

# NEURAL NETWORKS

# SPEED-UP AND COMPRESSION

**Lecture 2: Pruning**

# RECAP FROM LECTURE 1

We have considered

- Main neural network layers (linear, non-linear, normalization)
- Theoretical measurements of speed and memory
  - FLOP / MAC
  - number of parameters
- Empirical measurements of speed and memory
  - FlopCo
  - PyTorch tools
- Arithmetic intensity

# OUTLINE

- Neurons and connections in neural networks
  - Fully-connected & Convolutional layers
  - Pruning
- Unstructured pruning
- Structured pruning

# NEURONS AND CONNECTIONS

# FULLY-CONNECTED LAYER

Fully-connected layer

$$Y = WX$$

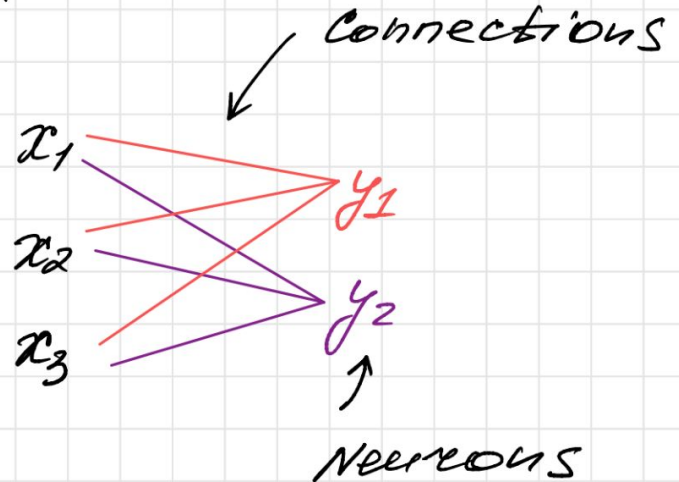
$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$

 $\cdot$ 

$x_1$
$x_2$
$x_3$

 $=$ 

$y_1$
$y_2$



# FULLY-CONNECTED LAYER: PRUNE

Fully-connected layer

$$Y = W^T X$$

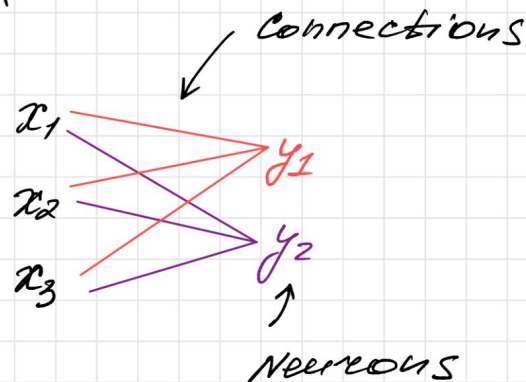
$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$

