# Lecture 16:
# Convolutional Networks II

# Administrative

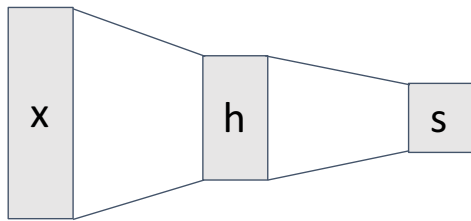HW4 Released, due Monday March 29, 11:59pm ET

Course Project:

- We will give ~6 suggested project descriptions

- Choose one, or propose your own

- We expect ~1 HW of work per person for project
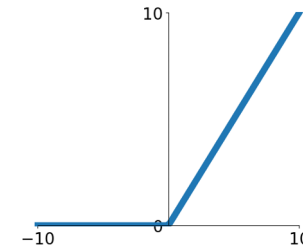
# Last Time: Convolutional Networks
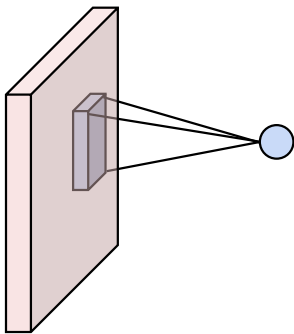
## Fully-Connected Layers
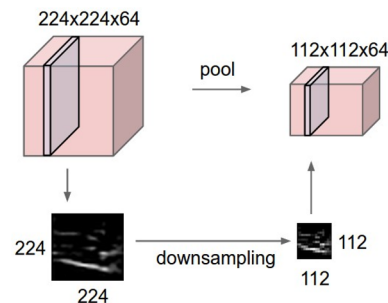


$$y = Wx + b$$

## Activation Function



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

## Convolution Layers



## Pooling Layers



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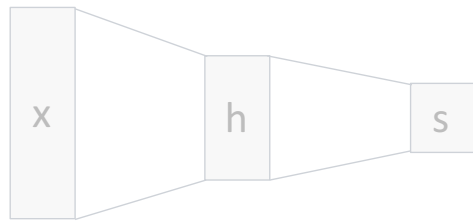
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Components of a Convolutional Network

Fully-Connected Layers

$$y = Wx + b$$

Activation Function

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

Convolution Layers

Pooling Layers

224x224x64

112x112x64

pool

224

downsampling

112

224

112

Normalization

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

# Batch Normalization

**Idea**: "Normalize" the outputs of each layer so they have zero mean and unit variance

Why? Helps reduce "internal covariate shift", improves optimization

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

# Batch Normalization

**Idea**: "Normalize" the outputs of each layer so they have zero mean and unit variance

Why? Helps reduce "internal covariate shift", improves optimization

We can normalize a batch of activations like this:

$$\hat{x} = \frac{x - E[x]}{\sqrt{Var[x]}}$$

This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

# Batch Normalization

**Input:** $x \in \mathbb{R}^{N \times D}$

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} \left(x_{i,j} - \mu_j\right)^2$$

Per-channel std, shape is D

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

Normalized x, Shape is N x D

N

X

D

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

# Batch Normalization

**Input:** $x \in \mathbb{R}^{N \times D}$

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

Per-channel mean, shape is D



N

X

D

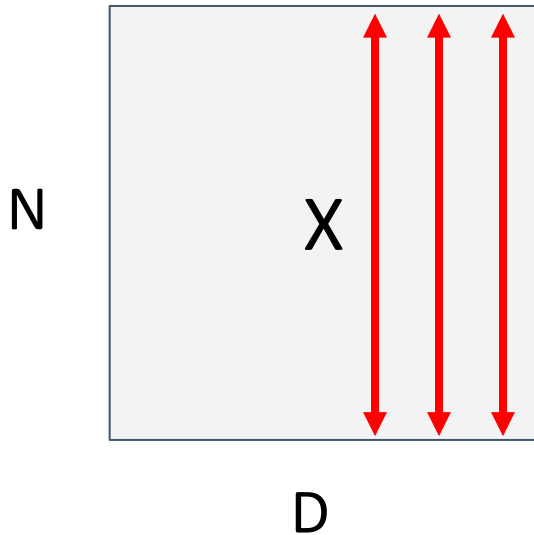$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} \left( x_{i,j} - \mu_j \right)^2$$

Per-channel std, shape is D

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

Normalized x, Shape is N x D

**Problem**: What if zero-mean, unit variance is too restrictive?

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

# Batch Normalization

**Input**: $x \in \mathbb{R}^{N \times D}$

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

Per-channel mean, shape is D

**Learnable scale and shift parameters:**

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} \left( x_{i,j} - \mu_j \right)^2$$

Per-channel std, shape is D

$$\gamma, \beta \in \mathbb{R}^D$$

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

Normalized x, Shape is N x D

Learning $\gamma = \sigma, \beta = \mu$ will recover the identity function (in expectation)

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output, Shape is N x D

# Batch Normalization

**Problem:** Estimates depend on minibatch; can't do this at test-time!

**Input**: $x \in \mathbb{R}^{N \times D}$

**Learnable scale and shift parameters:**

$$\gamma, \beta \in \mathbb{R}^D$$

Learning $\gamma = \sigma, \beta = \mu$ will recover the identity function (in expectation)

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} \left( x_{i,j} - \mu_j \right)^2$$

Per-channel std, shape is D

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

Normalized x, Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output, Shape is N x D

# Batch Normalization: Test-Time

**Input**: $x \in \mathbb{R}^{N \times D}$

$\mu_j = $ (Running) average of values seen during training

Per-channel mean, shape is D

**Learnable scale and shift parameters:**

$\gamma, \beta \in \mathbb{R}^{D}$

$\sigma_j^2 = $ (Running) average of values seen during training

Per-channel std, shape is D

Learning $\gamma = \sigma, \beta = \mu$ will recover the identity function (in expectation)

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

Normalized x, Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output, Shape is N x D

# Batch Normalization: Test-Time

**Input**: $x \in \mathbb{R}^{N \times D}$

$\mu_j =$ (Running) average of values seen during training — Per-channel mean, shape is D

**Learnable scale and shift parameters:**

$\sigma_j^2 =$ (Running) average of values seen during training — Per-channel std, shape is D

$$\gamma, \beta \in \mathbb{R}^D$$

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

Normalized x, Shape is N x D

During testing batchnorm becomes a linear operator!
Can be fused with the previous fully-connected or conv layer

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output, Shape is N x D

# Batch Normalization for ConvNets

Batch Normalization for
**fully-connected** networks

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

$$x : N \times D$$

Normalize

$$\mu, \sigma : 1 \times D$$
$$\gamma, \beta : 1 \times D$$
$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

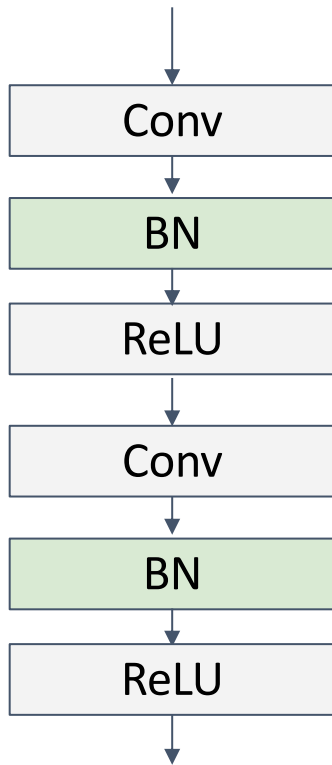$$x : N \times C \times H \times W$$

Normalize

$$\mu, \sigma : 1 \times C \times 1 \times 1$$
$$\gamma, \beta : 1 \times C \times 1 \times 1$$
$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

# Batch Normalization

Conv

BN

ReLU

Conv

BN

ReLU

Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x} = \frac{x - E[x]}{\sqrt{Var[x]}}$$

# Batch Normalization

Conv

BN

ReLU

Conv

BN

ReLU

- 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

Training iterations

# Batch Normalization

```
  ↓
┌─────────────┐
│    Conv     │
└─────────────┘
  ↓
┌─────────────┐
│     BN      │
└─────────────┘
  ↓
┌─────────────┐
│    ReLU     │
└─────────────┘
  ↓
┌─────────────┐
│    Conv     │
└─────────────┘
  ↓
┌─────────────┐
│     BN      │
└─────────────┘
  ↓
┌─────────────┐
│    ReLU     │
└─────────────┘
  ↓
```

- 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!
- Not well-understood theoretically
- Behaves differently during training and testing: this is a very common source of bugs!

# Convolutional Networks

### Fully-Connected Layers

$$y = Wx + b$$

### Activation Function

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

### Convolution Layers

### Pooling Layers

224x224x64

pool

112x112x64

224

224

downsampling

112

112

### Normalization

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

# Convolutional Networks

Fully-Connected Layers

Activation Function

Convolutio

tion

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

112x112x64

pool

224   112
downsampling   112

224

How can we combine these components into full architectures?

