# Convolutional Neural Networks V

CS7150, Spring 2025

Prof. Huaizu Jiang

Northeastern University

# Recap

# Batch Normalization for ConvNets
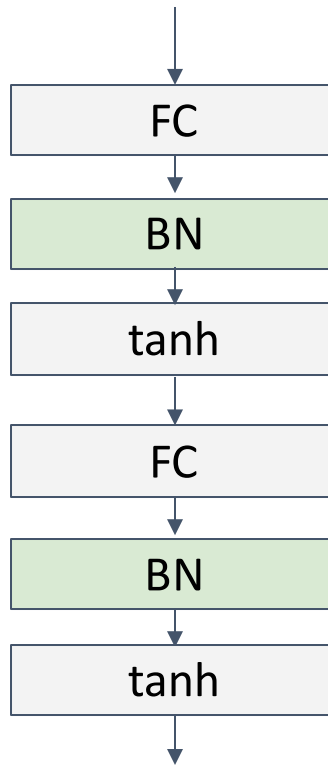
Batch Normalization for **fully-connected** networks

$$x : N \times D$$

Normalize $\downarrow$

$$\mu, \sigma : 1 \times D$$
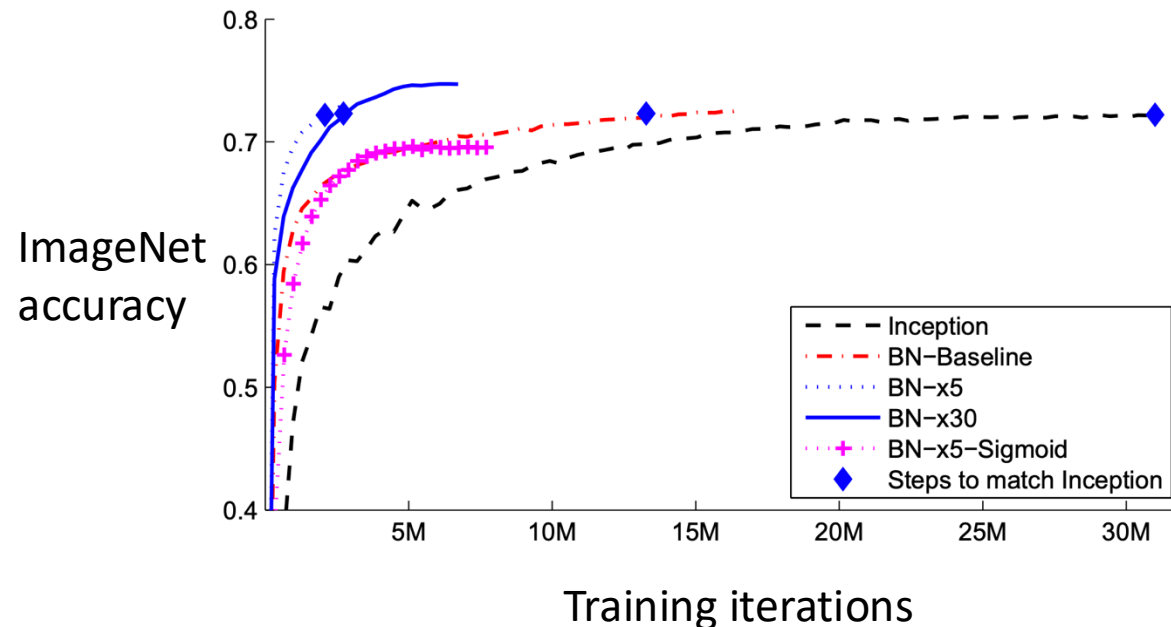$$\gamma, \beta : 1 \times D$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

Batch Normalization for **convolutional** networks
(Spatial Batchnorm, BatchNorm2D)

$$x : N \times C \times H \times W$$

Normalize $\downarrow \quad \downarrow \quad \downarrow$

$$\mu, \sigma$$
$$: 1 \times C \times 1 \times 1$$
$$\gamma, \beta$$
$$: 1 \times C \times 1 \times 1$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

# Batch Normalization
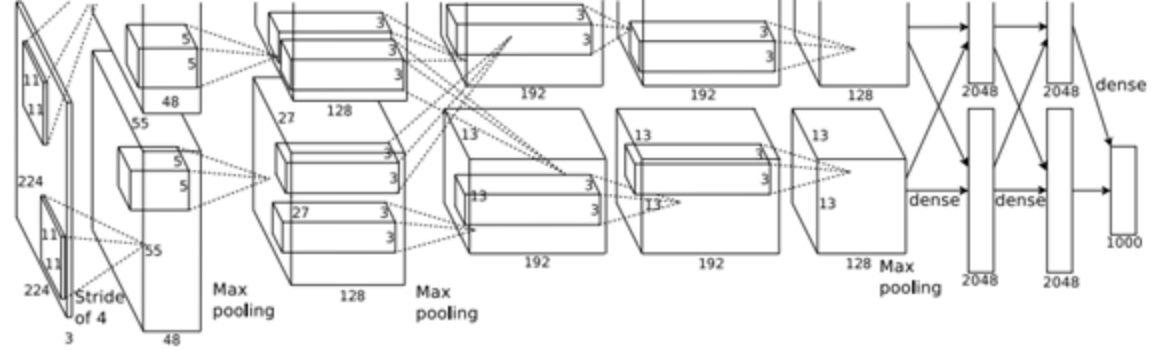
FC

BN

tanh

FC

BN

tanh

- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
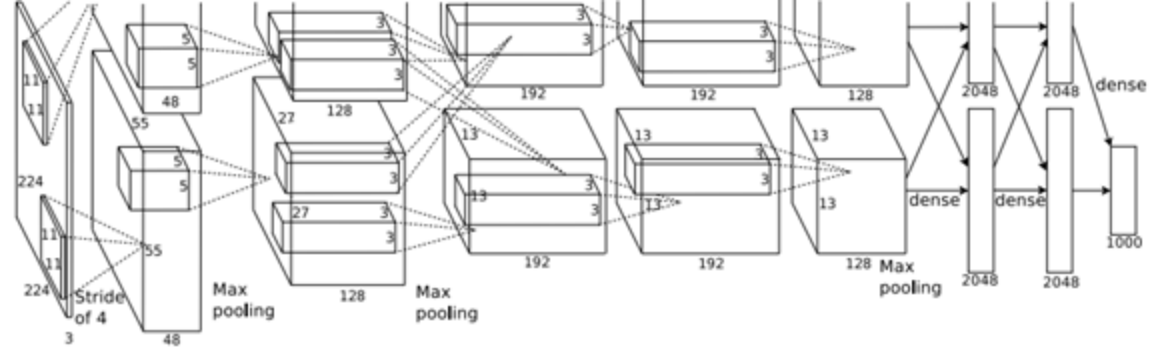- Free at test-time: can be fused with conv!