 $\cdot$ 

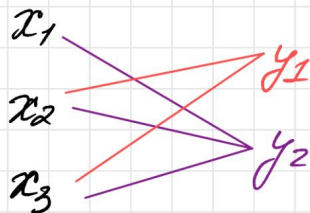
$x_1$
$x_2$
$x_3$

 $=$ 

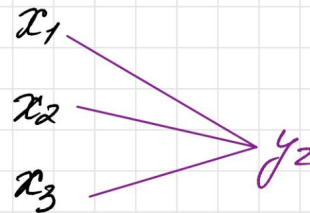
$y_1$
$y_2$



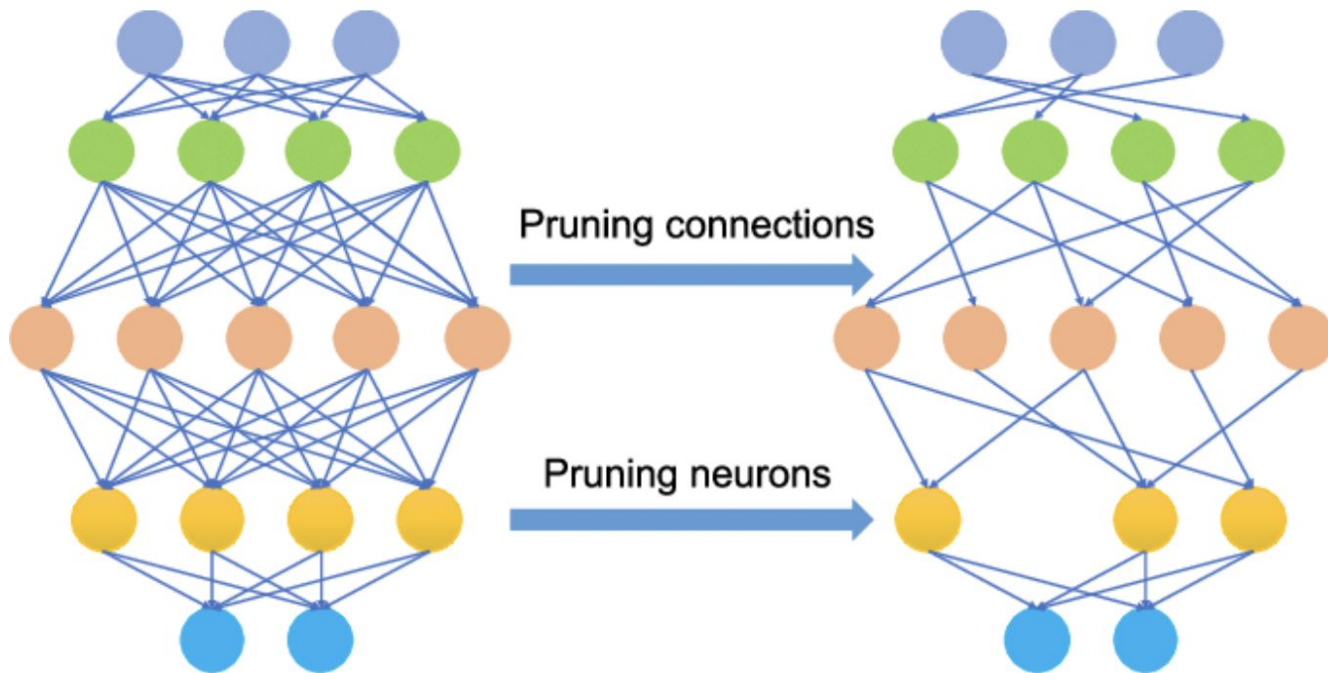
Prune connection  $w_{11}$



Prune neuron  $y_1$



# FULLY-CONNECTED NEURAL NETWORK



# CONVOLUTIONAL LAYER

Convolutional layer

$$Y = X * W$$

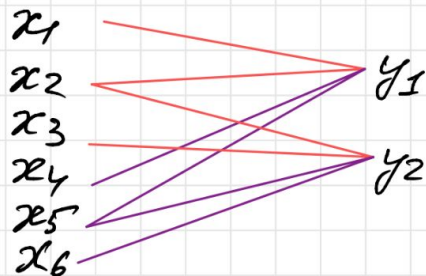
$x_1$	$x_2$	$x_3$
$x_4$	$x_5$	$x_6$
$x_7$	$x_8$	$x_9$

\*

$w_{11}$	$w_{12}$
$w_{21}$	$w_{22}$

=

$y_1 = w_{11}x_1 + w_{12}x_2 +$ $w_{21}x_4 + w_{22}x_5$	$y_2 = w_{11}x_2 + w_{12}x_3 +$ $w_{21}x_5 + w_{22}x_6$
$y_3 = w_{11}x_4 + w_{12}x_5 +$ $w_{21}x_7 + w_{22}x_8$	$y_4 = w_{11}x_5 + w_{12}x_6 +$ $w_{21}x_8 + w_{22}x_9$



In our example

$$C_{in} = C_{out} = 1$$

$$k_h = k_w = 2$$



# CONVOLUTIONAL LAYER: PRUNE SHAPE (ROW)

Convolutional layer  $Y = X * W$

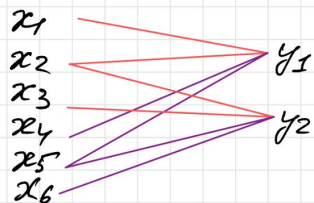
$x_1$	$x_2$	$x_3$
$x_4$	$x_5$	$x_6$
$x_7$	$x_8$	$x_9$

 $\times$ 

$w_{11}$	$w_{12}$
$w_{21}$	$w_{22}$

 $=$ 

$y_1 = w_{11}x_1 + w_{12}x_2 + w_{21}x_4 + w_{22}x_5$	$y_2 = w_{11}x_2 + w_{12}x_3 + w_{21}x_5 + w_{22}x_6$
$y_3 = w_{11}x_4 + w_{12}x_5 + w_{21}x_7 + w_{22}x_8$	$y_4 = w_{11}x_5 + w_{12}x_6 + w_{21}x_8 + w_{22}x_9$

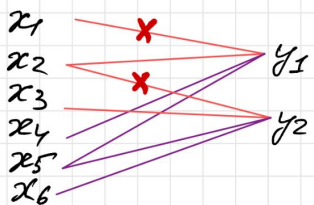


In our example

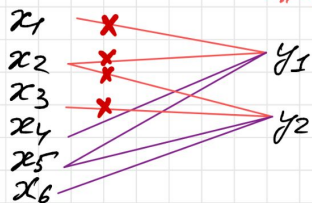
$$C_{in} = C_{out} = 1$$

$$k_h = k_w = 2$$

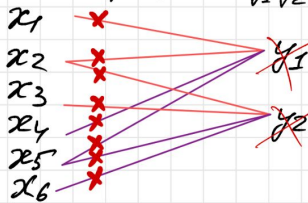
Prune weight  $w_{11}$



Prune weight row  $w_{11}, w_{12}$



Prune feature map row  $y_1, y_2$

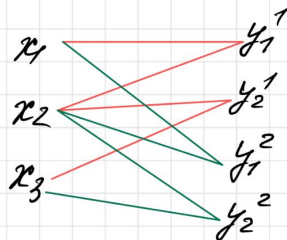


# CONVOLUTIONAL LAYER: PRUNE FILTER

Convolutional layer

$$Y = X * W$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} * \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{11}^1 & w_{12}^1 \end{bmatrix} = \begin{bmatrix} y_1^2 = w_{11}^2 x_1 + w_{12}^2 x_2 & y_2^2 = w_{11}^2 x_2 + w_{12}^2 x_3 \\ y_1^1 = w_{11}^1 x_1 + w_{12}^1 x_2 & y_2^1 = w_{11}^1 x_2 + w_{12}^1 x_3 \end{bmatrix}$$

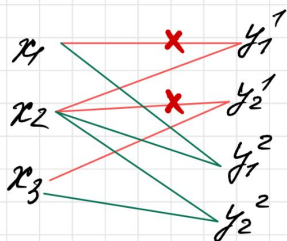


In our example:

