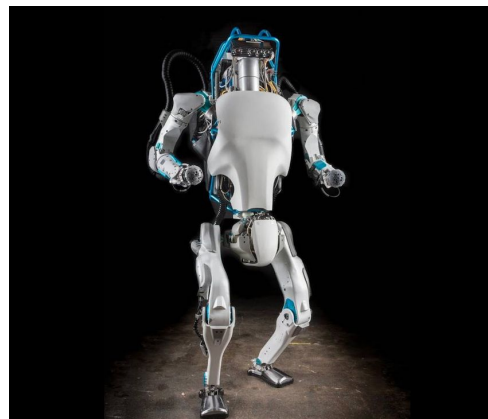# Neural Network Applications

## Self-Driving Car



This image is in the public domain

## Robots



This image is in the public domain
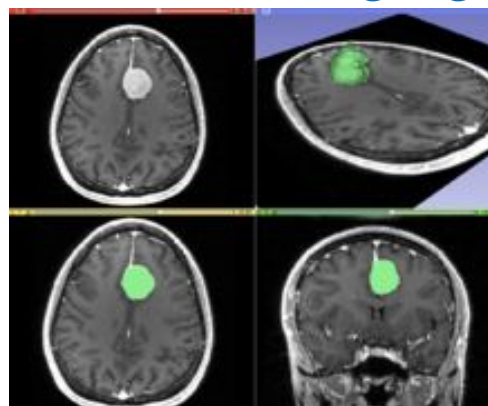
## Image processing



This image is in the public domain
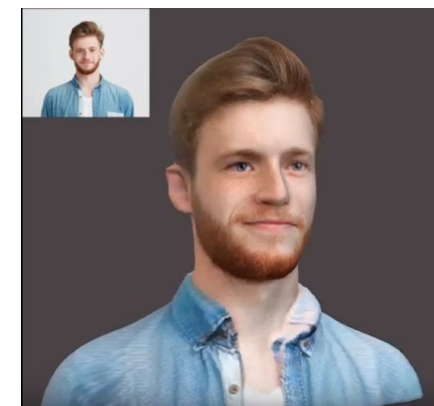
## Machine Translation



This image is in the public domain

## Medical imaging



This image is in the public domain

## 3D scanning



This image is in the public domain

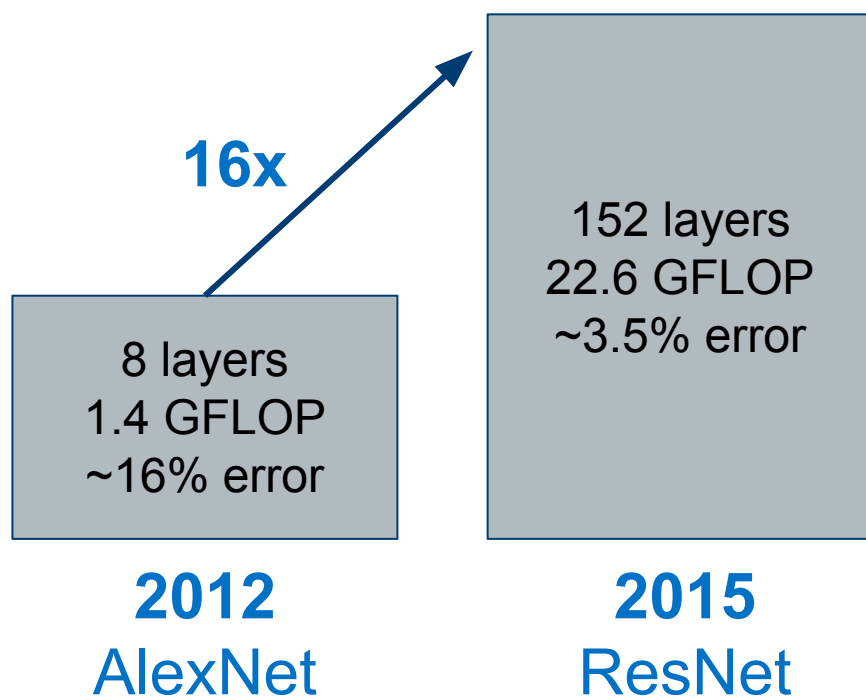# Where does Inference of Neural Networks Compute?

## Standalone

## Client-Server



All images are in the public domain

(intel)

# Models are Getting Larger

## Image Classification



16x

2012
AlexNet

8 layers
1.4 GFLOP
~16% error

2015
ResNet

152 layers
22.6 GFLOP
~3.5% error

## Speech Recognition



16x

2014
Deep Speech 1

80 GFLOP
~8% error

2015
Deep Speech 2

465 GFLOP
~5% error

intel
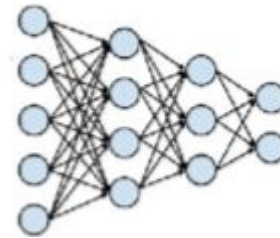
# The First Challenge: Model Size

- Hard to distribute large models through over-the-air update
- The first run is slow due to loading weights.

# The Second Challenge: Speed

Compression tool

Model(FP32)

Hardware

CPU, GPU

VPU, FPGA, ASIC

Accuracy

Speed

Trade off between accuracy and performance

# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
7. Winograd Transformation

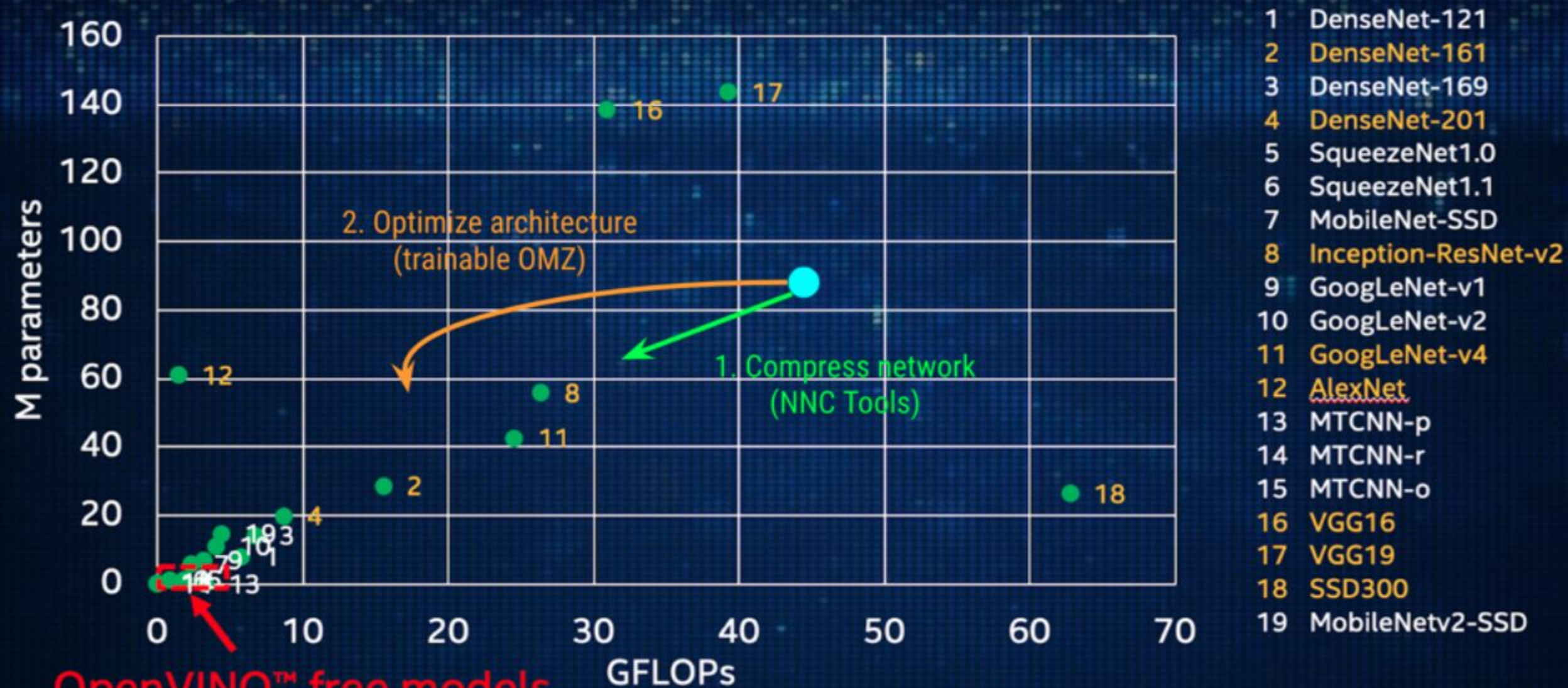OpenVINO™ PUBLIC AND FREE MODELS

1   DenseNet-121
2   DenseNet-161
3   DenseNet-169
4   DenseNet-201
5   SqueezeNet1.0
6   SqueezeNet1.1
7   MobileNet-SSD
8   Inception-ResNet-v2
9   GoogLeNet-v1
10  GoogLeNet-v2
11  GoogLeNet-v4
12  AlexNet
13  MTCNN-p
14  MTCNN-r
15  MTCNN-o
16  VGG16
17  VGG19
18  SSD300
19  MobileNetv2-SSD

