# Speeding up and compression of Convolutional Neural Networks
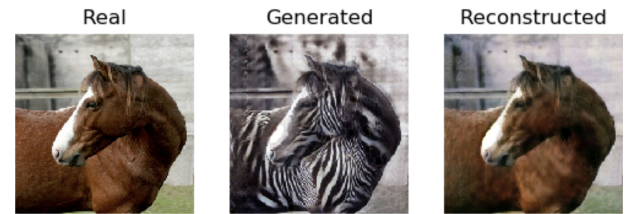
Konstantin Sobolev, Skoltech

# Outline

1. Introduction
   a. Motivation
   b. Preliminaries
   c. CNN reminder
2. Efficient architecture design
   a. MobileNet
   b. ShuffleNet
   c. GhostNet
3. Model compression
   a. Weight decomposition
   b. Fine-grained pruning
   c. Structural pruning
   d. Quantization
   e. Knowledge Distillation

# Motivation

Increasing importance and number of practical applications of CNN applications:
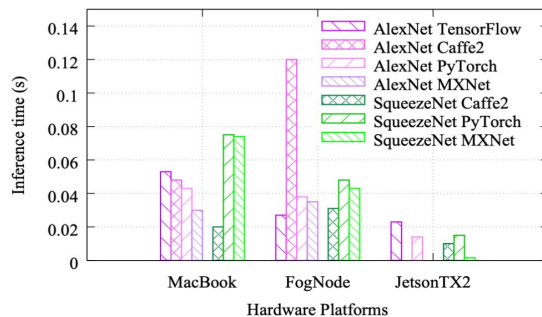
# Motivation

DL model limitations:

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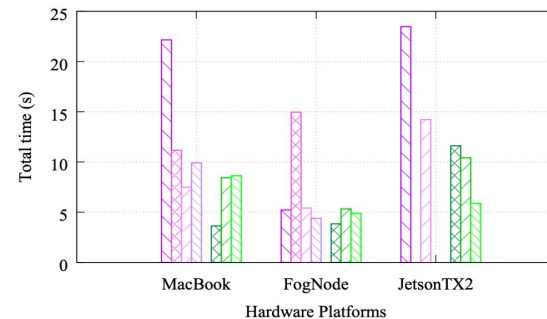- Great power consumption

Difficult to deploy on portable devices

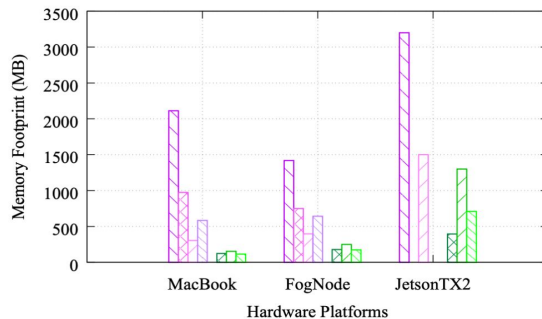(e.g. laptops and smartphones)

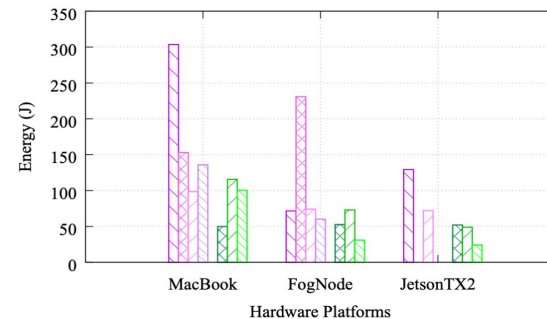Efficient architecture design is required



(a) Inference time

(b) Total time

(c) Memory footprint

(d) Energy

Comparison between cumbersome (AlexNet) and light-weight (SqueezeNet) CNN architectures on different edge platforms (MacBook, FogNode and JetsonTX2) and frameworks (TensorFlow, Caffe2, PyTorch and MXNet)

# Metrics to optimize

**Real metrics:**

- Inference time
- Memory consumption
- Battery consumption

**Proxy metrics:**

- **FLOP**s - number of computational operations required for inference
- **N** parameters - number of trained parameters
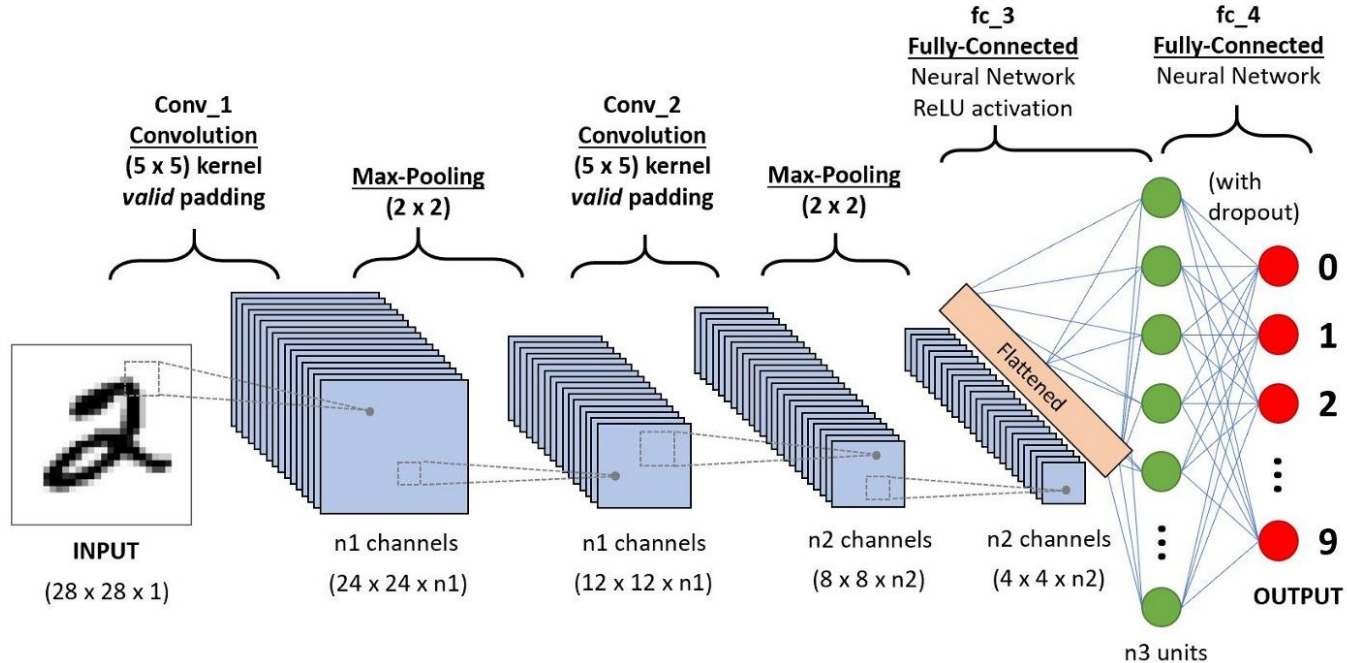
# Efficient architecture design methods