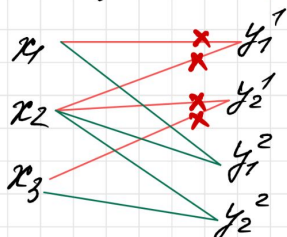
$$C_{in} = 1, \quad C_{out} = 2$$

$$k_h = 1, \quad k_w = 2$$

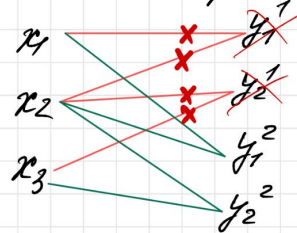
Prune weight  $w_{11}^1$



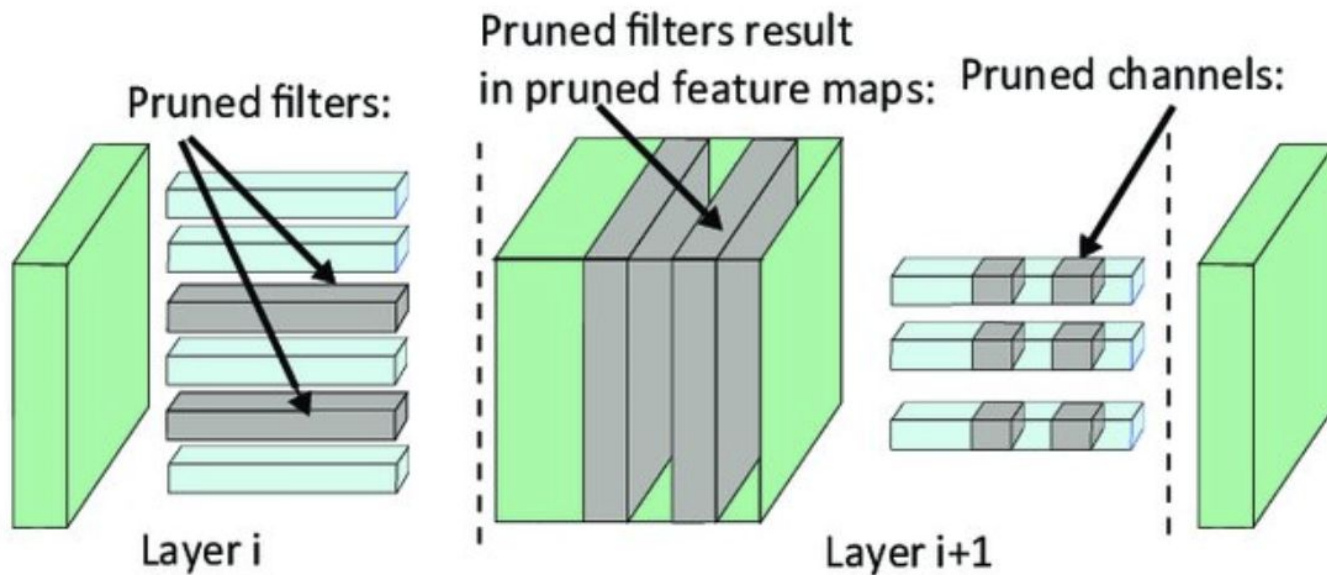
Prune weight filter  $w_{11}^1, w_{12}^1$



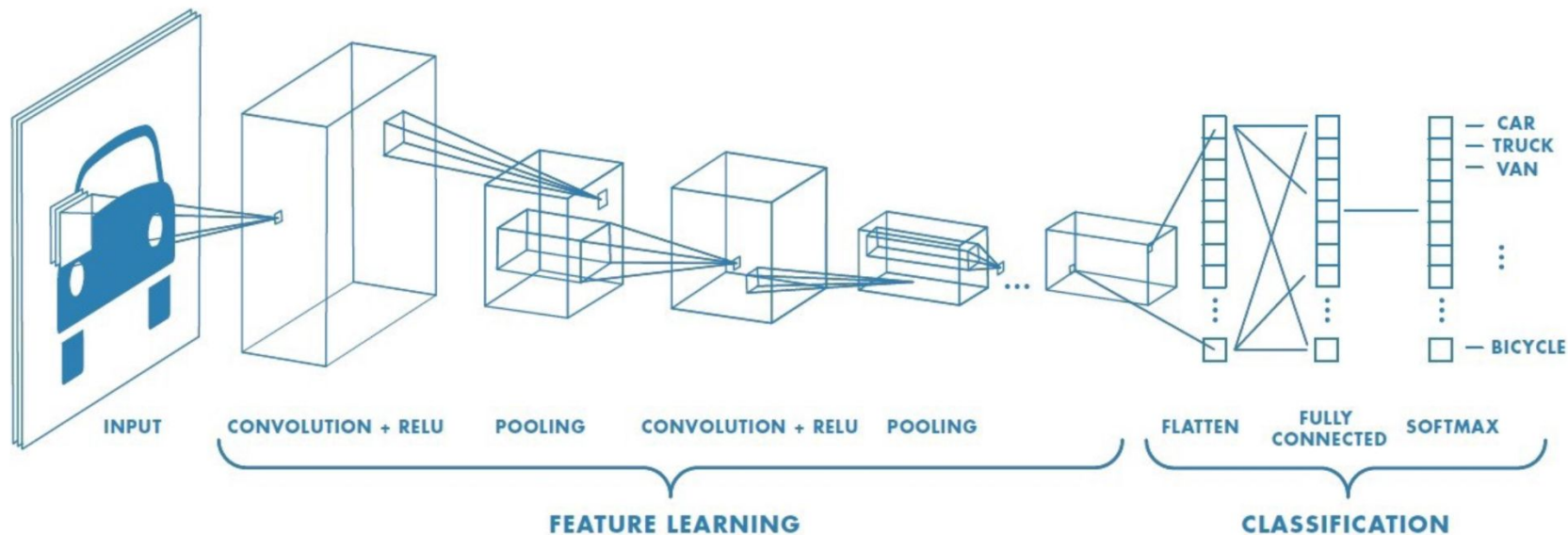
Prune feature map



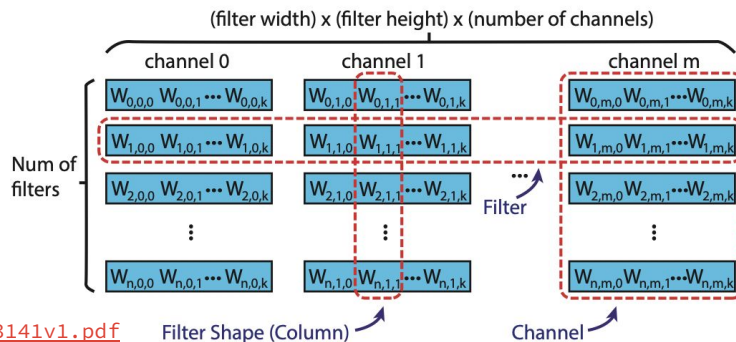
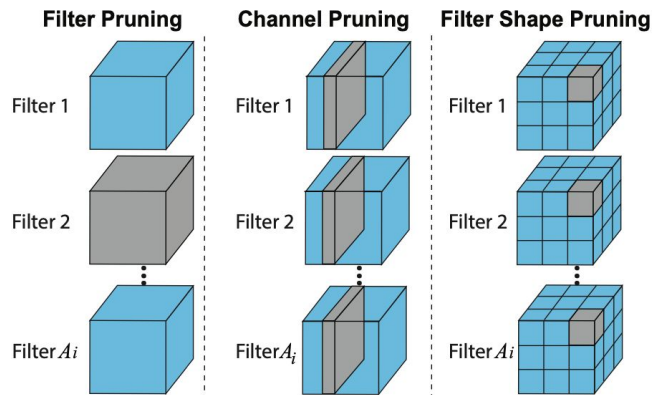
# CONVOLUTIONAL NEURAL NETWORK



