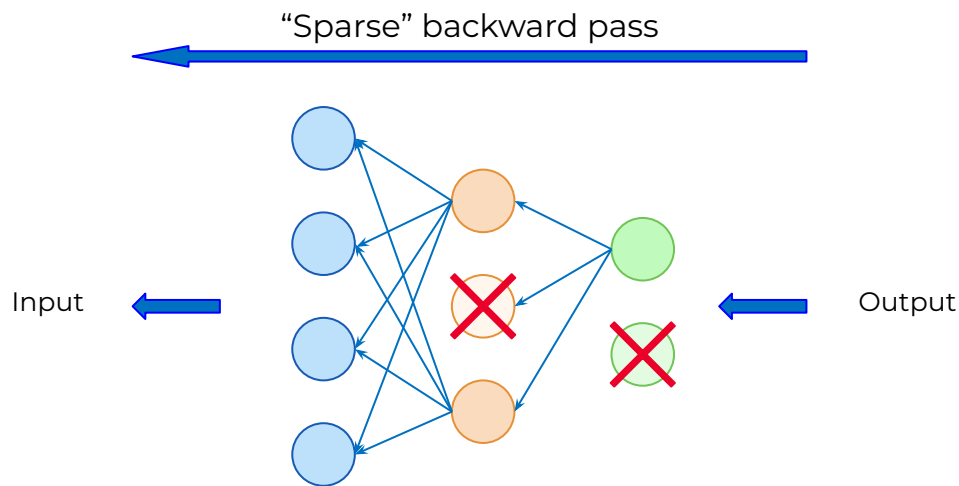
# Reducing the Training Cost
## The Idea



Backpropagation is the more computational-heavy part of the training.

We can reduce the training cost by slimming the backpropagation.

How?

# Reducing the Training Cost
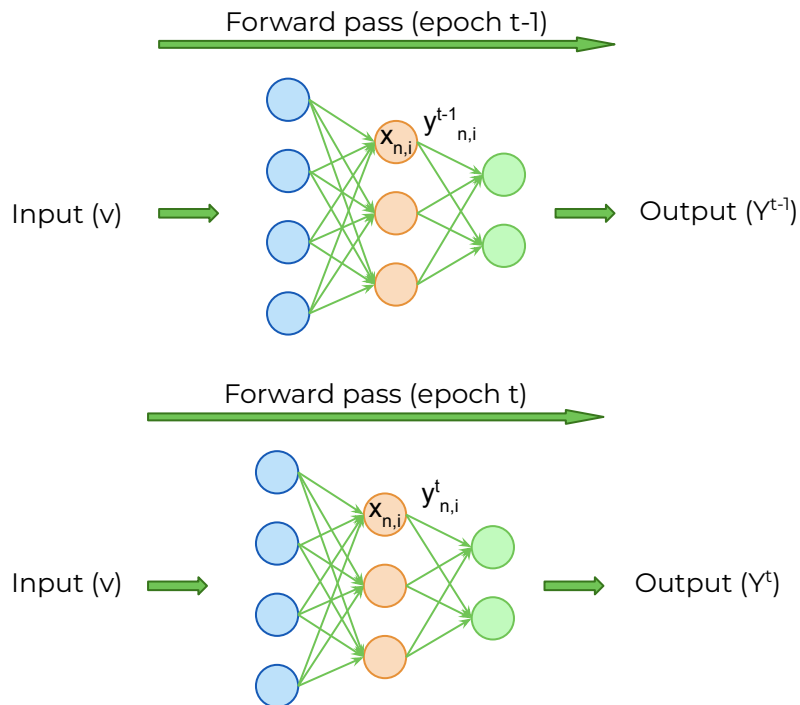## The Idea



"Sparse" backward pass

Input

Output

Are there neurons that "converge" before the end of the training? If so, we could disable their backpropagation.