Efficient neural architecture design

- MobileNet
- ShuffleNet
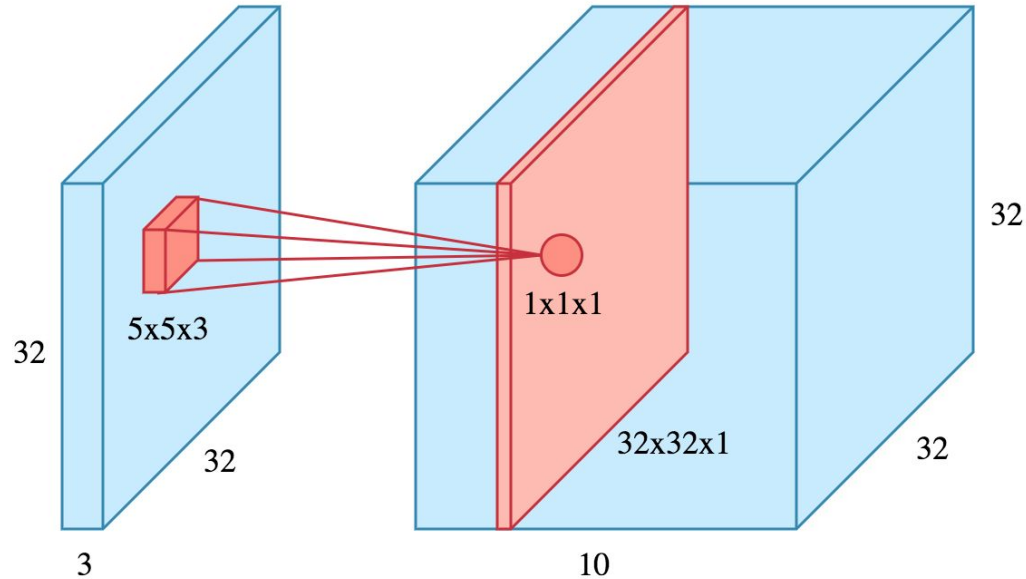- GhostNet

Neural network compression

- Weight Decomposition
- Unstructured Pruning
- Structured Pruning
- Quantization
- **Knowledge Distillation**

# Convolutional Neural Network reminder

Source: https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53

# Convolutional Neural Network reminder

Most computations are concentrated in convolutional layer

# Design efficient models from scratch

# MobileNet

Idea: Replace ordinary convolution by depth-wise separable convolution

Original convolution complexity:

$$D_k * D_k * C_{in} * D_f * D_f * C_{out}$$

Depth-wise separable convolution complexity:
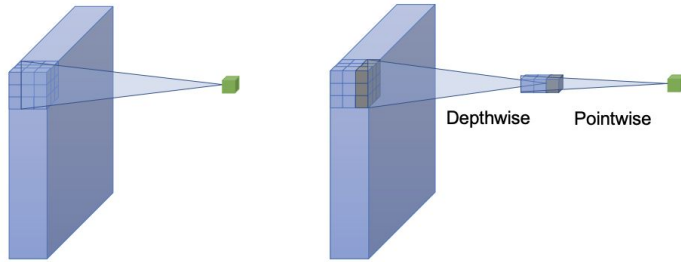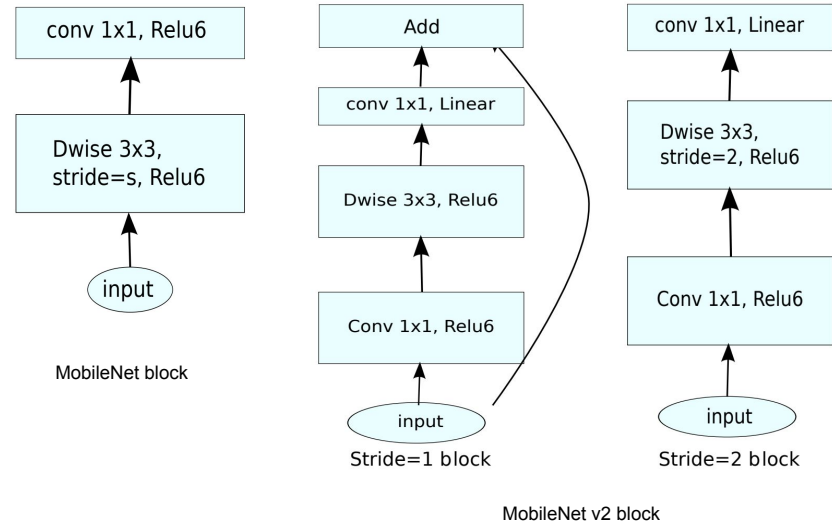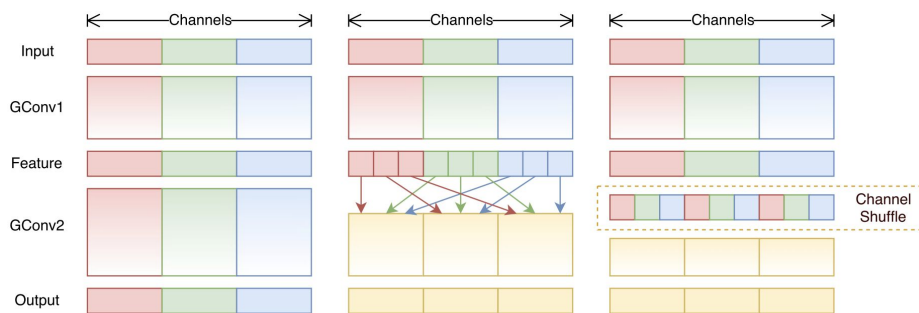
$$(D_k * D_k + C_{out}) * C_{in} * D_f * D_f$$

Depthwise     Pointwise

Figure 3: Standard convolution and depthwise separable convolution.