# CONVOLUTIONAL NEURAL NETWORK



# CONVOLUTIONAL LAYER: PRUNE FILTER/CHANNEL/SHAPE



# WHAT'S THE PROBLEM?

DL model limitations:

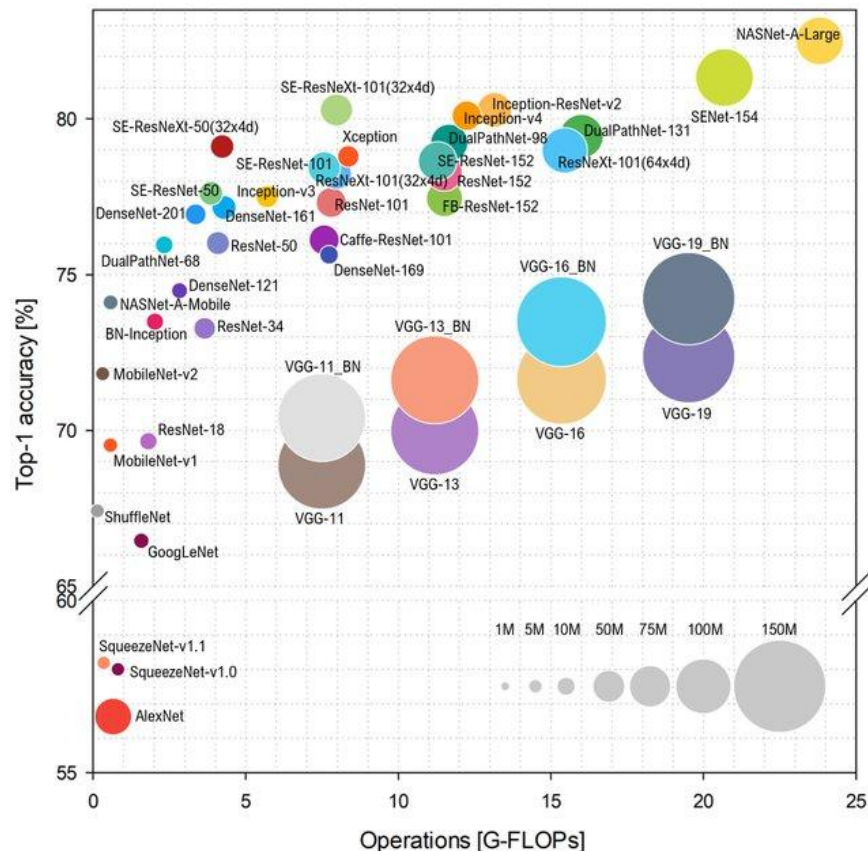
- High memory consumption
- Huge computational requirements
- Great power consumption



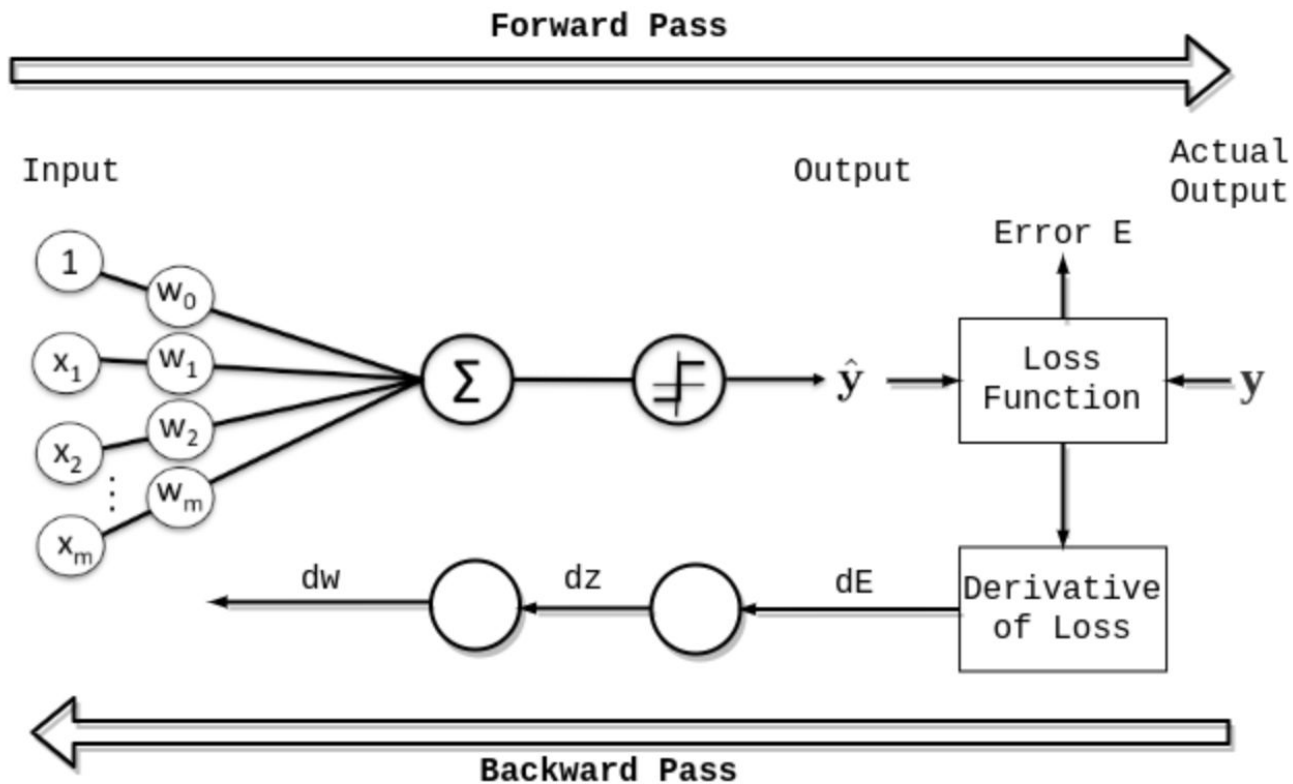
Difficult to deploy on portable devices  
(e.g. laptops and smartphones)