ImageNet accuracy vs Training iterations

Legend:
- Inception (black dashed)
- BN−Baseline (red dash-dot)
- BN−x5 (blue dotted)
- BN−x30 (blue solid)
- BN−x5−Sigmoid (magenta +)
- Steps to match Inception (blue diamond)

# AlexNet



How to choose this?
Trial and error =(

| Layer | Input size | | Layer | | | | Output size | | memory (KB) | params (k) | flop (M) |
|-------|-----------|-----|---------|--------|--------|-----|-------------|-----|-------------|------------|----------|
| | C | H / W | filters | kernel | stride | pad | C | H / W | | | |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 | 784 | 23 | 73 |
| pool1 | 64 | 56 | | 3 | 2 | 0 | 64 | 27 | 182 | 0 | 0 |
| conv2 | 64 | 27 | 192 | 5 | 1 | 2 | 192 | 27 | 547 | 307 | 224 |
| pool2 | 192 | 27 | | 3 | 2 | 0 | 192 | 13 | 127 | 0 | 0 |
| conv3 | 192 | 13 | 384 | 3 | 1 | 1 | 384 | 13 | 254 | 664 | 112 |
| conv4 | 384 | 13 | 256 | 3 | 1 | 1 | 256 | 13 | 169 | 885 | 145 |
| conv5 | 256 | 13 | 256 | 3 | 1 | 1 | 256 | 13 | 169 | 590 | 100 |
| pool5 | 256 | 13 | | 3 | 2 | 0 | 256 | 6 | 36 | 0 | 0 |
| flatten | 256 | 6 | | | | | 9216 | | 36 | 0 | 0 |
| fc6 | 9216 | | 4096 | | | | 4096 | | 16 | 37,749 | 38 |
| fc7 | 4096 | | 4096 | | | | 4096 | | 16 | 16,777 | 17 |
| fc8 | 4096 | | 1000 | | | | 1000 | | 4 | 4,096 | 4 |

# AlexNet



Interesting trends here!

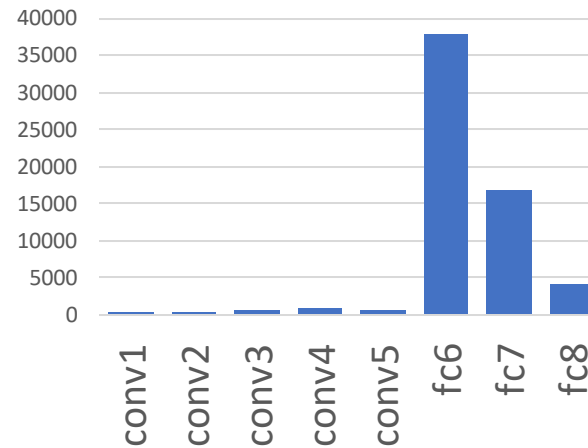| Layer | Input size | | Layer | | | | Output size | | Interesting trends | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Layer | C | H / W | filters | kernel | stride | pad | C | H / W | memory (KB) | params (k) | flop (M) |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 | 784 | 23 | 73 |
| pool1 | 64 | 56 | | 3 | 2 | 0 | 64 | 27 | 182 | 0 | 0 |
| conv2 | 64 | 27 | 192 | 5 | 1 | 2 | 192 | 27 | 547 | 307 | 224 |
| pool2 | 192 | 27 | | 3 | 2 | 0 | 192 | 13 | 127 | 0 | 0 |
| conv3 | 192 | 13 | 384 | 3 | 1 | 1 | 384 | 13 | 254 | 664 | 112 |
| conv4 | 384 | 13 | 256 | 3 | 1 | 1 | 256 | 13 | 169 | 885 | 145 |
| conv5 | 256 | 13 | 256 | 3 | 1 | 1 | 256 | 13 | 169 | 590 | 100 |
| pool5 | 256 | 13 | | 3 | 2 | 0 | 256 | 6 | 36 | 0 | 0 |
| flatten | 256 | 6 | | | | | 9216 | | 36 | 0 | 0 |
| fc6 | 9216 | | 4096 | | | | 4096 | | 16 | 37,749 | 38 |
| fc7 | 4096 | | 4096 | | | | 4096 | | 16 | 16,777 | 17 |
| fc8 | 4096 | | 1000 | | | | 1000 | | 4 | 4,096 | 4 |

# AlexNet



Most of the **memory usage** is in the early convolution layers
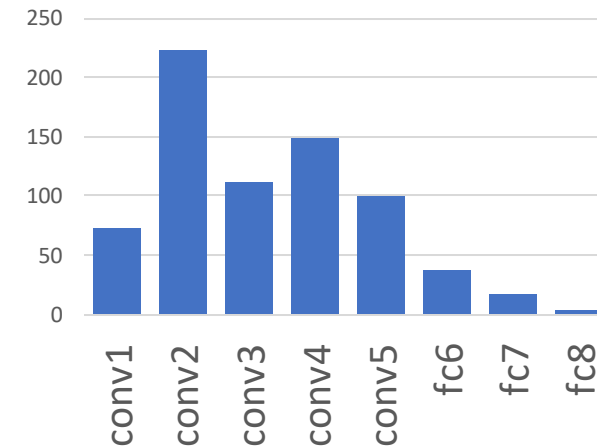
## Memory (KB)



Nearly all **parameters** are in the fully-connected layers

## Params (K)



Most **floating-point ops** occur in the convolution layers

## MFLOP

# VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1

**All max pool are 2x2 stride 2**

**After pool, double #channels**

Conv layers at each spatial resolution take the same amount of computation!