2. Optimize architecture (trainable OMZ)

1. Compress network (NNC Tools)
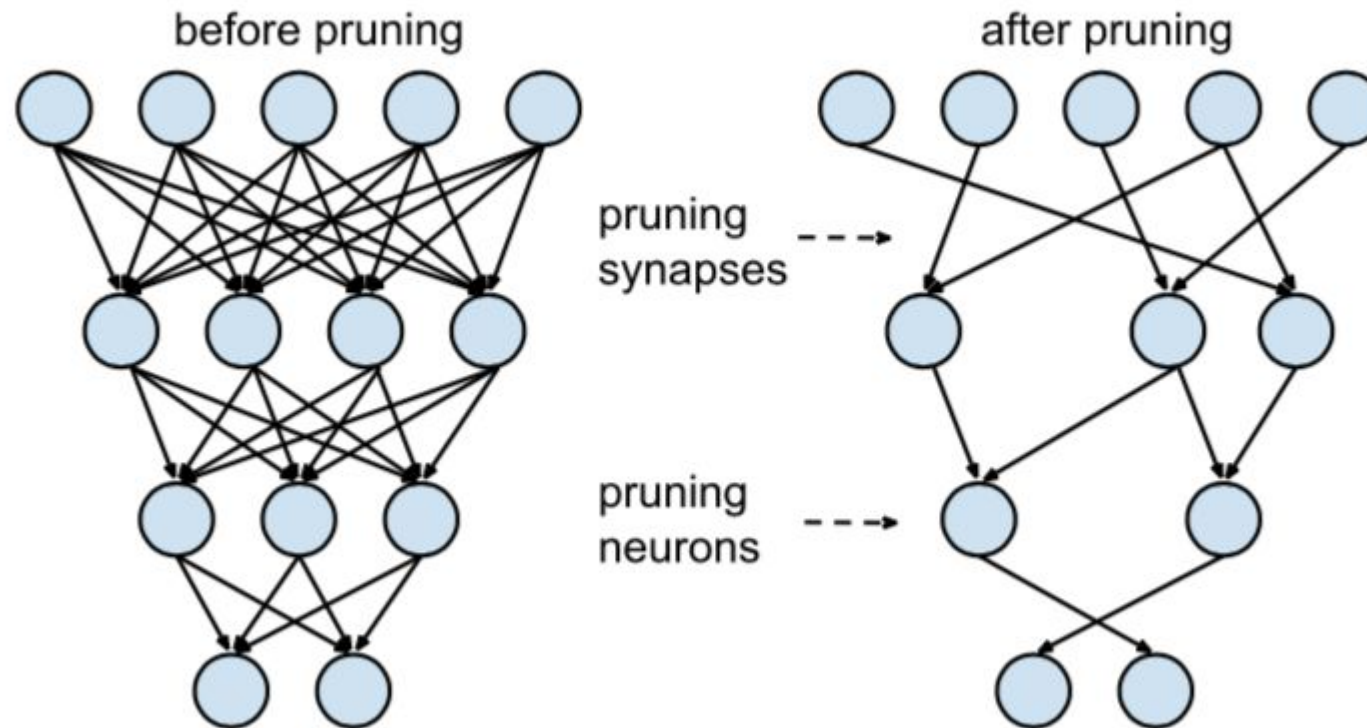
OpenVINO™ free models
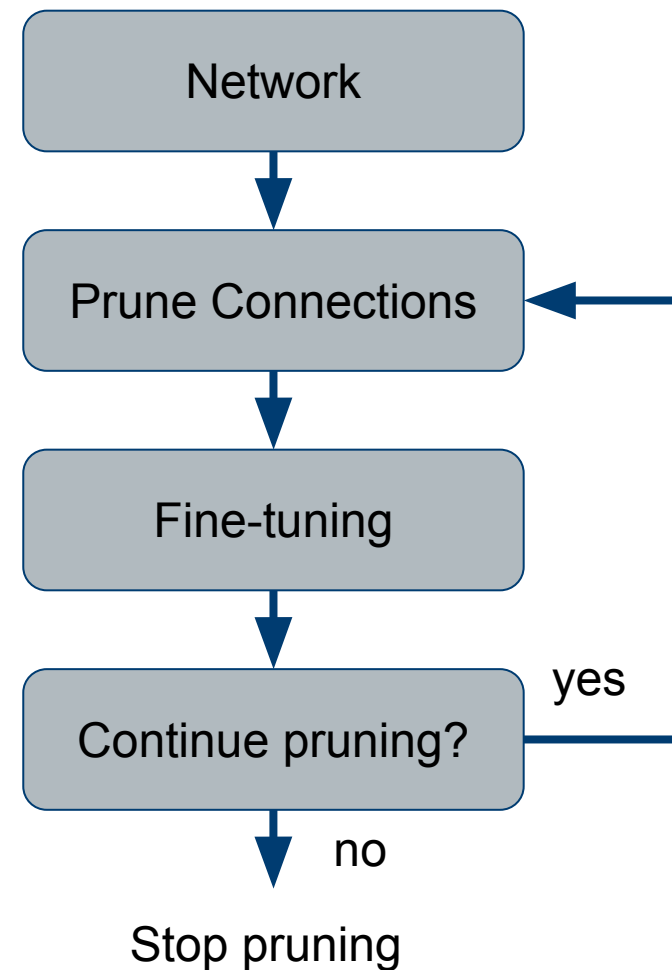
# Algorithms for Efficient Inference

1. **Pruning**
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
7. Winograd Transformation

# Pruning Neural Networks



before pruning                                    after pruning

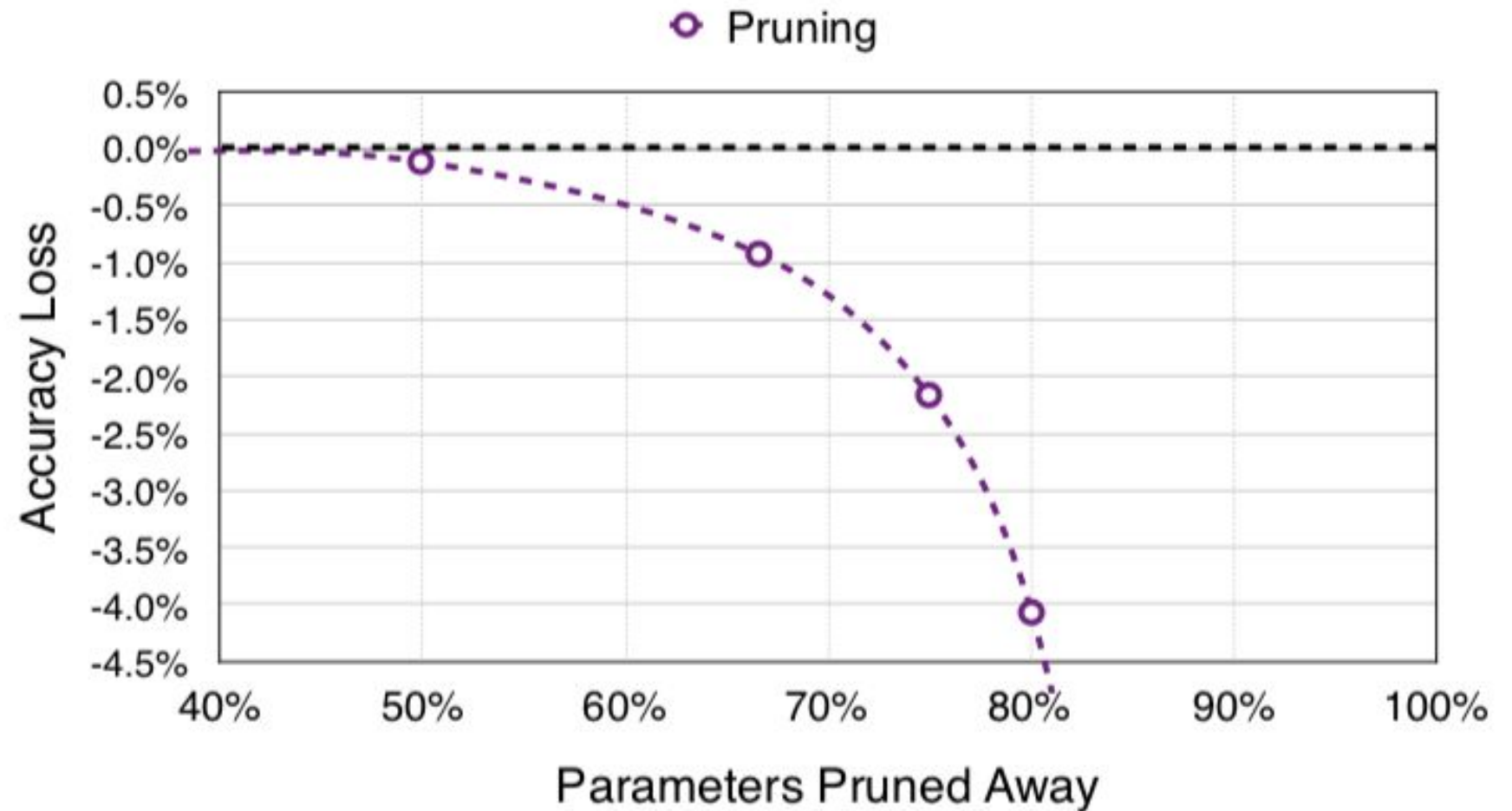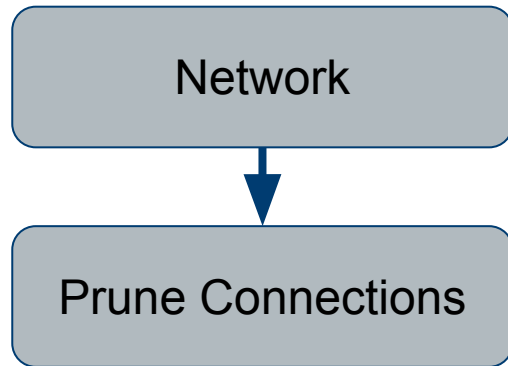pruning synapses - - ->

pruning neurons - - ->
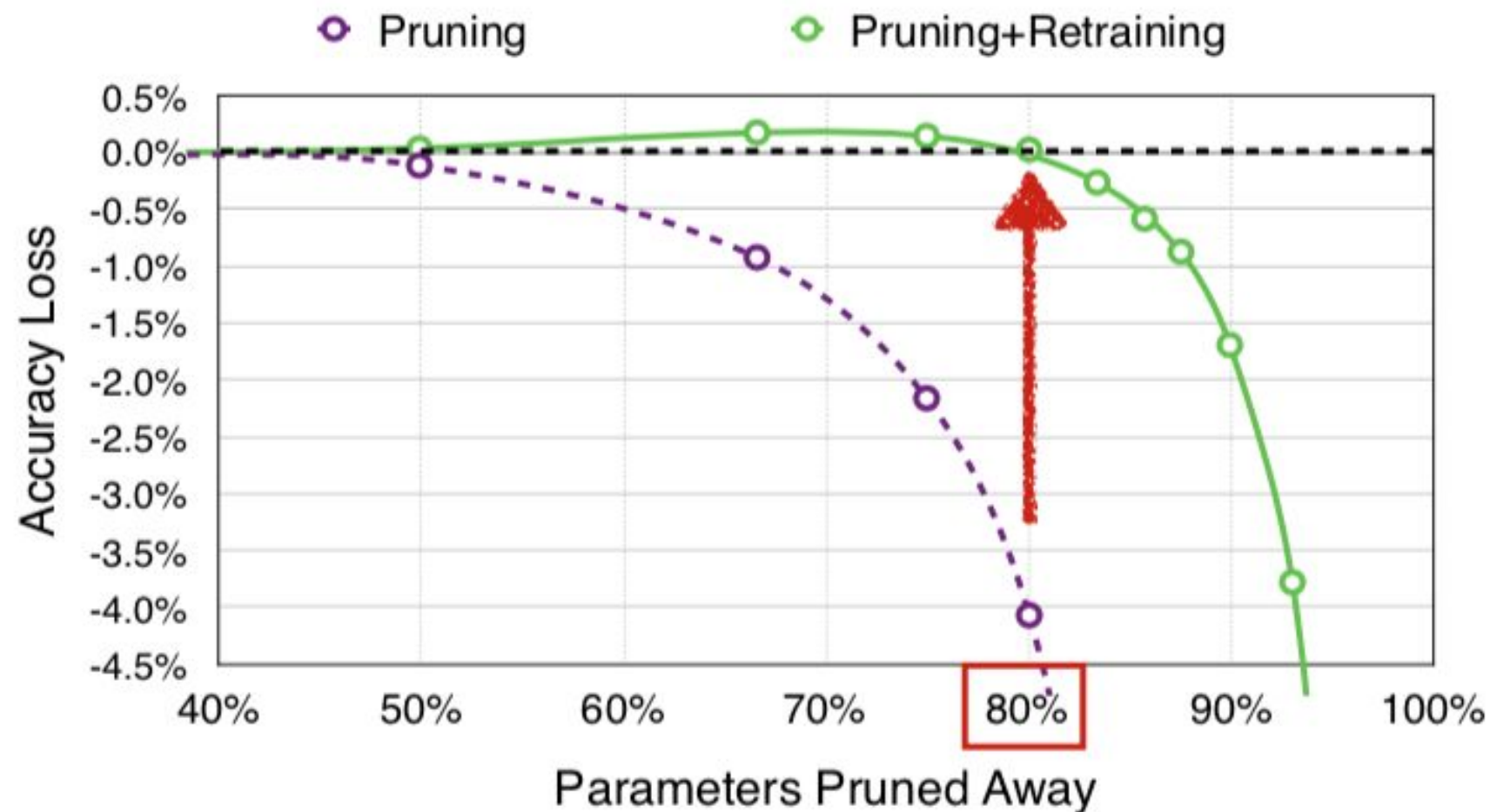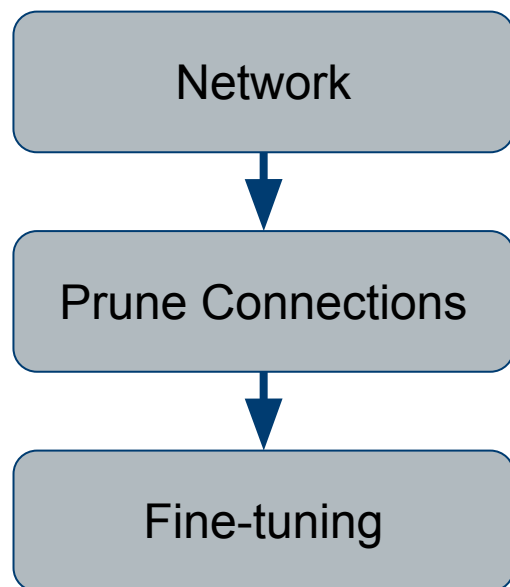
# Pruning Neural Networks