# ImageNet Classification Challenge

# AlexNet



227 x 227 inputs

5 Convolutional layers

Max pooling

3 fully-connected layers

ReLU nonlinearities

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# AlexNet



227 x 227 inputs
5 Convolutional layers
Max pooling
3 fully-connected layers
ReLU nonlinearities

Used "Local response normalization";
Not used anymore

Trained on two GTX 580 GPUs – only
3GB of memory each! Model split
over two GPUs

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# AlexNet



Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

# AlexNet



| | Input size | | Layer | | | | Output size | |
| Layer | C | H / W | filters | kernel | stride | pad | C | H / W |
|---|---|---|---|---|---|---|---|---|
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | ? | |

# AlexNet



| | Input size | | Layer | | | | Output size | | |
|---|---|---|---|---|---|---|---|---|---|
| Layer | C | H / W | filters | kernel | stride | pad | C | H / W | |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | ? | |

Recall: Output channels = number of filters

# AlexNet



| | Input size | | Layer | | | | Output size | |
|---|---|---|---|---|---|---|---|---|
| **Layer** | C | H / W | filters | kernel | stride | pad | C | H / W |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 |

Recall: W' = (W − K + 2P) / S + 1
= 227 − 11 + 2*2) / 4 + 1
= 220/4 + 1 = 56

# AlexNet



| Layer | Input size | | Layer | | | | Output size | | memory (KB) |
|---|---|---|---|---|---|---|---|---|---|
| | C | H / W | filters | kernel | stride | pad | C | H / W | |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 | ? |

# AlexNet



| | Input size | | Layer | | | | Output size | | |
|---|---|---|---|---|---|---|---|---|---|
| **Layer** | C | H / W | filters | kernel | stride | pad | C | H / W | memory (KB) |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 | 784 |

Number of output elements = C * H' * W'
$$= 64*56*56 = 200{,}704$$

Bytes per element = 4 (for 32-bit floating point)

KB = (number of elements) * (bytes per elem) / 1024
    = 200704 * 4 / 1024
    = **784**

# AlexNet



| Layer | Input size | | Layer | | | | Output size | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | C | H / W | filters | kernel | stride | pad | C | H / W | memory (KB) | params (k) |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 | 784 | ? |

# AlexNet

| | Input size | | | Layer | | | | Output size | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Layer | C | H / W | filters | kernel | stride | pad | C | H / W | memory (KB) | params (k) |
| conv1 | 3 | 227 | 64 | 11 | 4 | 2 | 64 | 56 | 784 | 23 |

Weight shape = $C_{out}$ x $C_{in}$ x K x K

= 64 x 3 x 11 x 11

Bias shape = $C_{out}$ = 64

Number of weights = 64*3*11*11 + 64

= **23,296**

# AlexNet



| Layer | Input size | | Layer | | | | Output size | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 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 | ? |

# AlexNet



| | Input size | | | Layer | | | | Output size | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 |

Number of floating point operations (multiply+add)
= (number of output elements) * (ops per output elem)
= ($C_{out}$ x H' x W') * ($C_{in}$ x K x K)
= (64 * 56 * 56) * (3 * 11 * 11)
= 200,704 * 363
= **72,855,552**

# AlexNet



| Layer | Input size | | Layer | | | | Output size | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 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 | ? | | | | |

# AlexNet



| | Input size | | | Layer | | | | Output size | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | | | | |

For pooling layer:

#output channels = #input channels = 64

W' = floor((W − K) / S + 1)
= floor(53 / 2 + 1) = floor(27.5) = **27**

# AlexNet



| | Input size | | | Layer | | | | Output size | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **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 | ? | |

#output elems = $C_{out}$ x H' x W'
Bytes per elem = 4
KB = $C_{out}$ * H' * W' * 4 / 1024
    = 64 * 27 * 27 * 4 / 1024
    = **182.25**

# AlexNet



| Layer | Input size | | | Layer | | | | Output size | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 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 | ? |

Pooling layers have no learnable parameters!

# AlexNet



| | Input size | | | Layer | | | | Output size | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 |

Floating-point ops for pooling layer
= (number of output positions) * (flops per output position)
= $(C_{out} * H' * W') * (K * K)$
= (64 * 27 * 27) * (3 * 3)
= 419,904
= **0.4 MFLOP**

# AlexNet



| Layer | Input size | | Layer | | | | Output size | | memory (KB) | params (k) | flop (M) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **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 |

Flatten output size = $C_{in}$ x H x W
= 256 * 6 * 6
= **9216**

# AlexNet



| 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 |

FC params = $C_{in} * C_{out} + C_{out}$

= 9216 * 4096 + 4096

= 37,725,832

FC flops = $C_{in} * C_{out}$

= 9216 * 4096

= 37,748,736

# AlexNet



| Layer | Input size | | Layer | | | | Output size | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **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



How to choose this?
Trial and error =(

| Layer | Input size | | Layer | | | | Output size | | memory (KB) | params (k) | flop (M) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **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



Interesting trends here!

| | Input size | | | Layer | | | | Output size | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **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

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

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

### Memory (KB)



### Params (K)



### MFLOP

# ImageNet Classification Challenge

# ImageNet Classification Challenge

# 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



AlexNet    VGG16    VGG19

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

# 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

Network has 5 convolutional **stages**:

Stage 1: conv-conv-pool

Stage 2: conv-conv-pool

Stage 3: conv-conv-pool

Stage 4: conv-conv-conv-[conv]-pool

Stage 5: conv-conv-conv-[conv]-pool

(VGG-19 has 4 conv in stages 4 and 5)



AlexNet      VGG16      VGG19

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

# 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

Option 1:

Conv(5x5, C -> C)

Params: $25C^2$

FLOPs: $25C^2HW$

**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

# 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

Option 1:

Conv(5x5, C -> C)

Option 2:

Conv(3x3, C -> C)

Conv(3x3, C -> C)

Params: $25C^2$

FLOPs: $25C^2HW$

Params: $18C^2$

FLOPs: $18C^2HW$

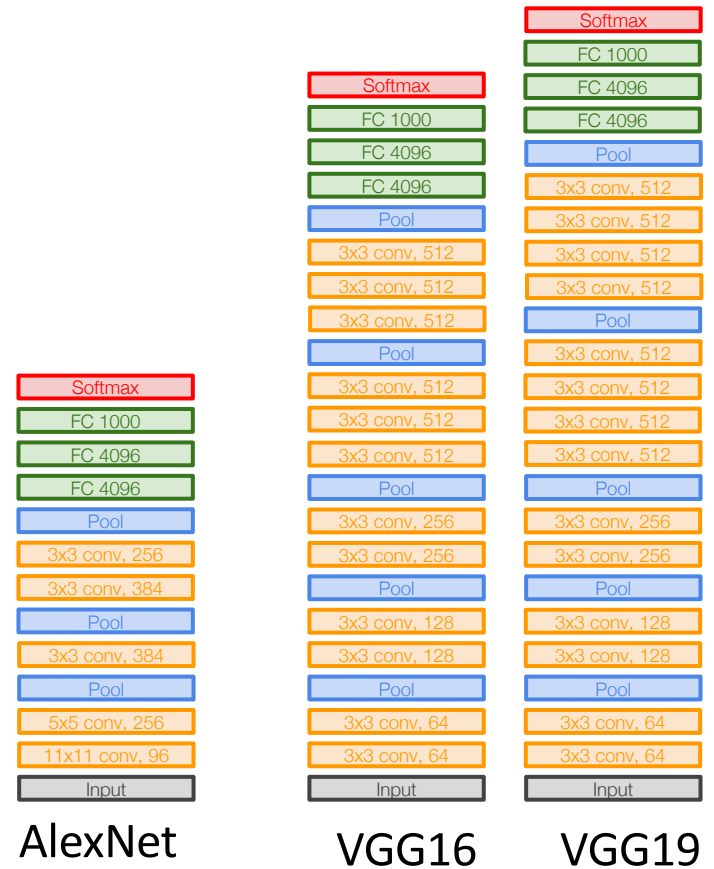| AlexNet |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 384 |
| Pool |
| 3x3 conv, 384 |
| Pool |
| 5x5 conv, 256 |
| 11x11 conv, 96 |
| Input |

| 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

# VGG: Deeper Networks, Regular Design

**Two 3x3 conv has same receptive field as a single 5x5 conv, but has fewer parameters and takes less computation!**

VGG Design rules:
**All conv are 3x3 stride 1 pad 1**
All max pool are 2x2 stride 2
After pool, double #channels

Option 1:
Conv(5x5, C -> C)