Input: C x 2H x 2W

Layer: Conv(3x3, C->C)

Memory: 4HWC

Params: $9C^2$

FLOPs: $36HWC^2$

Input: 2C x H x W

Conv(3x3, 2C -> 2C)

Memory: 2HWC

Params: $36C^2$

FLOPs: $36HWC^2$
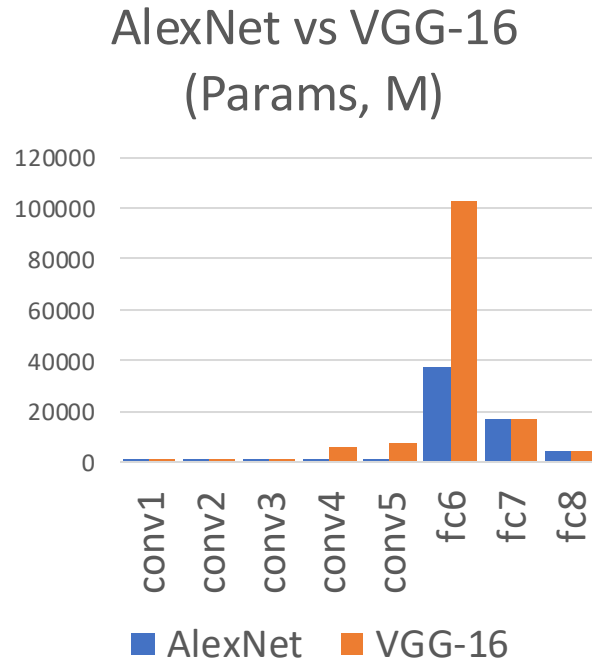
## AlexNet

| |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 384 |
| Pool |
| 3x3 conv, 384 |
| Pool |
| 5x5 conv, 256 |
| 11x11 conv, 96 |
| Input |

## VGG16

| |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

## VGG19

| |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

Slide credit: D Fouhey & J Johnson

# AlexNet vs VGG-16: Much Bigger!



AlexNet vs VGG-16 (Memory, KB)

AlexNet total: 1.9 MB
VGG-16 total: 48.6 MB (25x)

AlexNet vs VGG-16 (Params, M)

AlexNet total: 61M
VGG-16 total: 138M (2.3x)
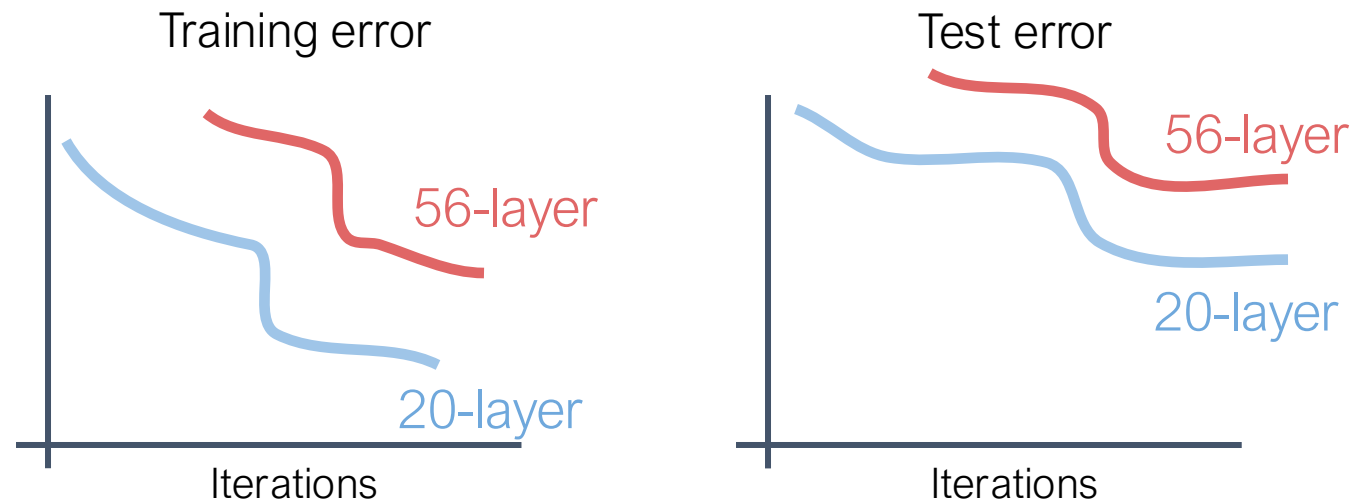
AlexNet vs VGG-16 (MFLOPs)

AlexNet total: 0.7 GFLOP
VGG-16 total: 13.6 GFLOP (19.4x)

# Residual Networks

Once we have Batch Normalization, we can train networks with 10+ layers. What happens as we go deeper?
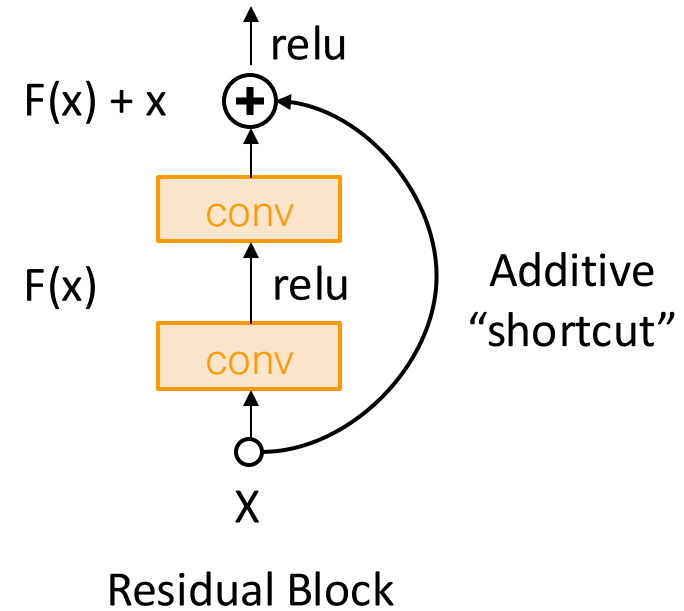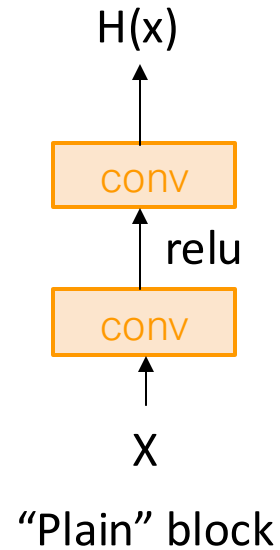


In fact the deep model seems to be **underfitting** since it also performs worse than the shallow model on the training set! It is actually **underfitting**

