

Training CNNs and CNN Architectures

Lecture Overview – Two Broad Sets of Topics

How to build CNNs?

[
 Layers in CNNs
 Activation Functions
 CNN Architectures
 Weight Initialization

How to train CNNs?

[
 Data Preprocessing
 Data augmentation
 Transfer Learning
 Hyperparameter Selection

Lecture Overview – Two Broad Sets of Topics

How to build CNNs?

Layers in CNNs
Activation Functions
CNN Architectures
Weight Initialization

How to train CNNs?

Data Preprocessing
Data augmentation
Transfer Learning
Hyperparameter Selection

Lecture Overview – Two Broad Sets of Topics

How to build CNNs?

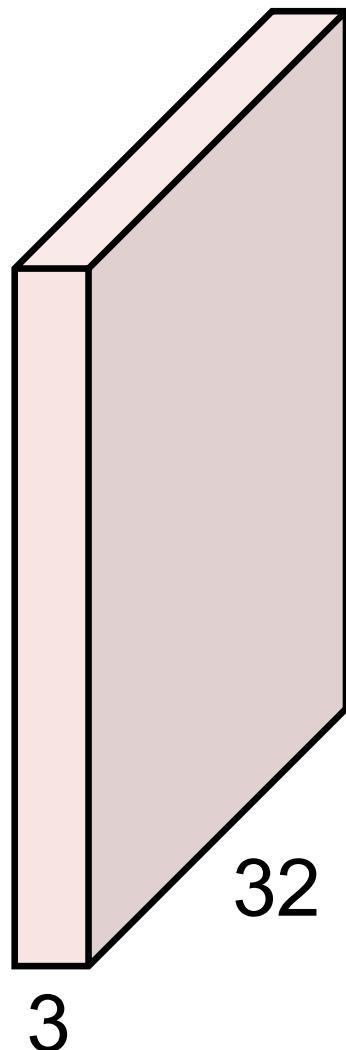
Layers in CNNs
Activation Functions
CNN Architectures
Weight Initialization

How to train CNNs?

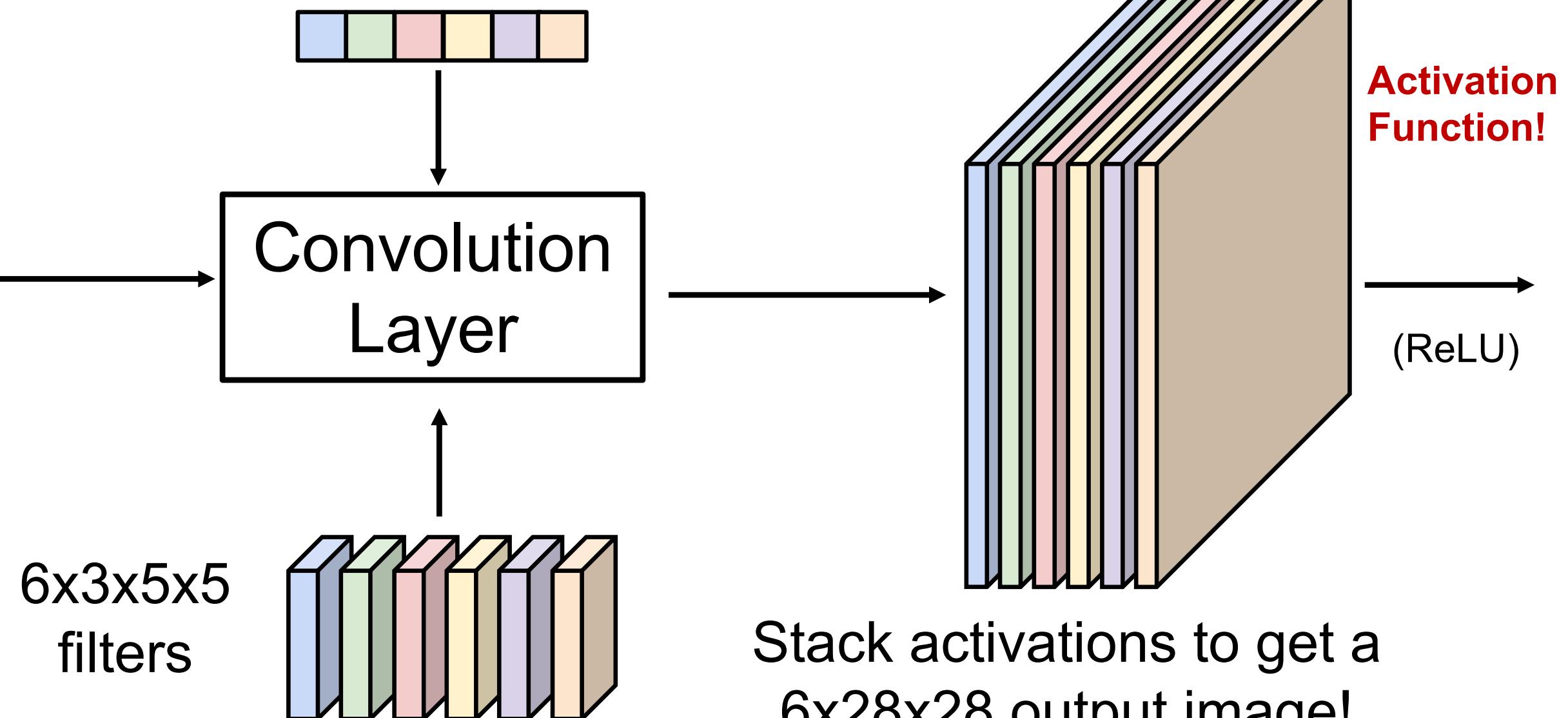
Data Preprocessing
Data augmentation
Transfer Learning
Hyperparameter Selection