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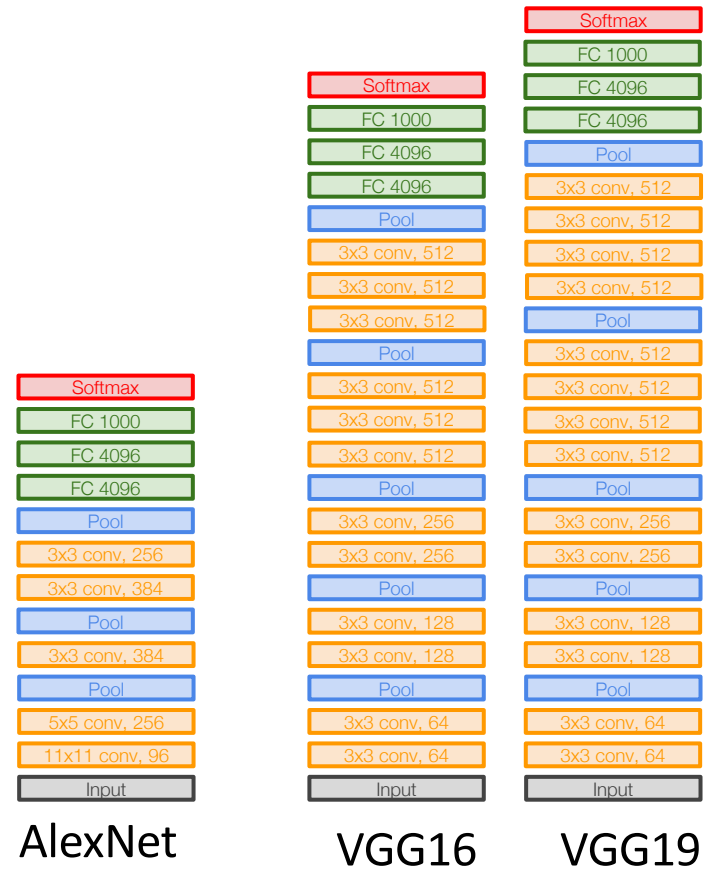Conv(3x3, C -> C)

Params: $25C^2$
FLOPs: $25C^2HW$

Params: $18C^2$
FLOPs: $18C^2HW$



AlexNet

VGG16

VGG19

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

# 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**

Input: C x 2H x 2W

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

Memory: 4HWC

Params: $9C^2$

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

# VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1

**All max pool are 2x2 stride 2**
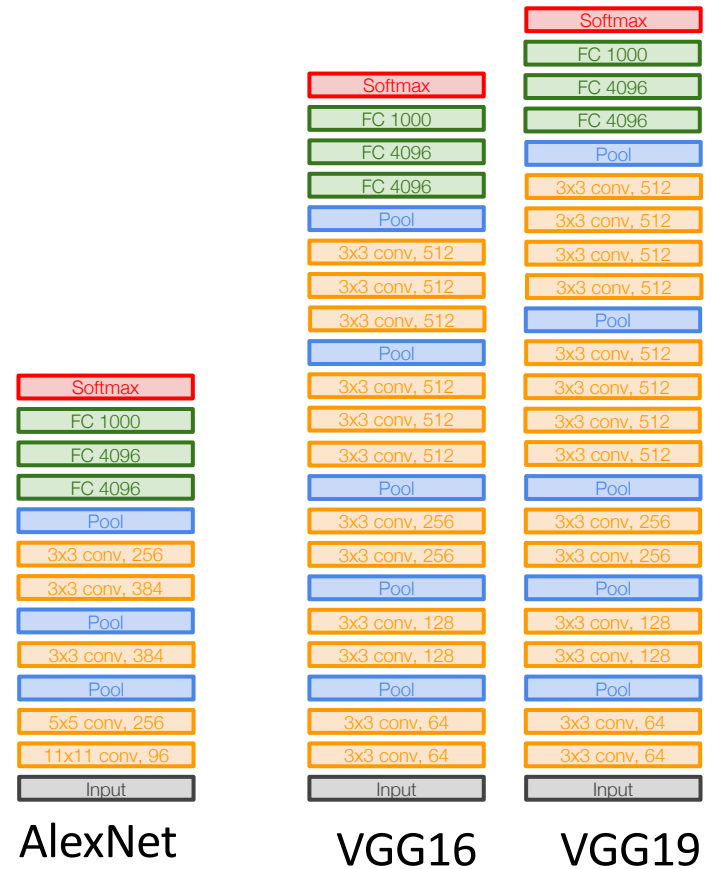
**After pool, double #channels**

| | |
|---|---|
| Input: C x 2H x 2W | Input: 2C x H x W |
| Layer: Conv(3x3, C->C) | Conv(3x3, 2C -> 2C) |
| | |
| Memory: 4HWC | Memory: 2HWC |
| Params: $9C^2$ | Params: $36C^2$ |
| FLOPs: $36HWC^2$ | FLOPs: $36HWC^2$ |



AlexNet          VGG16          VGG19

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

# 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)

Input: 2C x H x W

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

Memory: 4HWC

Params: $9C^2$

FLOPs: $36HWC^2$

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 |

# 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)

# ImageNet Classification Challenge

# ImageNet Classification Challenge

# Residual Networks

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

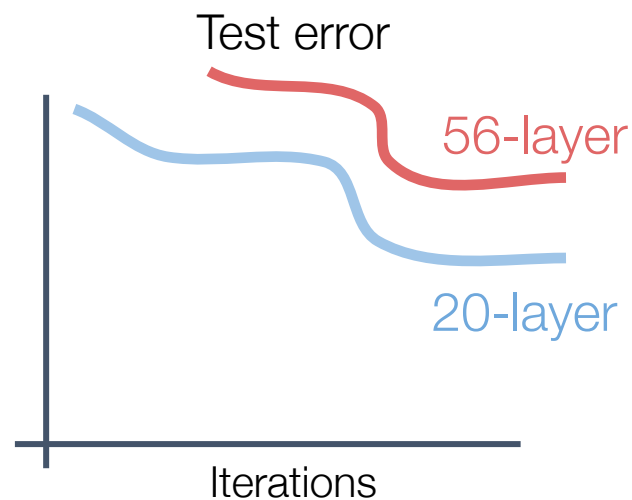He et al, "Deep Residual Learning for Image Recognition", CVPR 2016

# Residual Networks

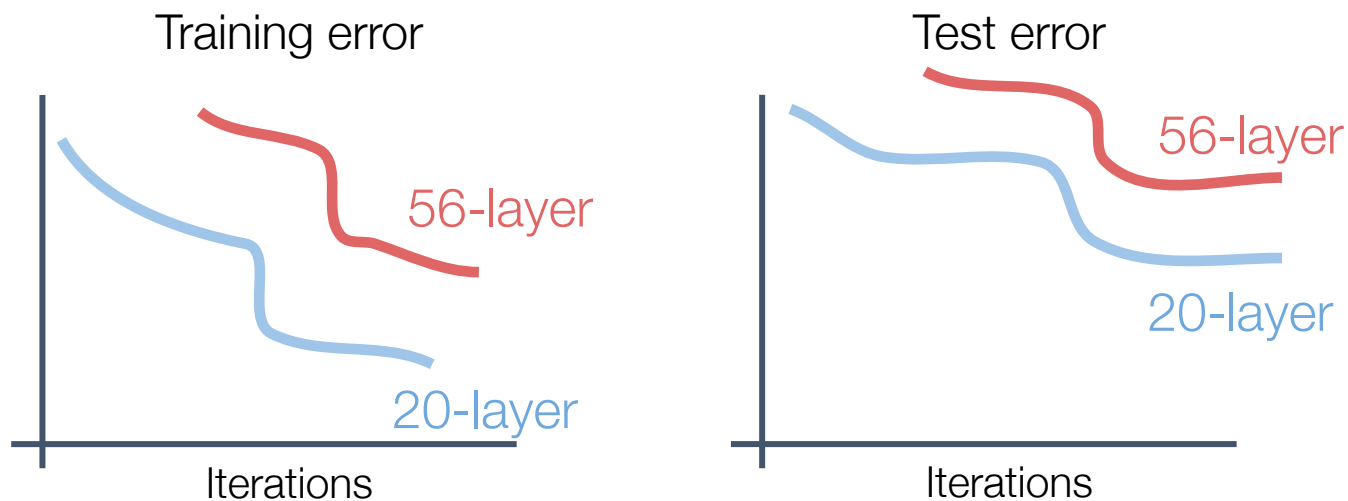Once we have Batch Normalization, we can train networks with 10+ layers. What happens as we go deeper?

Deeper model does worse than shallow model!