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016
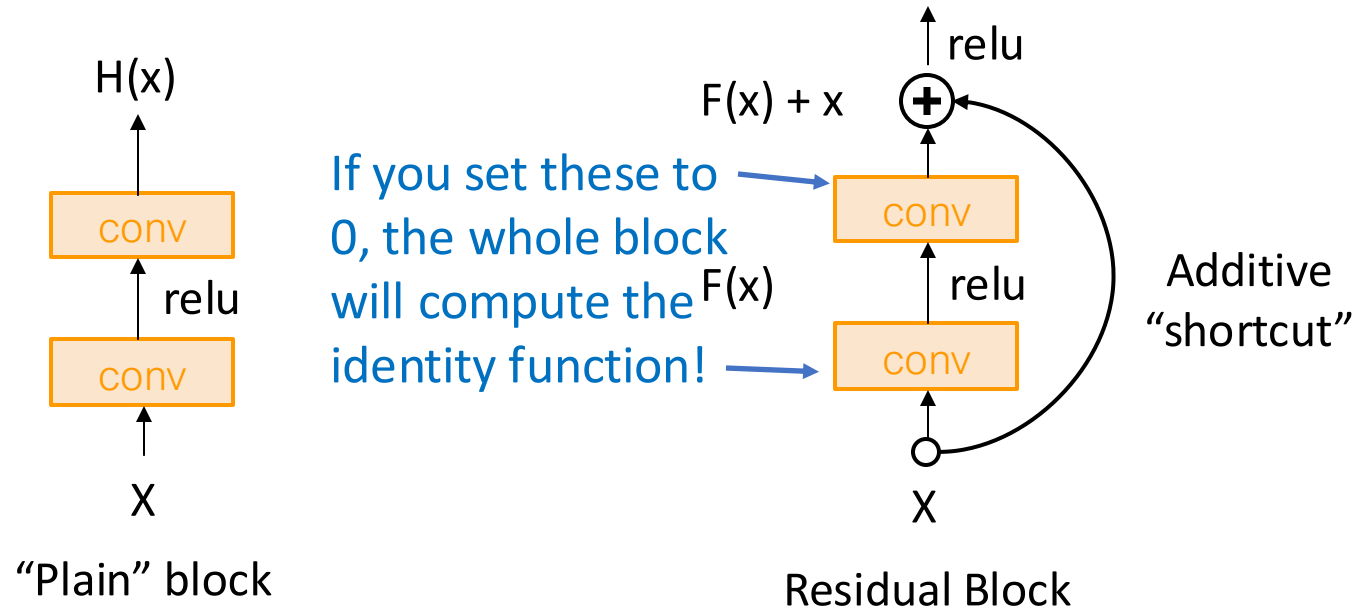
# Residual Networks

**Solution**: Change the network so learning identity functions with extra layers is easy!



H(x)

conv

relu

conv

X

"Plain" block

relu

F(x) + x

conv

F(x)    relu

conv

X

Additive "shortcut"

Residual Block

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016

# Residual Networks

**Solution**: Change the network so learning identity functions with extra layers is easy!



H(x)

conv

relu

conv

X

"Plain" block

If you set these to 0, the whole block will compute the identity function!

relu

F(x) + x ⊕

conv

F(x)

relu

conv

X

Additive "shortcut"

Residual Block

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016
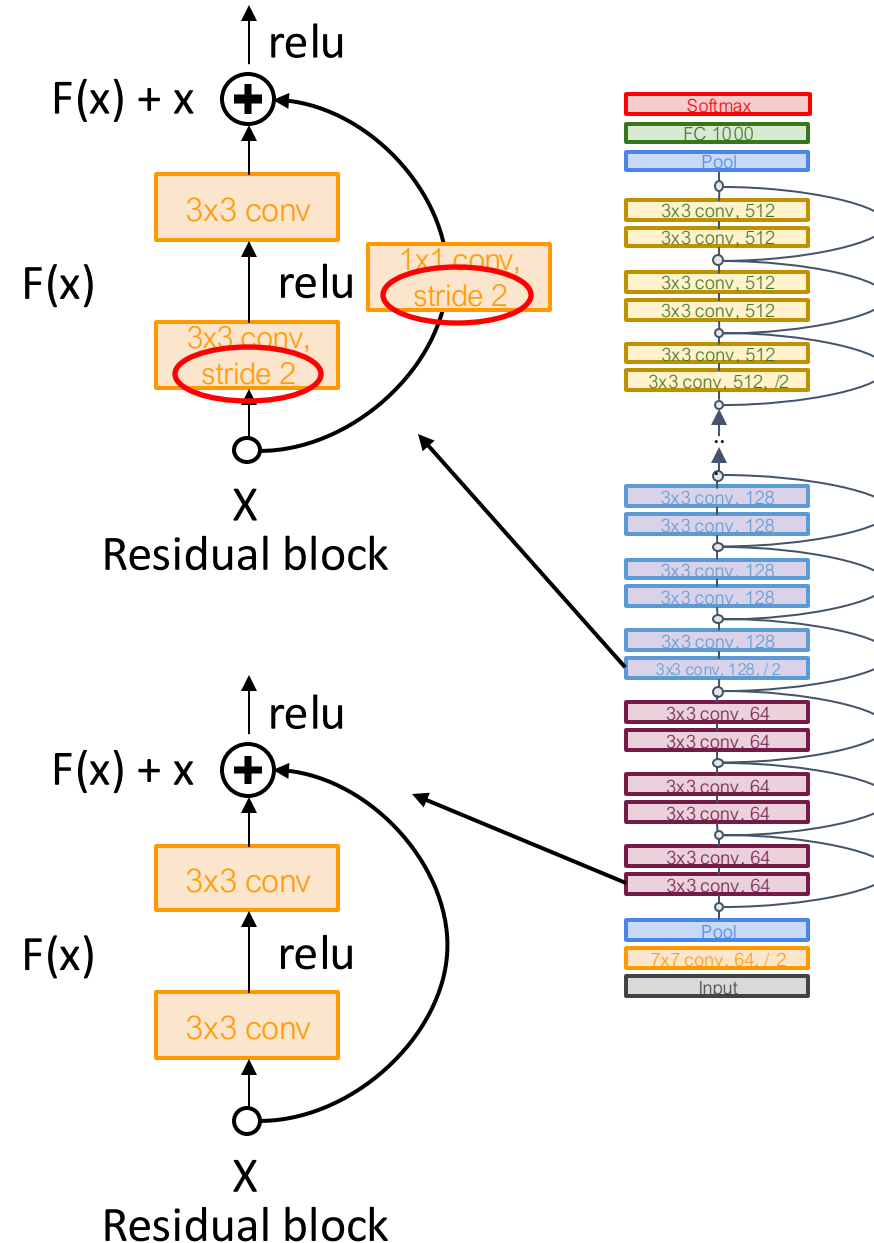
# Residual Networks