conv 1x1, Relu6

Dwise 3x3, stride=s, Relu6

input

MobileNet block

Add

conv 1x1, Linear

Dwise 3x3, Relu6

Conv 1x1, Relu6

input

Stride=1 block

conv 1x1, Linear

Dwise 3x3, stride=2, Relu6

Conv 1x1, Relu6

input

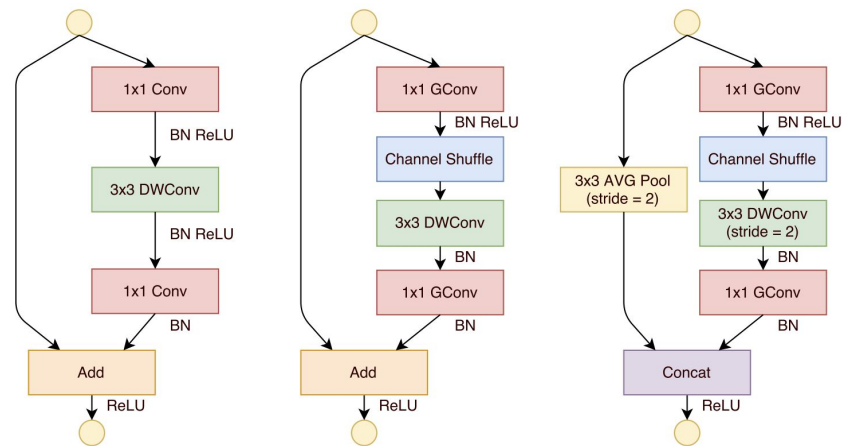Stride=2 block

MobileNet v2 block

10

# ShuffleNet

Idea: 1x1 Convolution is still an expensive operation. It can be accelerated by channel shuffling operation.
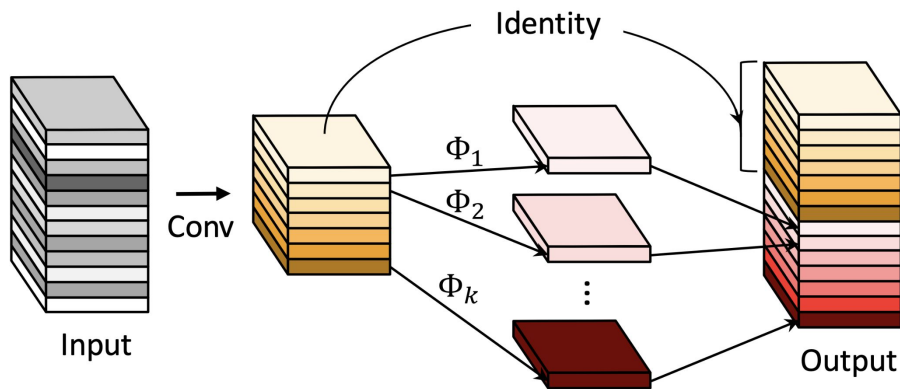


ShuffleNet operation visualization

MobileNet block

Shufflenet block

ShuffleNet block

# GhostNet

Idea: Generate redundancy in featuremap by cheap operations
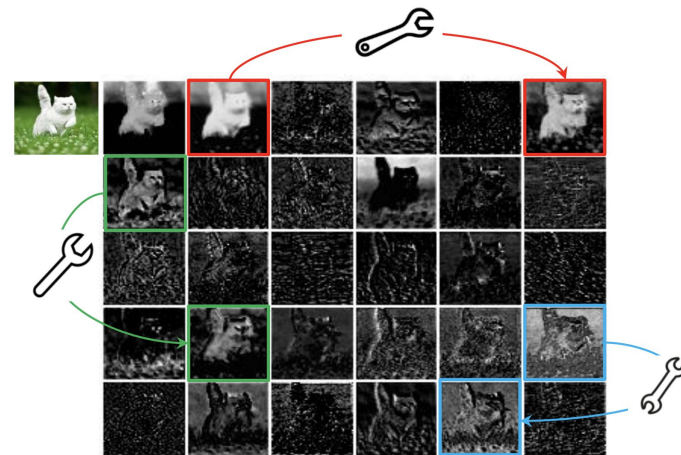


(b) The Ghost module.



Figure 1. Visualization of some feature maps generated by the first residual group in ResNet-50, where three similar feature map pair examples are annotated with boxes of the same color. One feature map in the pair can be approximately obtained by transforming the other one through cheap operations (denoted by spanners).

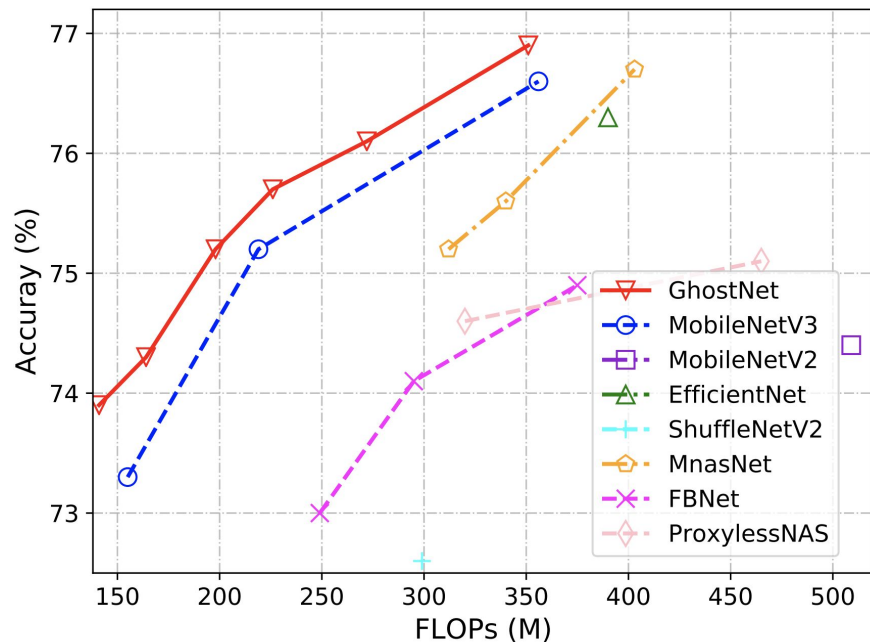Source: https://arxiv.org/pdf/1911.11907.pdf