Initial guess: Deep model is **overfitting** since it is much bigger than the other model

Test error

56-layer

20-layer

Iterations

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

# Residual Networks

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



Training error

56-layer

20-layer

Iterations

Test error

56-layer

20-layer

Iterations

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

A deeper model can <u>emulate</u> a shallower model: copy layers from shallower model, set extra layers to identity

Thus deeper models should do at least as good as shallow models

**Hypothesis**: This is an <u>optimization</u> problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models

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

# Residual Networks

A deeper model can <u>emulate</u> a shallower model: copy layers from shallower model, set extra layers to identity
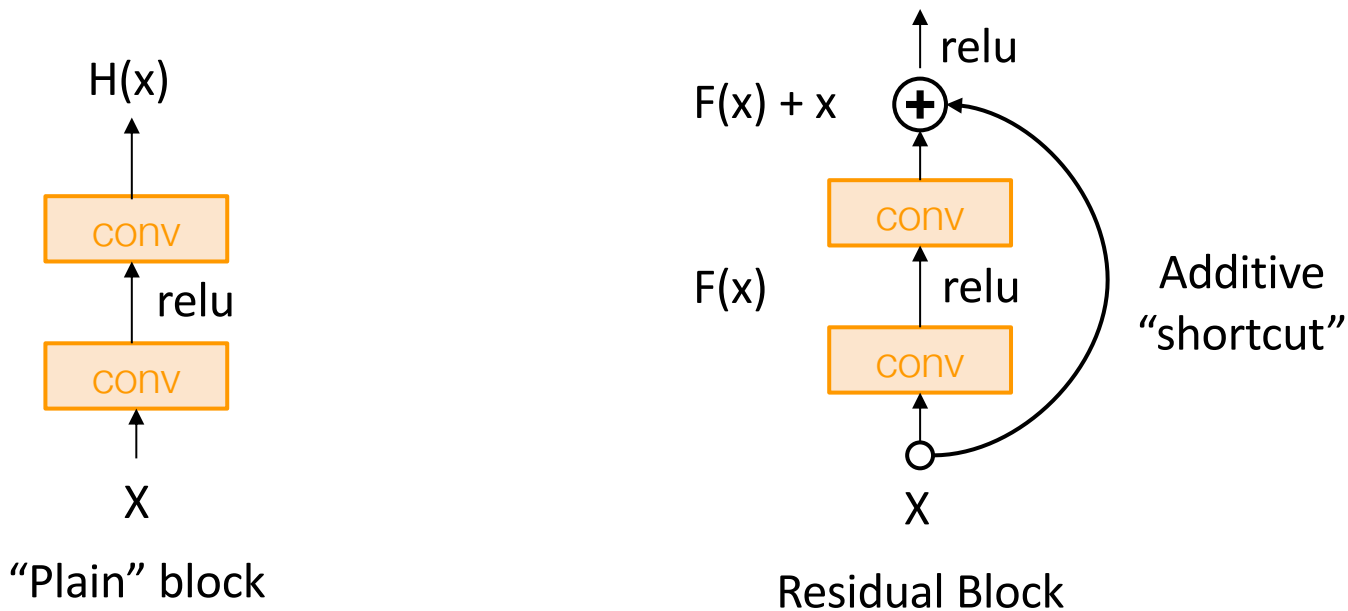
Thus deeper models should do at least as good as shallow models

**Hypothesis**: This is an <u>optimization</u> problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models

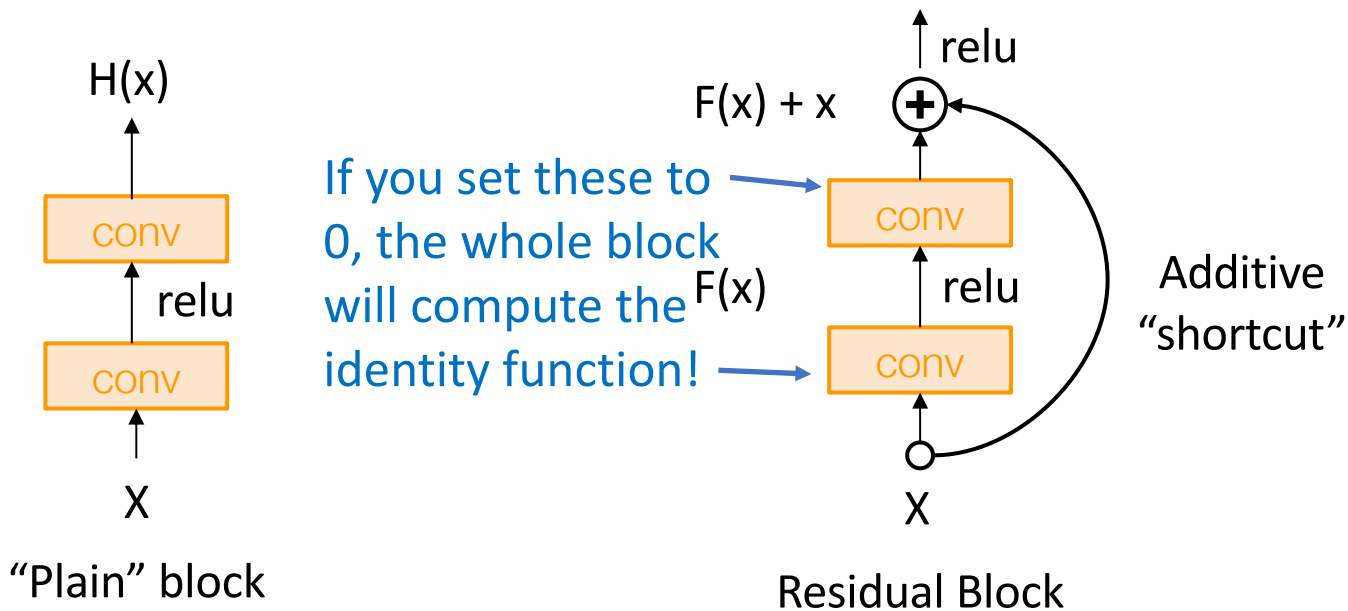**Solution**: Change the network so learning identity functions with extra layers is easy!

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
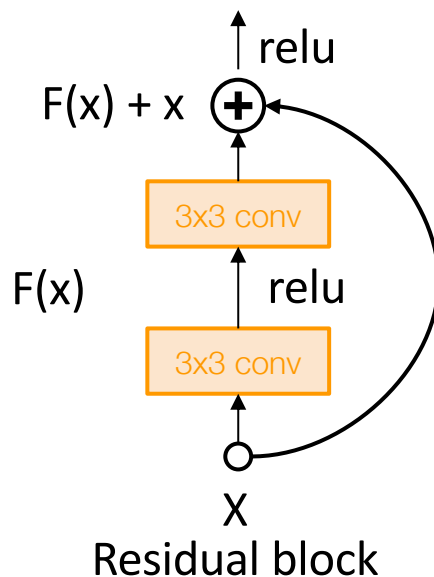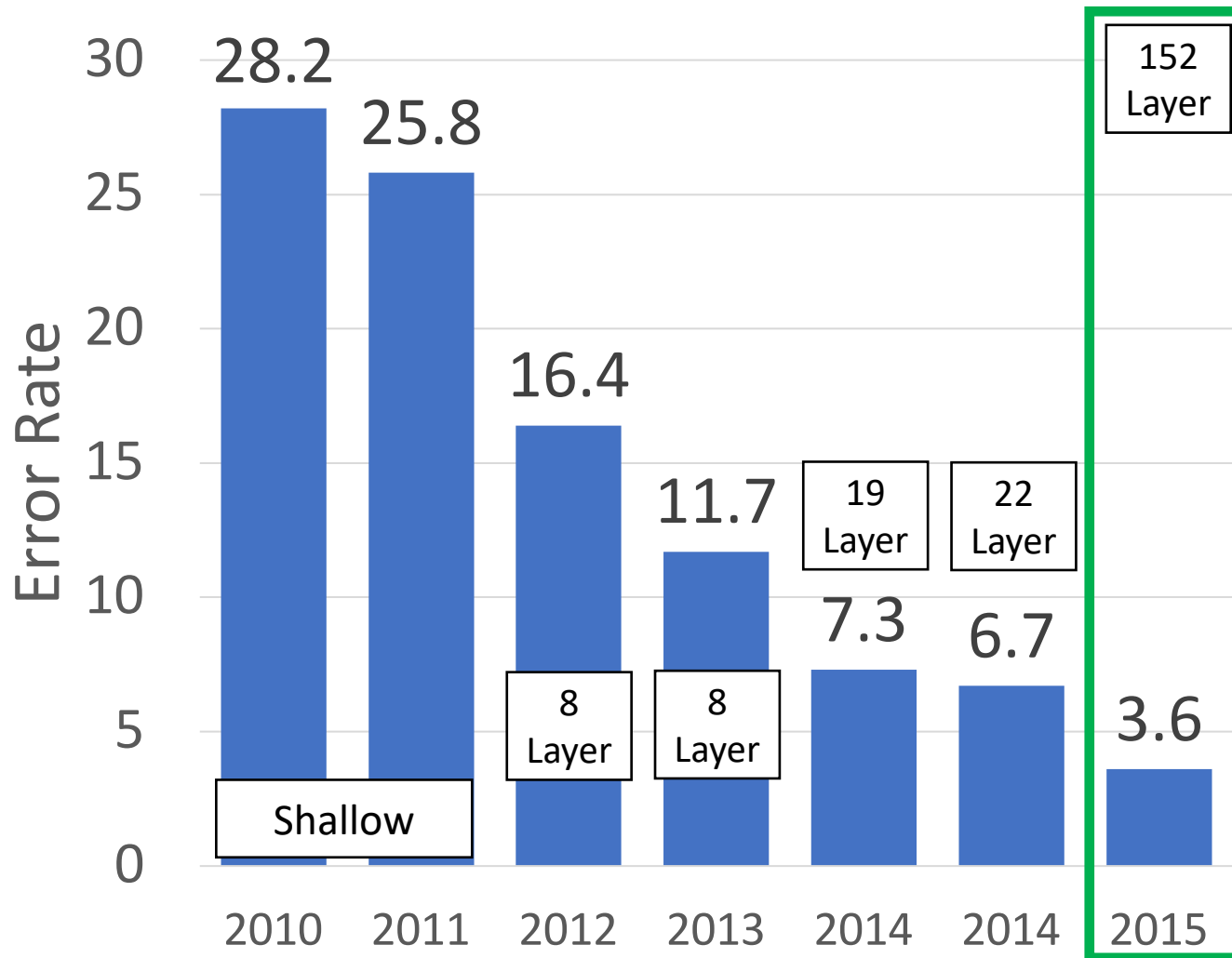
# Residual Networks