A residual network is a stack of many residual blocks

Regular design, like VGG: each residual block has two 3x3 conv

Network is divided into **stages**: the first block of each stage halves the resolution (with stride-2 conv) and doubles the number of channels

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016



Adapted from: D Fouhey & J Johnson

# Tiny Networks for Mobile Devices



**Object Detection**

Photo by Juanedc (CC BY 2.0)

**Finegrain Classification**

Photo by HarshLight (CC BY 2.0)

**MobileNets**

**Face Attributes**

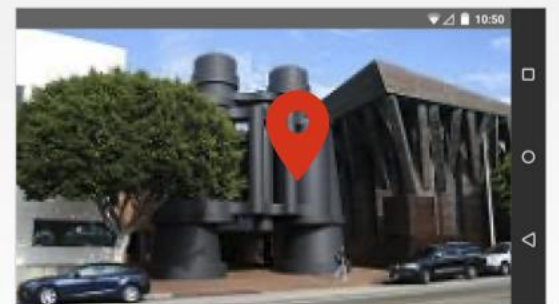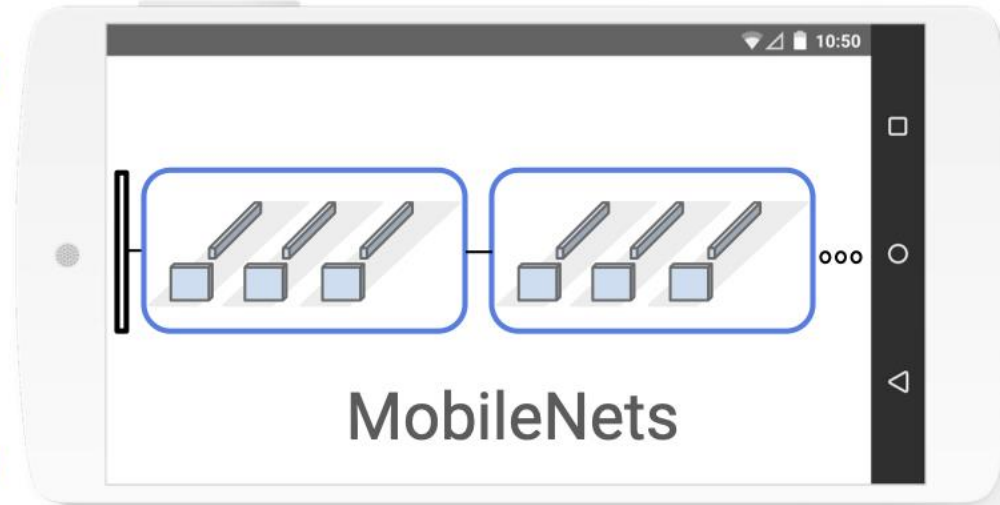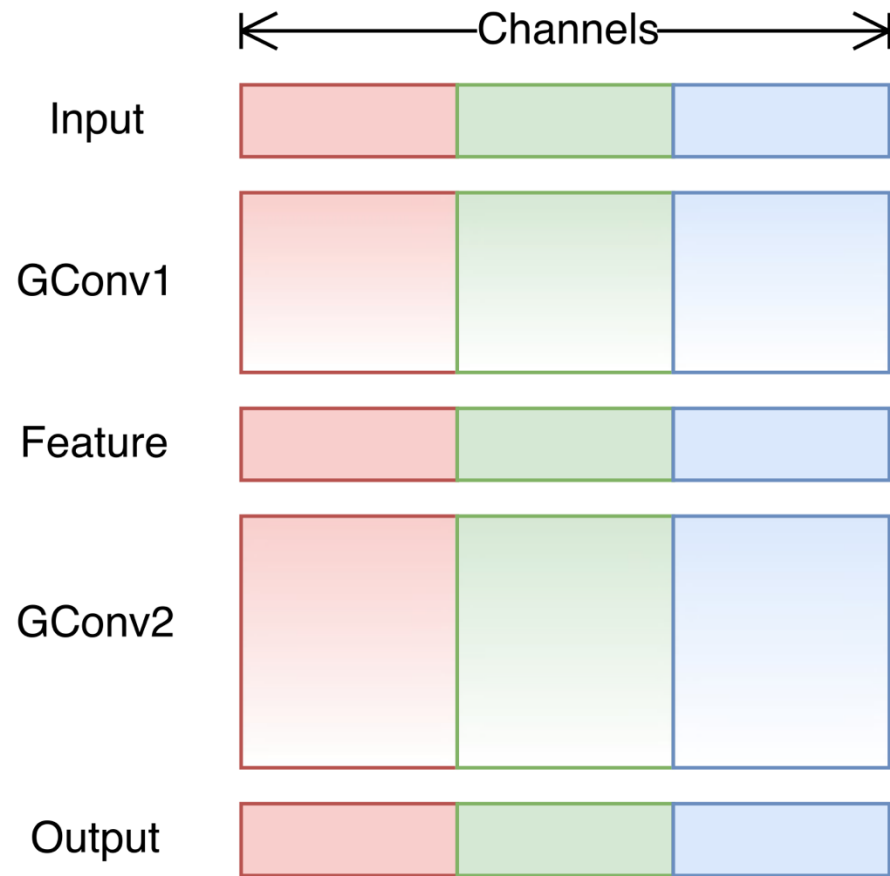Google Doodle by Sarah Harrison

**Landmark Recognition**

Photo by Sharon VanderKaay (CC BY 2.0)

[Howard et al., MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv 2017]