Recap: Convolution Layer

3x32x32 image

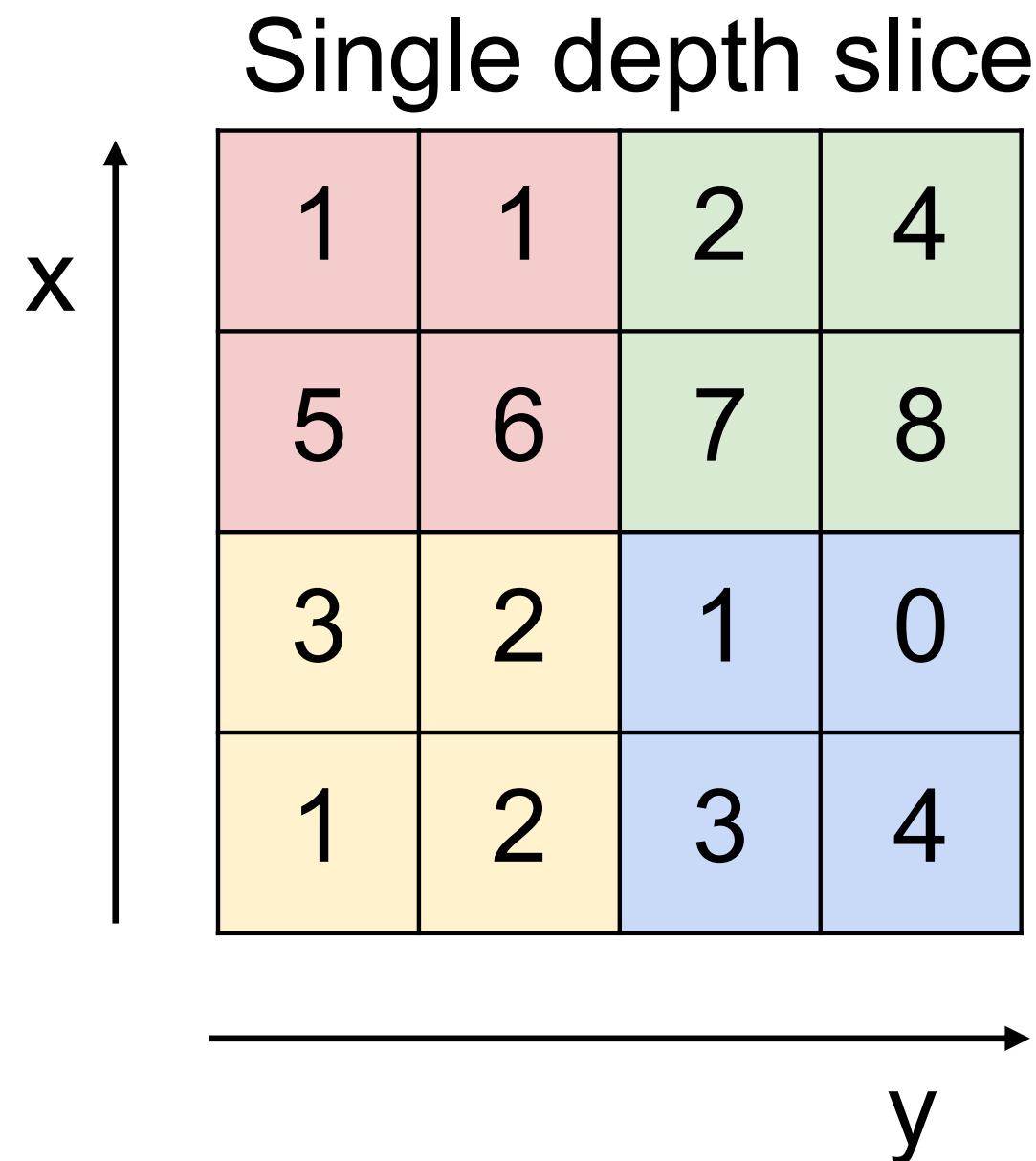


Don't forget bias terms!