- Criteria for pruning
  - Connections with low weights
  - Neurons or filters with low impact
- External limitations
  - Spatial structure of pruned weights

```
┌─────────────────┐
│     Network     │
└─────────────────┘
         │
         ▼
┌─────────────────┐◄──────┐
│ Prune Connections│       │
└─────────────────┘       │
         │                │
         ▼                │
┌─────────────────┐       │
│   Fine-tuning   │       │
└─────────────────┘       │
         │            yes │
         ▼                │
┌─────────────────┐───────┘
│ Continue pruning?│
└─────────────────┘
         │ no
         ▼
    Stop pruning
```
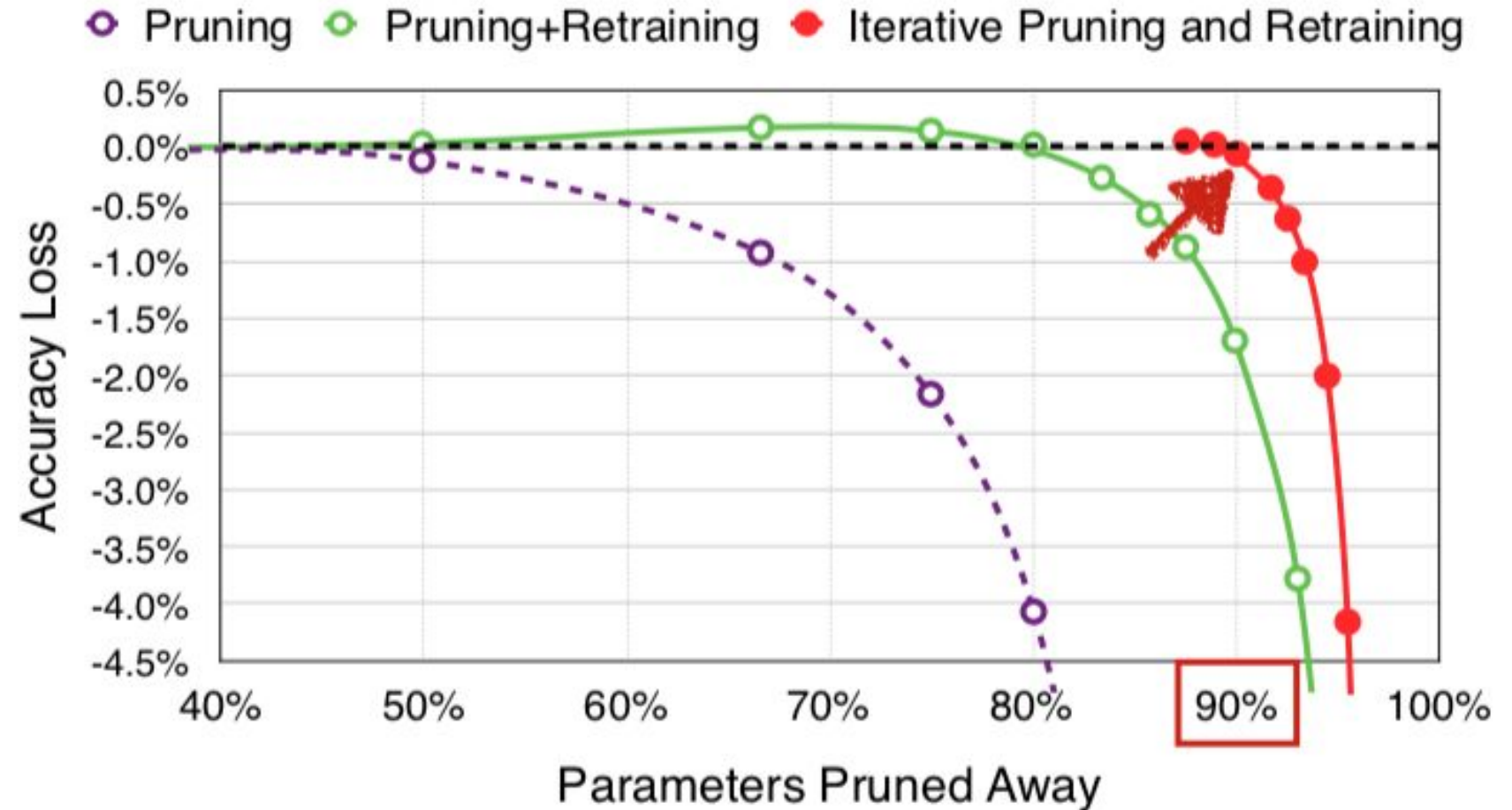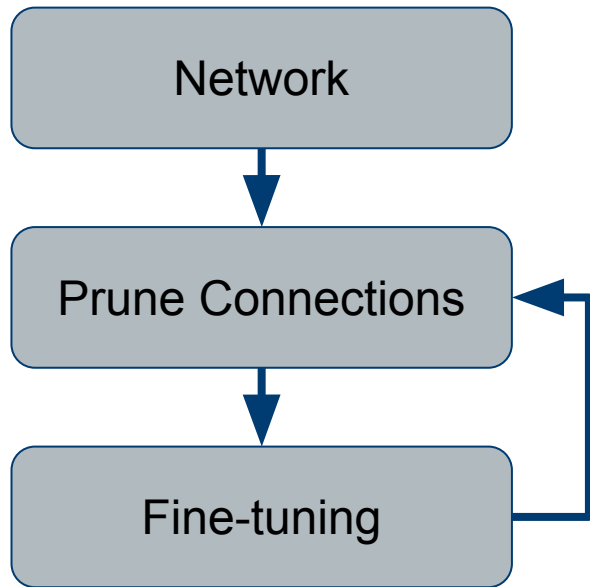
# Pruning Neural Networks

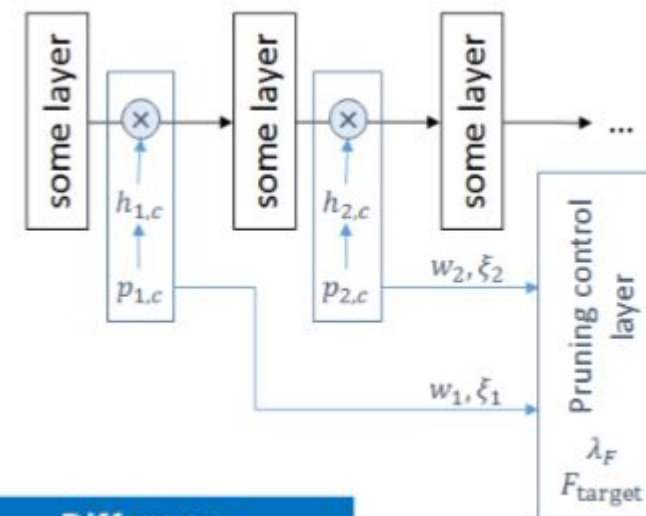# Retrain to Recover Accuracy

# Iteratively Retrain to Recover Accuracy

# Channel Pruning

- Drop some channels by Dropout layers with learnable dropout probabilities.



| Task | Original network | Pruned network | Difference |
|---|---|---|---|
| VGG-16 ILSVRC2012 classification (1000 classes) | 15.47 GFLOP 68.4% top-1 acc. 88.4% top-5 acc. | 3.87 GFLOP 67.4% top-1 acc. 88.0% top-5 acc. | −75% (FLOP) −1.0% top-1 acc. −0.4% top-5 acc. |
| MobileNet + SSD Vehicle detection (1 class) | 2.26 GFLOP 25 ms@CPU 87.7% AP | 1.33 GFLOP 19 ms@CPU 86.2% AP | −41% (FLOP) −1.5% AP |
| AlexNet ILSVRC2012 classification (1000 classes) | 724 MFLOP 56.8% top-1 acc. 79.9% top-1 acc. | 411 MFLOP 55.8% top-1 acc. 79.1% top-5 acc. | −43% (FLOP) −1.0% top-1 acc. −0.8% top-5 acc. |

# Sparsification

- The sparsification algorithm based on probabilistic approach and loss regularization.

Inception v3

Baseline: 77.46% top1 acc.

Sparsity Rate(SR) 50%: 77.25% top1 acc.

SR 92%: 76.6% top1 acc.

MobileNet v2

Baseline: 71.6% top1 acc.

SR 50%: 71.2% top1 acc.

SR 78%: 69.98% top1 acc.

