



# NETWORK COMPRESSION

Hung-yi Lee 李宏毅

# Smaller Model

Less parameters



Deploying ML models in resource-constrained environments



Lower latency, Privacy, etc.



# Outline

- Network Pruning
- Knowledge Distillation
- Parameter Quantization
- Architecture Design
- Dynamic Computation

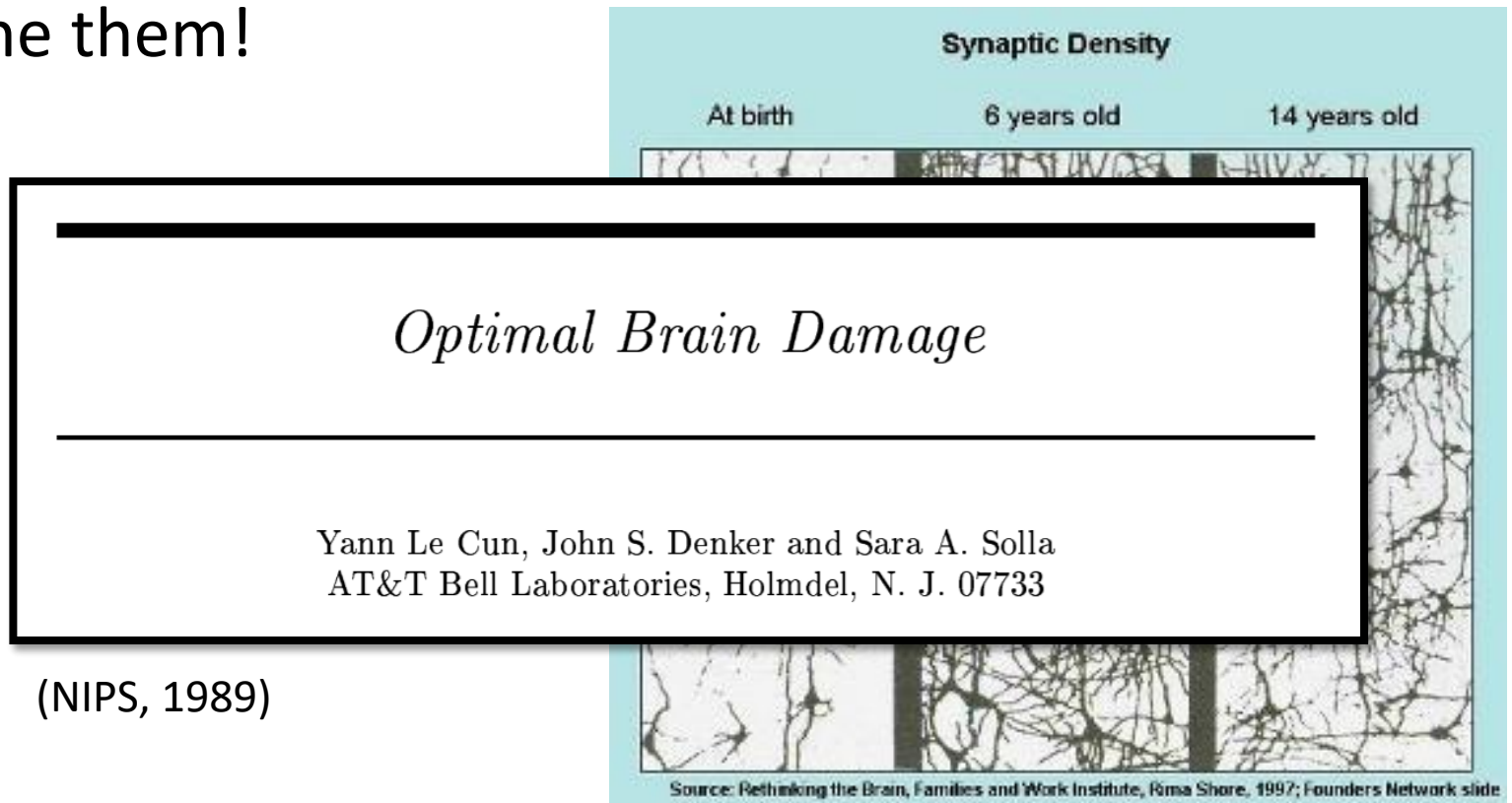
We will not talk about hard-ware solution today.

# Network Pruning



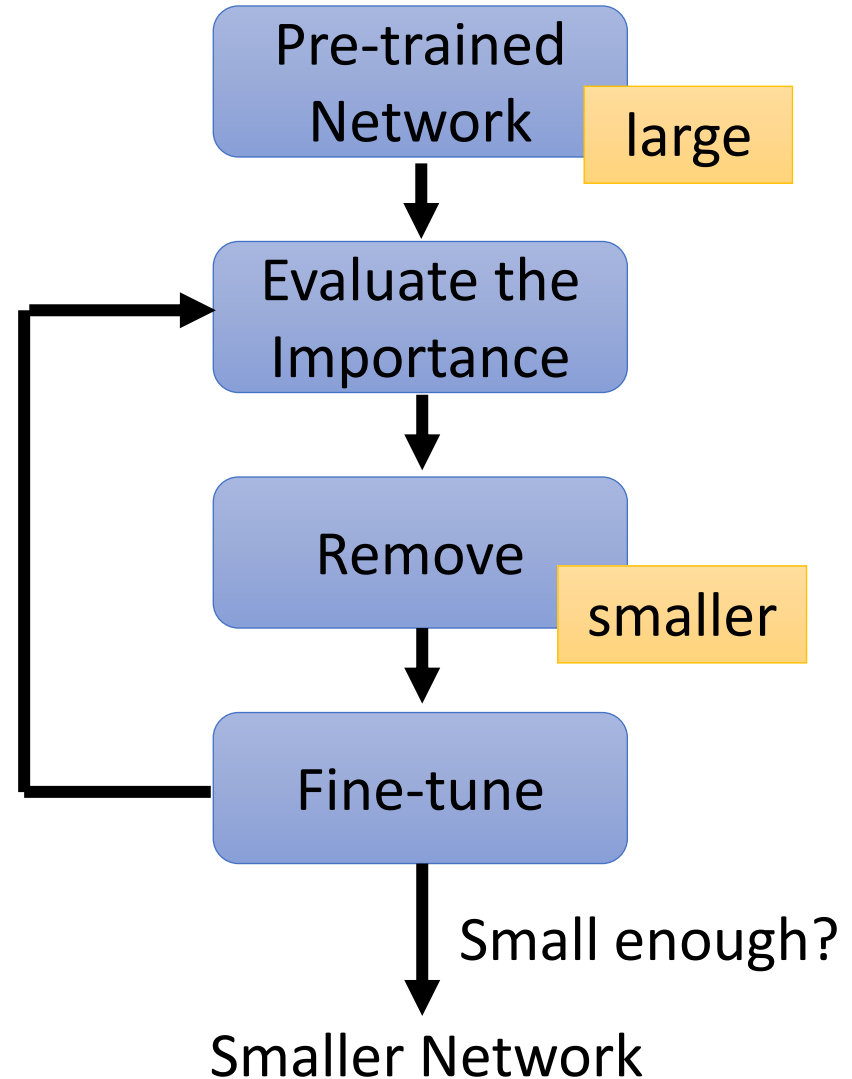
# Network can be pruned

- Networks are typically over-parameterized (there is significant redundant weights or neurons)
- Prune them!



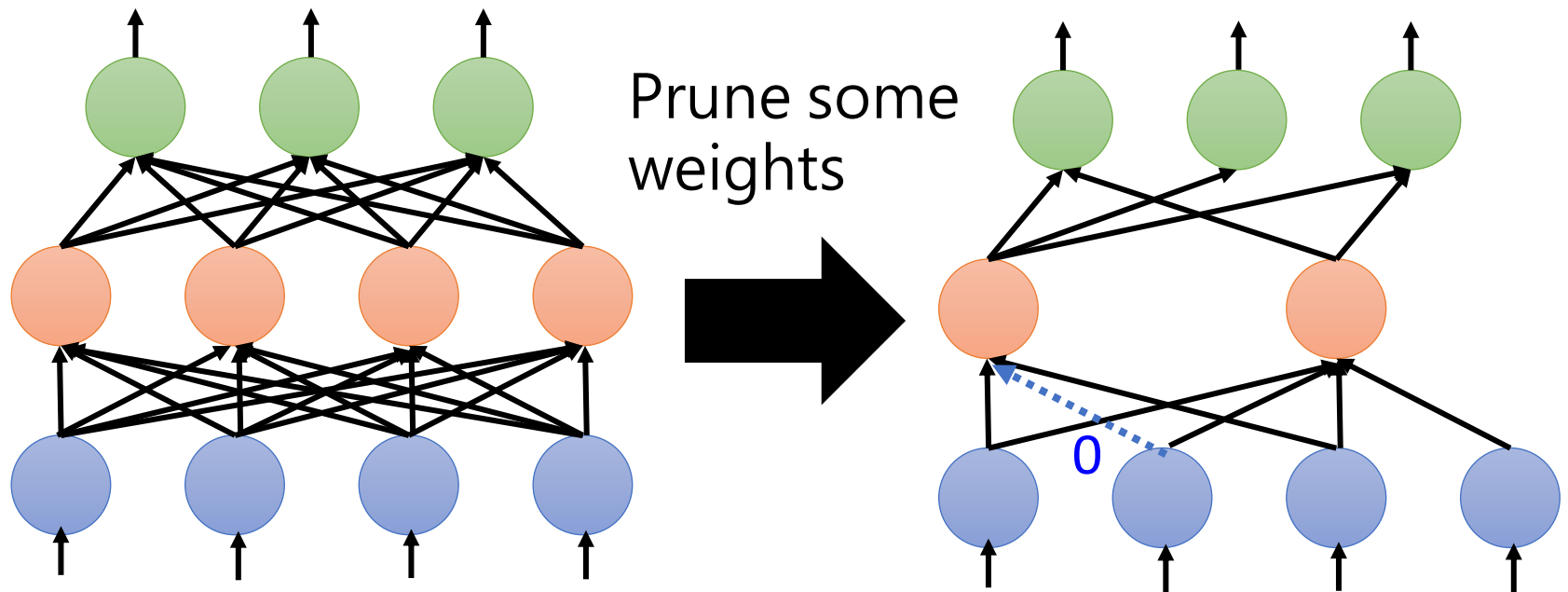
# Network Pruning

- Importance of a weight:  
absolute values, life long ...
- Importance of a neuron:  
the number of times it wasn't zero on a given data set .....
- After pruning, the accuracy will drop (hopefully not too much)
- Fine-tuning on training data for recover
- Don't prune too much at once, or the network won't recover.