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
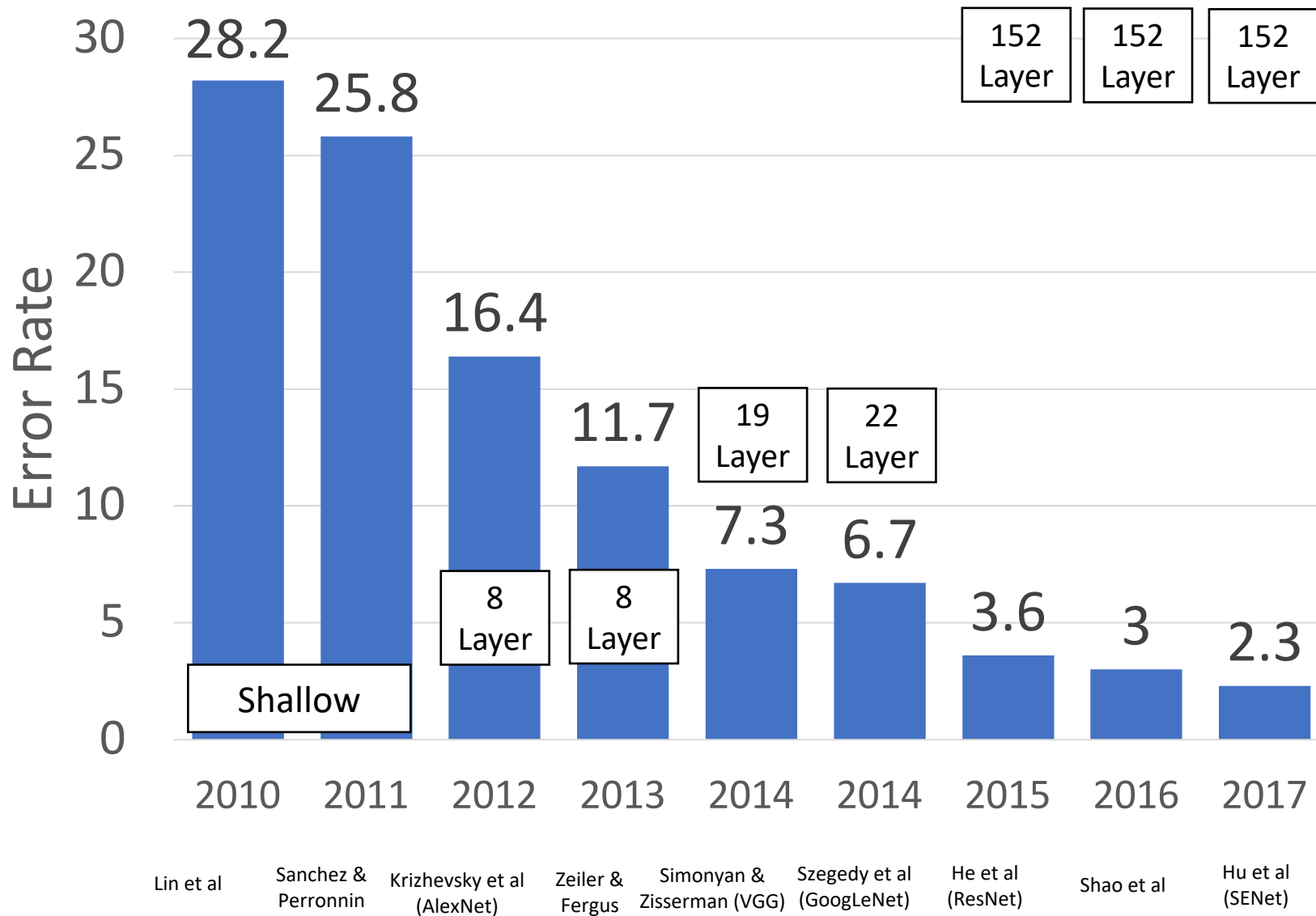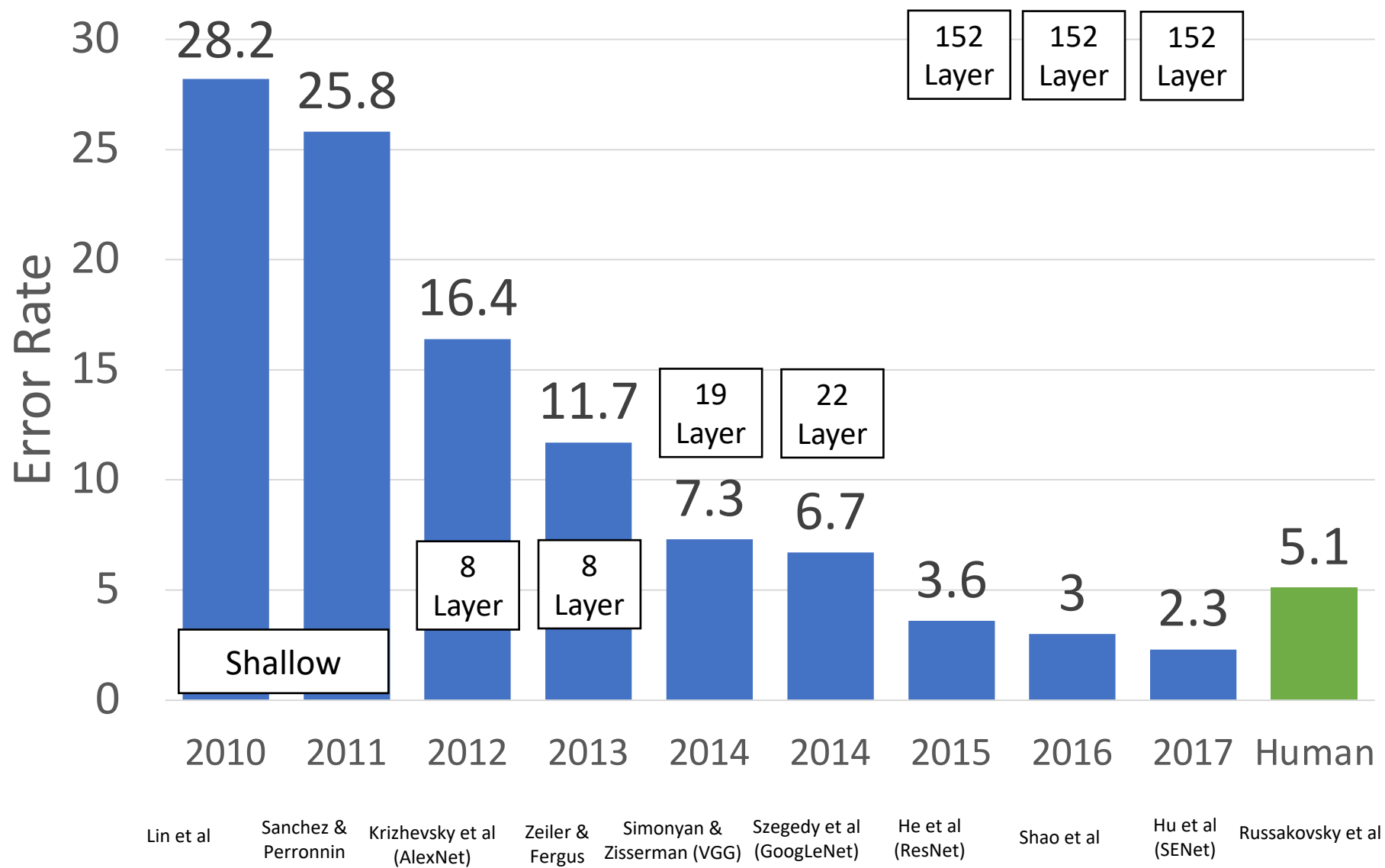