Ordinary convolution

$$output = conv(x, w)$$

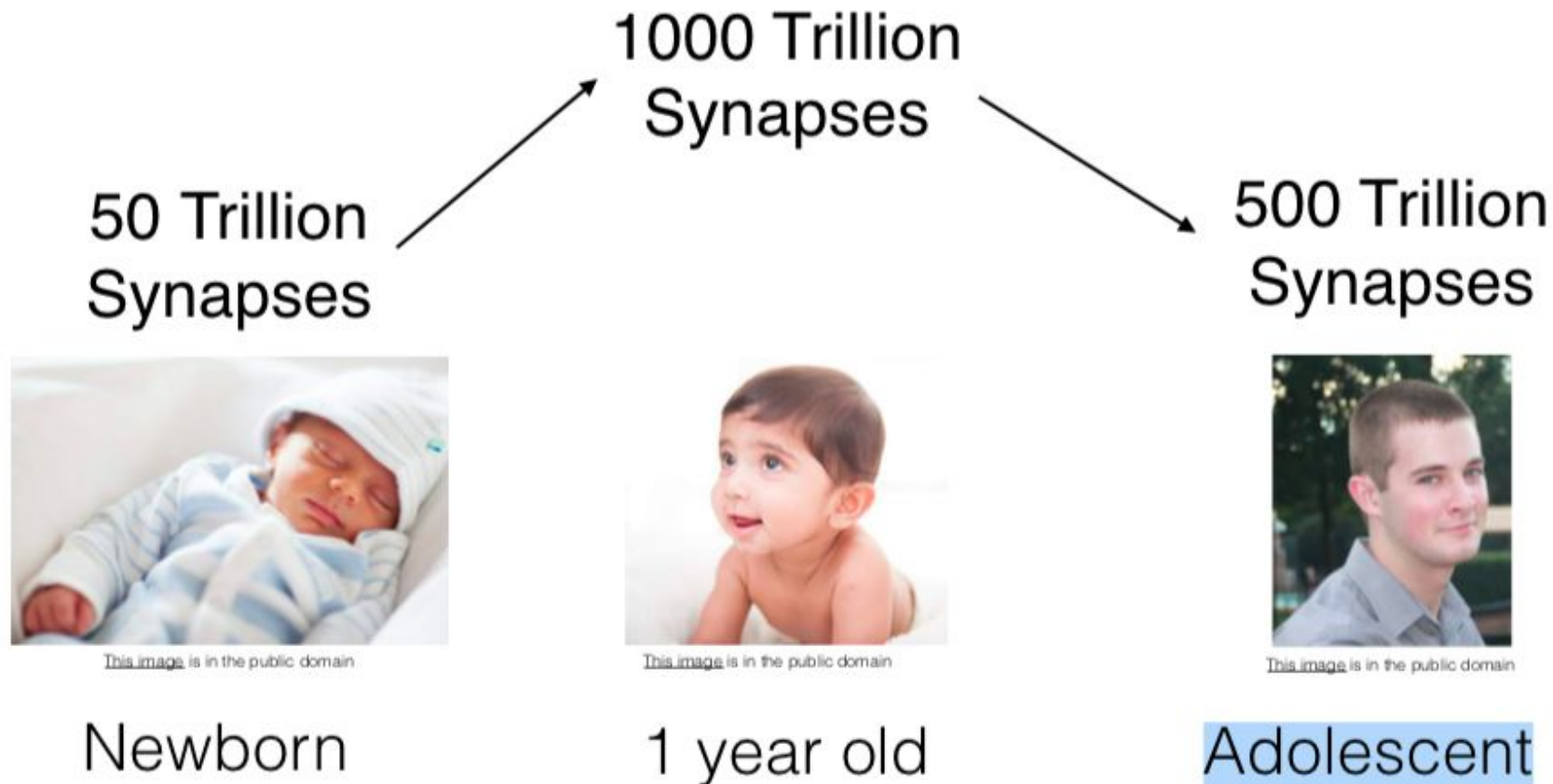Sparsifing weights we reparametrize weights as:

$$\bar{w} = w * z$$

Where:

$$w - weights, w \in R$$
$$z - binary\ mask, z \in [0, 1]$$

We train these masks using modificated loss

$$Loss = Loss_{task} + \alpha \left( \sum_{l}^{Layers} ||z|| - target \right)^2$$

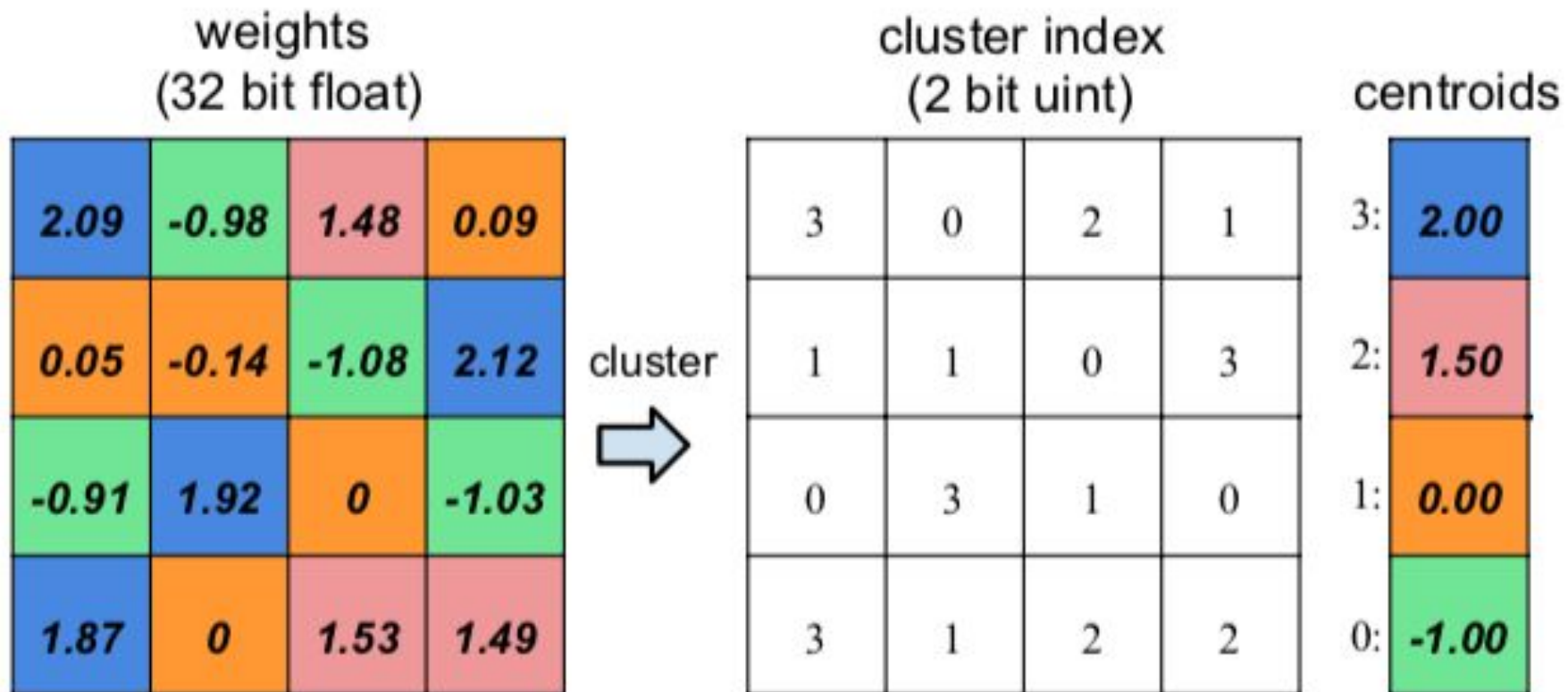# Pruning Happens in Human Brain



Christopher A Walsh. Peter Huttenlocher (1931-2013). Nature, 502(7470):172–172, 2013.

Slide credits by Song Han

# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
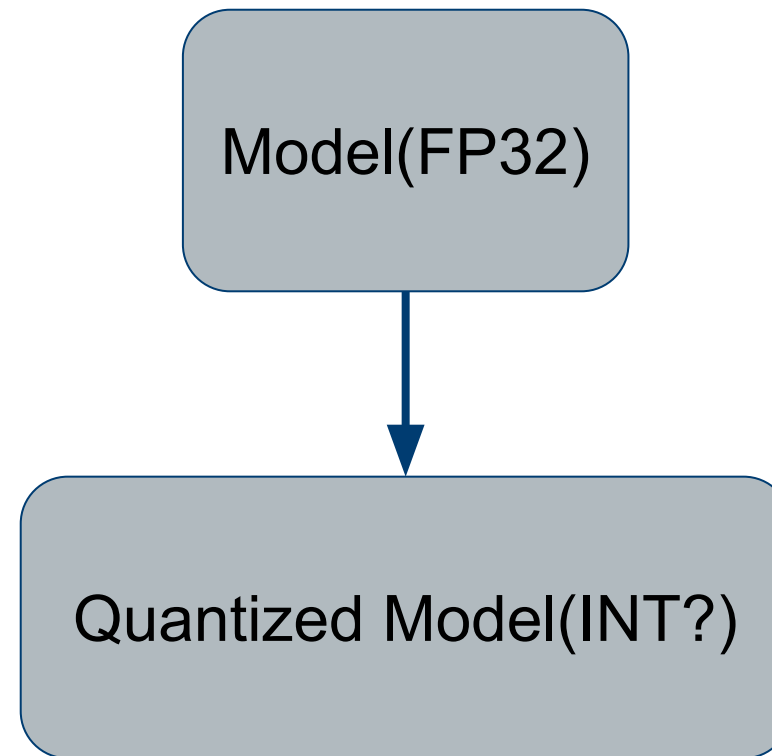6. Low Rank Approximation
7. Winograd Transformation