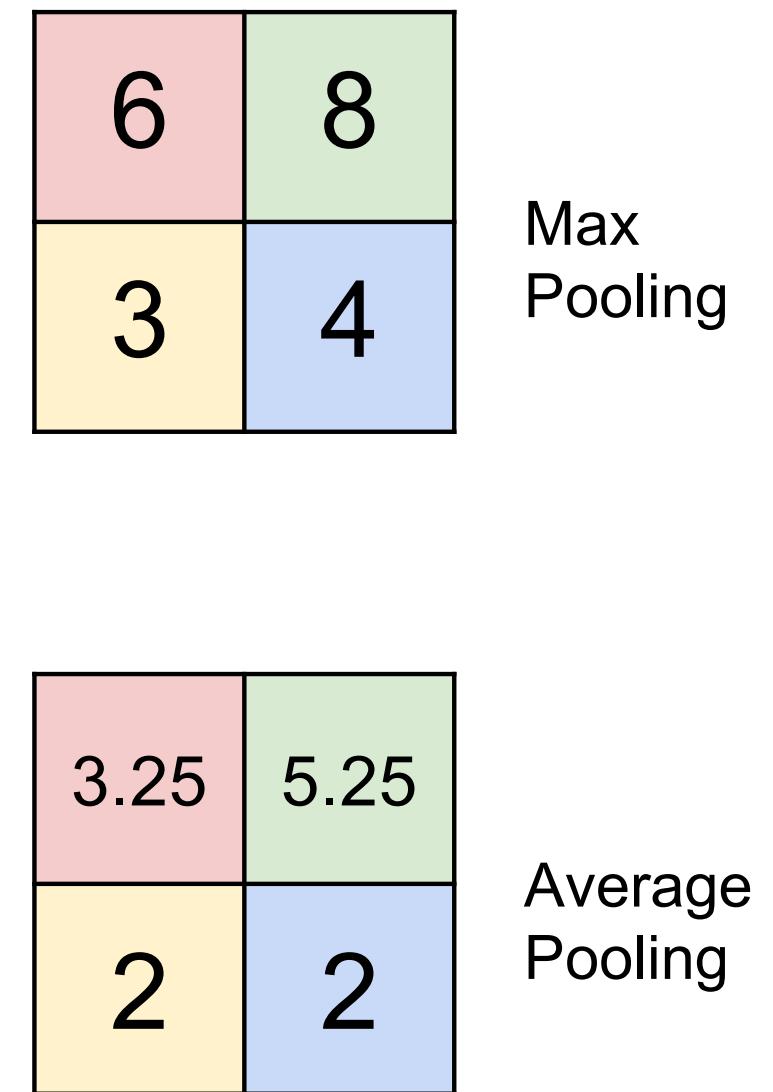


Slide inspiration: Justin Johnson

Recap: Pooling Layer

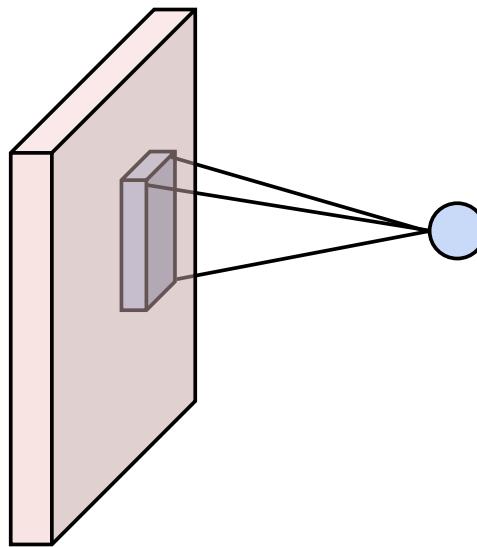


pool with 2x2 filters and
stride 2

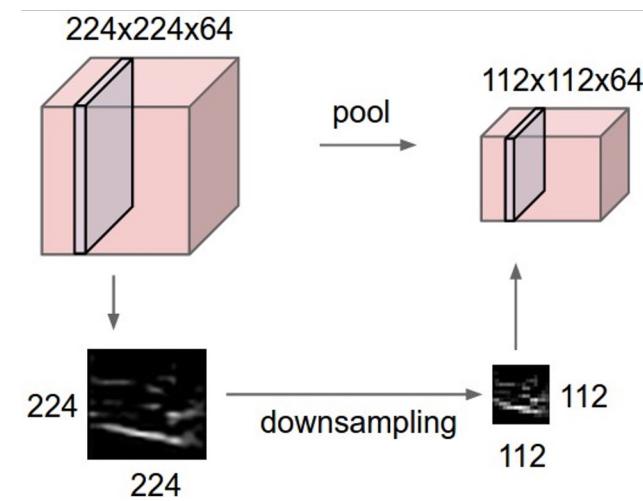


Components of CNNs

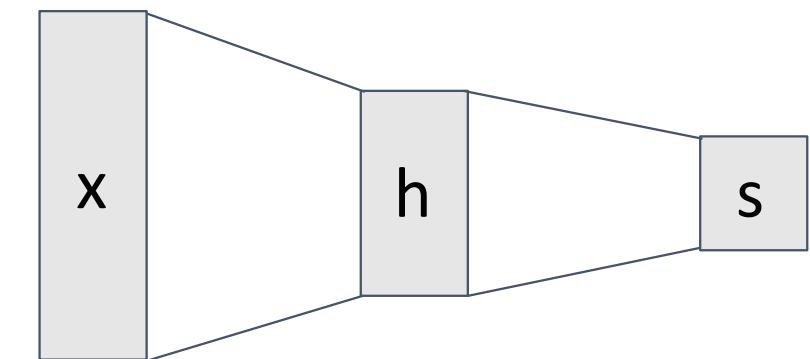