Residual block

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

# ImageNet Classification Challenge

# ImageNet Classification Challenge

# ImageNet Classification Challenge

# 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

# Weight Initialization: Activation Statistics

**Forward pass for a 6-layer net with hidden size 4096**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

# Weight Initialization: Activation Statistics

**Forward pass for a 6-layer net with hidden size 4096**

```
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```
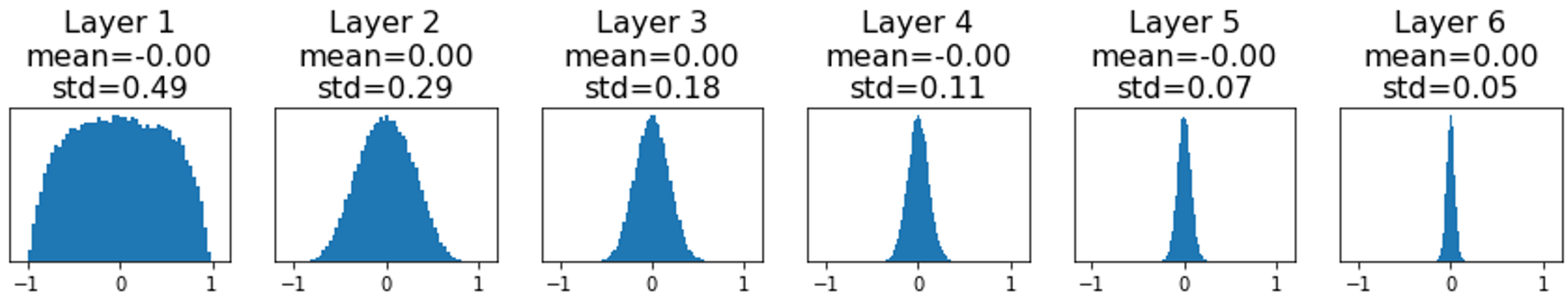
All activations tend to zero for deeper network layers

**Q**: What do the gradients dL/dW look like?



| Layer 1 | Layer 2 | Layer 3 | Layer 4 | Layer 5 | Layer 6 |
|---|---|---|---|---|---|
| mean=-0.00 std=0.49 | mean=0.00 std=0.29 | mean=0.00 std=0.18 | mean=-0.00 std=0.11 | mean=-0.00 std=0.07 | mean=0.00 std=0.05 |

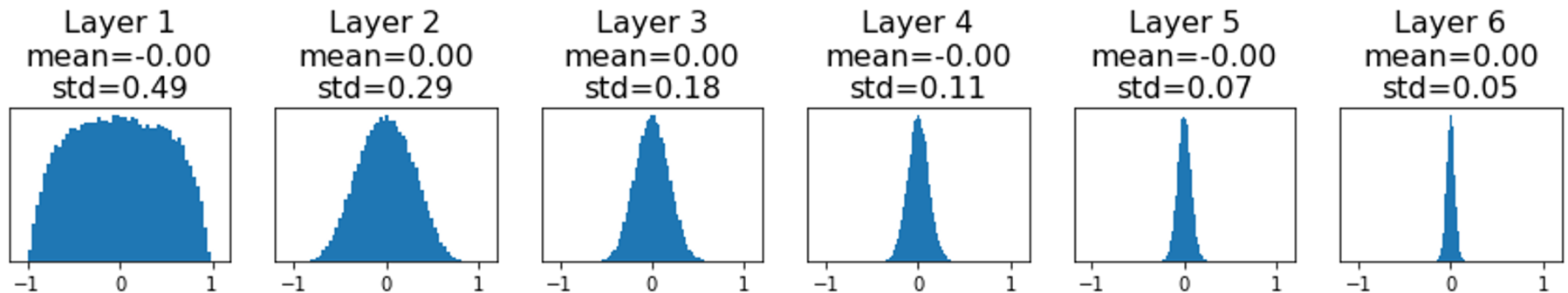# Weight Initialization: Activation Statistics

**Forward pass for a 6-layer net with hidden size 4096**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations tend to zero for deeper network layers

**Q**: What do the gradients dL/dW look like?

**A**: All zero, no learning =(



| Layer 1 mean=-0.00 std=0.49 | Layer 2 mean=0.00 std=0.29 | Layer 3 mean=0.00 std=0.18 | Layer 4 mean=-0.00 std=0.11 | Layer 5 mean=-0.00 std=0.07 | Layer 6 mean=0.00 std=0.05 |

Weights are **too small** at initialization!

# Weight Initialization: Activation Statistics

**Increase scale of weights at initialization 0.01 -> 0.05**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

# Weight Initialization: Activation Statistics

**Increase scale of weights at initialization 0.01 -> 0.05**

```
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```
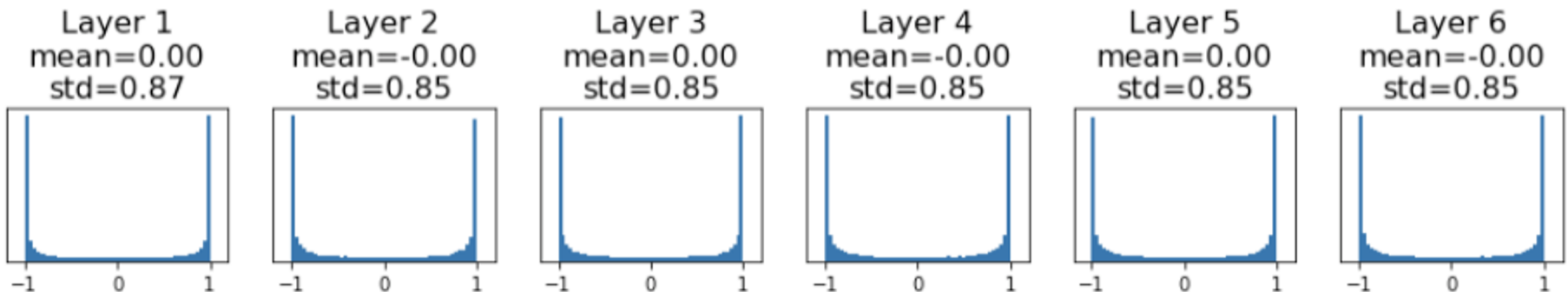
All activations saturate

**Q**: What do the gradients look like?



Weights are **too big** at initialization!
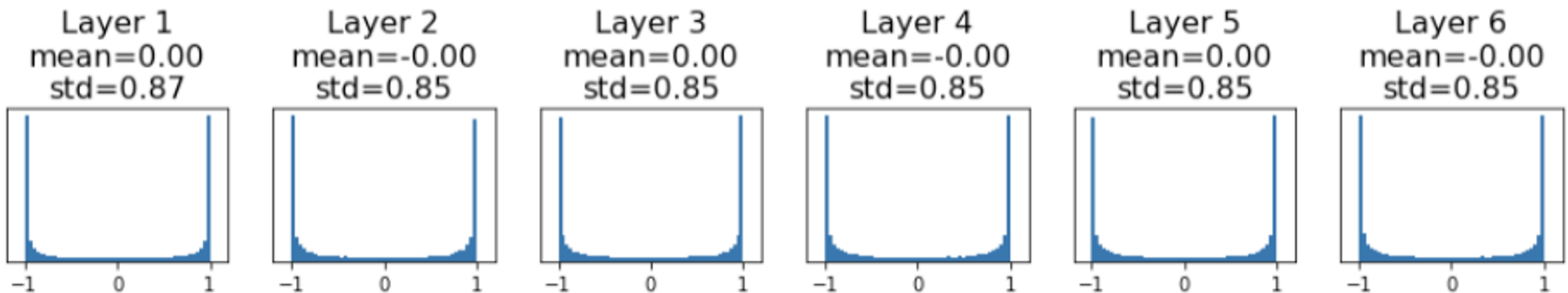
# Weight Initialization: Activation Statistics

**Increase scale of weights at initialization 0.01 -> 0.05**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations saturate

**Q**: What do the gradients look like?

**A**: Local gradients all zero, no learning =(



| Layer 1 | Layer 2 | Layer 3 | Layer 4 | Layer 5 | Layer 6 |
|---------|---------|---------|---------|---------|---------|
| mean=0.00 std=0.87 | mean=-0.00 std=0.85 | mean=0.00 std=0.85 | mean=-0.00 std=0.85 | mean=0.00 std=0.85 | mean=-0.00 std=0.85 |

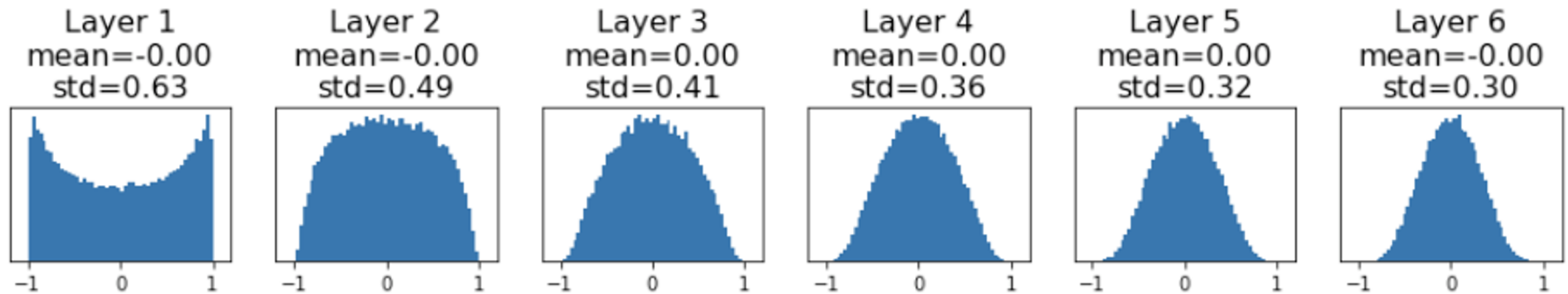# Weights are **too big** at initialization!