# Weight sharing

# Pruning + Weight sharing + Huffman Encoding

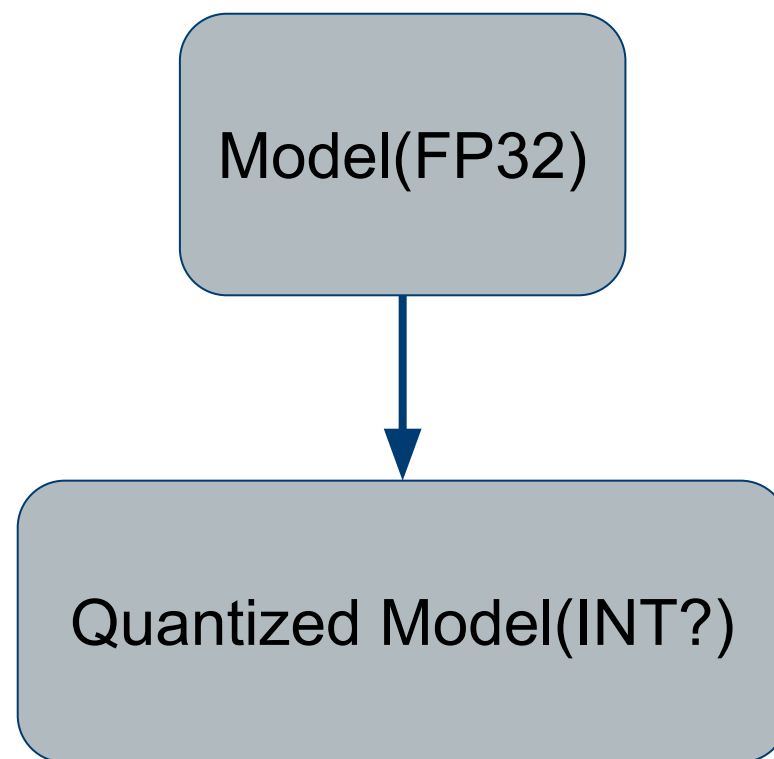| Network | Original Size | Compressed Size | Compression Ratio | Original Accuracy | Compressed Accuracy |
|---|---|---|---|---|---|
| LeNet-300 | 1070KB → 27KB | | 40x | 98.36% → 98.42% | |
| LeNet-5 | 1720KB → 44KB | | 39x | 99.20% → 99.26% | |
| AlexNet | 240MB → 6.9MB | | 35x | 80.27% → 80.30% | |
| VGGNet | 550MB → 11.3MB | | 49x | 88.68% → 89.09% | |
| GoogleNet | 28MB → 2.8MB | | 10x | 88.90% → 88.92% | |
| ResNet-18 | 44.6MB → 4.0MB | | 11x | 89.24% → 89.28% | |

(intel)

# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
7. Winograd Transformation

# Quantization

```
┌─────────────────────┐
│                     │
│    Model(FP32)      │
│                     │
└─────────────────────┘
           │
           ▼
┌─────────────────────┐
│                     │
│ Quantized Model(INT?)│
│                     │
└─────────────────────┘
```

# Quantization

- Level 1: Post-Training Quantization w/o dataset

Model(FP32)

Quantized Model(INT?)

- Int8
- BatchNorms

# Quantization

- Level 2: Post-Training Quantization w dataset



- Int8, Int4?

# Quantization



- Level 3: Quantization via training

| Dataset with Labels | Model(FP32) | Training pipeline |

Quantized Model(INT?)

- Int8, Int4, Binary

# Uniform Affine Quantization

- Quantization:

$$x_{int} = round\left(\frac{x}{\Delta}\right) + z$$

$$x_Q = clamp(0, N_{levels} - 1, x_{int})$$

- $\Delta$ specifies the step size of the quantizer and floating point zero maps to zero-point.
- $z$ - zero-point.
- $N_{levels} = 256$ for 8-bits of precision

- De-quantization:

$$x_{float} = (x_Q - z)\Delta$$

# Uniform Affine Quantization

- 2D convolution between a weight and an activation:

$$y(k, l, n) = \Delta_w \Delta_x conv(w_Q(k, l, m; n) - z_w, x_Q(k, l, m) - z_x)$$

$$y(k, l, n) = conv(w_Q(k, l, m; n), x_Q(k, l, m)) - z_w \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} x_Q(k, l, m)$$

$$- z_x \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \sum_{m=0}^{N-1} w_Q(k, l, m; n) + z_x z_w$$

# Uniform Symmetric Quantization

- Quantization, zero-point = 0
  - Activations:

$$x_{int} = round(\frac{x}{\Delta})$$

$$x_Q = clamp(-N_{levels}/2, N_{levels}/2 - 1, x_{int}) \qquad \text{if signed}$$

$$x_Q = clamp(0, N_{levels} - 1, x_{int}) \qquad \text{if un-signed}$$

  - Weights

$$x_Q = clamp(-(N_{levels}/2 - 1), N_{levels}/2 - 1, x_{int}) \qquad \text{if signed}$$

$$x_Q = clamp(0, N_{levels} - 2, x_{int}) \qquad \text{if un-signed}$$

(intel)

# Uniform Symmetric Quantization

- MKL DNN Int8 Workflow:

$$X_{s32} = W_{s8} \times \alpha u8 + b_{s32} \approx Q_\alpha Q_w X_{f32}$$

$$\text{where } X_{f32} = W_{f32} \times \alpha_{f32} + b_{f32}$$

$Q_\alpha = \frac{255}{R_\alpha}$ is the quantization factor for activations with non-negative values.

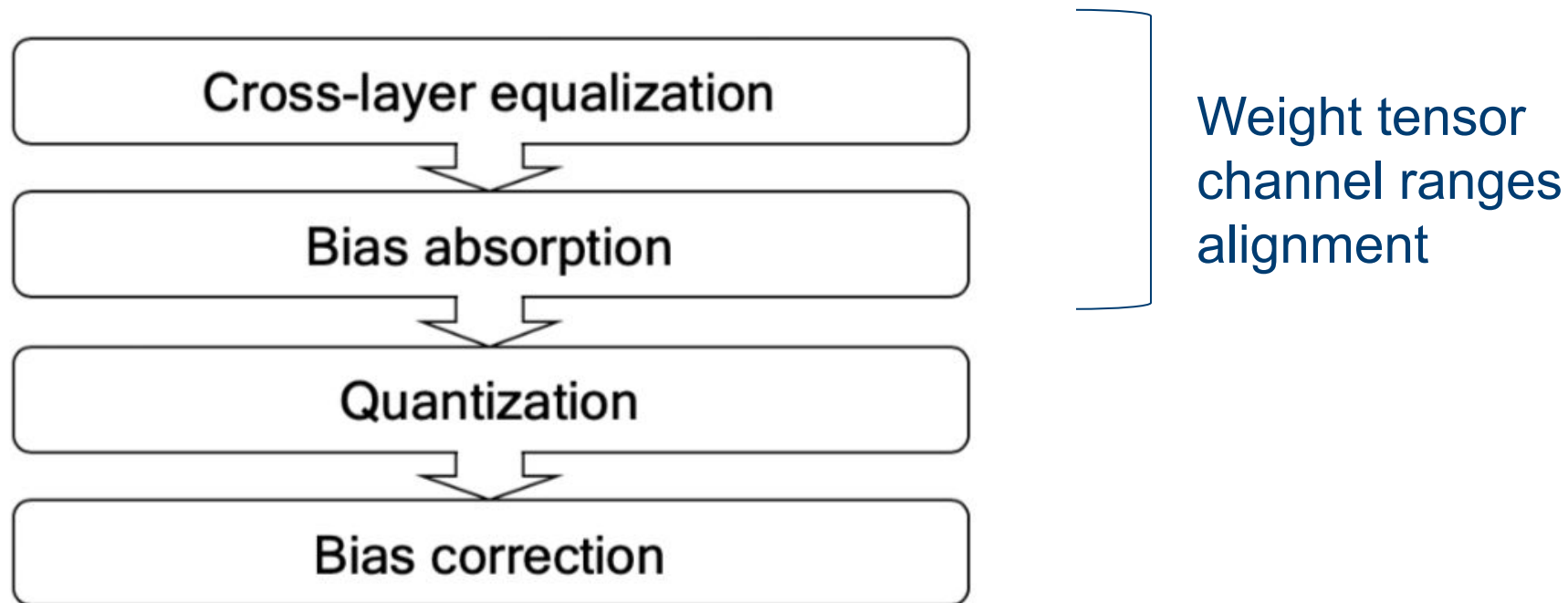$Q_w = \frac{127}{R_w}$ is the quantization factor for weights.

$$\alpha_{u8} = \lceil Q_\alpha \alpha_{f32} \rceil \in [0, 255]$$

$$W_{s8} = \lceil Q_w W_{f32} \rceil \in [-127, 127]$$

$$b_{s32} = \lceil Q_\alpha Q_w b_{f32} \rceil \in [-2^{31}, 2^{31} - 1]$$

(intel)

# Data-Free Quantization (Level 1)

- The basic idea is to use BatchNorm statistics to estimate the range of activations.
- Algorithm Flow:



Weight tensor channel ranges alignment

# Cross-Layer Equalization

Given two layers, $x_1 = f(W_1 x_0 + b_1)$ and $x_2 = f(W_2 x_1 + b_2)$ through scaling invariance we have that:

$$x_2 = f(W_2 f(W_1 x_0 + b_1) + b_2)$$
$$= f(W_2 S \hat{f}(S^{-1} W_1 x_0 + S^{-1} b_1) + b_2)$$
$$= f(\widehat{W}_2 \hat{f}(\widehat{W}_1 x_0 + \hat{b}_1) + b_2)$$

where:

$$s_i = \frac{1}{r_i^{(1)}} \sqrt{r_i^{(1)} r_i^{(2)}} \qquad ; r_i^{(1)} = \quad 2 \cdot \max_j |\widehat{W}_{ij}^{(1)}|$$

# Quantization

- Uniform Affine Quantization is used for weights and activations

    - (min(w), max(w)) - range for weights

    - (mean - 3 * var, mean + 3 * var) - range for activations

# Quantization bias correction

- x -input
- $\mathbf{W}$ - weights FP32
- $\widetilde{\mathbf{W}}$ - quantized weights

$$\widetilde{\mathbf{y}} = \mathbf{y} + \epsilon\mathbf{x}, \quad \epsilon = \widetilde{\mathbf{W}} - \mathbf{W}$$

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{y}] + \mathbb{E}[\epsilon\mathbf{x}] - \mathbb{E}[\epsilon\mathbf{x}]$$

$$= \mathbb{E}[\widetilde{\mathbf{y}}] - \mathbb{E}[\epsilon\mathbf{x}]$$

$$\mathbb{E}[\mathbf{x}_c] = \mathbb{E}\left[\mathrm{ReLU}\left(\mathbf{x}_c^{pre}\right)\right]$$

$$= \gamma_c \mathcal{N}\left(\frac{-\boldsymbol{\beta}_c}{\gamma_c}\right) + \boldsymbol{\beta}_c\left[1 - \Phi\left(\frac{-\boldsymbol{\beta}_c}{\gamma_c}\right)\right]$$