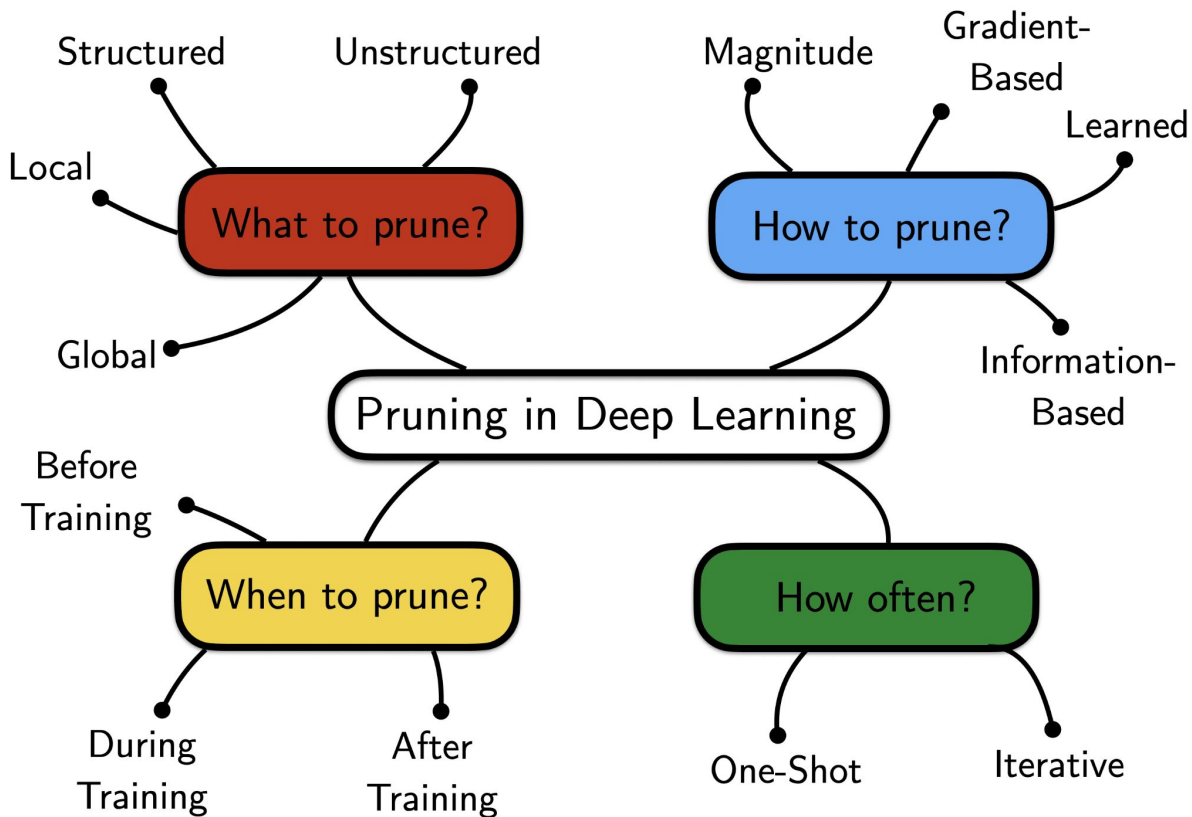
Efficient architecture design is required



# RECAP: EPOCH OF NEURAL NETWORK TRAINING



# PRUNING

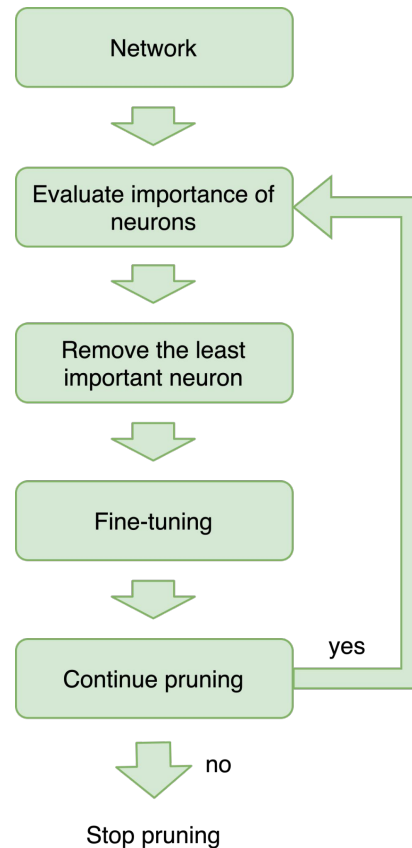
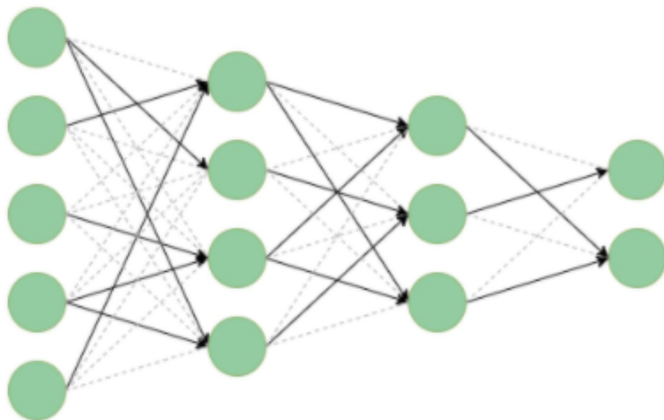