How can we find these neurons?
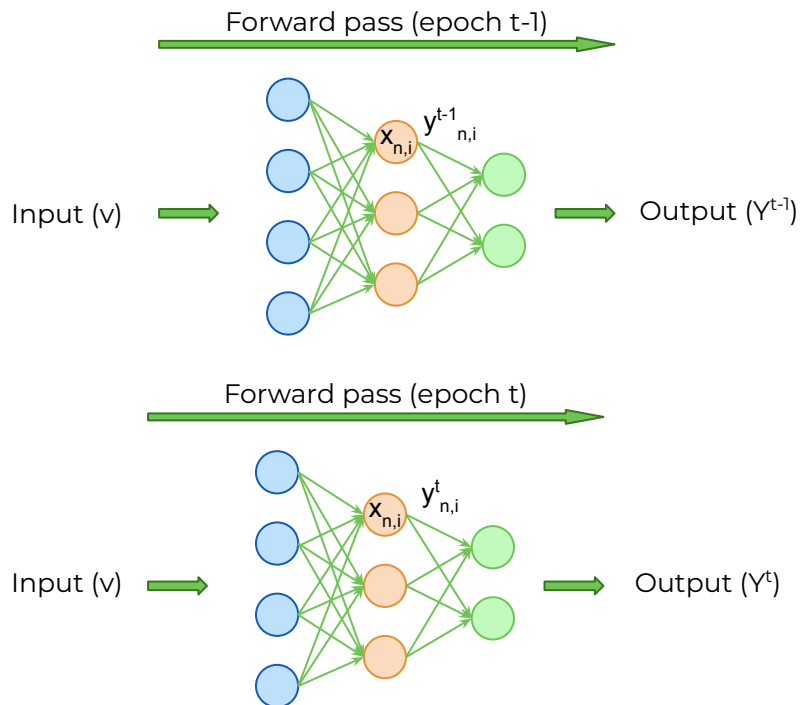
# Reducing the Training Cost
## Equilibrium evaluation



We evaluate a neuron's state and disable the update for neurons that reached equilibrium.

We consider the neuron's output in 2 adjacent epochs and evaluate the similarity.

$$\phi_{n,i}^{t} = \sum_{v}^{V} \sum_{m=1}^{M_i} y_{n,i,m,v}^{t} \cdot y_{n,i,m,v}^{t-1}$$

# Reducing the Training Cost

Equilibrium evaluation



To asses the convergence to equilibrium, we evaluate the variation of similarities.

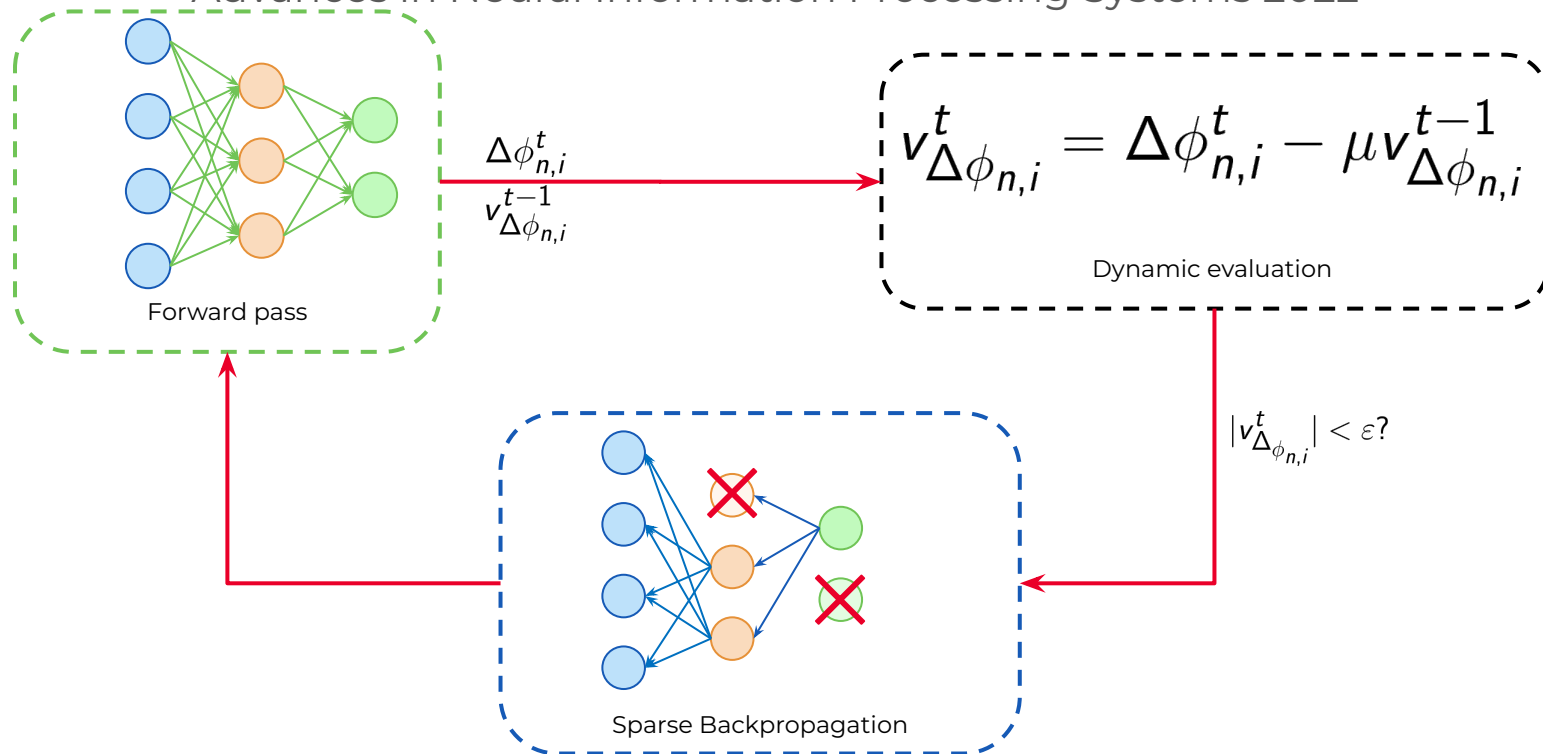$$\Delta\phi_{n,i}^{t} = \phi_{n,i}^{t} - \phi_{n,i}^{t-1}$$