# Group-based Convolution
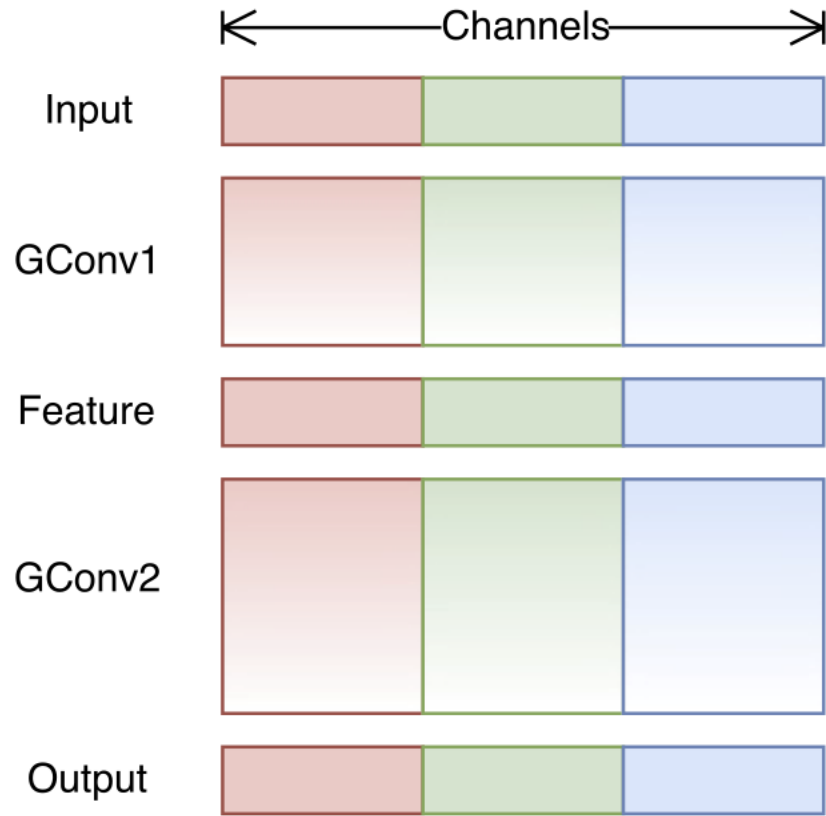


**Input**: $C_{in}$ x H x W
**Hyperparameters**:
- **Kernel size**: $K_H$ x $K_W$
- **Number filters**: $C_{out}$
- **Padding**: P
- **Stride**: S
- **Groups:** G

**Weight matrix**: $C_{out}$ /G x $C_{in}$ /G x $K_H$ x $K_W$ x G

**Bias vector**: $C_{out}$ /G

**FLOPS:** $C_{out}$ /G x $C_{in}$ /G x $K_H$ x $K_W$ x G x H x W

# ShuffleNet



[Zhang et al., ShuffleNet: An Extremely Efficient Convolutional Neural Network for Mobile Devices. CVPR 2018]

# ShuffleNet Units



[Zhang et al., ShuffleNet: An Extremely Efficient Convolutional Neural Network for Mobile Devices. CVPR 2018]
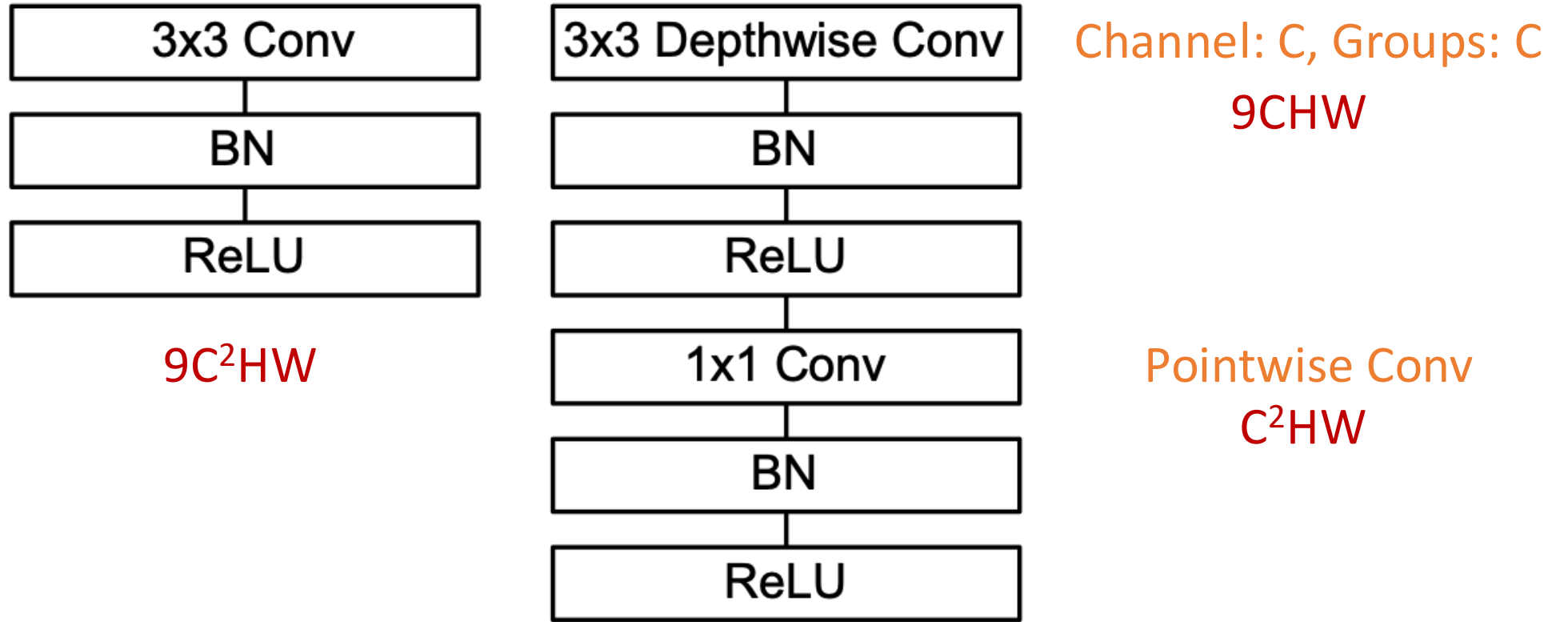
# Today's Class

- Lightweight convolutional neural networks
- Tips of training deep convolutional neural networks
- PyTorch tutorial

# MobileNet



| 3x3 Conv | 3x3 Depthwise Conv | Channel: C, Groups: C |
| BN | BN | 9CHW |
| ReLU | ReLU | |
| 9C$^2$HW | 1x1 Conv | Pointwise Conv |
| | BN | C$^2$HW |
| | ReLU | |

3x3 Conv → BN → ReLU

9C$^2$HW

3x3 Depthwise Conv → BN → ReLU → 1x1 Conv → BN → ReLU

Channel: C, Groups: C
9CHW

Pointwise Conv
C$^2$HW

Computation reduction: $9C^2HW/(9CHW + C^2HW) = 9C/(9+C)$

[Howard et al., MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv 2017]

# Training Convolutional Networks

1. Download big datasets
2. Design CNN architecture
3. Initialize Weights
4. For t = 1 to T:
    1. Form minibatch
    2. Compute loss + gradient
    3. Update Weights
5. Apply trained model to task

# Training Convolutional Networks

1. Download big datasets

2. Design CNN architecture

**3. Initialize Weights**

4. For t = 1 to T:

    1. Form minibatch

    2. Compute loss + gradient

    3. Update Weights

5. Apply trained model to task

# Training Convolutional Networks

1. Download big datasets
2. Design CNN architecture
3. Initialize Weights

4. For t = 1 to T:

    1. Form minibatch

    2. Compute loss + gradient

    3. Update Weights

5. Apply trained model to task

If the model is big, won't we overfit?