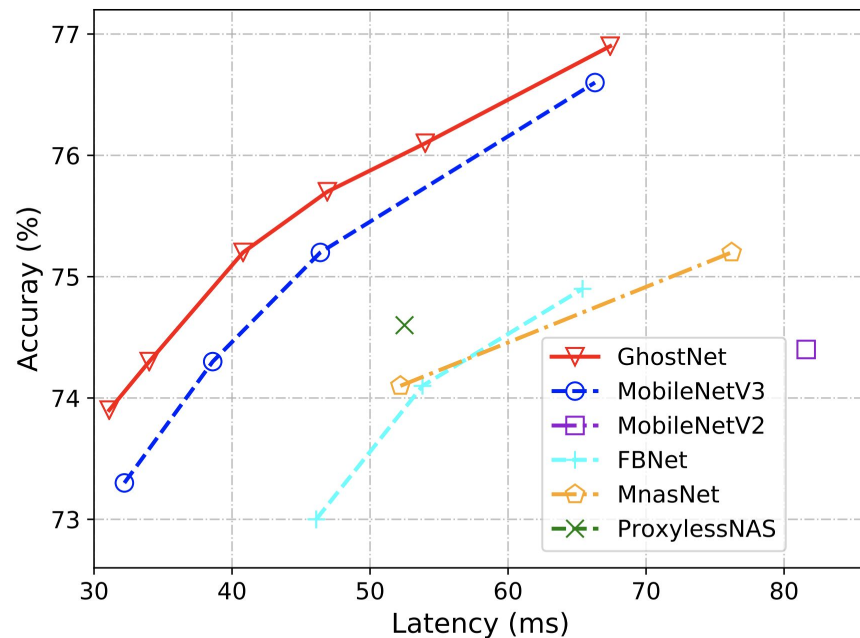# Results



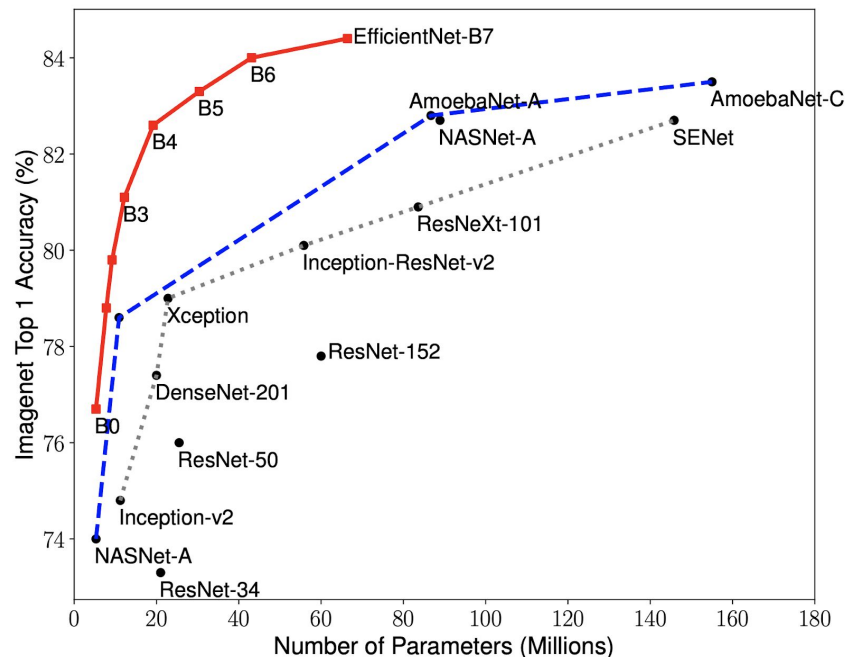Figure 6. Top-1 accuracy *v.s.* FLOPs on ImageNet dataset.



Figure 7. Top-1 accuracy *v.s.* latency on ImageNet dataset.

# Model compression

# Model compression

- Most neural networks are redundant

- Redundancy helps to achieve a better performance

- Redundancy helps during training but harms during inference

# Weight Decomposition: Preliminaries

| Scalar | Vector | Matrix | Tensor |
|--------|--------|--------|--------|

| Scalar-level | Matrix-level | Block-Matrix level | Tensor level | Level of thinking in algebra field |
|---|---|---|---|---|
| 1960's | 1980's | 2000's | Now | |

# Weight Decomposition: Preliminaries

Matrix decomposition:



Tensor decomposition:

Source: https://www.researchgate.net/publication/359367892_electronics-11-009451/figures?lo=1&utm_source=google&utm_medium=organic

# Weight Decomposition

**Pipeline:**
1. Extract convolutional kernel
2. Decompose it into factors: A, G, B
3. Replace initial layer by sequence of layers with factors as kernels
4. Fine-tune network

**Result:**
1. Faster inference
2. Lower memory consumption
3. Longer battery life

# Weight Decomposition



(a) Weight SVD

(b) Spatial SVD

(c) CP-decomposition

(d) Tucker decomposition

(e) Tensor-train decomposition



Imagenet (Resnet18)

- uncompressed model
- Spatial SVD
- Weight SVD
- Tucker decomposition
- Tensor-train decomposition
- CP-decomposition

# Weight Decomposition

Pros:

- Achieves high theoretical speedup with low performance degradation
- Does not require additional hardware support

Cons:

- Poorly works with non-standard convolutions (PW, DW, Group)
- Requires careful rank selection

# Pruning: Preliminaries

Fully-connected layer     $Y = W^T X$

$$\begin{array}{|c|c|c|} \hline W_{11} & W_{12} & W_{13} \\ \hline W_{21} & W_{22} & W_{23} \\ \hline \end{array} \cdot \begin{array}{|c|} \hline x_1 \\ \hline x_2 \\ \hline x_3 \\ \hline \end{array} = \begin{array}{|c|} \hline y_1 \\ \hline y_2 \\ \hline \end{array}$$

Connections

$x_1$
$x_2$
$x_3$

$y_1$
$y_2$

Neurons

Prune connection $W_{11}$

$x_1$
$x_2$
$x_3$

$y_1$
$y_2$

Prune neuron $y_1$

$x_1$
$x_2$
$x_3$

$y_2$

# Unstructured Pruning

aka Fine-Grained Pruning aka Weight Sparsification
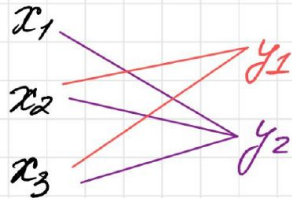


Network

↓

Evaluate importance of neurons

↓

Remove the least important neuron

↓

Fine-tuning

↓

Continue pruning — yes

↓ no

Stop pruning

Source: https://arxiv.org/pdf/1506.02626.pdf

# Unstructured Pruning

Source:    https://arxiv.org/pdf/1506.02626.pdf

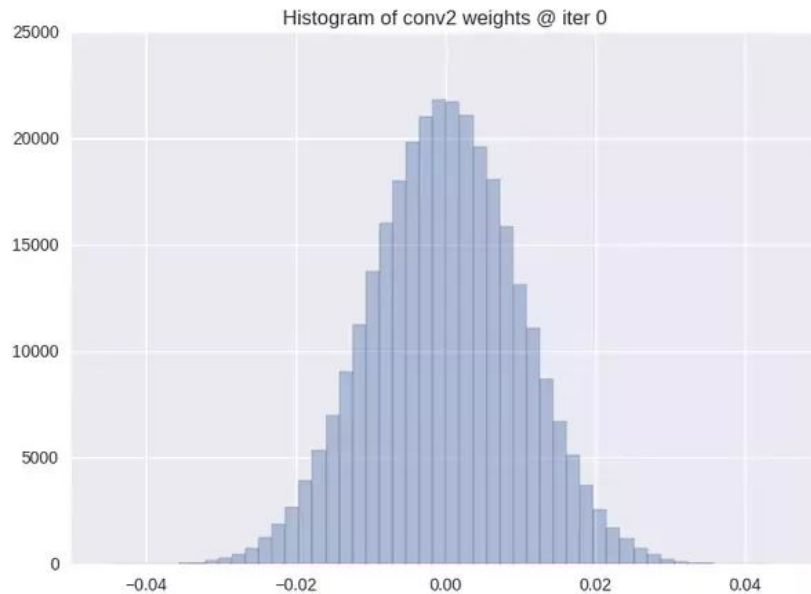# Unstructured Pruning

List of possible criterions:

- Weight-based criteria (L1/L2 norm)

- Gradient-based criteria



Histogram of conv2 weights @ iter 0

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

$$g_i = \frac{\partial E}{\partial u_i} \quad \text{and} \quad h_{ij} = \frac{\partial^2 E}{\partial u_i \partial u_j}$$