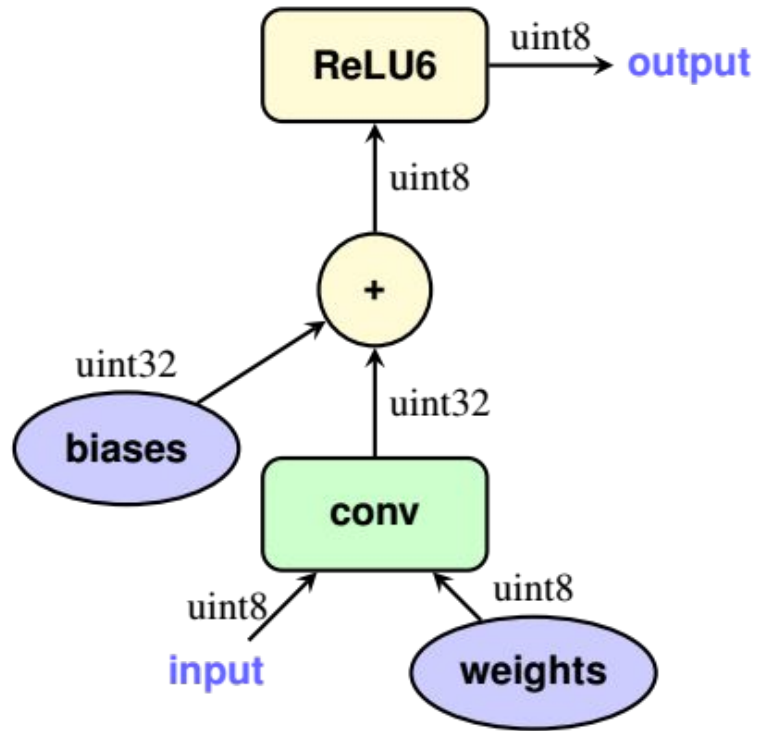# Quantization with represented dataset (Level 2)

- Run N examples through the FP32 model and collect for each layer the per-channel pre-activation statistics:
    - moving average of the minimum and maximum values across batches to determine the quantizer parameters for activations.
    - mean activation for quantization bias correction
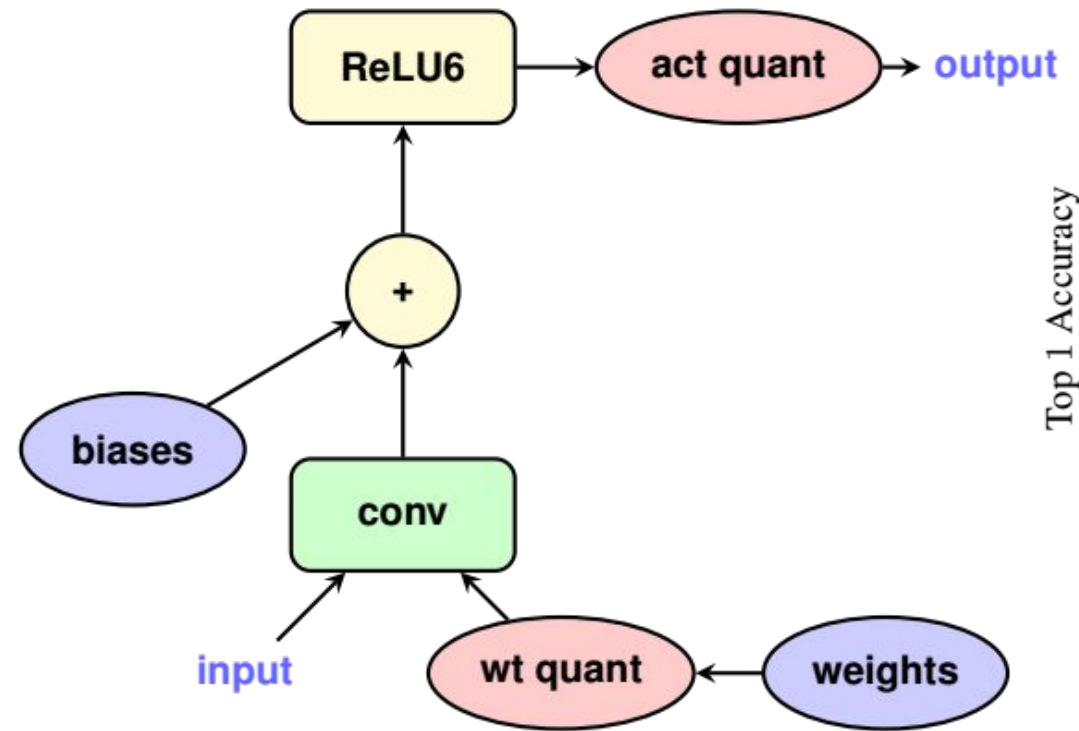
# Quantization via training (Level 3)

- Simulated quantization layer (Fake Quantization Layer)

$$\text{clamp}(r; a, b) := \min\left(\max(x, a), b\right)$$

$$s(a, b, n) := \frac{b - a}{n - 1}$$

$$q(r; a, b, n) := \left\lfloor \frac{\text{clamp}(r; a, b) - a}{s(a, b, n)} \right\rceil s(a, b, n) + a, \tag{12}$$

# Quantization via training (Level 3)



(a) Integer-arithmetic-only inference

(b) Training with simulated quantization

# Quantization via training (Level 3)

Algorithm:

1. Create a training graph of the floating-point model.
2. Insert fake quantization operations in locations where tensors will be downcasted to fewer bits during inference.
3. Train in simulated quantized mode until convergence.
4. Create and optimize the inference graph for running in a low bit inference engine.
5. Run inference using the quantized inference graph.
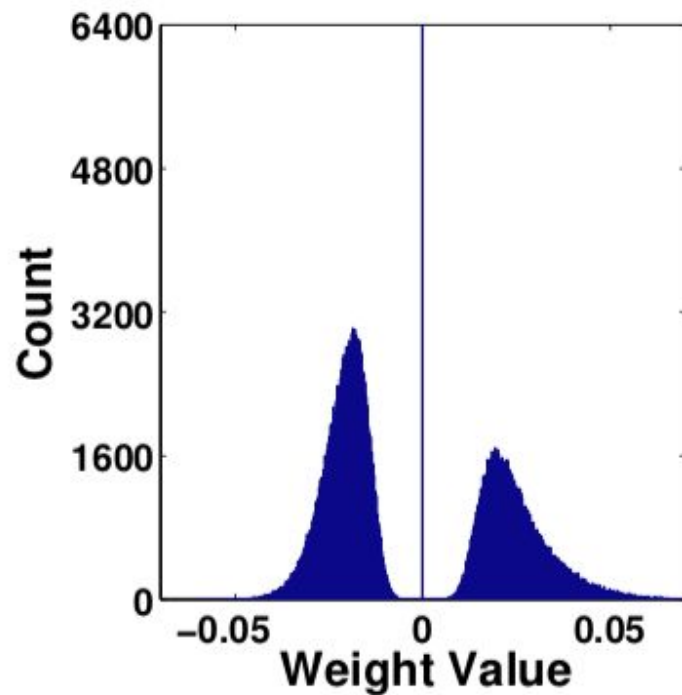
# Summary: Quantization approaches

| | ~D | ~BP | ~AC | MobileNetV2 | | MobileNetV1 | | ResNet18 | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | FP32 | INT8 | FP32 | INT8 | FP32 | INT8 | INT6 |
| DFQ (ours) | ✓ | ✓ | ✓ | 71.7% | **71.2%** | 70.8% | **70.5%** | 69.7% | **69.7%** | 66.3% |
| Per-layer [18] | ✓ | ✓ | ✓ | 71.9% | 0.1% | 70.9% | 0.1% | 69.7% | 69.2%* | 63.8%* |
| Per-channel [18] | ✓ | ✓ | ✓ | 71.9% | 69.7% | 70.9% | 70.3% | 69.7% | 69.6%* | **67.5%*** |
| QT [16] ^ | ✗ | ✗ | ✓ | 71.9% | 70.9% | 70.9% | 70.0% | - | **70.3%[†]** | 67.3%[†] |
| SR+DR[†] | ✗ | ✗ | ✓ | - | - | - | 71.3% | - | 68.2% | 59.3% |
| QMN [31] | ✗ | ✗ | ✗ | - | - | 70.8% | 68.0% | - | - | - |
| RQ [21] | ✗ | ✗ | ✗ | - | - | - | **70.4%** | - | 69.9% | **68.6%** |

Table 5. Top1 ImageNet validation results for different models and quantization approaches. The top half compares level 1 approaches (~D: data free, ~BP: backpropagation-free, ~AC: Architecture change free) whereas in the second half we also compare to higher level approaches in literature. Results with * indicates our own implementation since results are not provided, ^ results provided by [18] and [†] results from table 2 in [21].
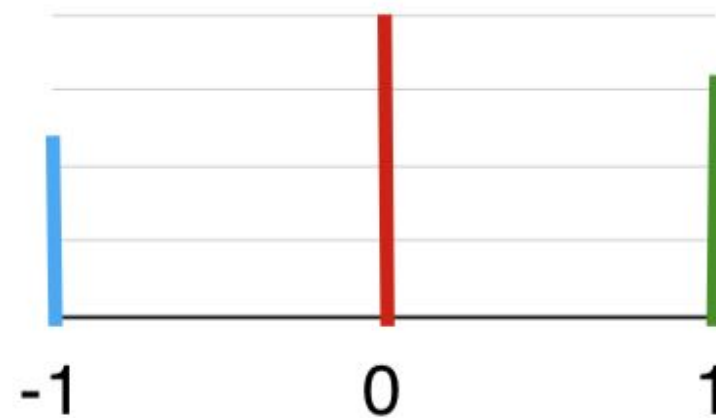
# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
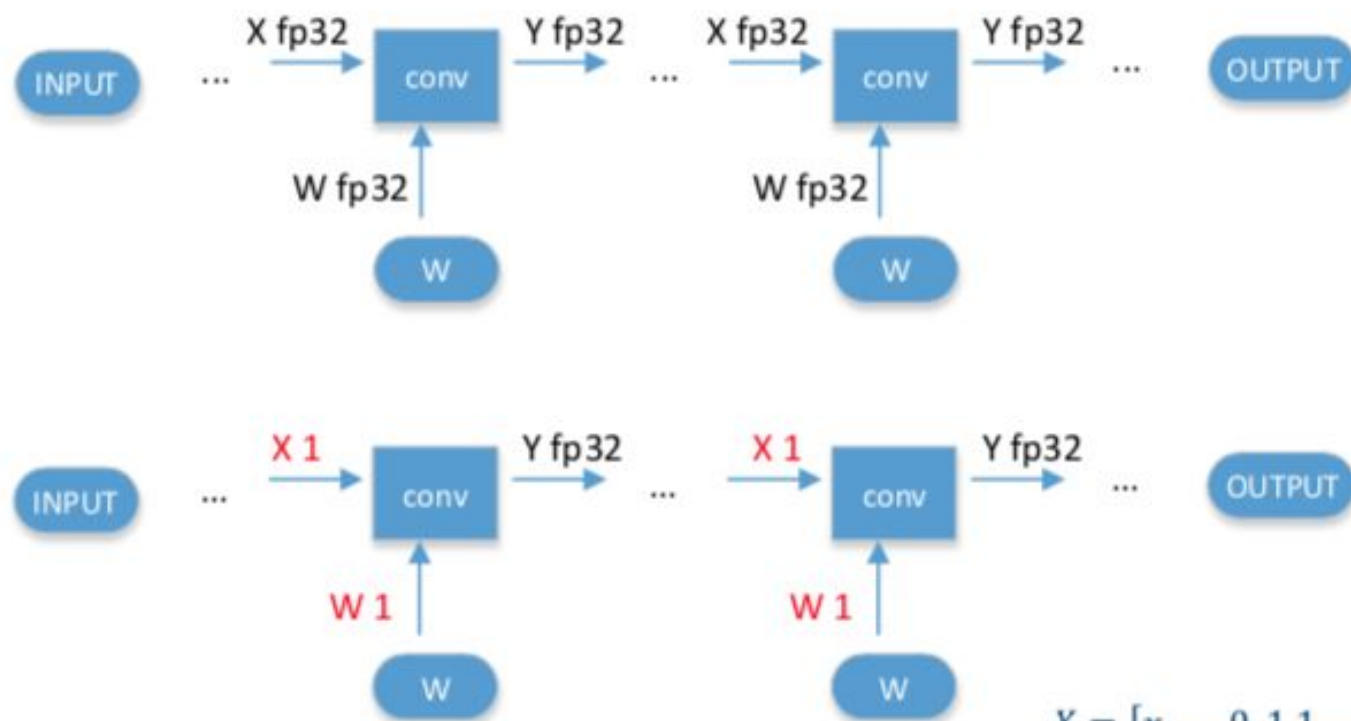7. Winograd Transformation