# Regularizing CNNs: Weight Decay

$$L_{reg} = \frac{1}{2} \sum_{\ell} \|W_\ell\|^2 \qquad \frac{\partial L_{reg}}{\partial W_\ell} = W_\ell$$

Add L2 regularization term $L_{reg}$ to the loss penalizing large weight matrices

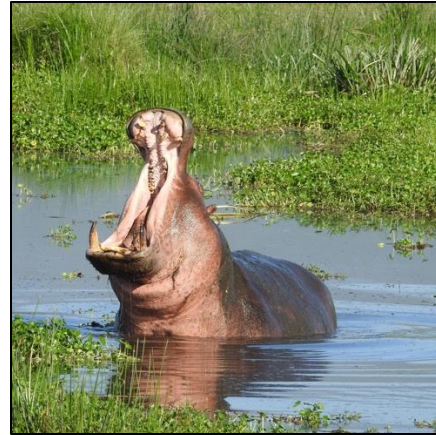Usually don't regularize bias terms, or BatchNorm scale / shift params

*Technical note: Adding an explicit term to the loss is "L2 Regularization"; "Weight decay" adds a term to the gradient. They are equivalent for SGD, but not quite the same for other optimizers like Adam

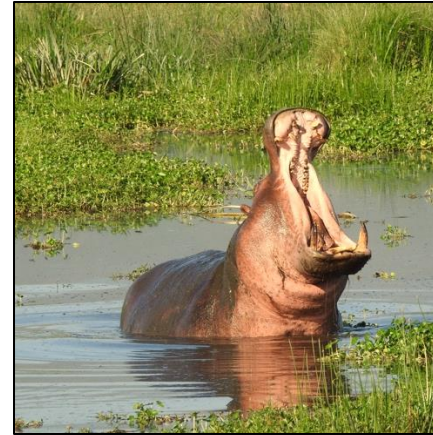# Regularizing CNNs: Data Augmentation

Hippo    Hippo?    Hippo?    Hippo?



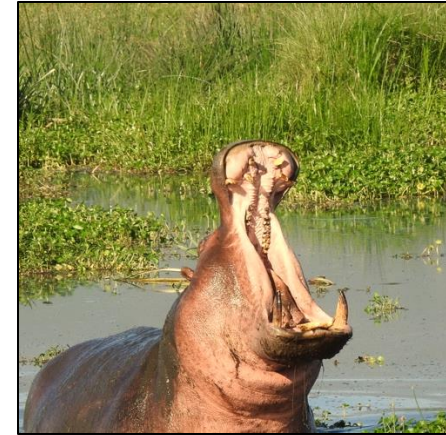Horizontal Flip    Color Jitter    Image Cropping

# Regularizing CNNs: Data Augmentation

Apply random transformations to input images during training
Artificially "inflate" the size of your dataset



Training
Image

Augmented
images

Hippo

Hippo

# Training Convolutional Networks

1. Download big datasets
2. Design CNN architecture
3. Initialize Weights
4. For t = 1 to T:
   1. Form minibatch
   2. Compute loss + gradient
   3. Update Weights
5. Apply trained model to task

If the model is big, won't we overfit?

# Training Convolutional Networks

1. **Download big datasets**

   <span style="color:darkred">What if we can't find one?</span>

2. Design CNN architecture
3. Initialize Weights
4. For t = 1 to T:
   1. Form minibatch
   2. Compute loss + gradient
   3. Update Weights
5. Apply trained model to task

# Transfer Learning: Feature Extraction

1. Train on ImageNet

| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
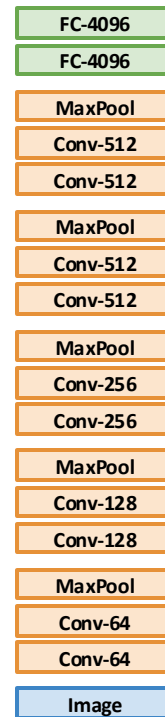
# Transfer Learning: Feature Extraction



1. Train on ImageNet

2. CNN as feature extractor

Remove last layer

Freeze these

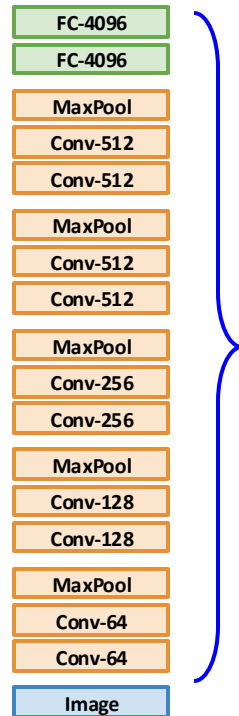Use your small dataset to train a **linear classifier** on top of pretrained CNN features

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014

Slide credit: D Fouhey & J Johnson

# Transfer Learning: Fine-Tuning

**1. Train on ImageNet**

| |
|---|
| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

**2. CNN as feature extractor**

| |
|---|
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

Remove last layer

Freeze these

**3. Bigger dataset: Fine-Tuning**

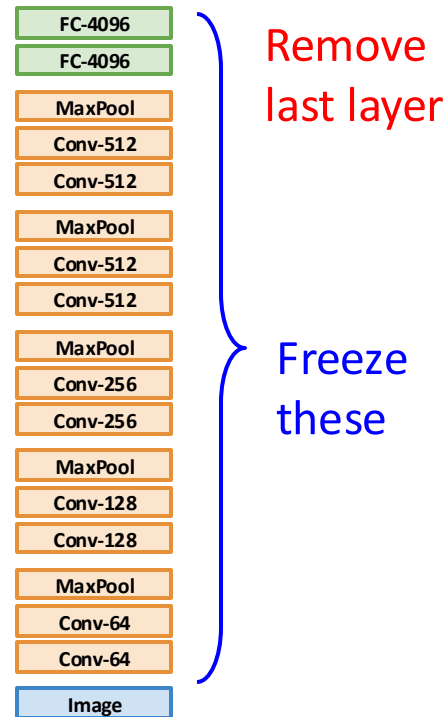| |
|---|
| FC |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