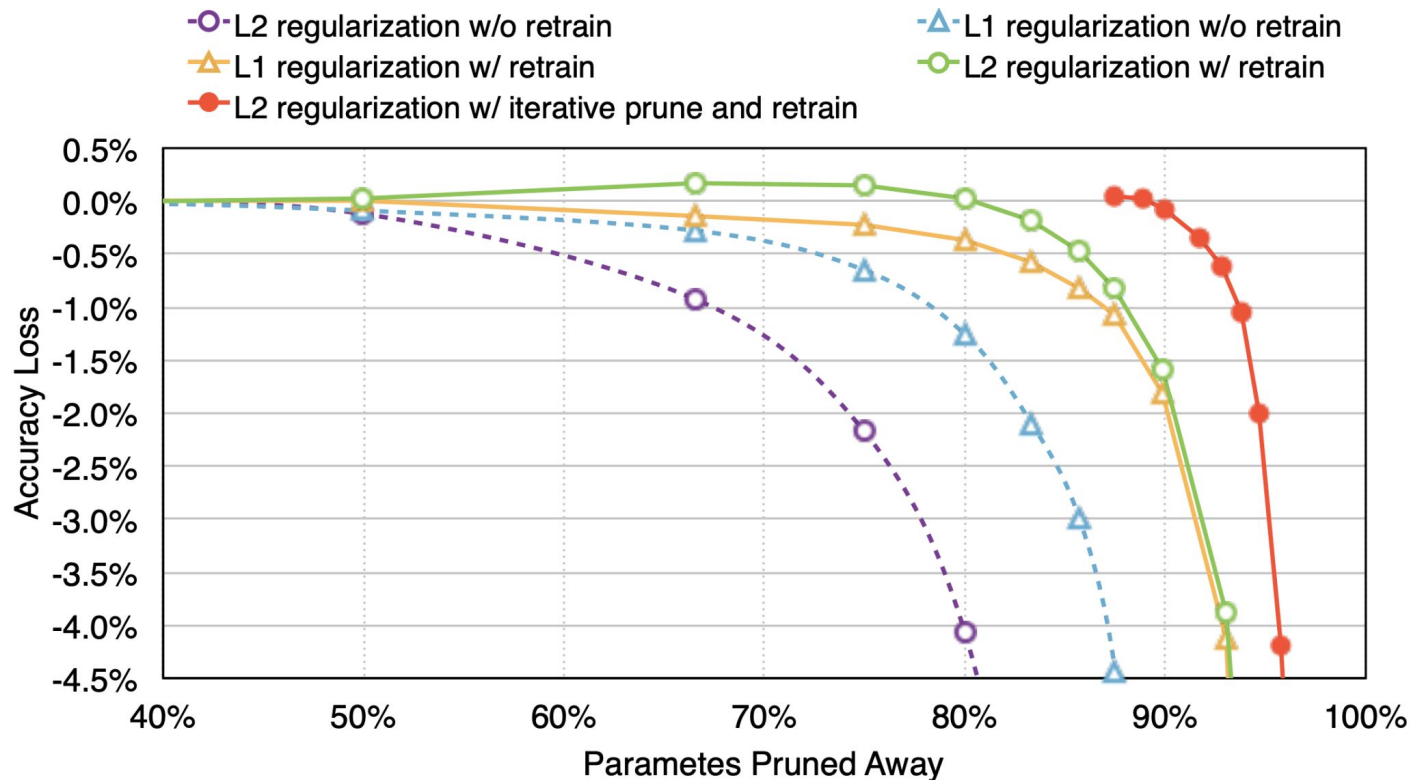
# UNSTRUCTURED PRUNING

# UNSTRUCTURED PRUNING

or Fine-Grained Pruning or Weight Sparsification



# UNSTRUCTURED PRUNING

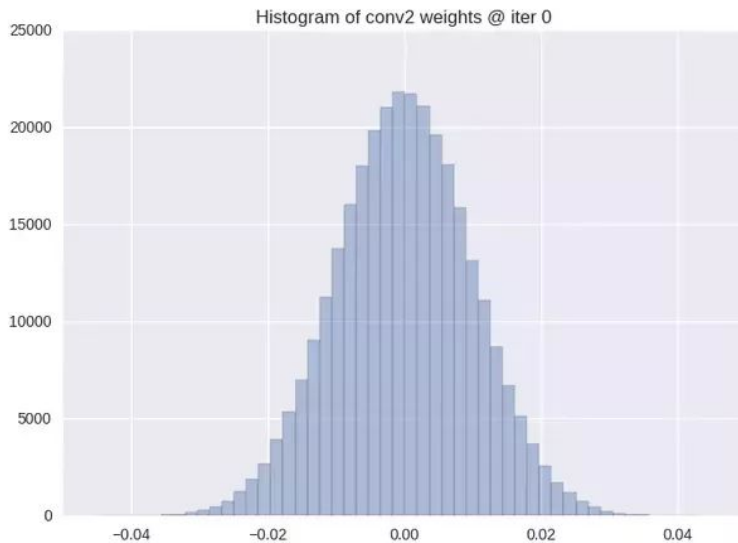


# UNSTRUCTURED PRUNING

List of possible criteria:

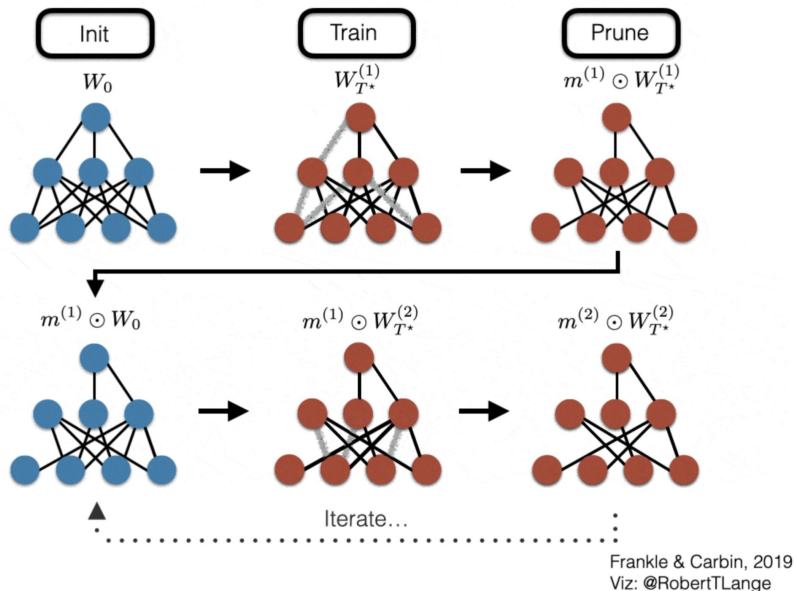
- Weight-based criteria (L1/L2 norm)
- Gradient-based criteria