# Weight Initialization: Xavier

**"Xavier" initialization: std = 1 / sqrt(Din)**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010

# Weight Initialization: Xavier

**"Xavier" initialization: std = 1 / sqrt(Din)**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.tanh(x.dot(W))
    hs.append(x)
```
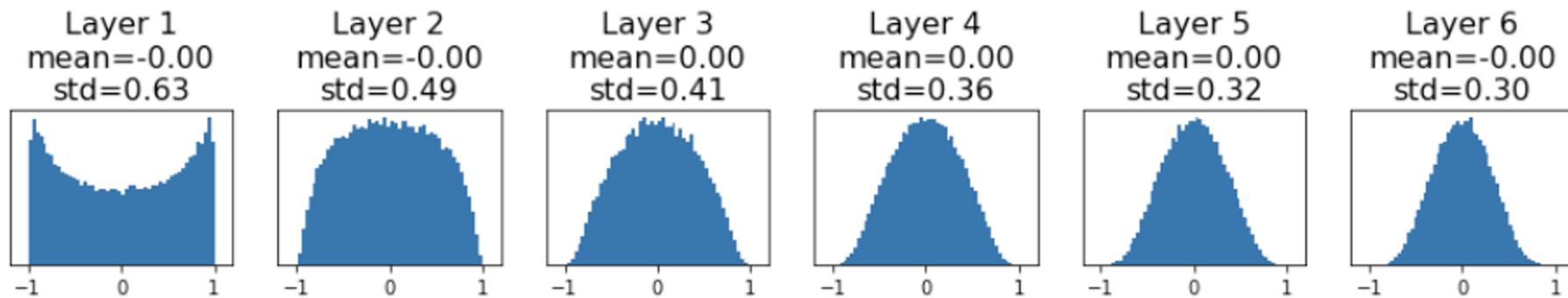


| Layer 1 mean=-0.00 std=0.63 | Layer 2 mean=-0.00 std=0.49 | Layer 3 mean=0.00 std=0.41 | Layer 4 mean=0.00 std=0.36 | Layer 5 mean=0.00 std=0.32 | Layer 6 mean=-0.00 std=0.30 |

## Weights are **just right** at initialization!

Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010

# Weight Initialization: Xavier

**"Xavier" initialization: std = 1 / sqrt(Din)**

```
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

For conv layers, Din is kernel_size$^2$ * input_channels



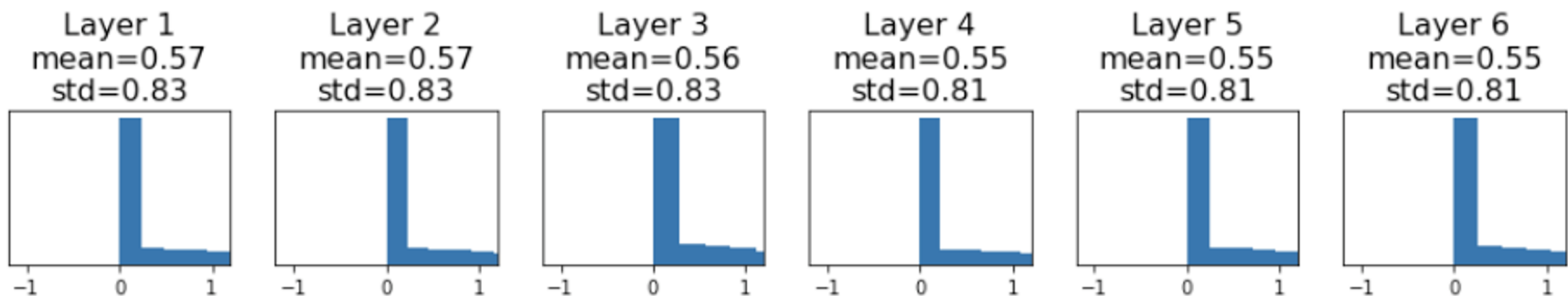| Layer 1 mean=-0.00 std=0.63 | Layer 2 mean=-0.00 std=0.49 | Layer 3 mean=0.00 std=0.41 | Layer 4 mean=0.00 std=0.36 | Layer 5 mean=0.00 std=0.32 | Layer 6 mean=-0.00 std=0.30 |