# Network Pruning - Practical Issue

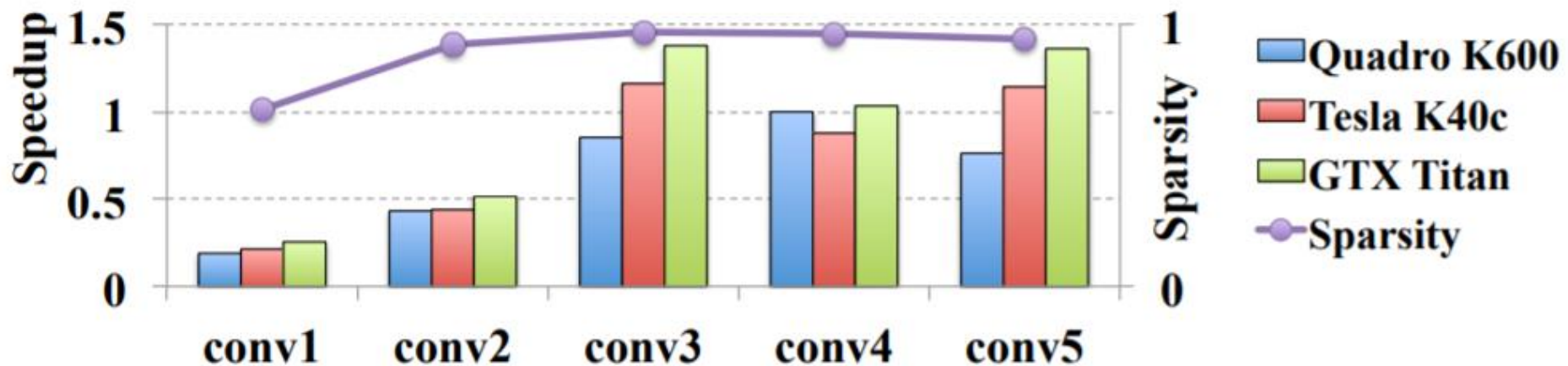
- Weight pruning



Hard to implement, hard to speedup .....

# Network Pruning - Practical Issue

- Weight pruning



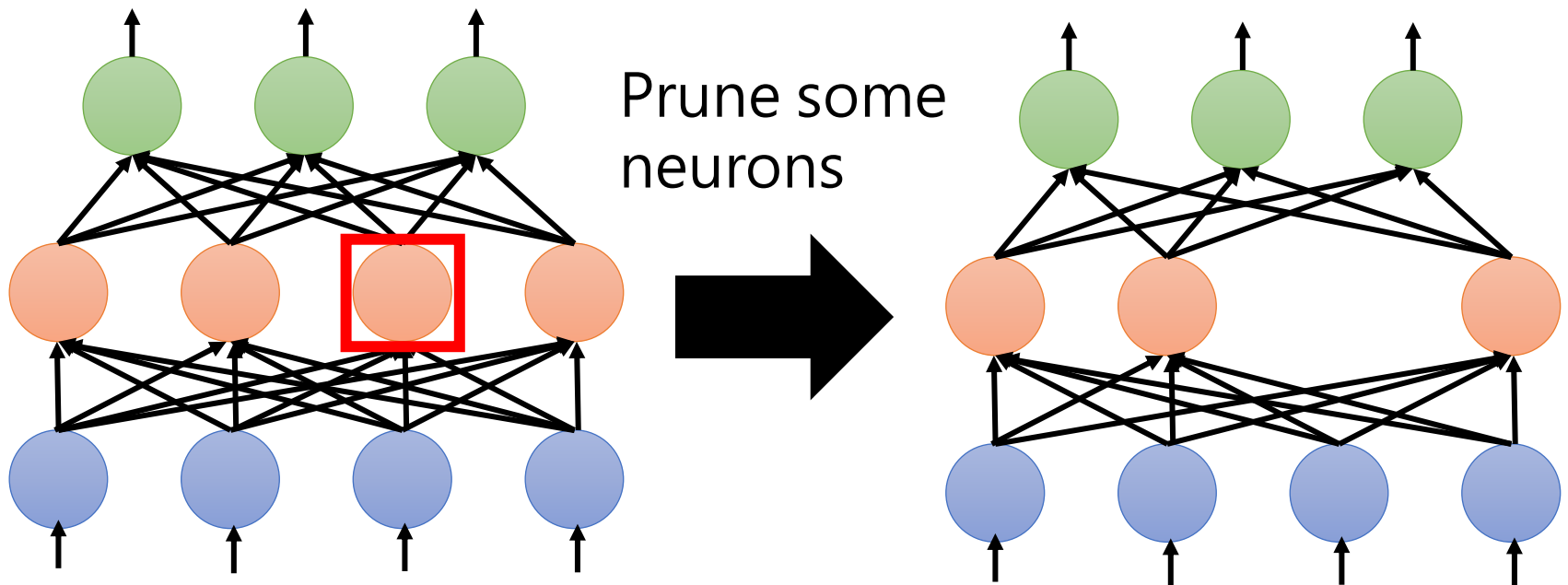
<https://arxiv.org/pdf/1608.03665.pdf>



# Network Pruning - Practical Issue

- Neuron pruning

The network architecture is regular.



Easy to implement, easy to speedup .....

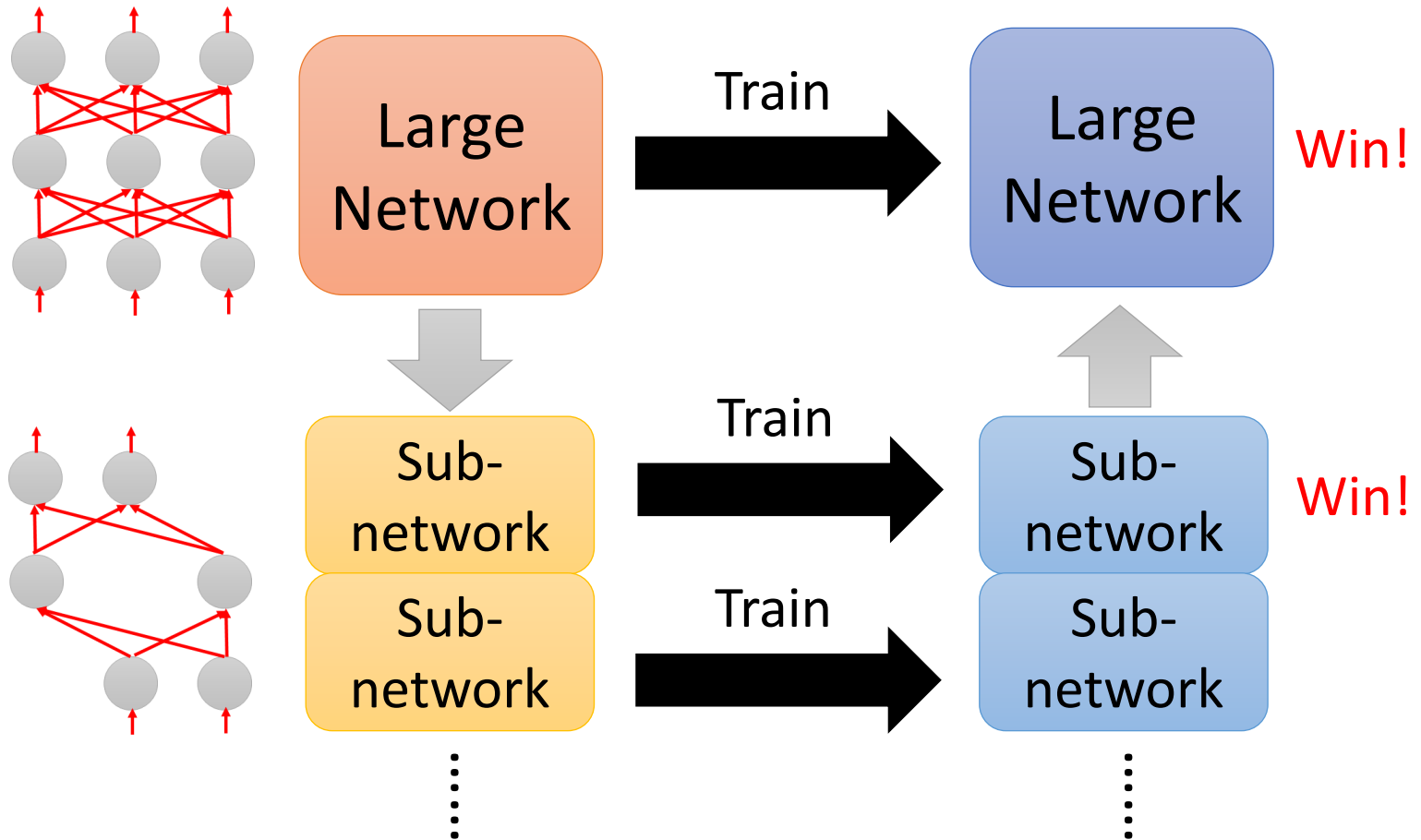
# Why Pruning?