Convolution Layers



Pooling Layers



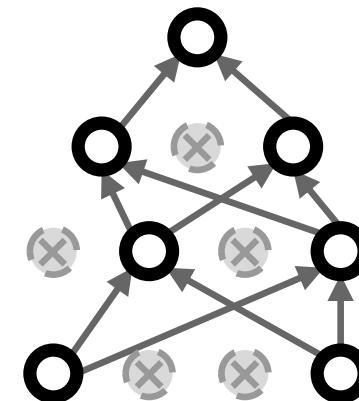
Fully-Connected Layers



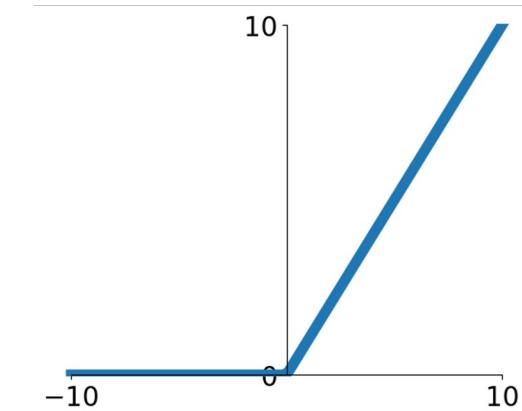
Normalization Layers

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

Dropout (sometimes)