Source:    https://openaccess.thecvf.com/content_CVPR_2020/papers/He_Learning_Filter_Pruning_Criteria_for_Deep_Convolutional_Neural_Networks_Acceleration_CVPR_2020_paper.pdf
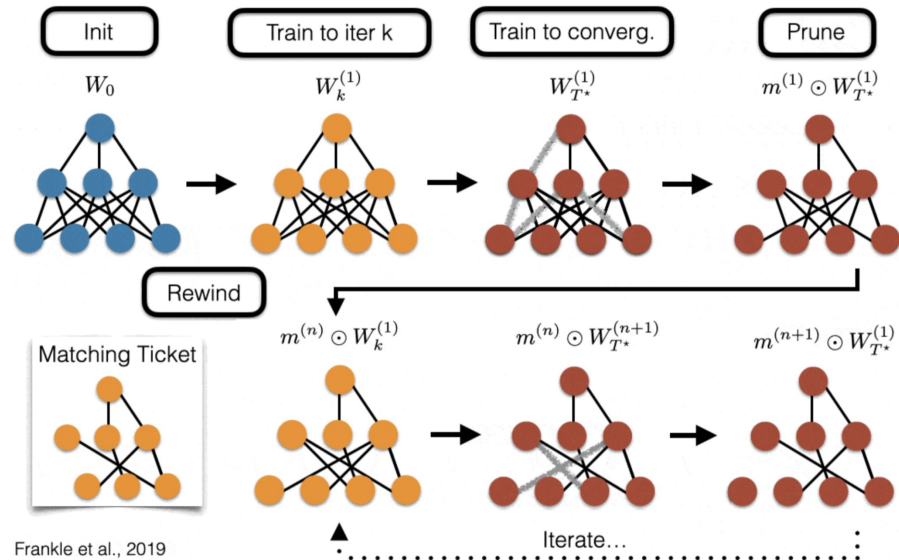
# Unstructured Pruning: Lottery Ticket Hypothesis



Searching for Tickets: Iterative Magnitude Pruning

Iterative Magnitude Pruning with Rewinding

Frankle & Carbin, 2019
Viz: @RobertTLange

Frankle et al., 2019
Viz: @RobertTLange

# Unstructured Pruning: Lottery Ticket Hypothesis



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

Resnet-18

- rewind to 0
- rewind to 500
- random reinit
- random reinit

**B** Rewinding Resnet-50 on ImageNet

Resnet-50 on ImageNet (Iterative)

- rewind to 6 (oneshot)
- rewind to 0 (iterative)
- rewind to 6 (iterative)
- random reinit

26

# Unstructured Pruning

Pros:

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

Cons:

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

# Structured Pruning

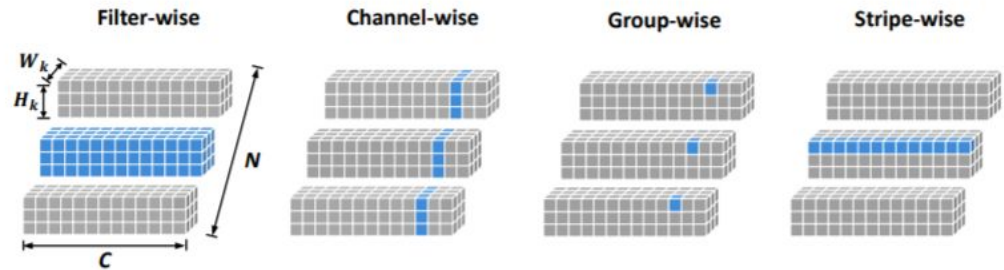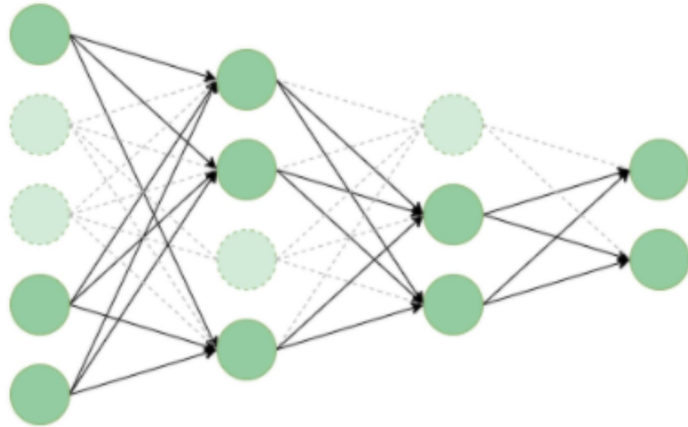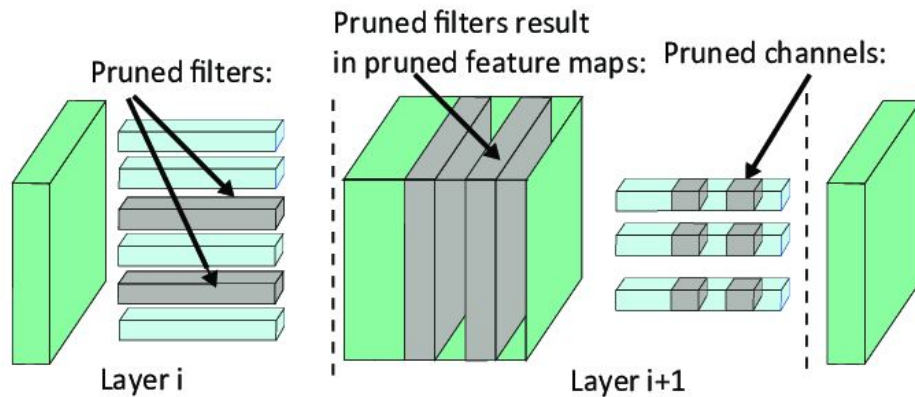Removing structural parts instead of individual weights



Filter-wise    Channel-wise    Group-wise    Stripe-wise

Figure 2: The visualization of different types of pruning.

Source:   https://heartbeat.fritz.ai/neural-network-pruning-research-review-2020-bc21a77f0295

# Structured Pruning

List of possible criterions:

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



Pruned filters:

Pruned filters result in pruned feature maps:

Pruned channels:

Layer i

Layer i+1

Source: https://openaccess.thecvf.com/content_CVPR_2020/papers/He_Learning_Filter_Pruning_Criteria_for_Deep_Convolutional_Neural_Networks_Acceleration_CVPR_2020_paper.pdf