- How about simply train a smaller network?
- It is widely known that smaller network is more difficult to learn successfully.
  - Larger network is easier to optimize?  
[https://www.youtube.com/watch?v=\\_VuWvQUMQVk](https://www.youtube.com/watch?v=_VuWvQUMQVk)
- Lottery Ticket Hypothesis  
<https://arxiv.org/abs/1803.03635>



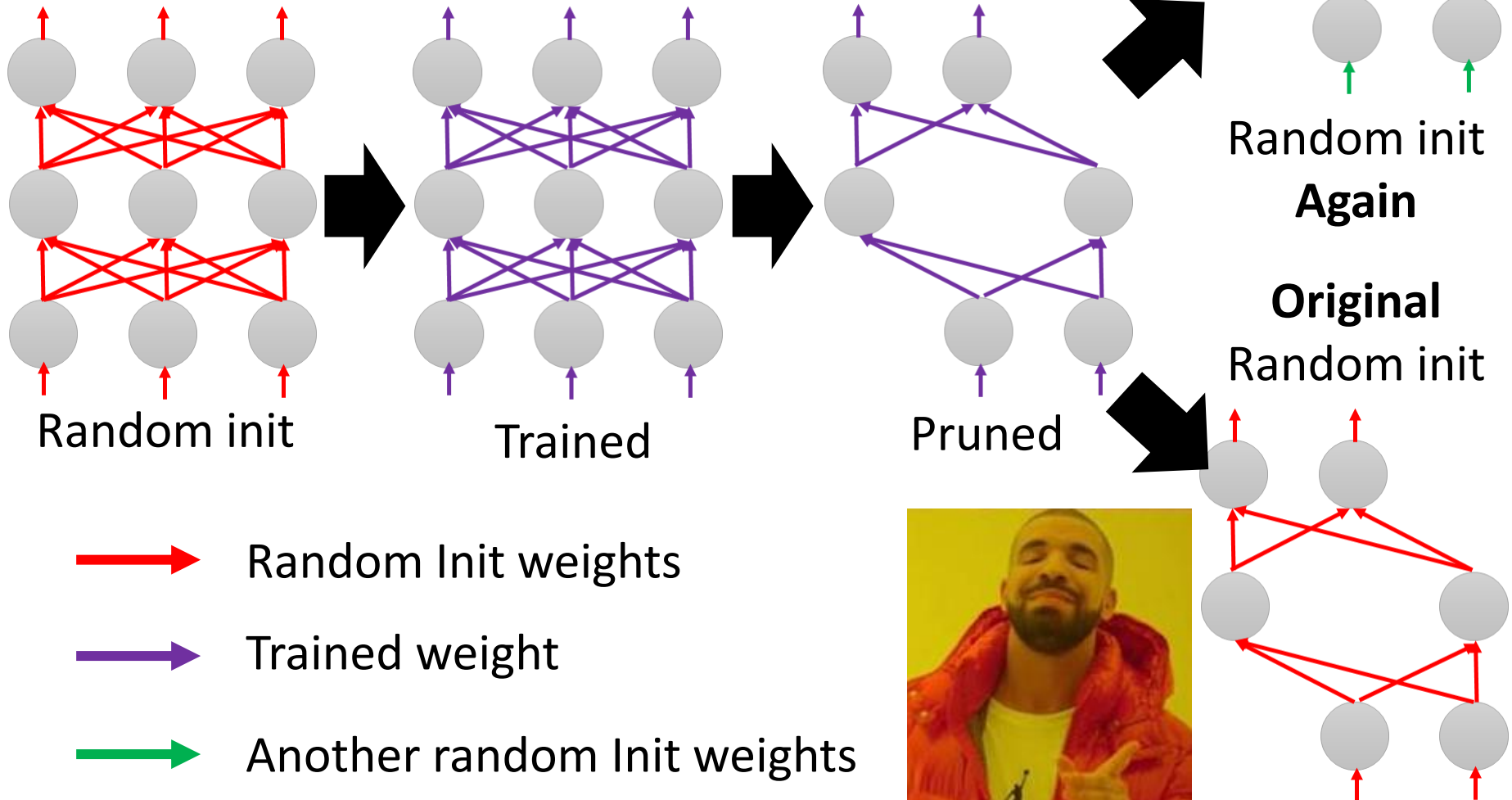
# Why Pruning?

## Lottery Ticket Hypothesis



# Why Pruning?

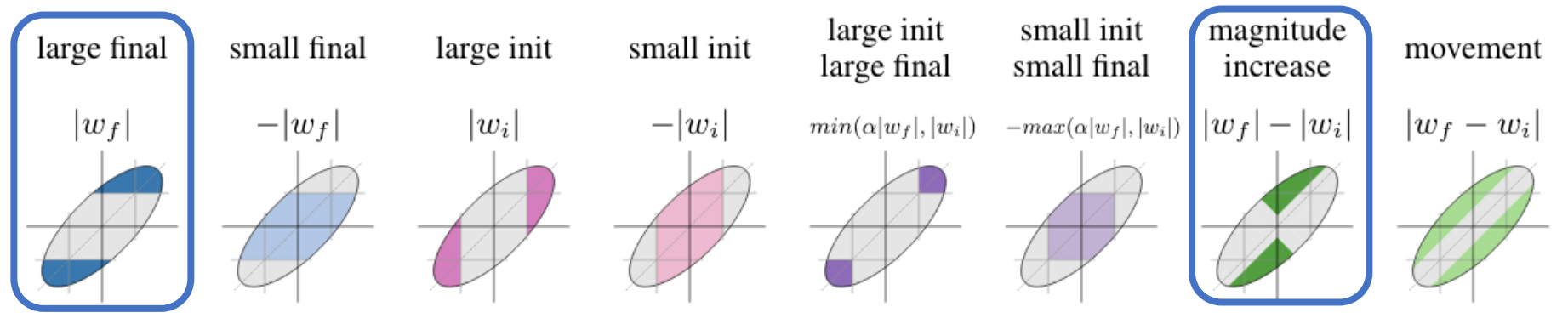
## Lottery Ticket Hypothesis



# Why Pruning?

## Lottery Ticket Hypothesis

- Different pruning strategy



- “sign-ificance” of initial weights: Keeping the sign is critical

0.9, 3.1, -9.1, 8.5 .....  $\longrightarrow$   $+\alpha, +\alpha, -\alpha, +\alpha$  .....

- Pruning weights from a network with random weights

# Why Pruning?

<https://arxiv.org/abs/1810.05270>

- Rethinking the Value of Network Pruning

Dataset	Model	Unpruned	Pruned Model	Fine-tuned	Scratch-E	Scratch-B
CIFAR-10	VGG-16	93.63 ( $\pm 0.16$ )	VGG-16-A	93.41 ( $\pm 0.12$ )	93.62 ( $\pm 0.11$ )	<b>93.78</b> ( $\pm 0.15$ )
	ResNet-56	93.14 ( $\pm 0.12$ )	ResNet-56-A	92.97 ( $\pm 0.17$ )	92.96 ( $\pm 0.26$ )	<b>93.09</b> ( $\pm 0.14$ )
			ResNet-56-B	92.67 ( $\pm 0.14$ )	92.54 ( $\pm 0.19$ )	<b>93.05</b> ( $\pm 0.18$ )
	ResNet-110	93.14 ( $\pm 0.24$ )	ResNet-110-A	93.14 ( $\pm 0.16$ )	<b>93.25</b> ( $\pm 0.29$ )	93.22 ( $\pm 0.22$ )
			ResNet-110-B	92.69 ( $\pm 0.09$ )	92.89 ( $\pm 0.43$ )	<b>93.60</b> ( $\pm 0.25$ )
ImageNet	ResNet-34	73.31	ResNet-34-A	72.56	72.77	<b>73.03</b>
			ResNet-34-B	72.29	72.55	<b>72.91</b>

- **New** random initialization, not **original** random initialization in “Lottery Ticket Hypothesis”
- Limitation of “Lottery Ticket Hypothesis” (small  $l_r$ , unstructured)