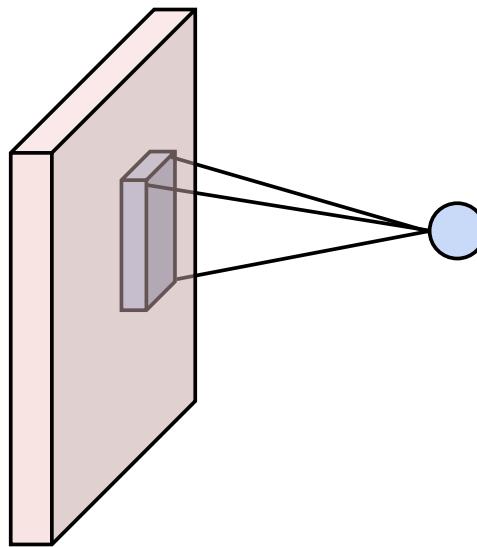


Activation Functions

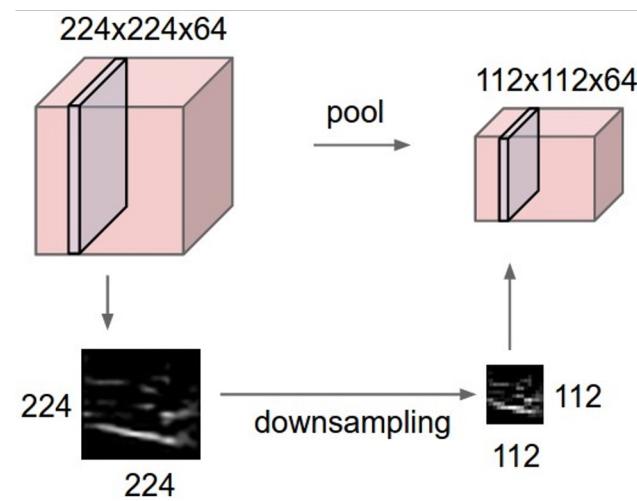


Components of CNNs

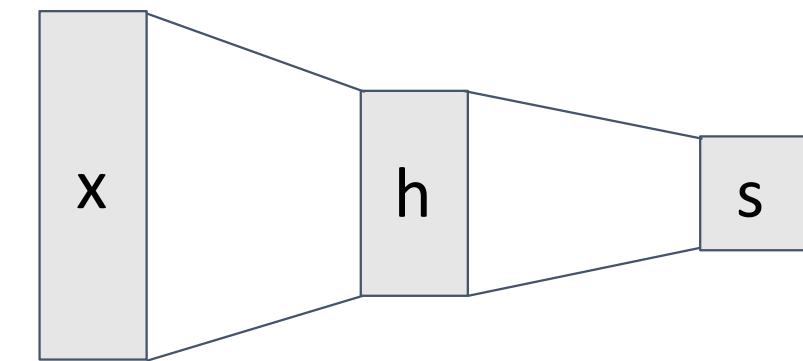
Convolution Layers



Pooling Layers



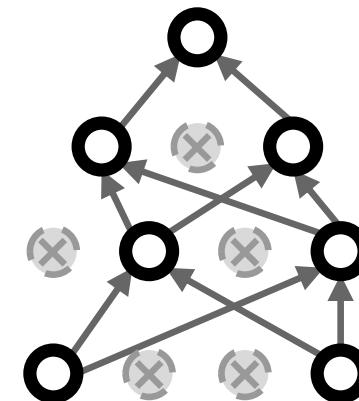
Fully-Connected Layers



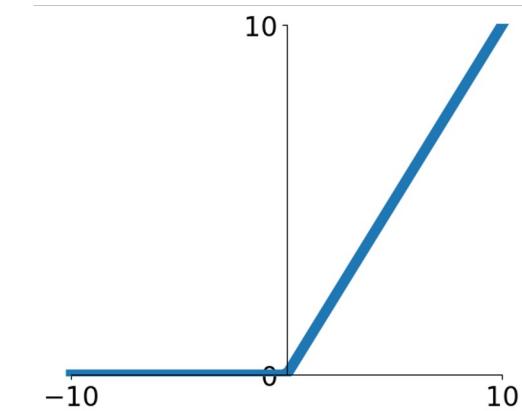
Normalization Layers

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

Dropout (sometimes)



Activation Functions



Example Normalization Layer: LayerNorm

High-level Idea: Learn parameters that let us **scale / shift the input data**

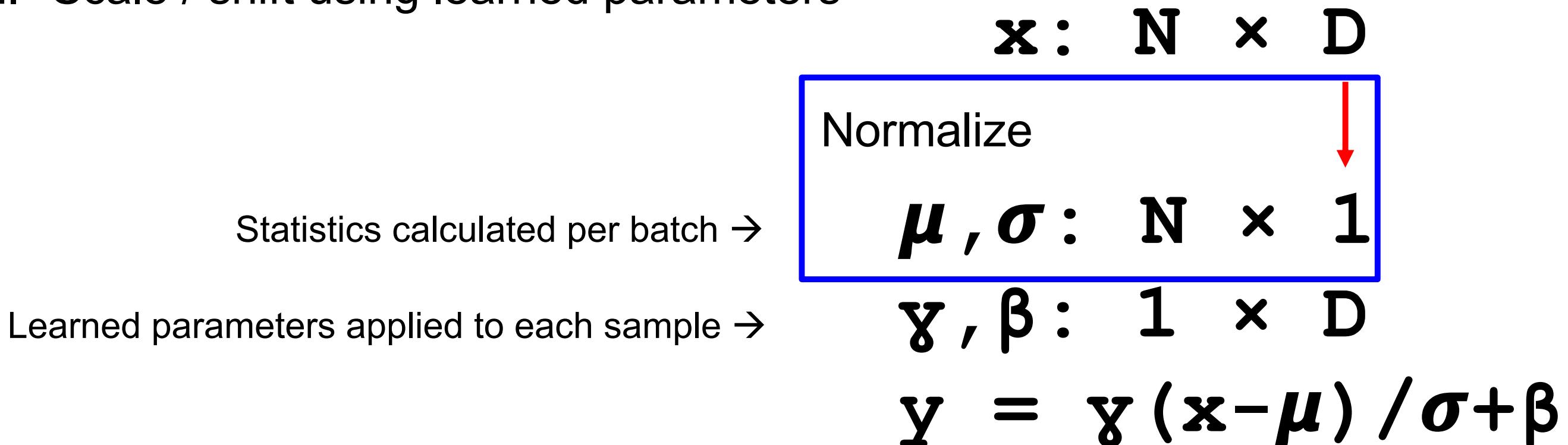
1. Normalize input data
2. Scale / shift using learned parameters