We say that we reach equilibrium when

$$\Delta\phi_{n,i}^{t} \to 0$$

**A round-trip journey in pruned artificial neural networks**

# To update or not to update? Neuron at equilibrium

Advances in Neural Information Processing Systems 2022



$$v^t_{\Delta\phi_{n,i}} = \Delta\phi^t_{n,i} - \mu v^{t-1}_{\Delta\phi_{n,i}}$$

Dynamic evaluation

Forward pass

$\Delta\phi^t_{n,i}$

$v^{t-1}_{\Delta\phi_{n,i}}$

$|v^t_{\Delta\phi_{n,i}}| < \varepsilon$?

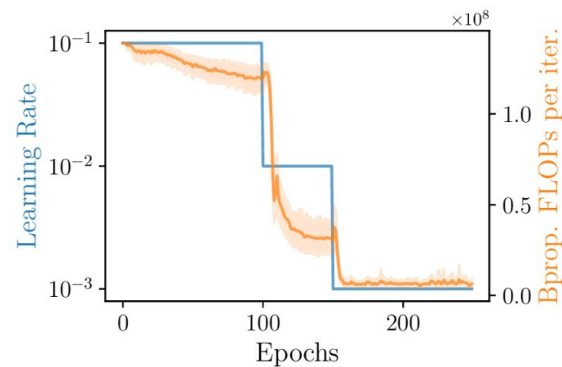Sparse Backpropagation

# NEq

## Some Neurons may Unfreeze



SGD



Adam

In the first phase of the train (high learning rate and stochastic noise) the amount of the trained neurons is higher.

Adam drives the neurons towards equilibrium faster.

At the first learning rate decay, for SGD, the number of updated neurons decreases and then increase, as SGD looks for large minima, preventing equilibrium in high learning rate regimes.

# NEq
## Experiments

We evaluate our approach on different combinations of architecture and dataset:

- ResNet-32 on CIFAR-10

- ResNet-18 on ImageNet

- Swin-B on Imagenet

- DeepLabV3 on COCO

All the learning policies used are borrowed from other works and are un-optimized to test the adaptability of NEq.

The pruning performance is evaluated according to multiple metrics:

- Average FLOPs per iteration at backpropagation.

- Final performance of the model evaluated on the test set (classification accuracy or IoU).
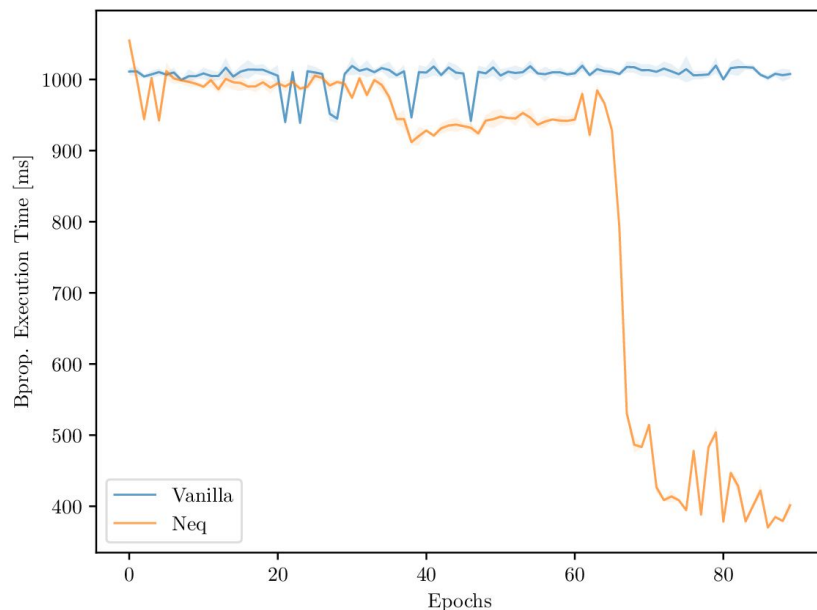
# NEq
## Experiments