# Knowledge Distillation



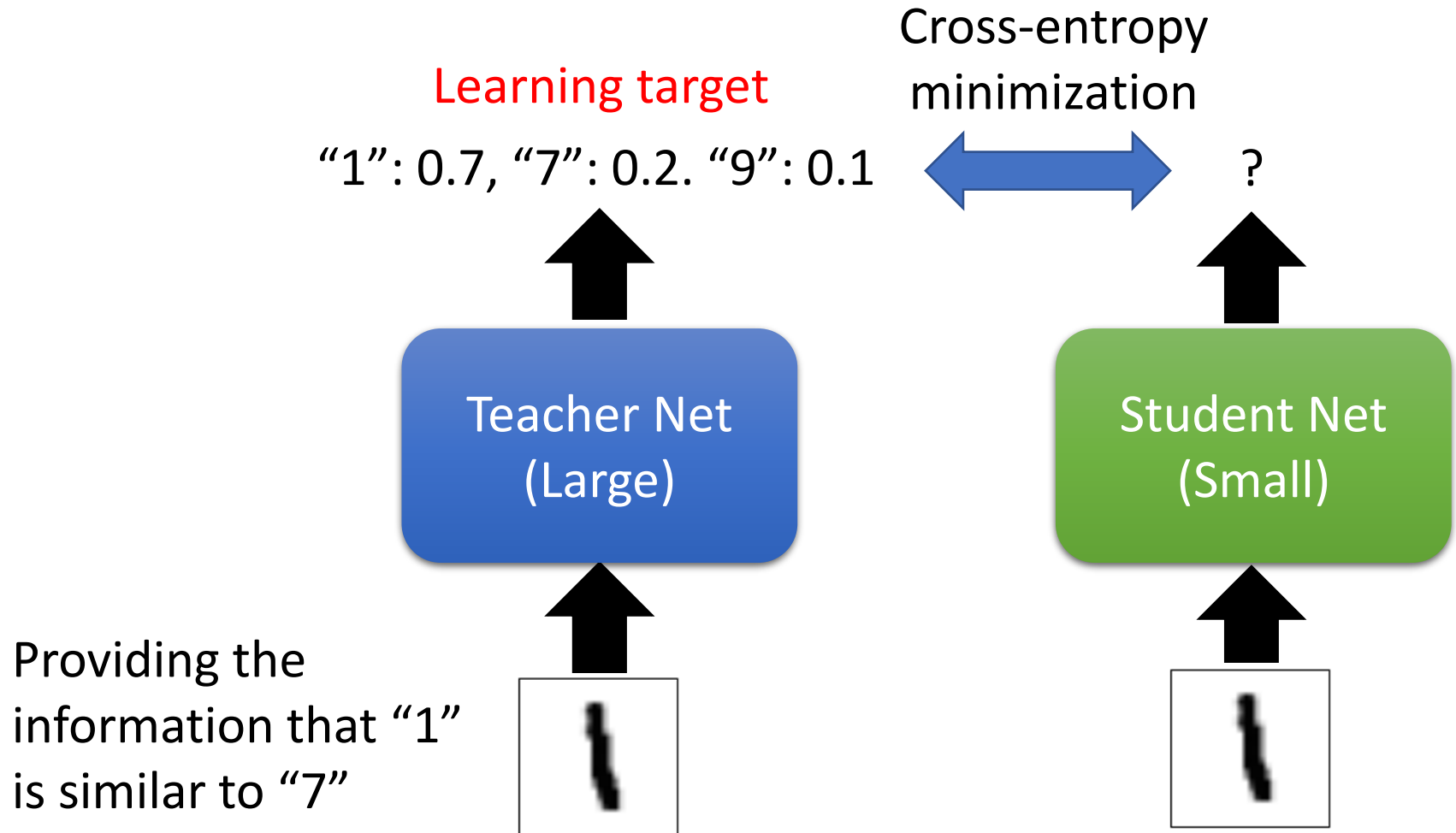
# Knowledge Distillation

Knowledge Distillation

<https://arxiv.org/pdf/1503.02531.pdf>

Do Deep Nets Really Need to be Deep?

<https://arxiv.org/pdf/1312.6184.pdf>





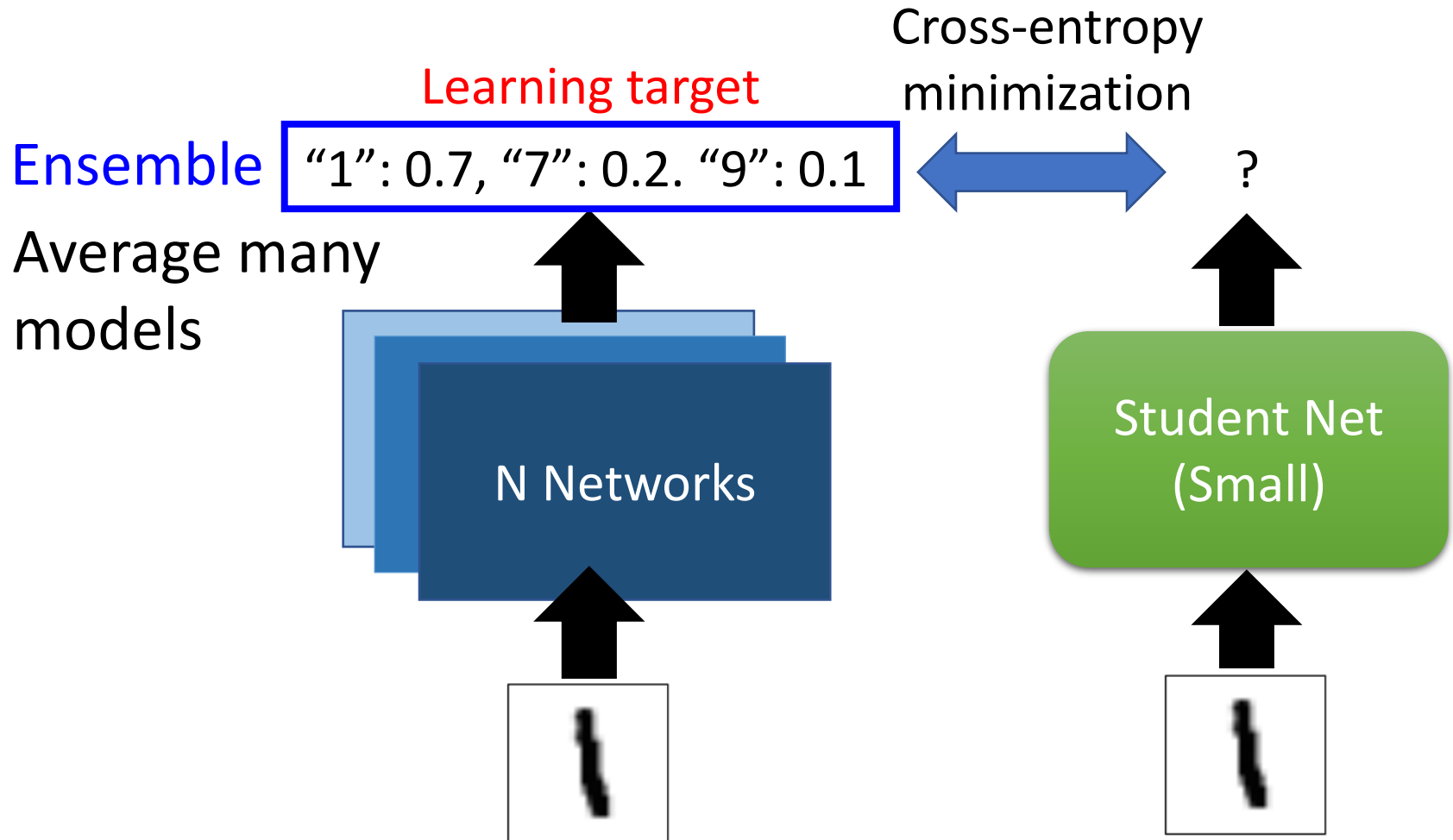
# Knowledge Distillation

Knowledge Distillation

<https://arxiv.org/pdf/1503.02531.pdf>

Do Deep Nets Really Need to be Deep?

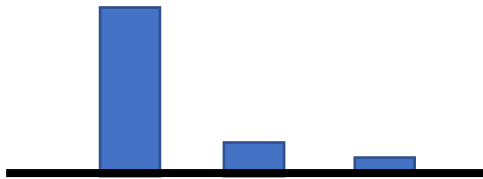
<https://arxiv.org/pdf/1312.6184.pdf>



# Knowledge Distillation

- Temperature for softmax

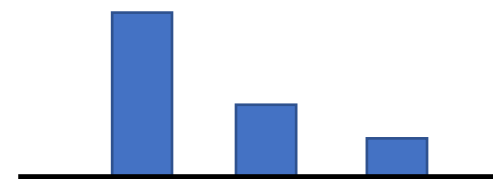
$$y'_i = \frac{\exp(y_i)}{\sum_j \exp(y_j)} \xrightarrow{T=100} y'_i = \frac{\exp(y_i/T)}{\sum_j \exp(y_j/T)}$$



$$y_1 = 100 \quad y'_1 = 1$$

$$y_2 = 10 \quad y'_2 \approx 0$$

$$y_3 = 1 \quad y'_3 \approx 0$$



$$y_1/T = 1 \quad y'_1 = 0.56$$

$$y_2/T = 0.1 \quad y'_2 = 0.23$$

$$y_3/T = 0.01 \quad y'_3 = 0.21$$

# Parameter Quantization



# Parameter Quantization

- 1. Using less bits to represent a value
- 2. Weight clustering

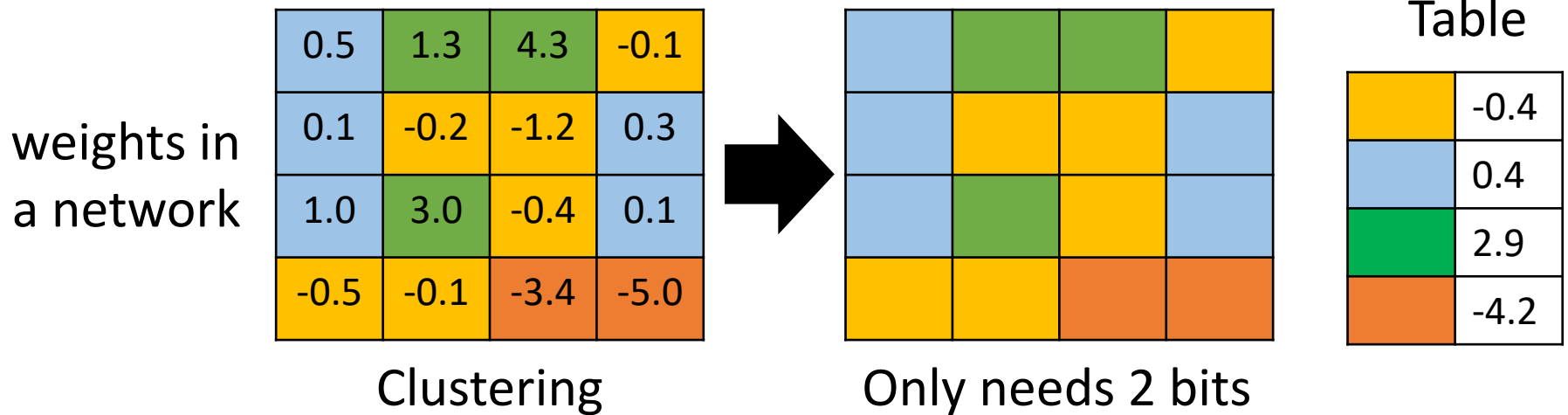
weights in  
a network

0.5	1.3	4.3	-0.1
0.1	-0.2	-1.2	0.3
1.0	3.0	-0.4	0.1
-0.5	-0.1	-3.4	-5.0

Clustering

# Parameter Quantization

- 1. Using less bits to represent a value
- 2. Weight clustering



- 3. Represent frequent clusters by less bits, represent rare clusters by more bits
  - e.g. Huffman encoding