Ba, Kiros, and Hinton, “Layer Normalization”, arXiv 2016

Example Normalization Layer: LayerNorm

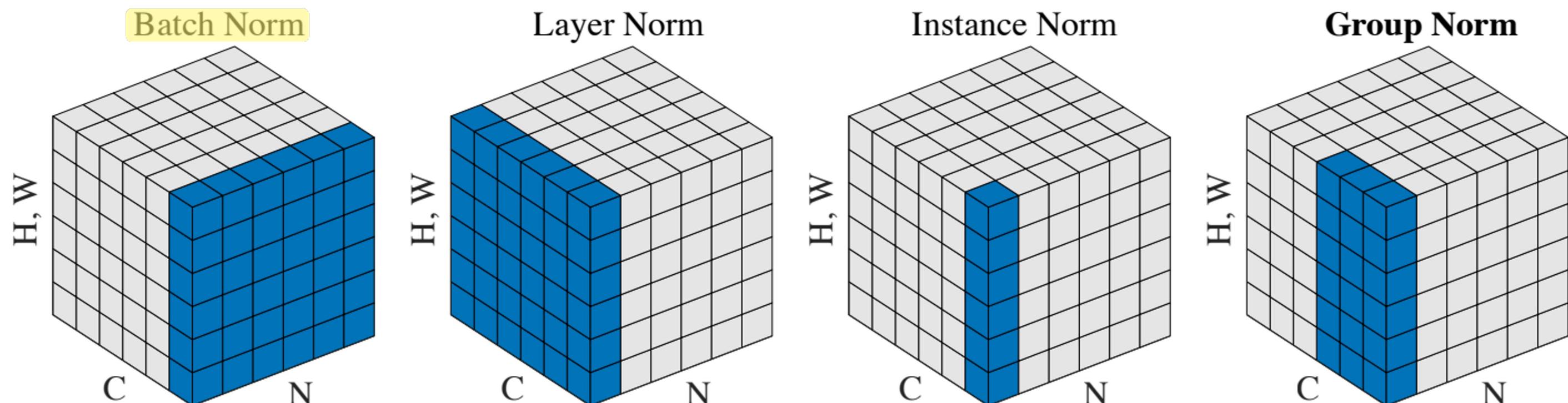
High-level Idea: Learn parameters that let us **scale / shift the input data**

1. Normalize input data
2. Scale / shift using learned parameters



Ba, Kiros, and Hinton, “Layer Normalization”, arXiv 2016

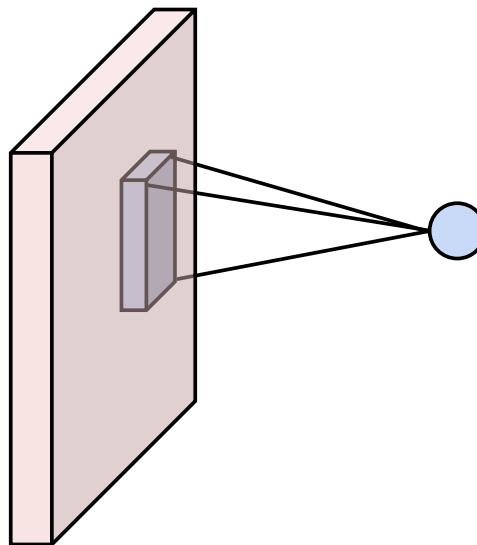
Other Normalization Layers



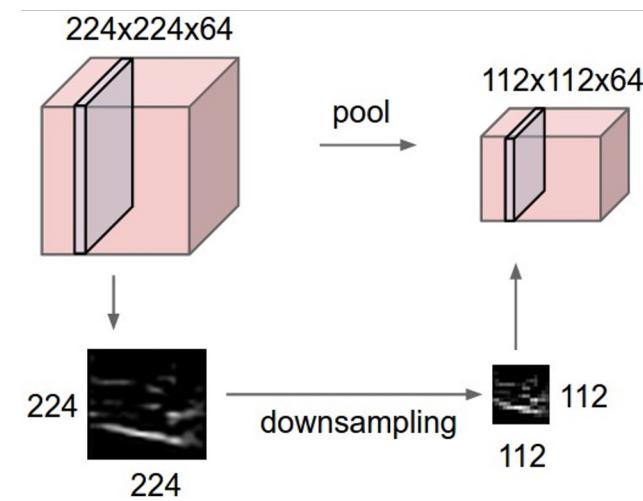
Wu and He, "Group Normalization", ECCV 2018