# Structured Pruning
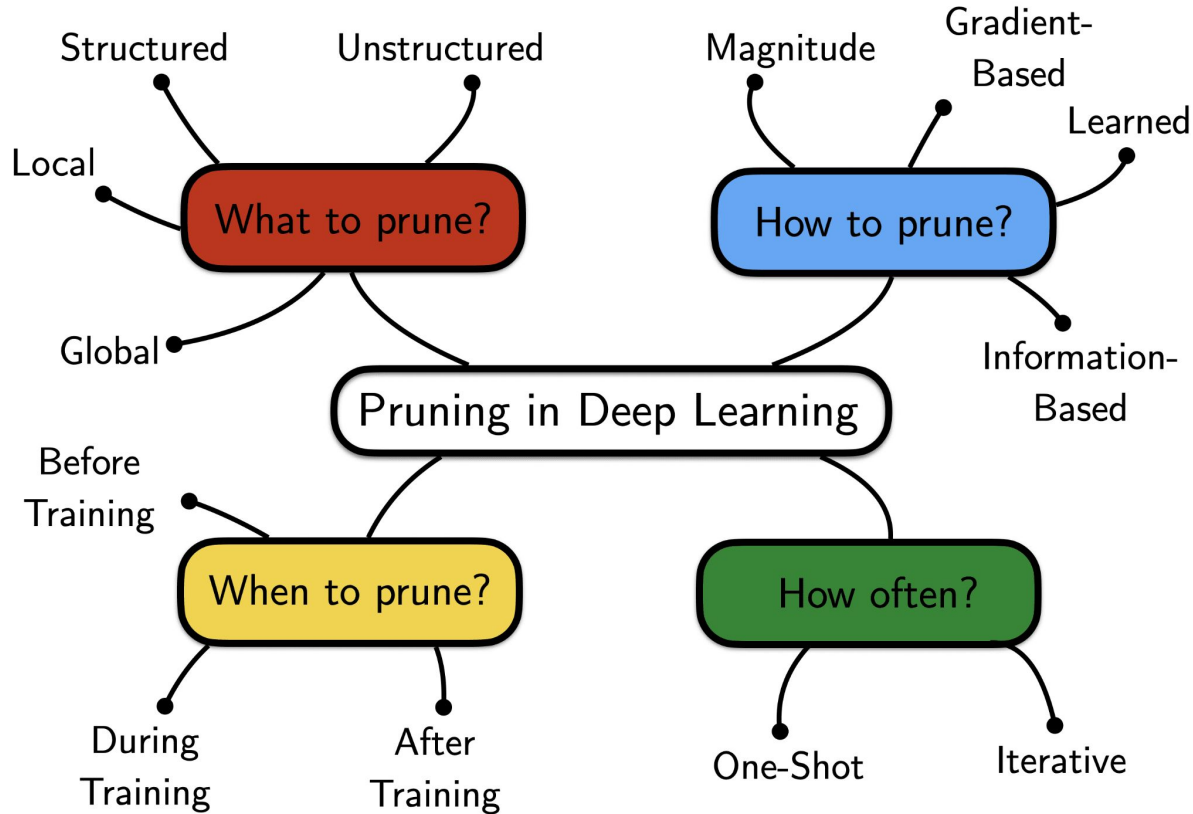
Pros:

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

Cons:

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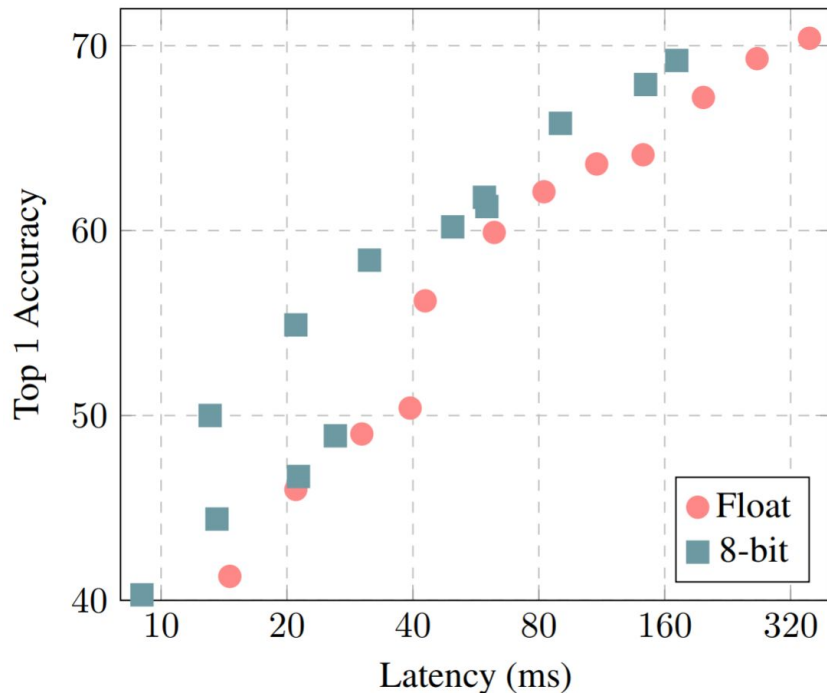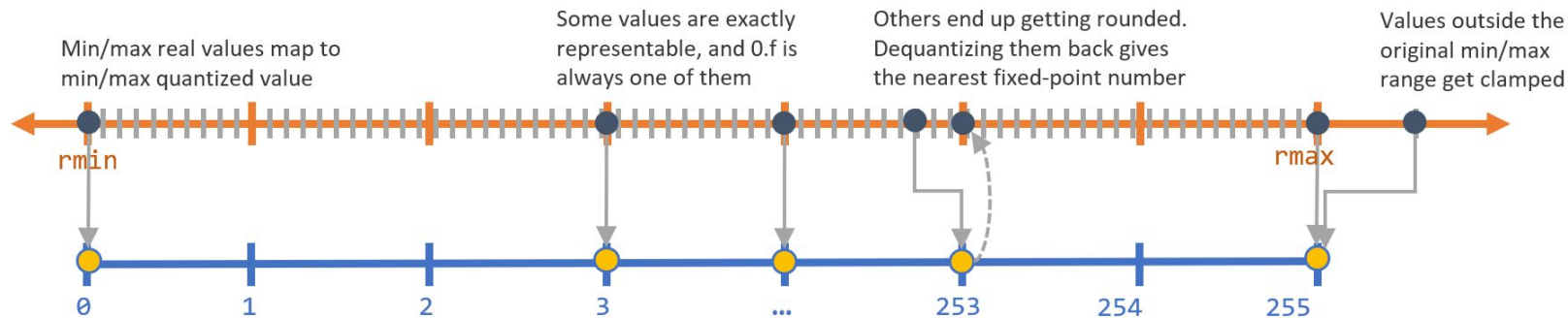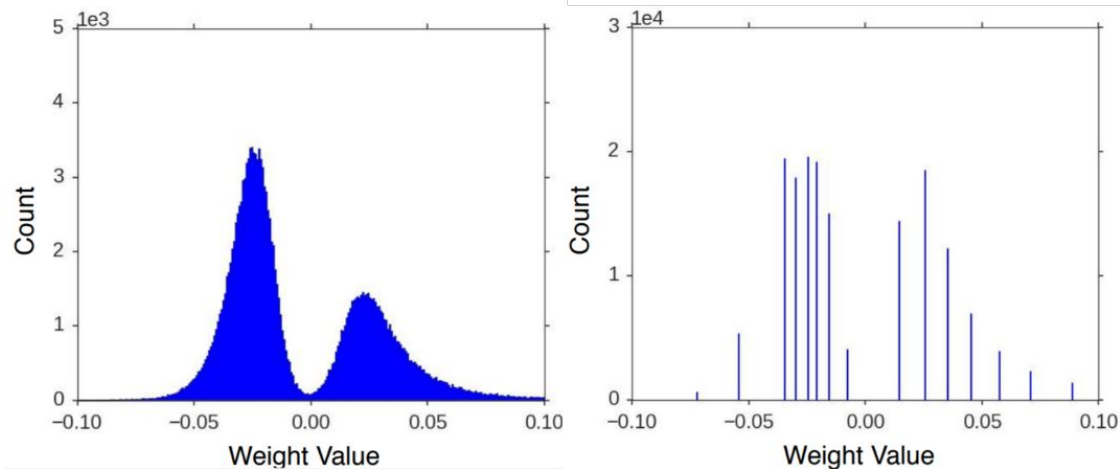
# Pruning

# Quantization

Reduce redundancy in weight numerical representation

- DNNs are known to be quite robust to noise and other small perturbations
- Weights and activations by a particular layer often tend to lie in a small range
- Arithmetic with lower bit-depth is faster
- Reduces memory consumption (e.g. in moving from 32-bits to 8-bits, we get (almost) 4x reduction in memory)