# Binary Weights

Your weights are always +1 or -1

- Binary Connect

Binary Connect:

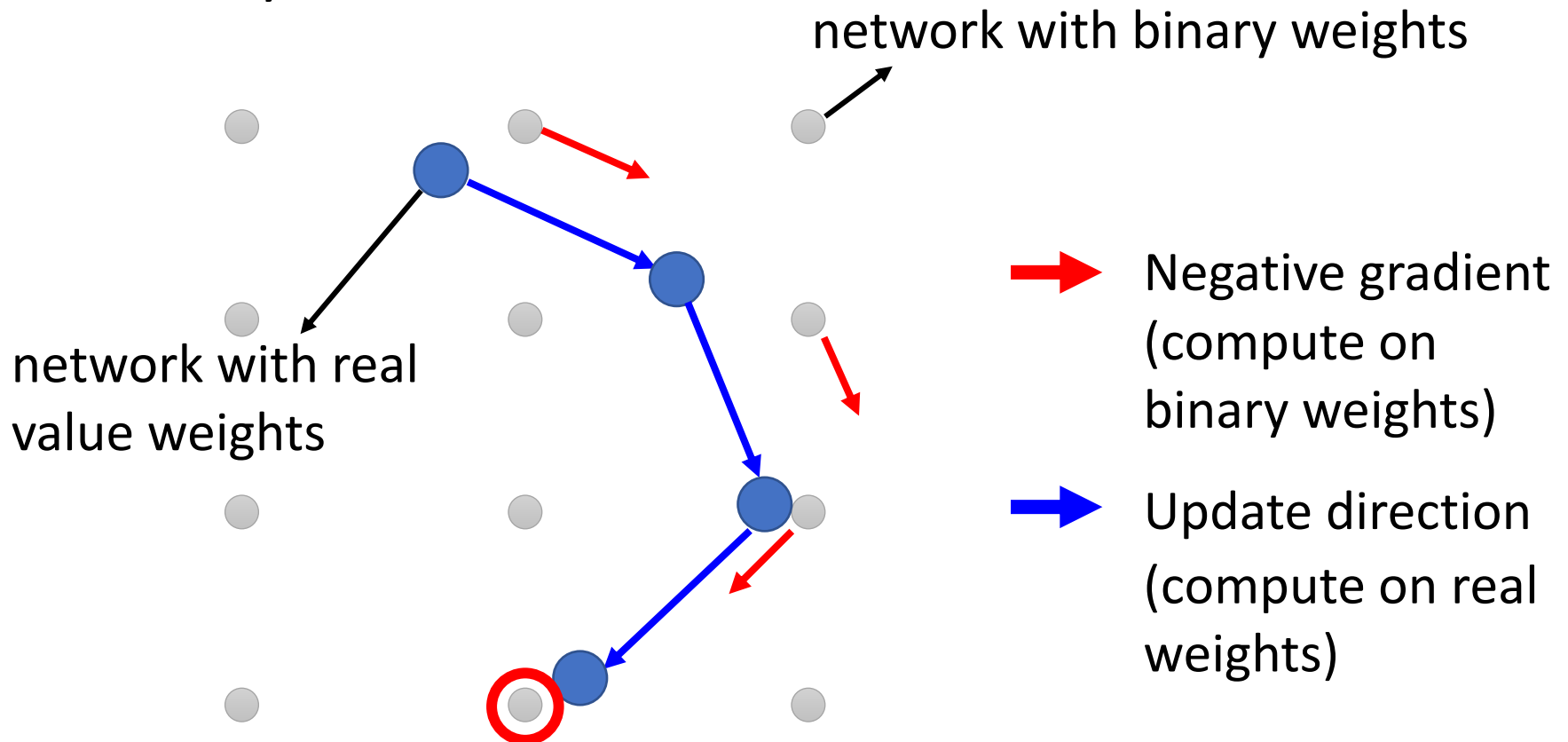
<https://arxiv.org/abs/1511.00363>

Binary Network:

<https://arxiv.org/abs/1602.02830>

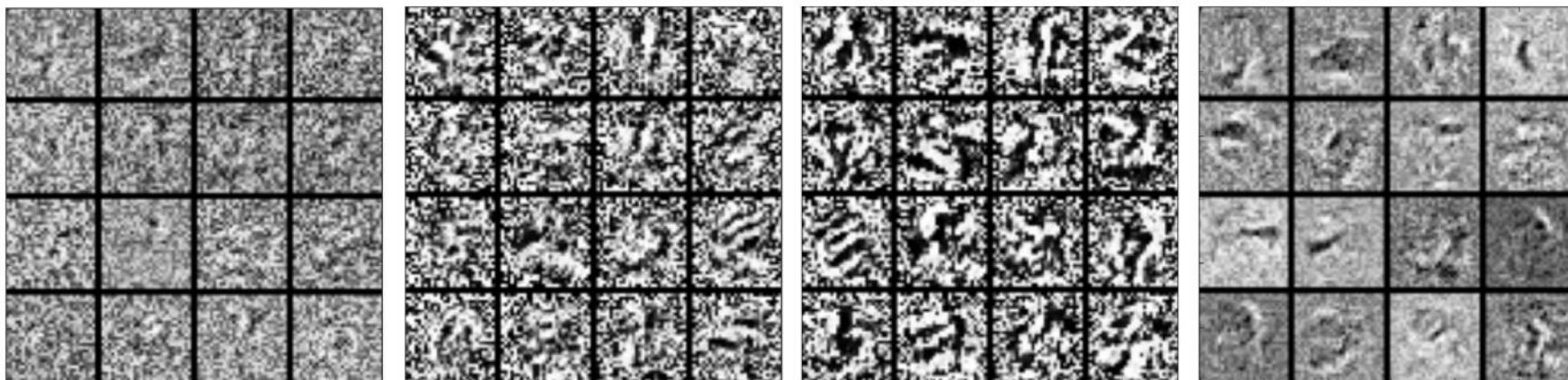
XNOR-net:

<https://arxiv.org/abs/1603.05279>



# Binary Connect

Method	MNIST	CIFAR-10	SVHN
No regularizer	$1.30 \pm 0.04\%$	10.64%	2.44%
BinaryConnect (det.)	$1.29 \pm 0.08\%$	9.90%	2.30%
BinaryConnect (stoch.)	$1.18 \pm 0.04\%$	<b>8.27%</b>	2.15%
50% Dropout	$1.01 \pm 0.04\%$		