# Binary / Ternary Net

# Binary Net



1 value inference in conv is dot product

$$y_1 = w_1 * x_1 + w_2 * x_2 + \cdots + w_{32} * x_{32}$$

**32 * and 32+ operations**

$$y_1 = \{-1, +1\} * \{-1, +1\} + w_2 * x_2 + \cdots$$

$$x_i * w_i = \begin{array}{c|cc} x\backslash w & -1 & +1 \\ \hline -1 & +1 & -1 \\ +1 & -1 & +1 \end{array} \quad \{-1, +1\},$$
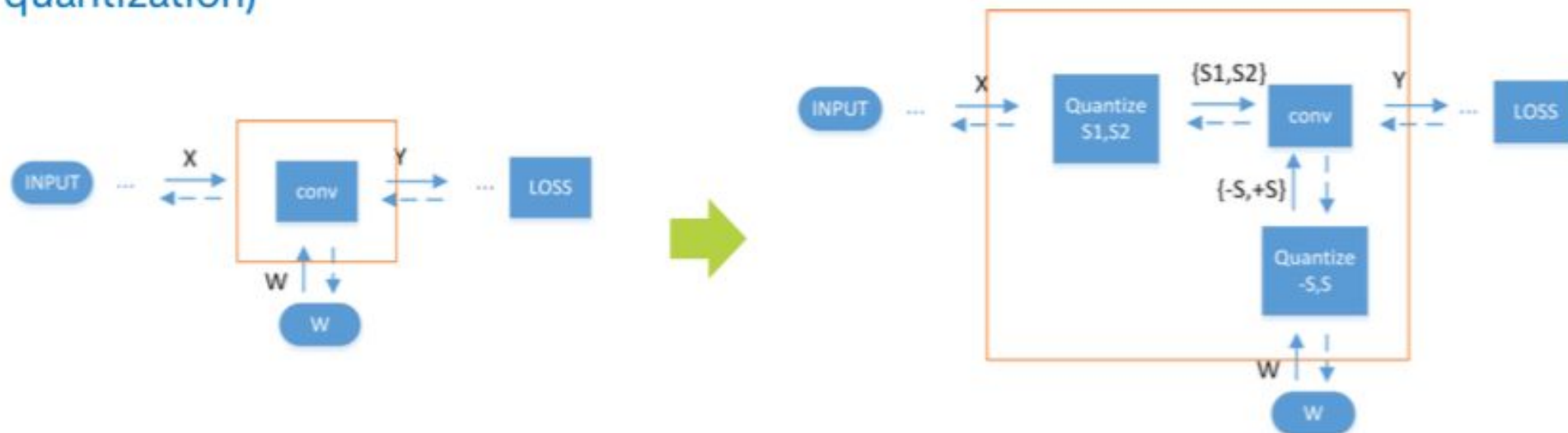
$$a \; XNOR \; b = \begin{array}{c|cc} a\backslash b & 0 & 1 \\ \hline 0 & 1 & 0 \\ 1 & 0 & 1 \end{array} \quad \{0, 1\},$$

$X = [x_1, \ldots 0, 1, 1 \ldots, x_{32}], W = [w_1, \ldots, w_{32}]$ -> int32 number

Calc $y = 2 * popcount(W \; XNOR \; X) - 32$

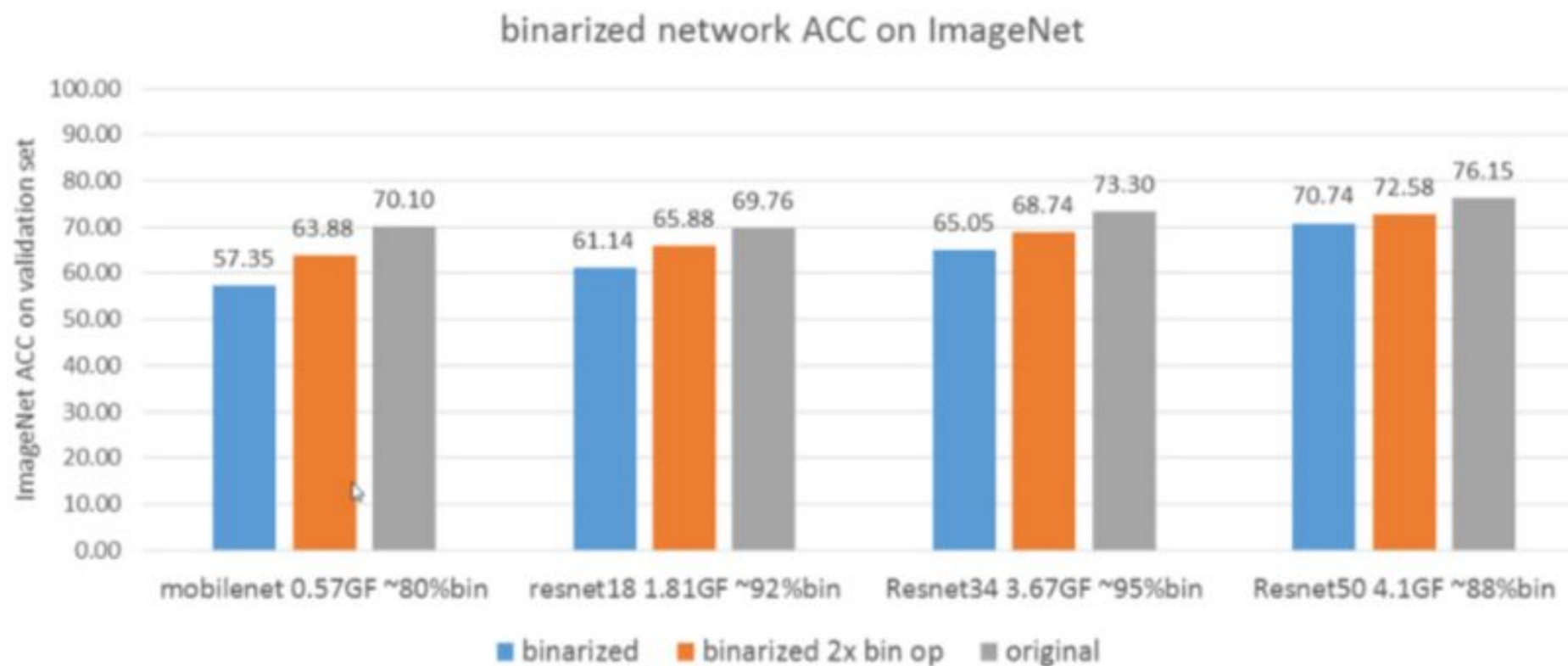**We have only1 XNOR and 1POPCOUNT operations**

# Binary Net

During binarization process selected convolutional layers of the original CNN are replaced with binary convolution alternatives.

Network is modified by inserting special quantization layer for input activations and weights that converts any full-precision value into two pretrained values (so-called "fake" quantization)

# Binary Net



binarized network ACC on ImageNet

# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
7. Winograd Transformation

# Distillation

- **L2 – Ba [14]**
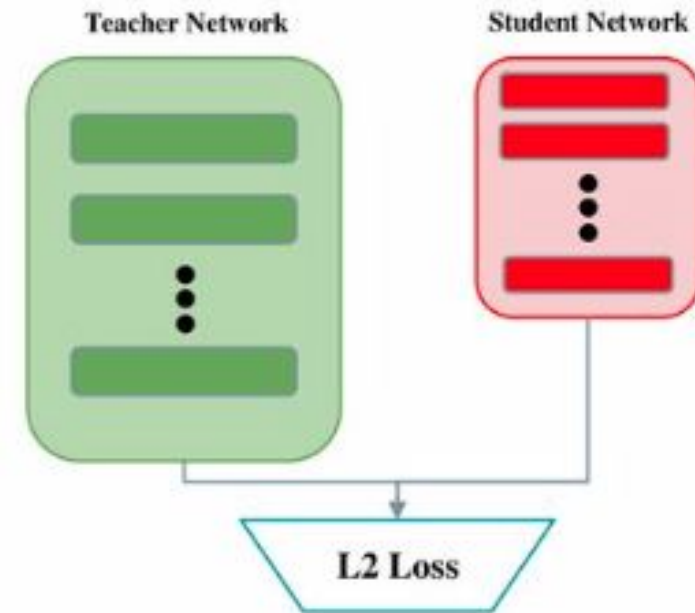  - L2 loss between teacher and student logits
  - No labels required



Figure 6. Teacher-Student model, L2

[14] J. Ba and R. Caruana, "Do deep nets really need to be deep?" In Advances in neural information processing systems, 2014, pp. 2654–2662.
[15] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," ArXiv preprint arXiv:1503.02531, 2015.
[16] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," ArXiv preprint arXiv:1412.6550, 2015.