$$\delta E = \sum_i g_i \delta u_i + \frac{1}{2} \sum_i h_{ii} \delta u_i^2 + \frac{1}{2} \sum_{i \neq j} h_{ij} \delta u_i \delta u_j + O(\|\delta \mathcal{U}\|^3)$$
$$g_i = \frac{\partial E}{\partial u_i} \quad \text{and} \quad h_{ij} = \frac{\partial^2 E}{\partial u_i \partial u_j}$$

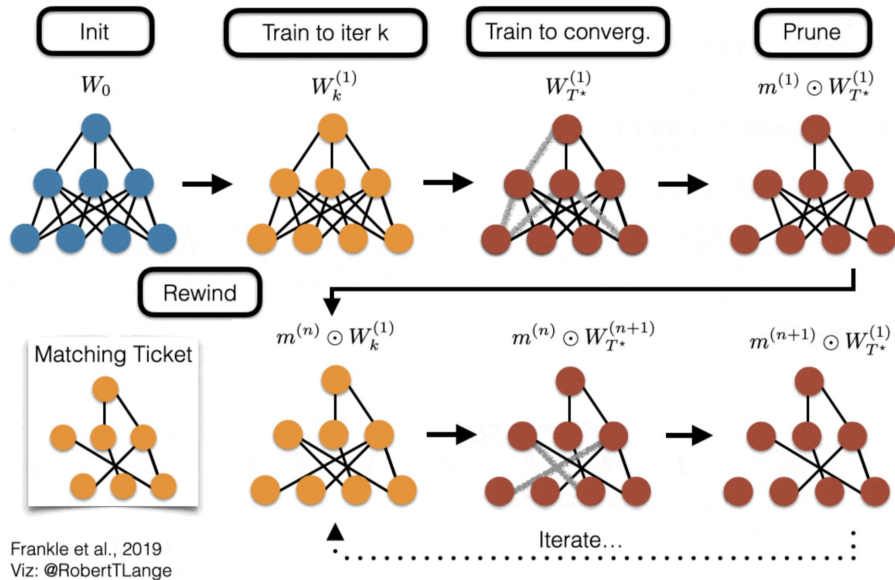


# UNSTRUCTURED PRUNING: LOTTERY TICKET HYPOTHESIS

## Searching for Tickets: Iterative Magnitude Pruning

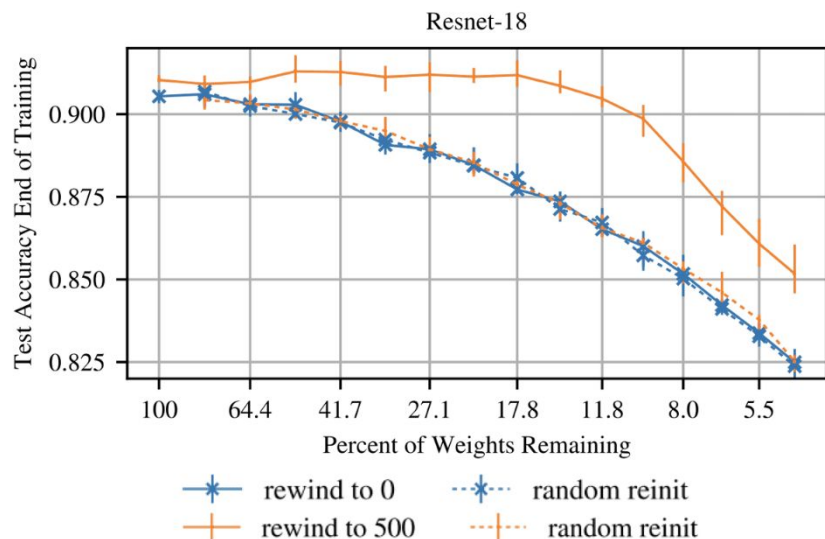


## Iterative Magnitude Pruning with Rewinding

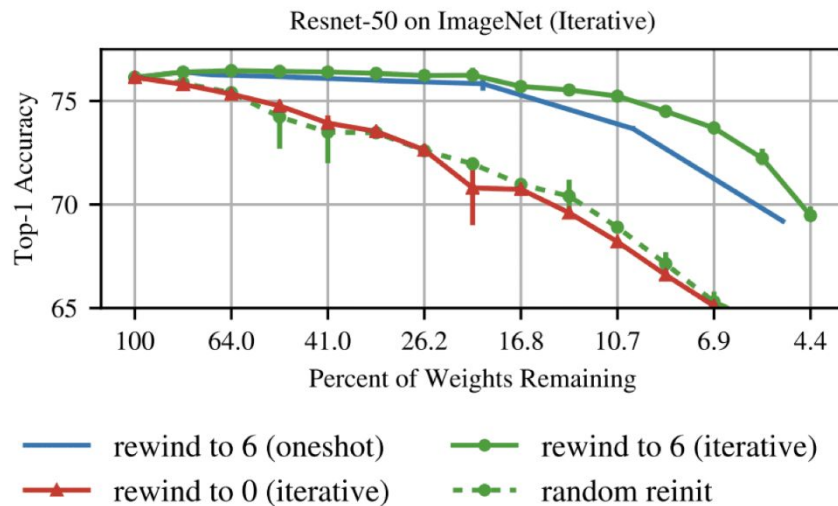


# UNSTRUCTURED PRUNING: LOTTERY TICKET HYPOTHESIS

**A** Rewinding Resnet-20 on CIFAR-10



**B** Rewinding Resnet-50 on ImageNet



# UNSTRUCTURED PRUNING

- Benefits

- Achieves high rates of weight reduction without acc. drop (~ 95% of weights can be removed)