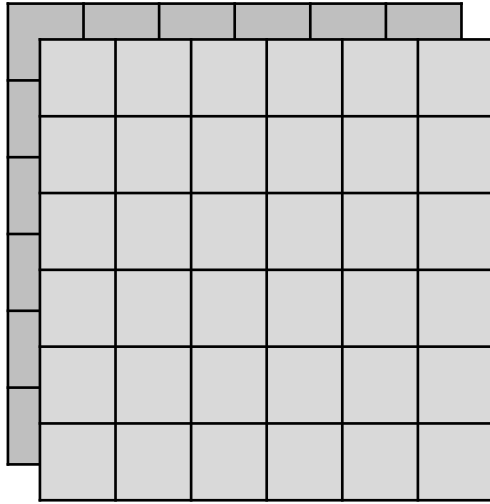
# Architecture Design

## *Depthwise Separable Convolution*

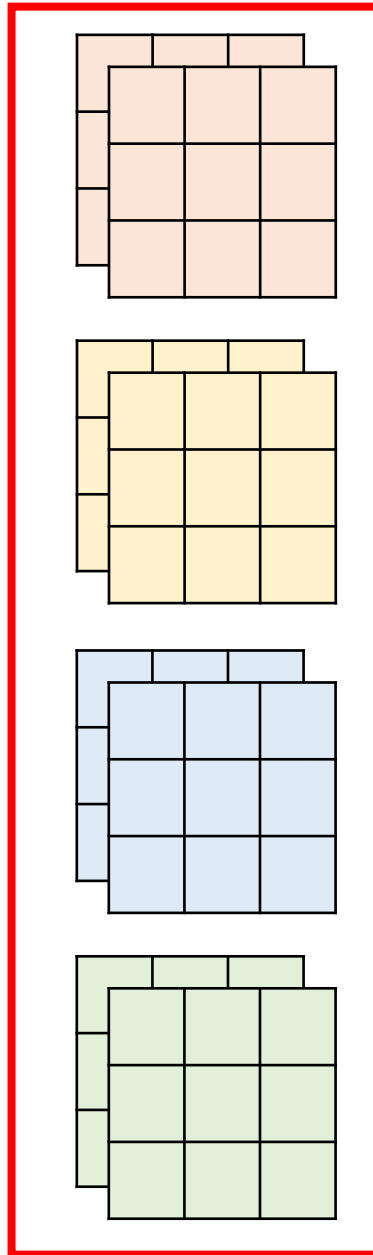


# Review: Standard CNN

Input feature map

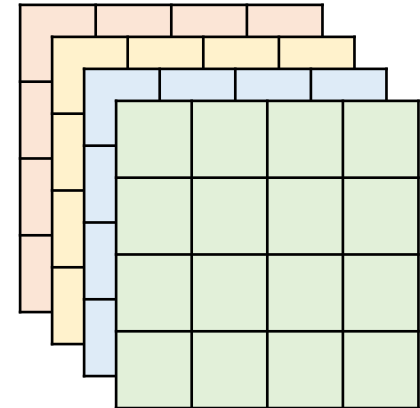


2 channels



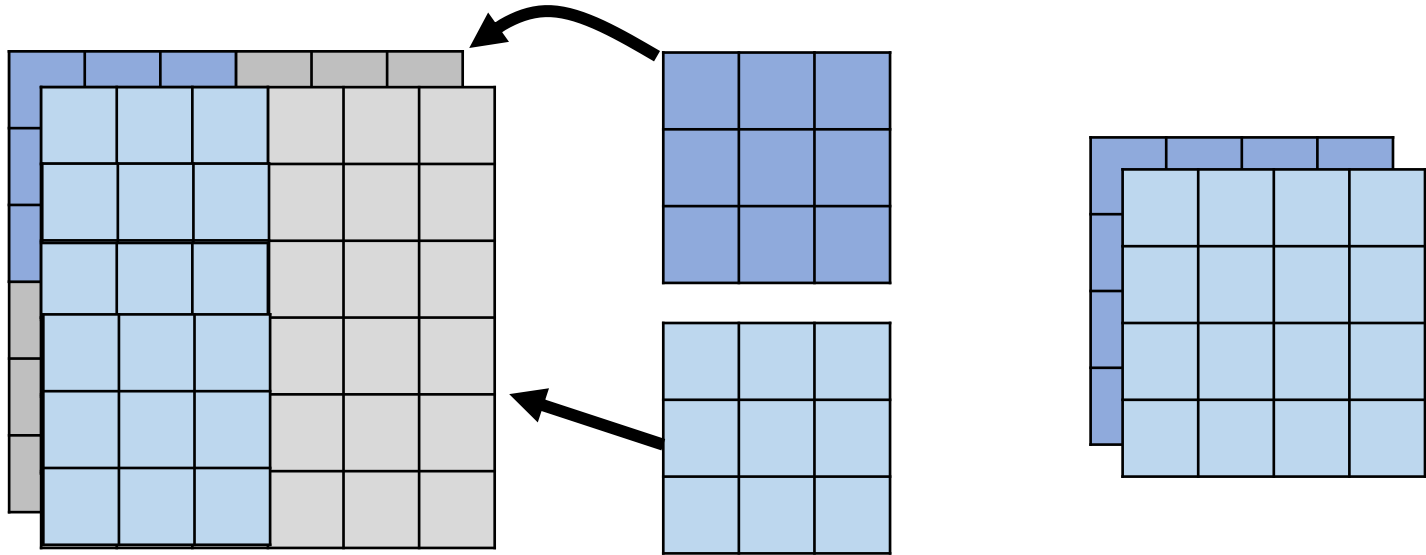
$$3 \times 3 \times 2 \times 4 = 72$$

parameters



# *Depthwise Separable Convolution*

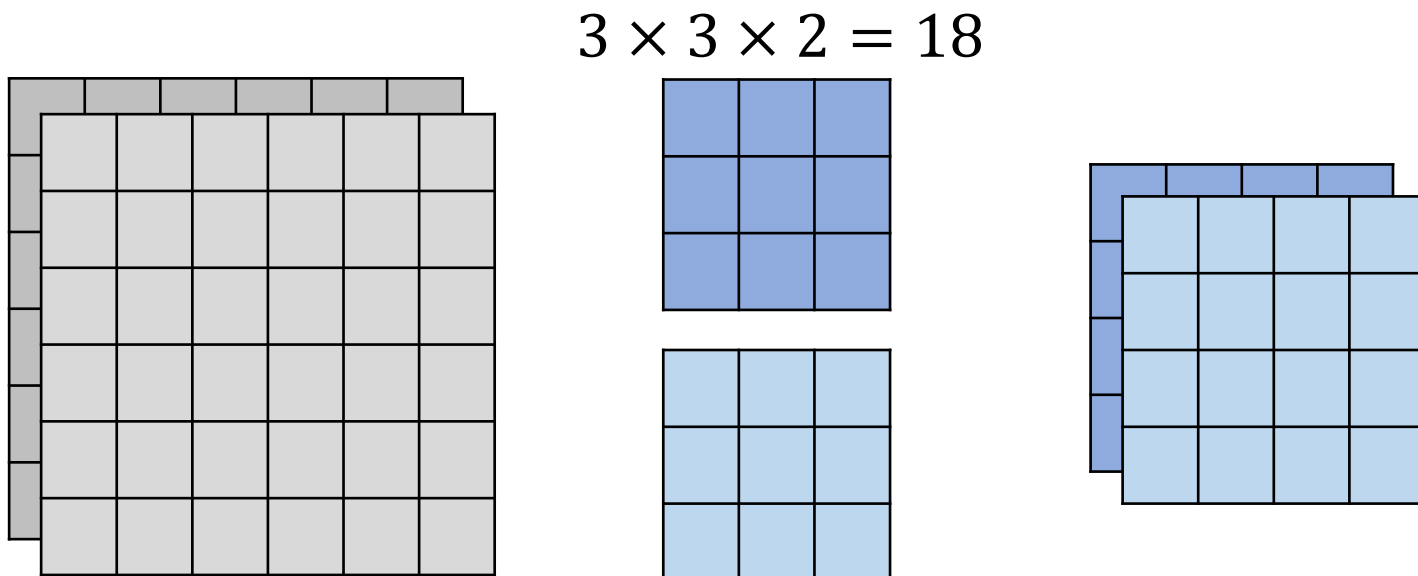
## 1. Depthwise Convolution



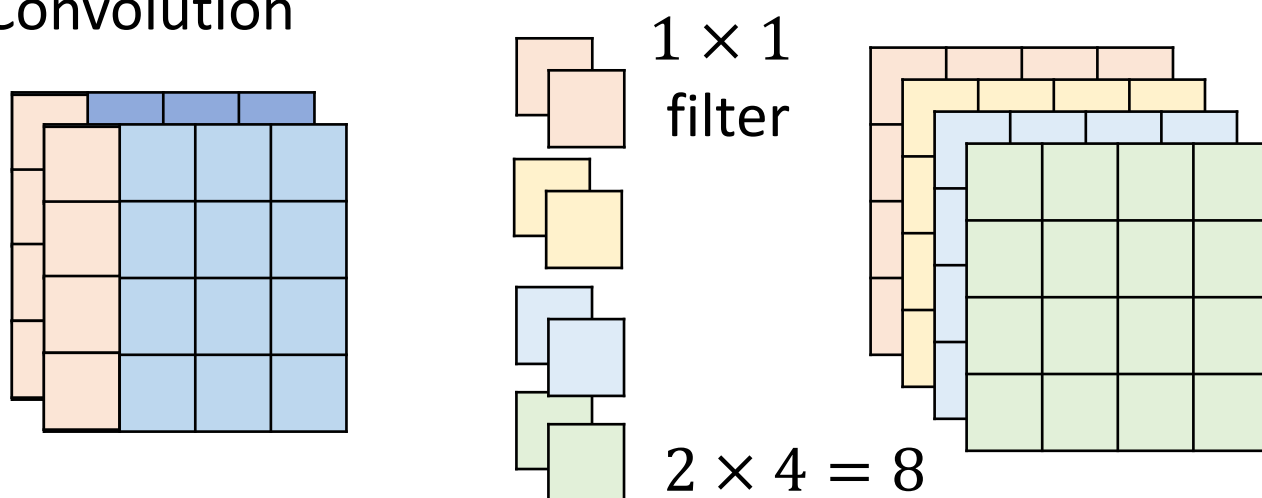
- Filter number = Input channel number
- Each filter only considers one channel.
- The filters are  $k \times k$  matrices
- There is no interaction between channels.

# *Depthwise Separable Convolution*

## 1. Depthwise Convolution



## 2. Pointwise Convolution



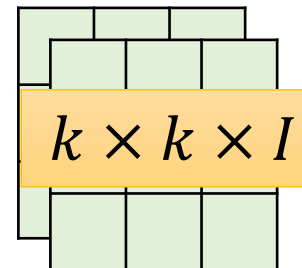
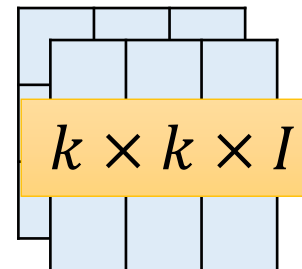
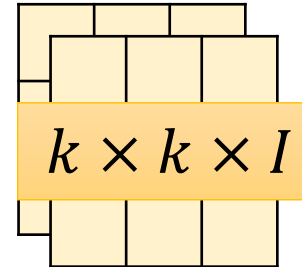
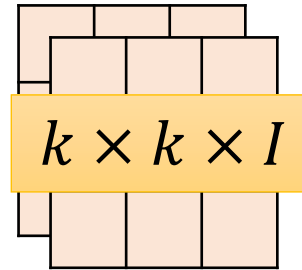
$I$ : number of input channels