# Distillation

- L2 – Ba [14]
  - L2 loss between teacher and student logits
  - No labels required
- **Knowledge Distillation [15]**
  - Soft target: softmax cross entropy with teacher logits
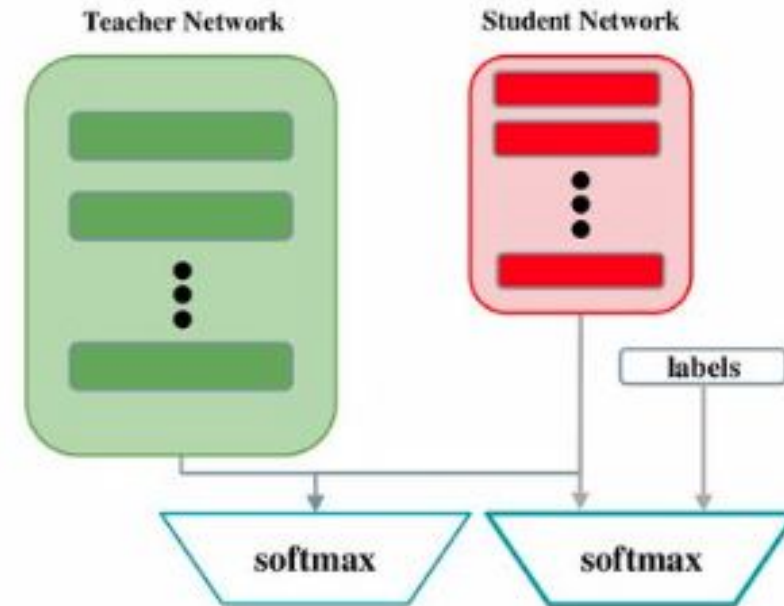  - Hard target: softmax cross entropy with correct labels



Figure 7. Teacher-Student model, KD

[14] J. Ba and R. Caruana, "Do deep nets really need to be deep?" In Advances in neural information processing systems, 2014, pp. 2654–2662.
[15] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," ArXiv preprint arXiv:1503.02531, 2015.
[16] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," ArXiv preprint arXiv:1412.6550, 2015.

# Distillation

- L2 – Ba [14]
  - L2 loss between teacher and student logits
  - No labels required
- Knowledge Distillation [15]
  - Soft target: softmax cross entropy with teacher logits
  - Hard target: softmax cross entropy with correct labels
- **FitNets [16]**
  - Knowledge Distillation with hints in the middle points of the network
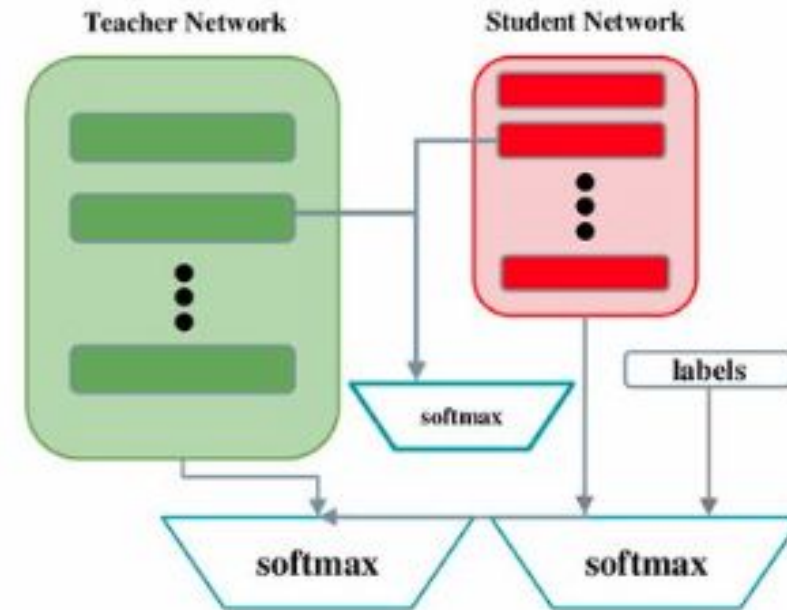  - Student is deeper than the teacher



Figure 8. Teacher-Student model, FitNets

[14] J. Ba and R. Caruana, "Do deep nets really need to be deep?" In Advances in neural information processing systems, 2014, pp. 2654–2662.

[15] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," ArXiv preprint arXiv:1503.02531, 2015.

[16] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," ArXiv preprint arXiv:1412.6550, 2015.

# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
7. Winograd Transformation

# Low Rank Approximation

- Basis filter set => Basis feature maps
- Final feature map = linear combination of basis feature maps
- Rank-1 basis filter => decomposed into a sequence of horizontal and vertical filters
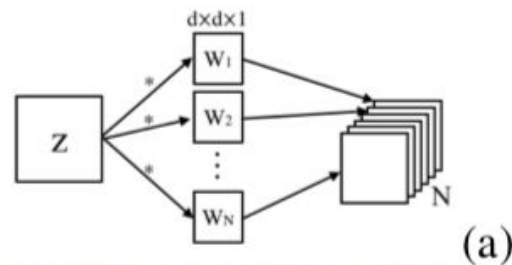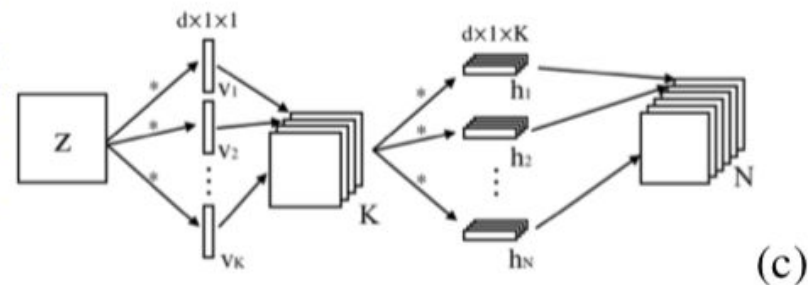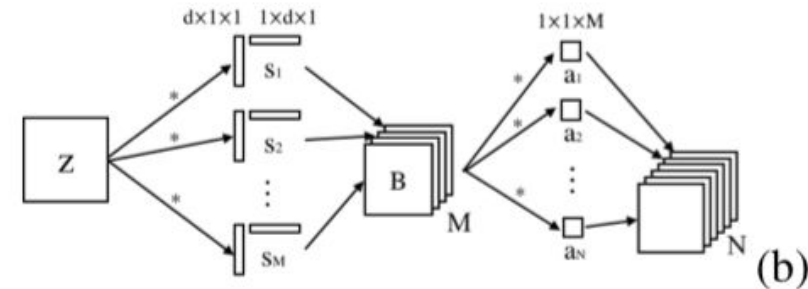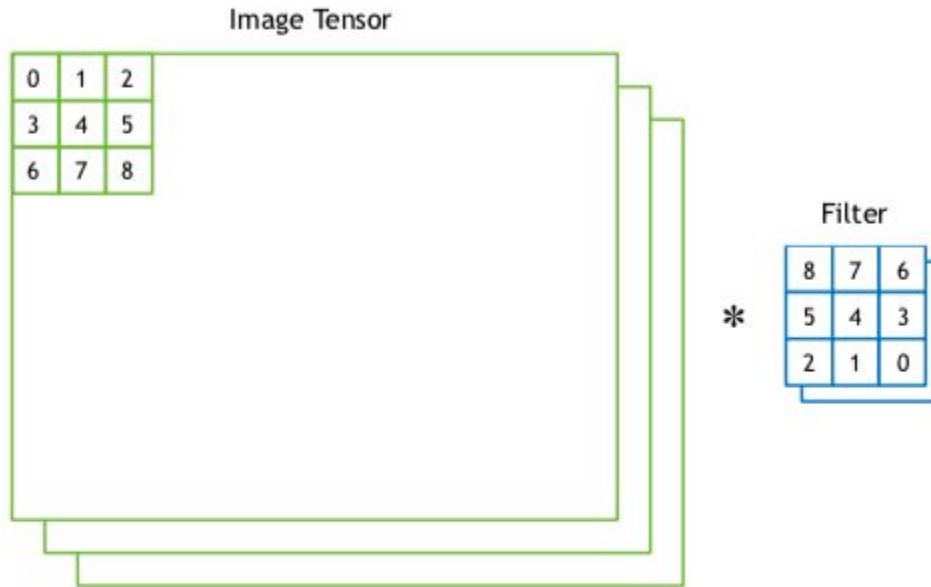- ~2.4x speedup, no performance drop



Figure 1: (a) The original convolutional layer acting on a single-channel input i.e. C=1. (b) The approximation to that layer using the method of Scheme 1. (c) The approximation to that layer using the method of Scheme 2. Individual filter dimensions are given above the filter layers.

# Algorithms for Efficient Inference

1. Pruning
2. Weight Sharing
3. Quantization
4. Binary / Ternary Net
5. Distillation
6. Low Rank Approximation
7. **Winograd Transformation**

# Winograd Transformation



Image Tensor

Filter

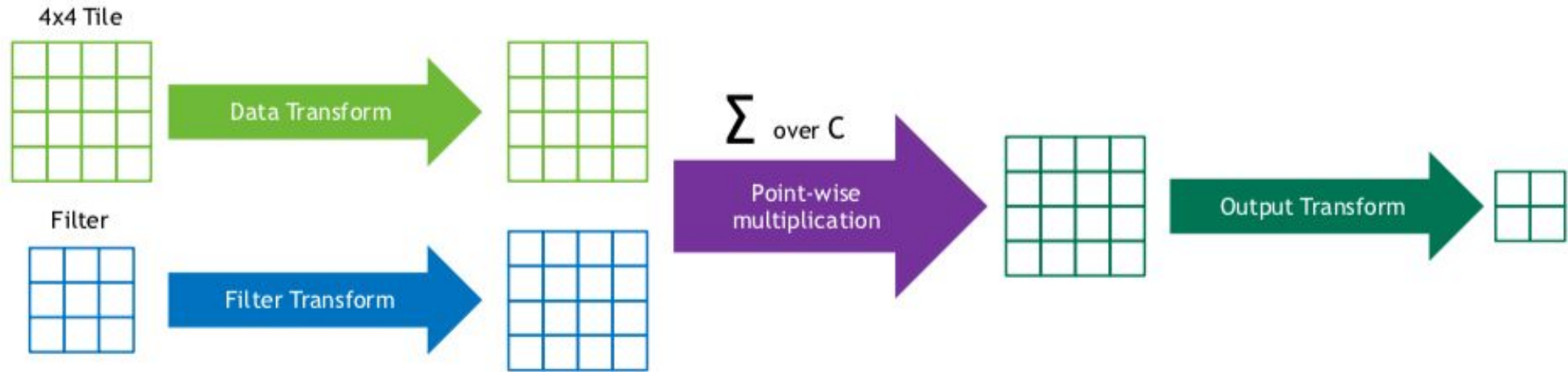9xC FMAs/Output: Math Intensive

9xK FMAs/Input: Good Data Reuse

Direct convolution: we need 9xCx4 = 36xC FMAs for 4 outputs

# Winograd Transformation



Direct convolution: we need 9xCx4 = 36xC FMAs for 4 outputs
Winograd convolution: we need 16xC FMAs for 4 outputs: **2.25x** fewer FMAs

# Summary

| Method | Advantages | Disadvantages |
|---|---|---|
| Binarization & Quantization | Low latency and memory usage | High loss of accuracy |
| Pruning | Prevents overfitting, the accuracy can increase | Slow |
| Factorization | Can achieve state of the art results while decreasing the computation cost | Dependent on framework |
| Distillation | Applicable to all architectures Doesn't change the network | only applicable to classification task |

# Thank you for your attention!