| Dataset | Model | Approach | Bprop. FLOPs per iteration | Performance |
|---|---|---|---|---|
| CIFAR-10 | ResNet-32 | Baseline | 138.94M ± 0.0M | 92.85% ± 0.23%[†] |
| | | Stochastic ($p = 0.2$) | 112.99M ± 0.00M (-18.68%) | 92.78% ± 0.19% (-0.07%)[†] |
| | | Stochastic ($p = 0.5$) | 69.75M ± 0.00M (-49.8%) | 91.88% ± 0.27% (-0.97%)[†] |
| | | Stochastic* | 86.34M ± 0.00M (-37.85%) | 92.23% ± 0.25% (-0.62%)[†] |
| | | Neq | 84.81M ± 0.63M (-38.96%) | 92.96% ± 0.21% (+0.11%)[†] |
| ImageNet-1K | ResNet-18 | Baseline | 3.64G ± 0.0G | 69.90% ± 0.04%[†] |
| | | Stochastic ($p = 0.2$) | 2.94G ± 0.00G (-19.26%) | 69.42% ± 0.16% (-0.48%)[†] |
| | | Stochastic ($p = 0.5$) | 1.85G ± 0.00G (-49.11%) | 69.18% ± 0.03% (-0.72%)[†] |
| | | Stochastic* | 2.82G ± 0.00G (-22.58%) | 69.45% ± 0.06% (-0.45%)[†] |
| | | Neq | 2.80G ± 0.03G (-23.08%) | 69.62% ± 0.06% (-0.28%)[†] |
| | Swin-B | Baseline | 30.28G ± 0.00G | 84.71% ± 0.04% [†] |
| | | Stochastic ($p = 0.2$) | 24.65G ± 0.00G (-18.6%) | 84.54% ± 0.04% (-0.83%)[†] |
| | | Stochastic ($p = 0.5$) | 16.15G ± 0.00G (-46.67%) | 84.40% ± 0.02% (-0.31%)[†] |
| | | Stochastic* | 11.02G ± 0.00G (-63.67%) | 84.27% ± 0.04% (-0.44%)[†] |
| | | Neq | 10.78G ± 0.02G (-64.39%) | 84.35% ± 0.02% (-0.36%)[†] |
| COCO | DeepLabv3 | Baseline | 305.06G ± 0.0G | 67.71% ± 0.02%[‡] |
| | | Stochastic ($p = 0.2$) | 248.69G ± 0.00G (-18.48%) | 67.11% ± 0.02% (-0.60%)[‡] |
| | | Stochastic ($p = 0.5$) | 163.42G ± 0.00G (-46.43%) | 66.91% ± 0.04% (-0.80%)[‡] |
| | | Stochastic* | 229.00G ± 0.00G (-24.93%) | 67.02% ± 0.03% (-0.69%)[‡] |
| | | Neq | 217.29G ± 0.04G (-28.77%) | 67.22% ± 0.04% (-0.49%)[‡] |

# NEq

## Faster Backpropagation



Backpropagation execution time for vanilla and NEq ResNet-18.

We observe a reduction in the wall-clock time of around -17.52%.

# Outline
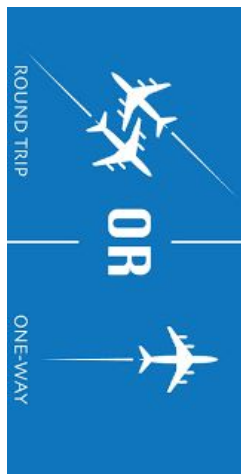
**A round-trip journey in pruned artificial neural networks**

# Conclusions

- We shared our recent experiences in **frugal deep learning**
  - **one way**: simplify the model via pruning for faster/lower memory footprint when deploying nets
  - **return**: prune the backward pass to reduce the training cost

- Future research
  - joint approaches including quantization and target device constraints
  - dataset pruning
  - efficient automatic hyper-parameter tuning