$O$ : number of output channels

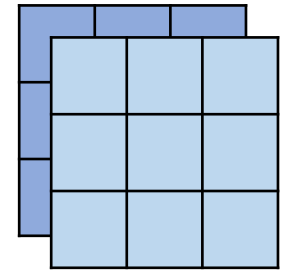
$k \times k$ : kernel size

$$\frac{k \times k \times I + I \times O}{k \times k \times I \times O}$$

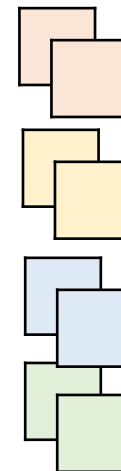
$$= \frac{1}{O} + \frac{1}{k \times k}$$



$$(k \times k \times I) \times O$$



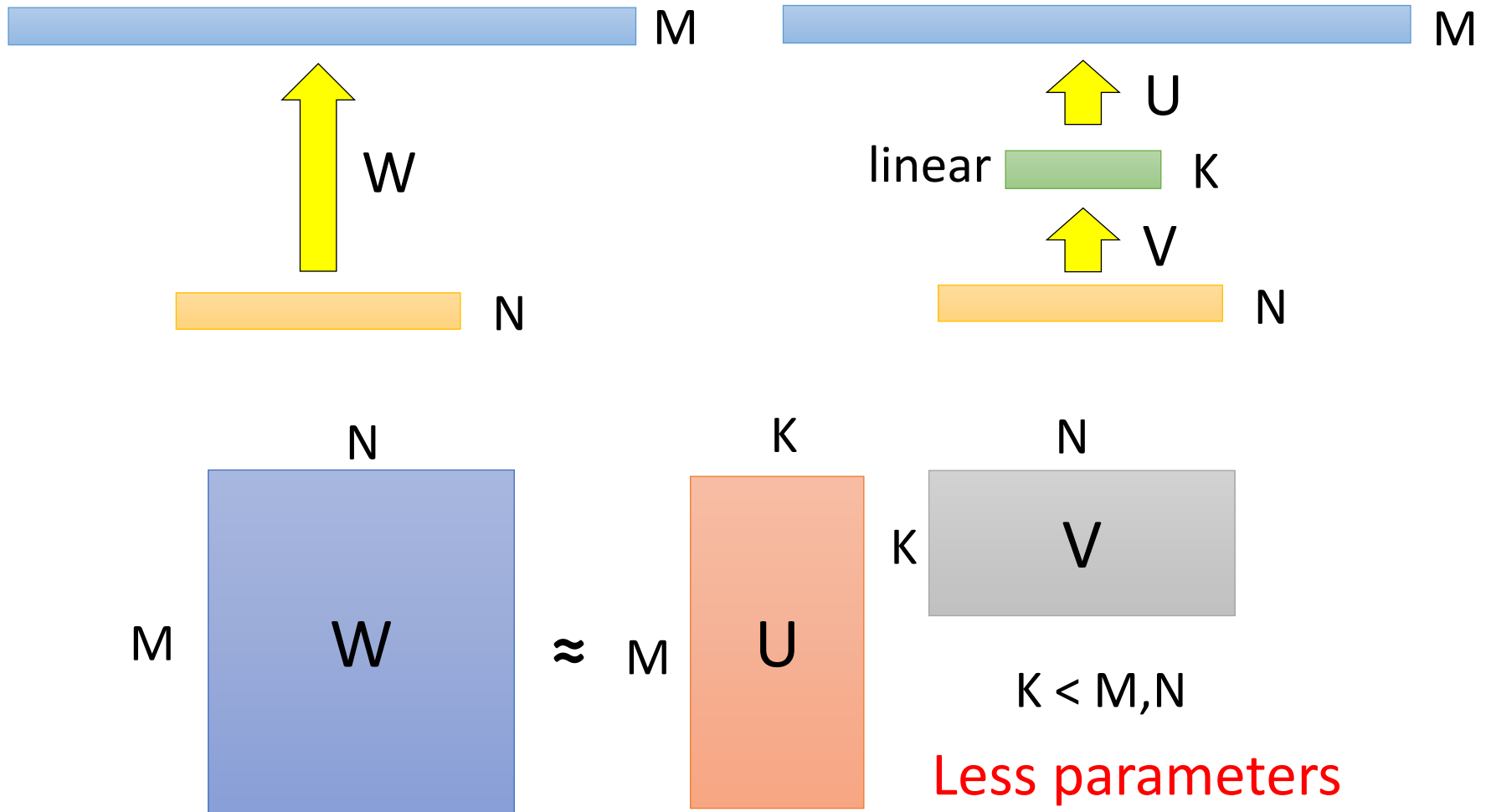
$$k \times k \times I$$

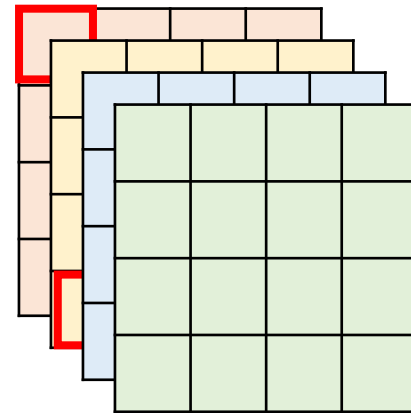
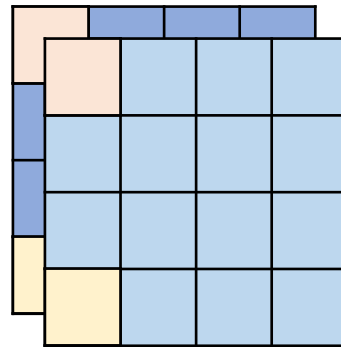
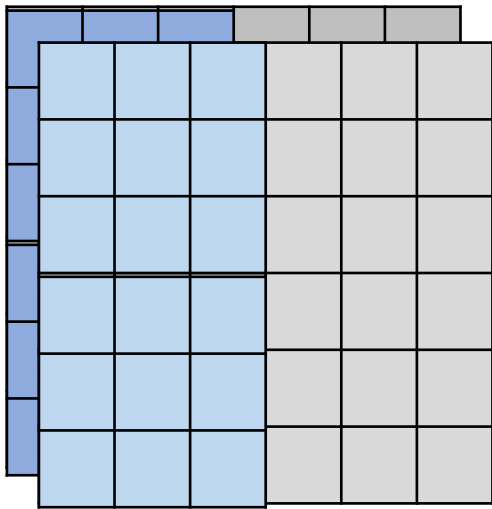
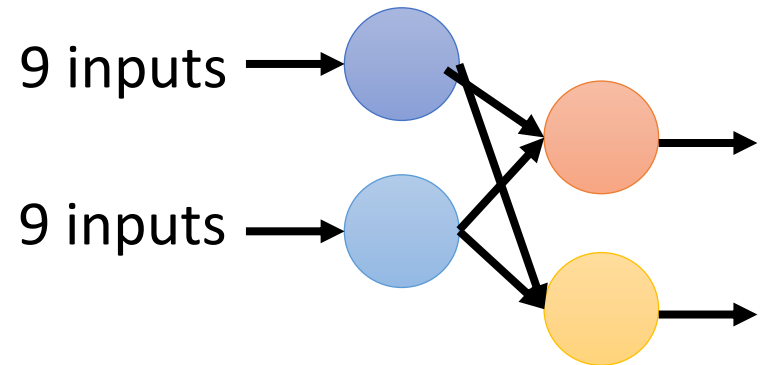
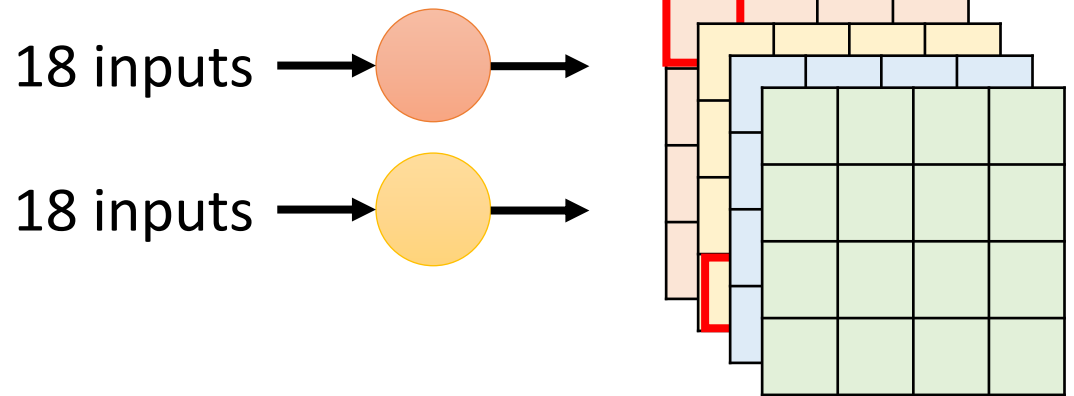
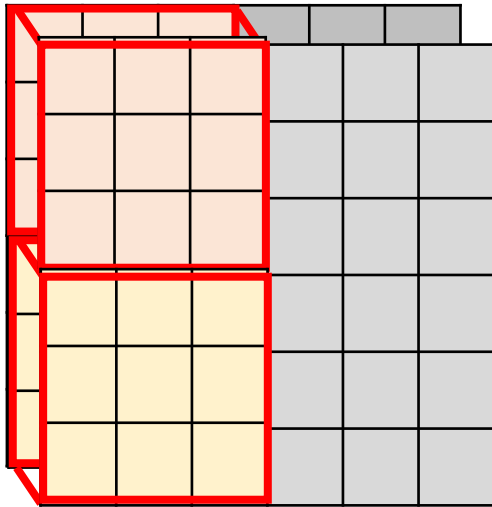


$$I \times O$$

$$k \times k \times I + I \times O$$

# Low rank approximation





# To learn more .....

- SqueezeNet
  - <https://arxiv.org/abs/1602.07360>
- MobileNet
  - <https://arxiv.org/abs/1704.04861>
- ShuffleNet
  - <https://arxiv.org/abs/1707.01083>
- Xception
  - <https://arxiv.org/abs/1610.02357>

# Dynamic Computation

The image features a dark grey background. On the right side, there is a complex geometric structure composed of several parallel lines. A prominent blue line runs diagonally from the bottom left towards the top right. To its right, a grey line follows a similar path but is slightly offset. Another blue line runs vertically from the top right towards the bottom left, intersecting the diagonal lines. The overall effect is a sense of depth and dynamic movement, consistent with the title 'Dynamic Computation'.



# Dynamic Computation

- Can network adjust the computation power it need?