Components of CNNs

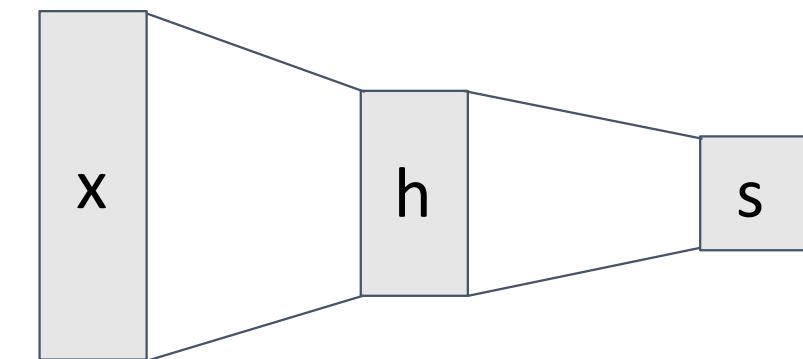
Convolution Layers



Pooling Layers



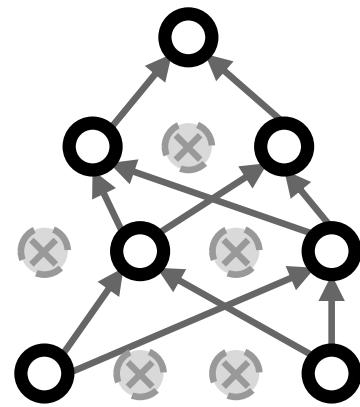
Fully-Connected Layers



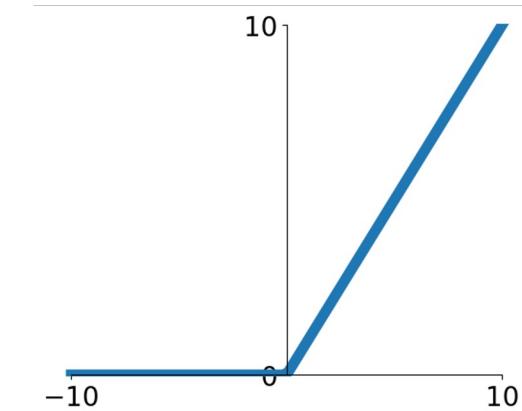
Normalization Layers

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

Dropout (sometimes)



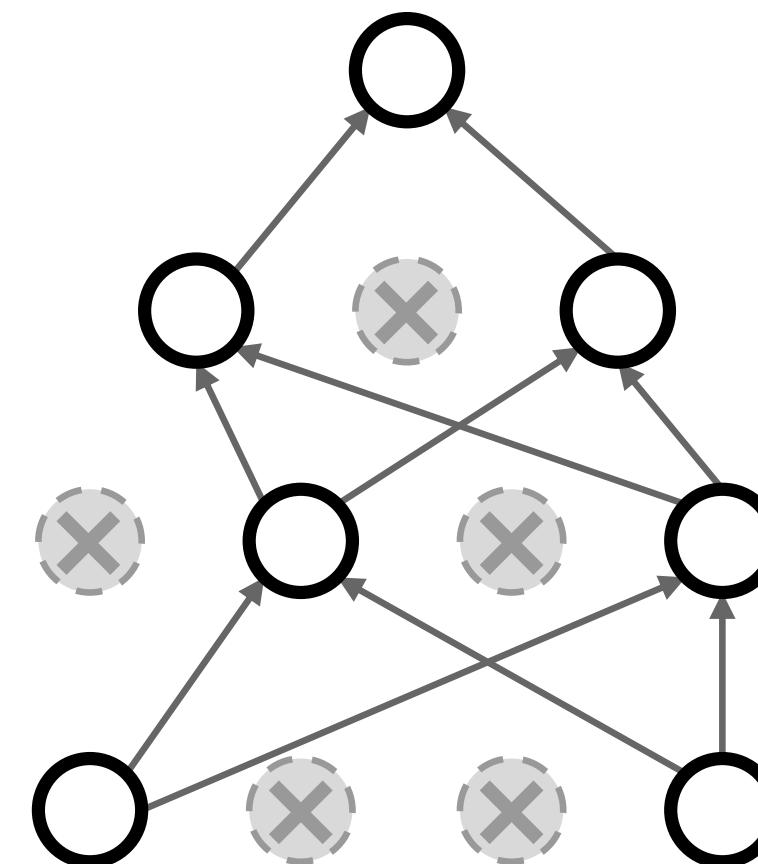
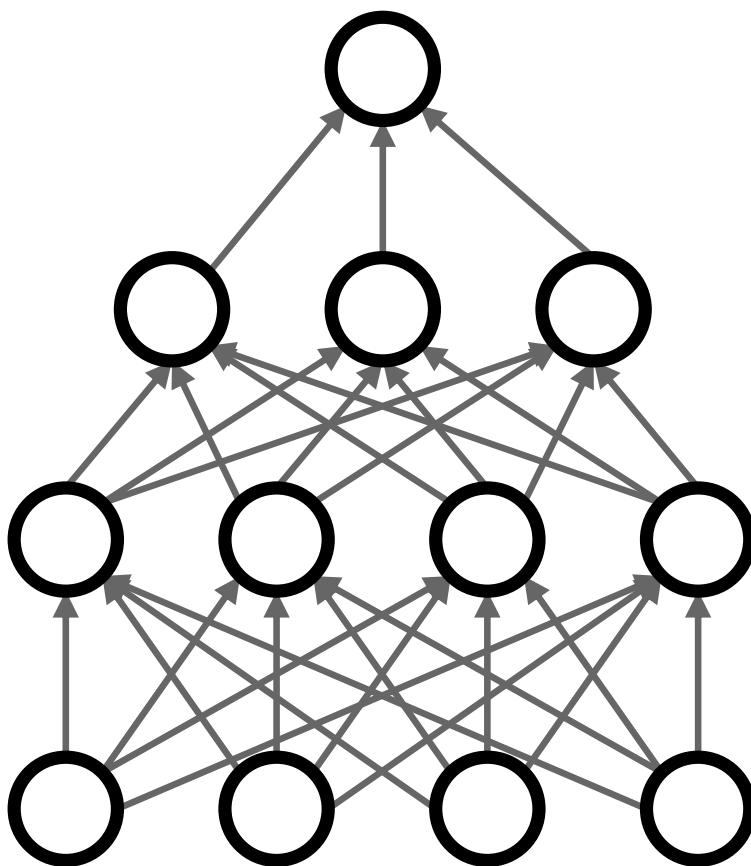
Activation Functions