## Weights are **just right** at initialization!

Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010

# Weight Initialization: MSRA

**For ReLU networks: std = 2 / sqrt(Din)**

```python
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

"Just right" – activations nicely scaled for all layers



Layer 1 mean=0.57 std=0.83 | Layer 2 mean=0.57 std=0.83 | Layer 3 mean=0.56 std=0.83 | Layer 4 mean=0.55 std=0.81 | Layer 5 mean=0.55 std=0.81 | Layer 6 mean=0.55 std=0.81

# Weights are **just right** at initialization!

He et al, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification", ICCV 2015

# 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

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

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

# 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?

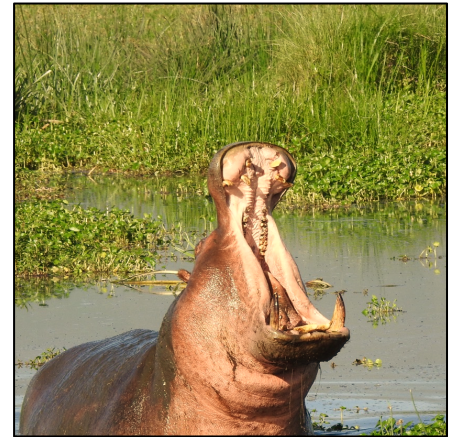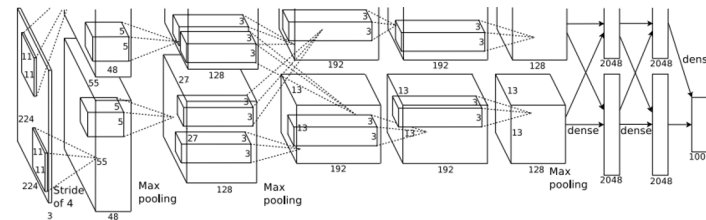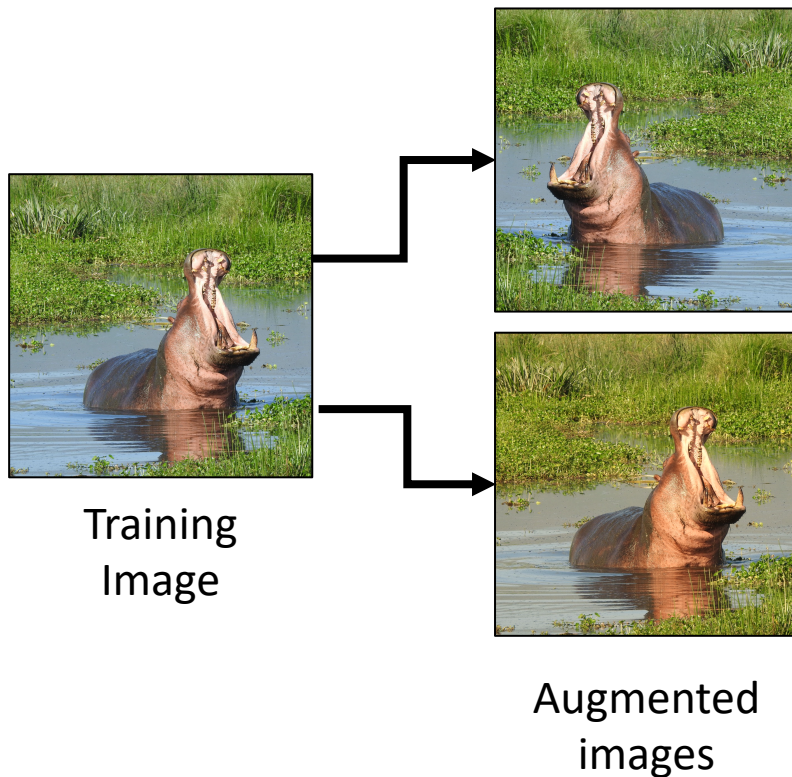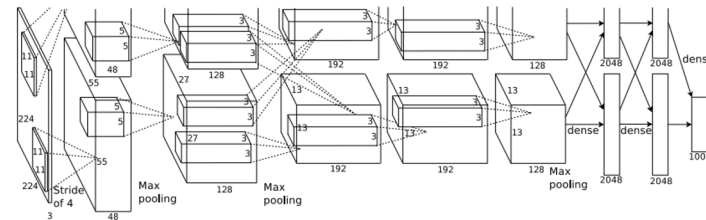Hippo?

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

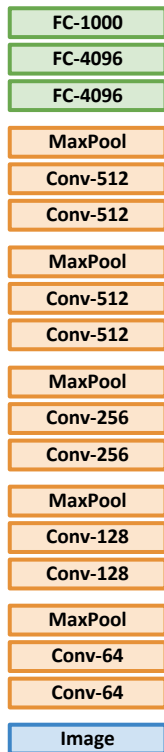Hippo

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

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

What if we can't find one?
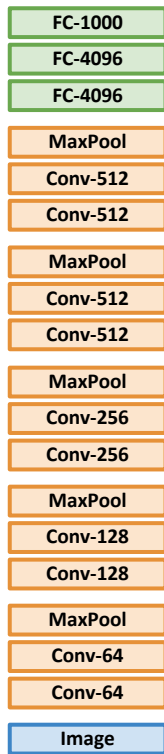
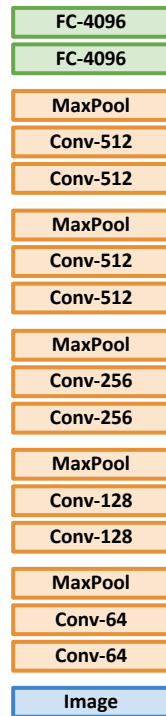# Transfer Learning: Feature Extraction

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**

| 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

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

# Transfer Learning: Fine-Tuning

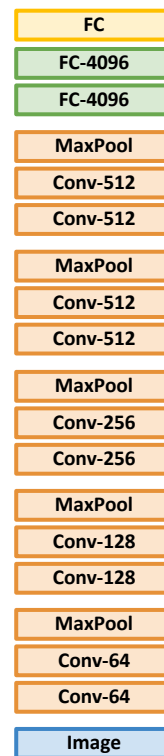**1. Train on ImageNet**



**2. CNN as feature extractor**



Remove last layer

Freeze these

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



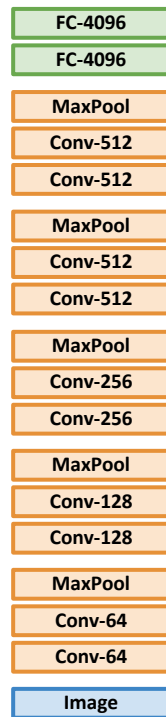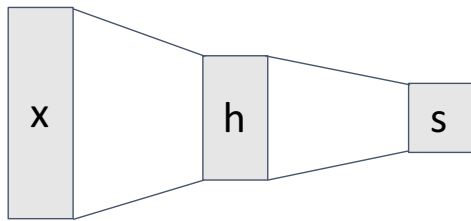Reinitialize last layer and continue training whole network on your dataset

# 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
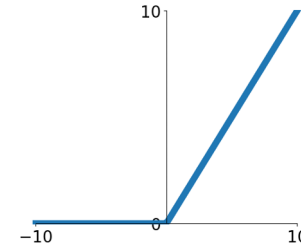
# Recap: Convolutional Networks

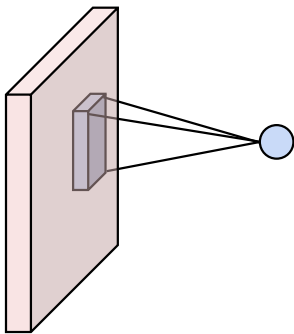## Fully-Connected Layers



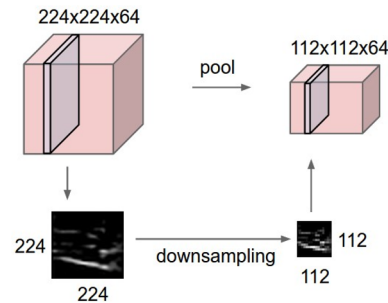$$y = Wx + b$$

## Activation Function



$$y = \max(0, x)$$
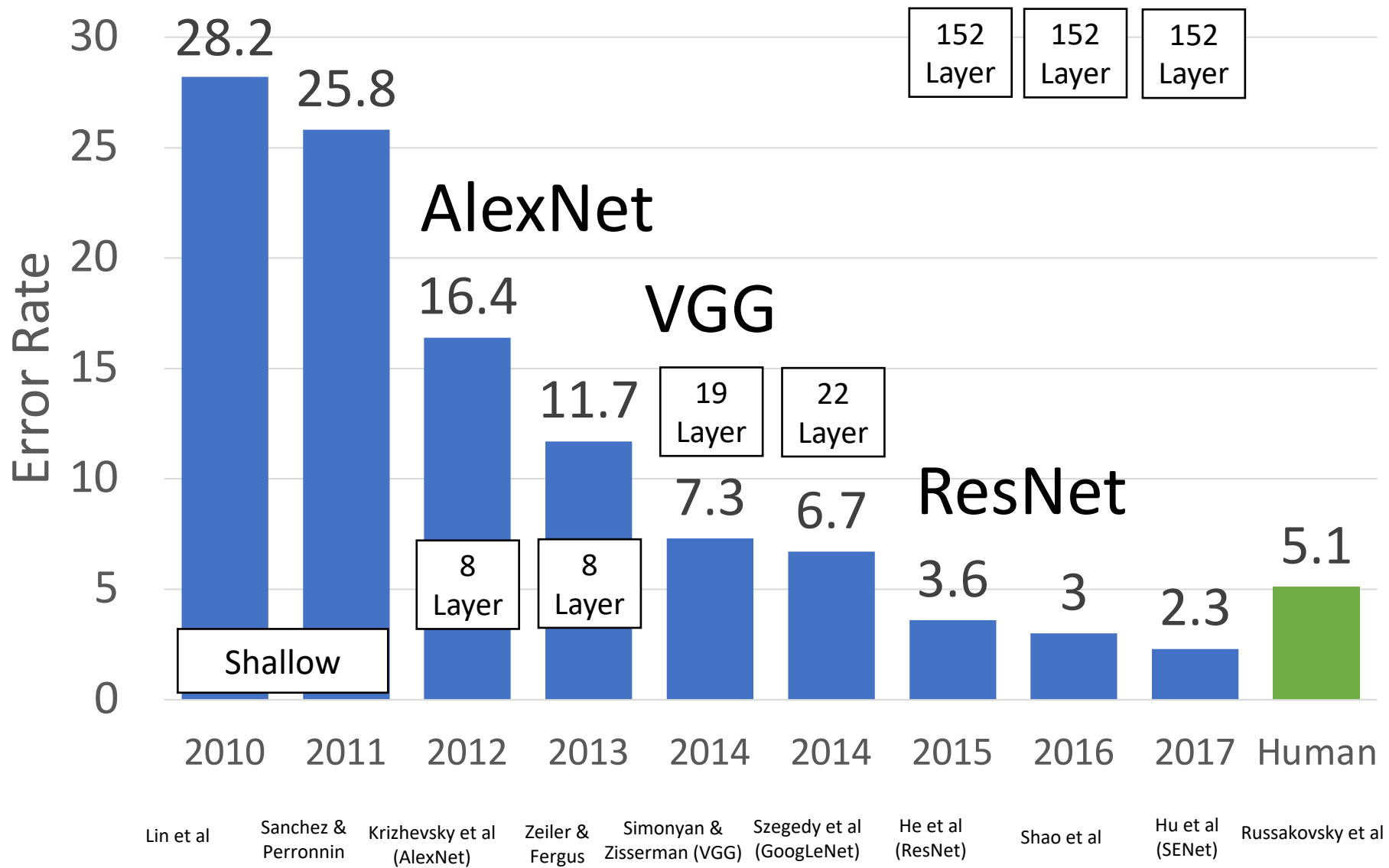
## Convolution Layers



## Pooling Layers



## Normalization

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

# Recap: CNN Architectures

# 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     Regularization + Data Augmentation
    3. Update Weights
5. Apply trained model to task

# So Far: Image Classification

 ⟶ Cat

# What about <u>Localizing</u> Objects?

# Next time:
# Detection + Segmentation