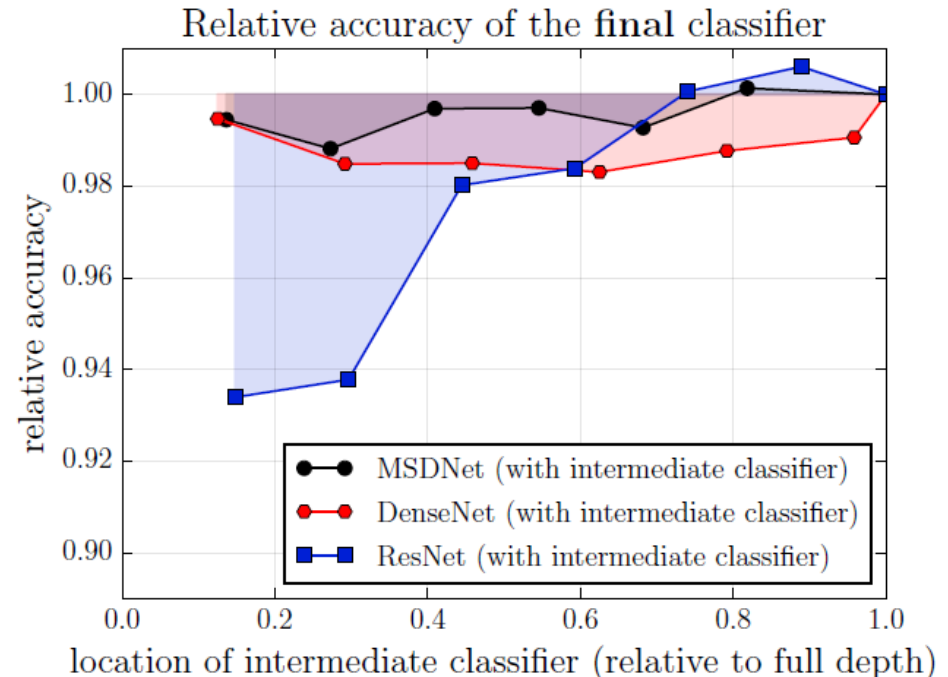
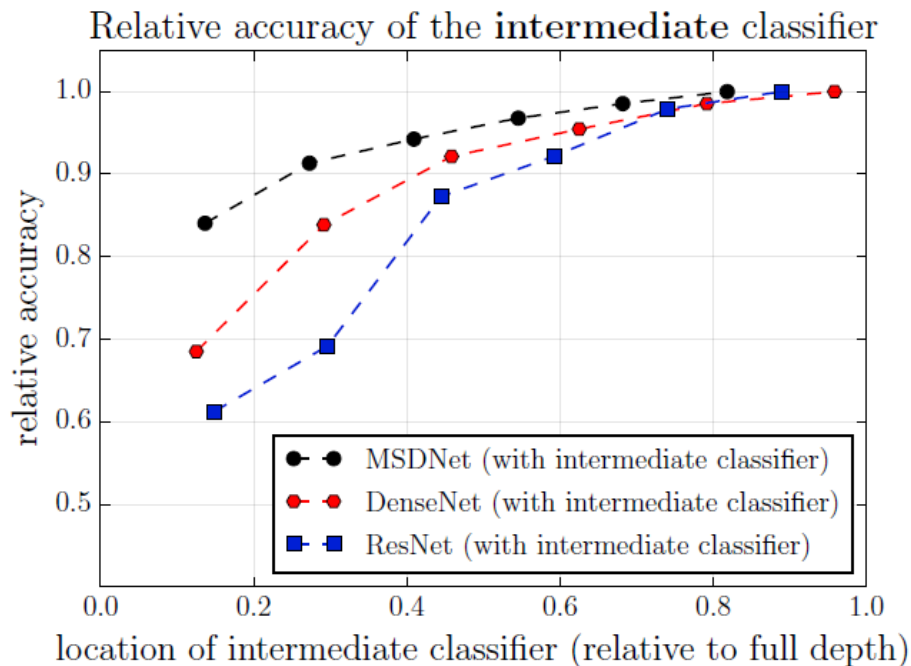
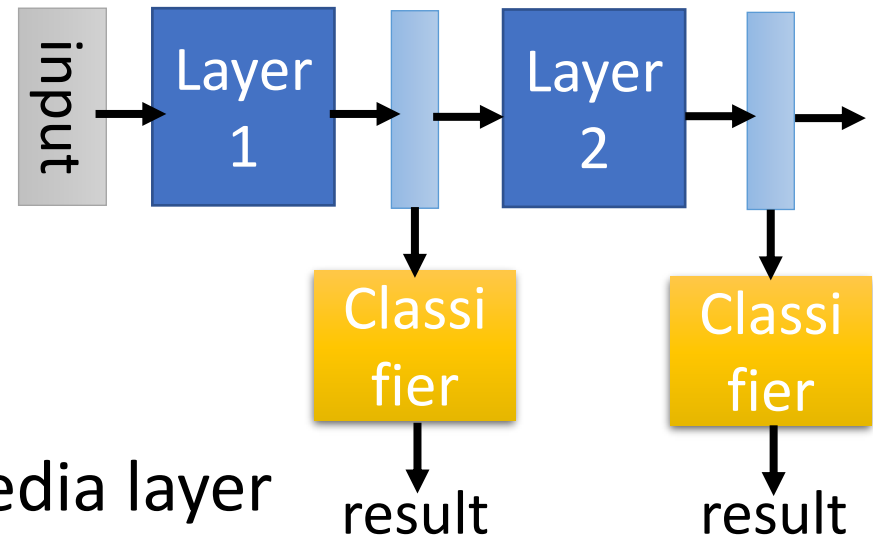
資源充足，那麼就做到最好

減少運算量，先求有再求好  
(但也不要太差)

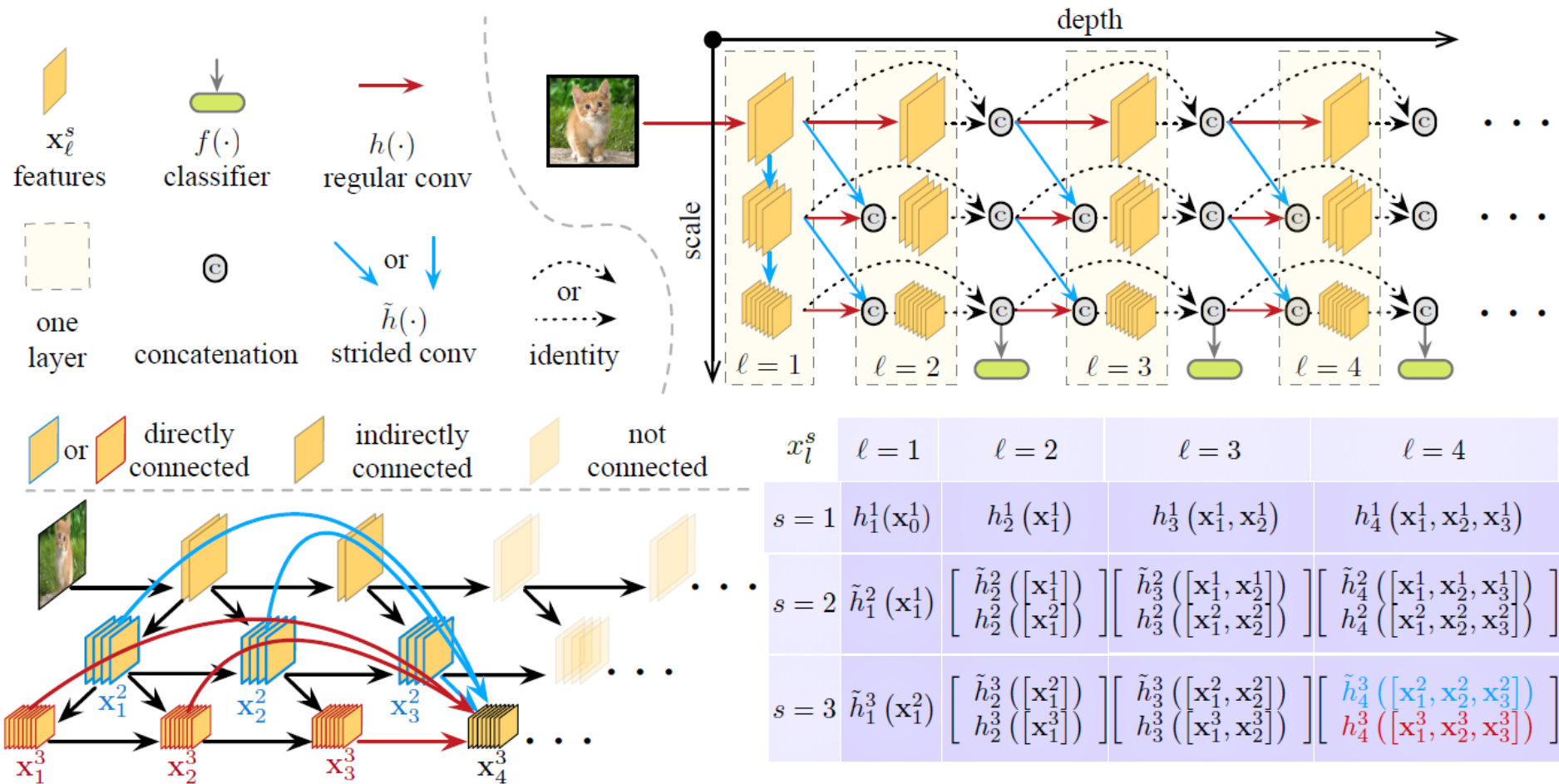


# Possible Solutions

- 1. Train multiple classifiers
- 2. Classifiers at the intermedia layer



# Multi-Scale Dense Networks



# Concluding Remarks

- Network Pruning
- Knowledge Distillation
- Parameter Quantization
- Architecture Design
- Dynamic Computation