# Quantization



Min/max real values map to min/max quantized value

Some values are exactly representable, and 0.f is always one of them

Others end up getting rounded. Dequantizing them back gives the nearest fixed-point number

Values outside the original min/max range get clamped

rmin
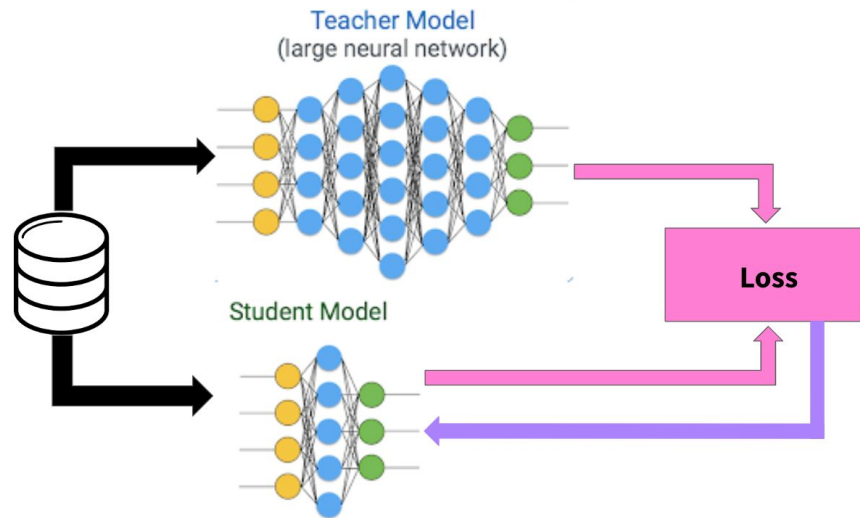
rmax

0   1   2   3   ...   253   254   255
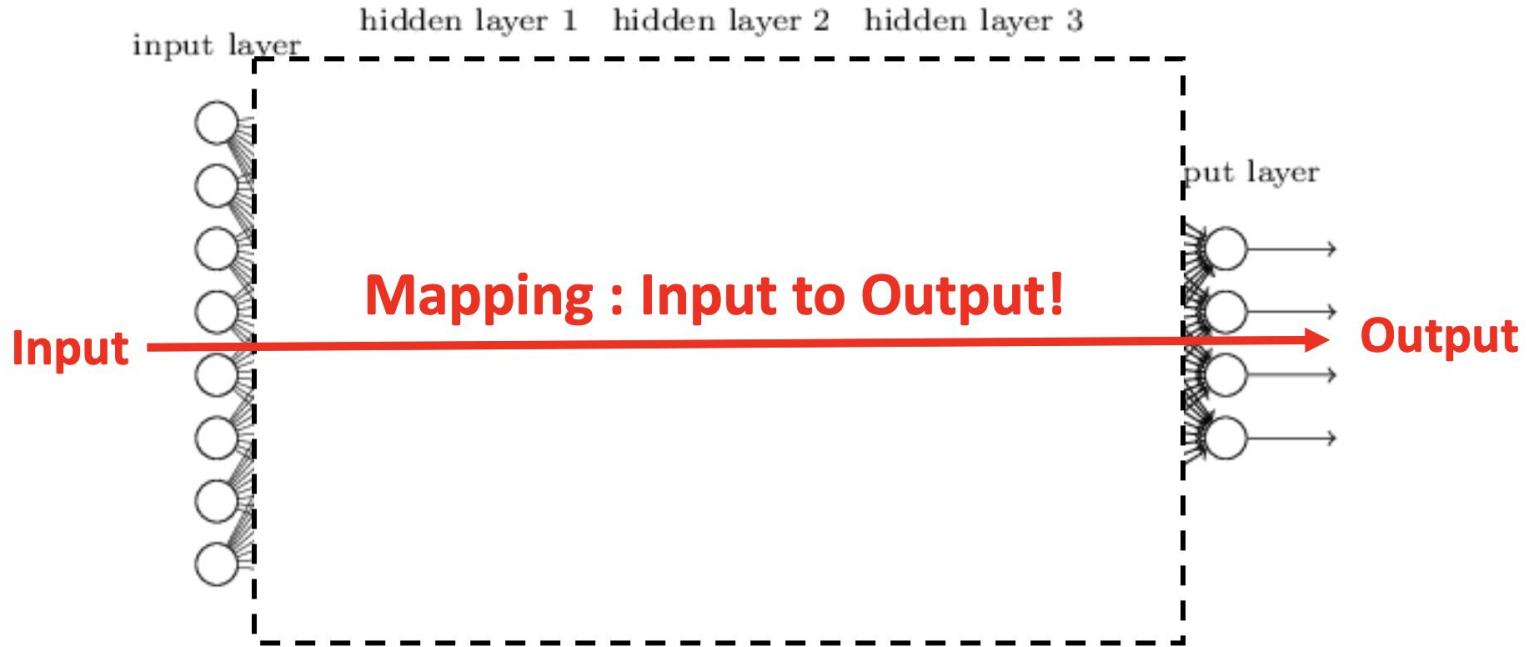
# Quantization

Pros:

- Easy combine with other methods (low-rank, pruning, etc.)
- Easy to apply (supported in modern NN frameworks: pytorch, tensorflow)
- Efficient model acceleration and compression
- Few quantization options => easy to find the best

# Knowledge Distillation

**Knowledge distillation** is a process of distilling or transferring the knowledge from a (set of) large, cumbersome model(s) to a lighter, easier-to-deploy single model
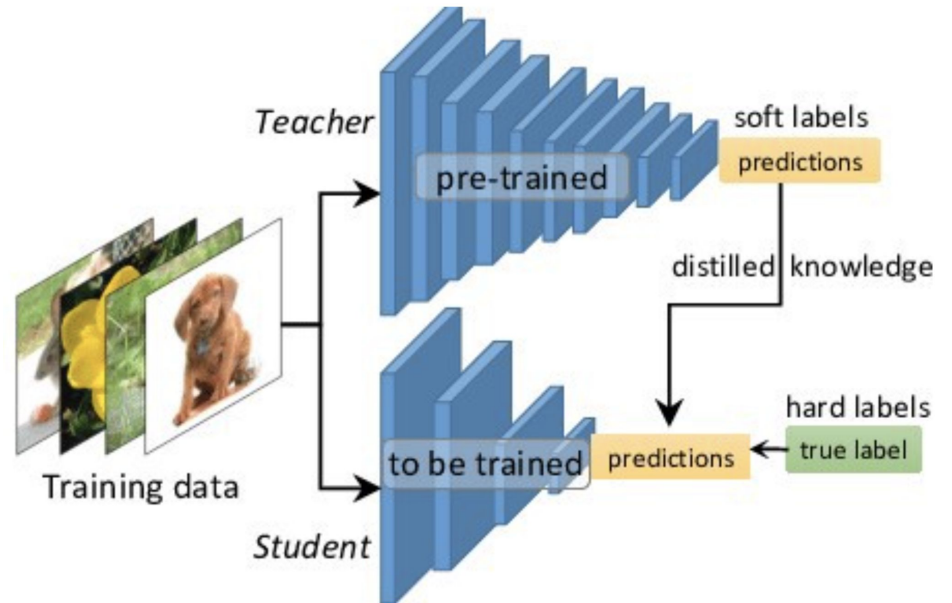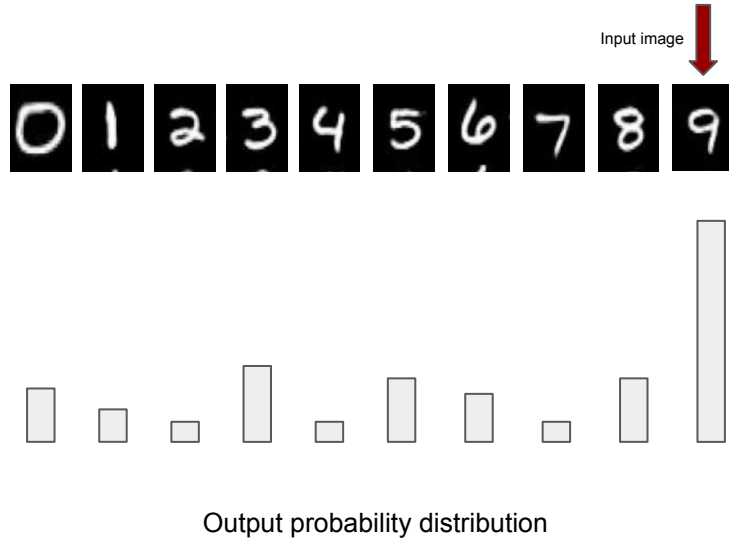
35