Reinitialize last layer and continue training whole network on your dataset

# Transfer Learning: Fine-Tuning

## 1. Train on ImageNet

| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

## 2. CNN as feature extractor

| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

Remove last layer

Freeze these

## 3. Bigger dataset: **Fine-Tuning**

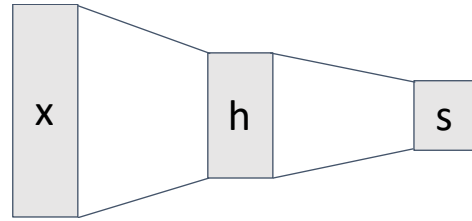| FC |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

Reinitialize last layer and continue training whole network on your dataset

Some tricks:
- Train with feature extraction first before fine-tuning
- Lower the learning rate: use ~1/10 of LR used in original training
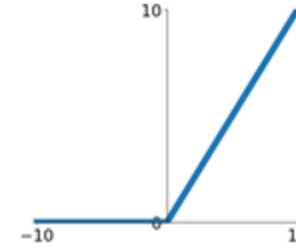- Sometimes freeze lower layers to save computation

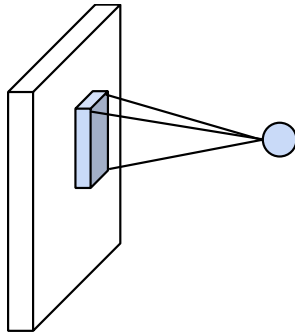# Recap: Convolutional Networks

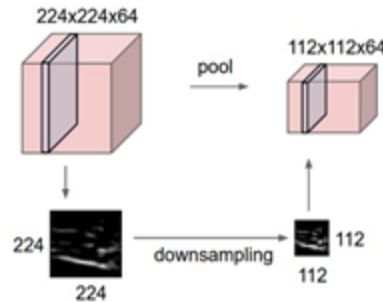### Fully-Connected Layers



$$y = Wx + b$$

### Activation Function
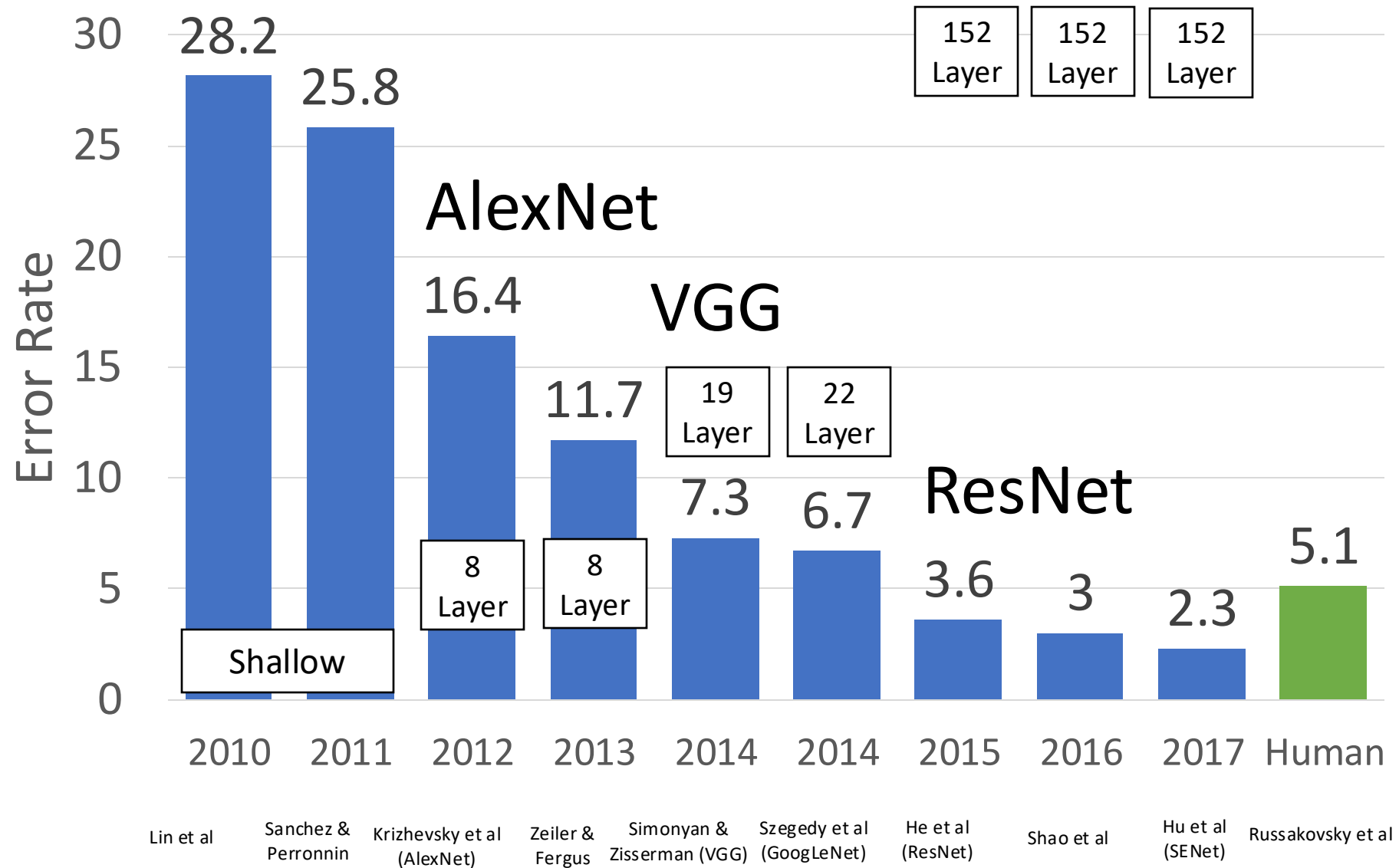


$$y = \max(0, x)$$

### Convolution Layers



### Pooling Layers



### Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Recap: CNN Architectures



Slide credit: D Fouhey & J Johnson

# Recap: Training CNNs

1. Download big datasets  *Transfer Learning*
2. Design CNN architecture
3. Initialize Weights  *Xavier / MSRA Init*
4. For t = 1 to T:
   1. Form minibatch
   2. Compute loss + gradient
   3. Update Weights
5. Apply trained model to task

*Regularization + Data Augmentation*

# Next Class

## More about
## Convolutional Neural Networks