- Drawbacks

- Requires hardware support for sparse computation speedup
- Hard to find sparsity level for all layers of the network

# STRUCTURED PRUNING



# STRUCTURED PRUNING

Removing structural parts instead of individual weights

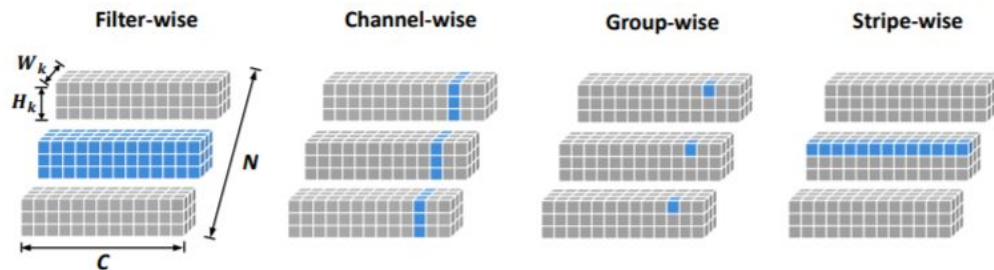
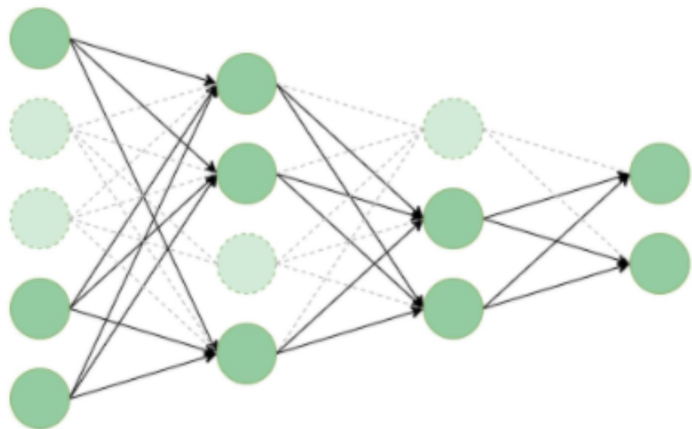
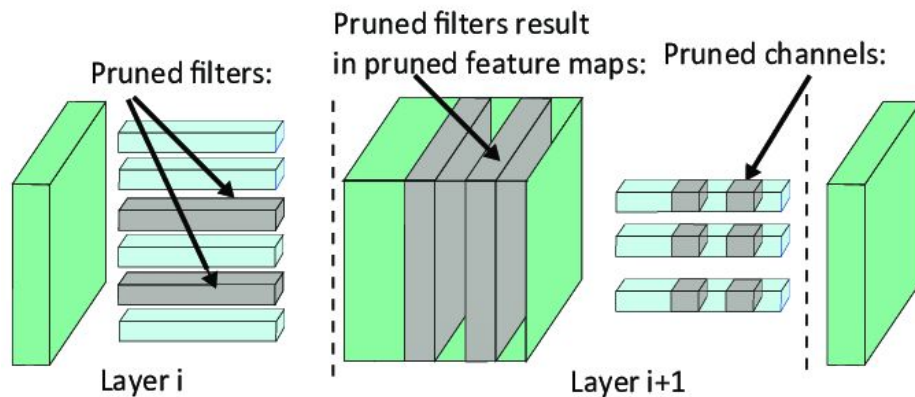


Figure 2: The visualization of different types of pruning.

# STRUCTURED PRUNING

List of possible criteria:

- Weight-based criteria:
  - L1/L2 norm
  - Scaling parameter in BN
- Activation-based criteria:
  - PCA of activations
- Gradient-based criteria
- Greedy and One-shot Pruning



# STRUCTURED PRUNING

- Benefits

- Efficiently accelerates model
- No need of hardware/software support (pruned model is structurally equivalent to initial model)

- Drawbacks

- Pruning channels/filters in one layer affects previous/subsequent layers
- Hard to find best filter configuration for the whole network

# RESUME

# RESUME

Current pipeline in pruning can be divided into two components

- Identifying the most promising neurons to be pruned
- Training and fine-tuning the pruned model to recover the base model's prediction performance

A successful pruning algorithm is an iterative progression of these components

```
Data: CNN model, training data
while Compression requirement not met or exceed budget
do
  | train model (to convergence);
  | compute pruning criteria;
  | prune parameters below threshold;
end
```

**Algorithm 1:** Workflow for model pruning.

# RESUME

Table 1: Saliency measurements used in pruning.

	Basic idea	Saliency Expression
• Data-agnostic	Minimize pruned deterioration	$H_{ii}w_i^2; w_i^2/H_{ii}^{-1}$
	Remove weights with small values	$ w $
	Weight similarity and redundancy	$  w_i - w_j  ;   w - \mathbb{E}_{\text{geo}}[w]  $
	Structured L1/2-norm penalty	Penalize $\sum_{c'} w_{c',c,i,j}^2; \sum_{c,i,j} w_{c',c,i,j}^2$
	Magnitude of Batch Norm	$\sum_{c'} w_{c',c,i,j}^2; \sum_{c,i,j}  w_{c',c,i,j} $
• Data-aware	Remove by filter similarity	-
	Remove inactive neurons (APoZ)	Geometric mean
	Remove activations with flat gradient	$\sum_i \mathbb{I}_{z_i=0}$
	Reconstruction error on channel pruning	$ \sum_i \frac{\mathcal{L}}{\partial z_i} z_i $
	Entropy	$  z_i - z_i^{\text{pruned}}  $
	Reconstruction error and L1 norm	$\sum_m P(z_i) \log P(z_i)$ , where $P(z_i)$ is probability of activations in bins
	Neural Importance Score	-
	Remove insensitive neurons	$s_k =  w_{i+1} ^T s_{k+1}$ Reset weights with small updates to initial value

# PRUNING