# Knowledge Distillation



input layer     hidden layer 1    hidden layer 2    hidden layer 3

put layer

**Mapping : Input to Output!**

**Input** → **Output**

A more abstract view of the **knowledge**, that frees it from any particular instantiation, is that it is a **learned mapping from input vectors to output vectors**.

# Knowledge Distillation

Idea: Train less redundant model by using knowledge from big models

Input image



Output probability distribution

# Knowledge Distillation

- Response-Based Knowledge



$$L_{ResD}(p(z_t, T), p(z_s, T)) = \mathcal{L}_R(p(z_t, T), p(z_s, T))$$

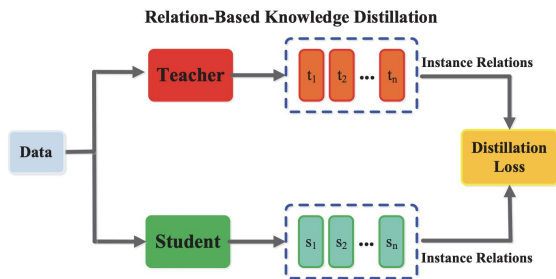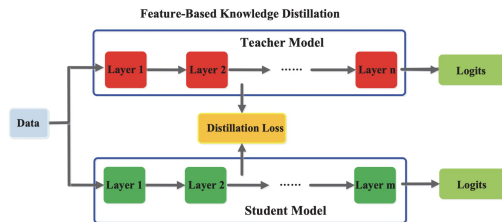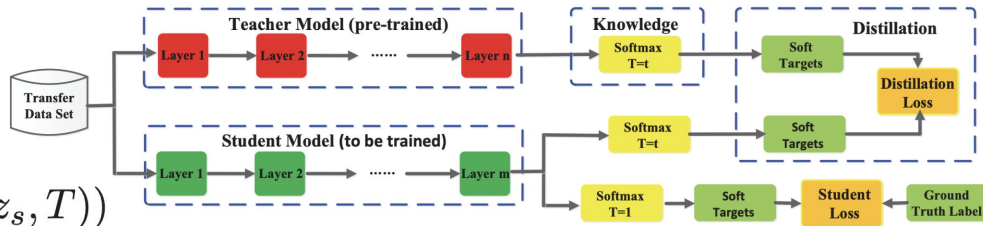$$p(z_i, T) = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

- Feature-Based Knowledge



$$L_{FeaD}(f_t(x), f_s(x)) = \mathcal{L}_F(\Phi_t(f_t(x)), \Phi_s(f_s(x)))$$
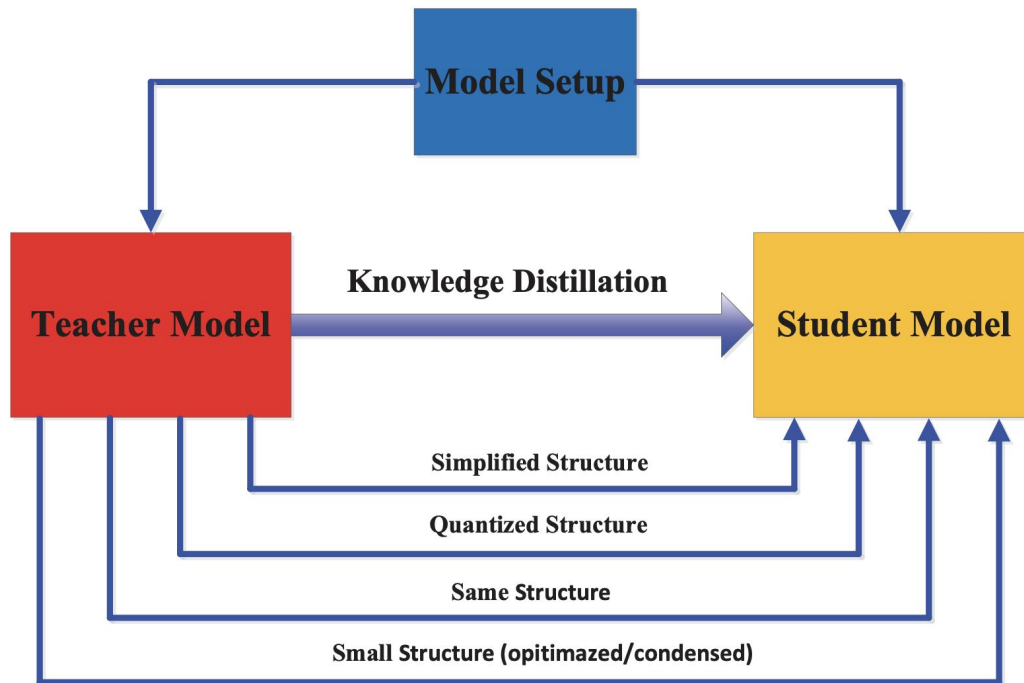
- Relation-Based Knowledge

$$L_{RelD}(f_t, f_s) = \mathcal{L}_{R^1}(\Psi_t(\hat{f}_t, \check{f}_t), \Psi_s(\hat{f}_s, \check{f}_s))$$

$$L_{RelD}(F_t, F_s) = \mathcal{L}_{R^2}(\psi_t(t_i, t_j), \psi_s(s_i, s_j))$$

# Knowledge Distillation

# Knowledge Distillation

Pros:

- Improves model performance

Cons:

- Capacity gap
- Controversy of the method

# Summary

Efficient neural architecture design

- MobileNet
- ShuffleNet
- GhostNet

Neural network compression

- Weight Decomposition
- Unstructured Pruning
- Structured Pruning
- Quantization
- **Knowledge Distillation**

# Thanks for your attention!