Regularization: Dropout

In each forward pass, randomly set some neurons to zero

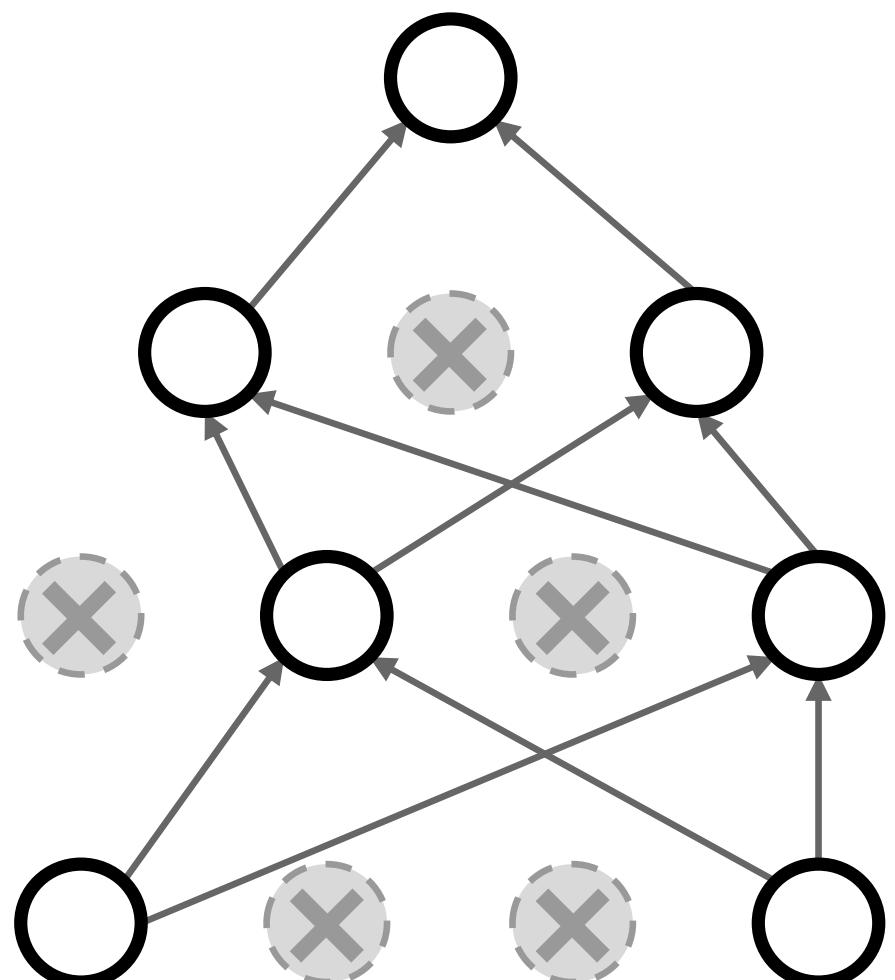
Probability of dropping is a hyperparameter; 0.5 is common



Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014

Regularization: Dropout

How can this possibly be a good idea?



Forces the network to have a redundant representation;
Prevents co-adaptation of features



Lecture Overview – Two Broad Sets of Topics

How to build CNNs?

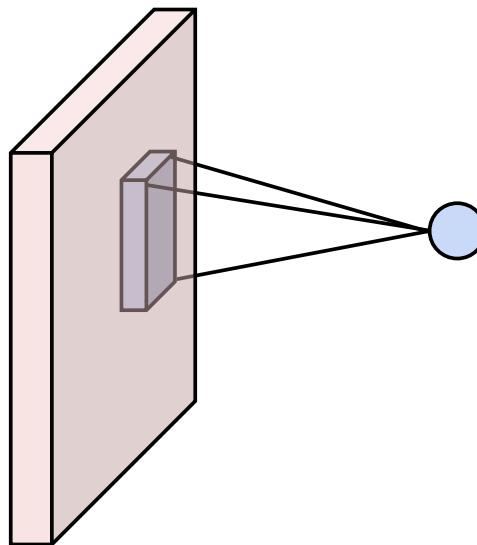
Layers in CNNs
Activation Functions
CNN Architectures
Weight Initialization

How to train CNNs?

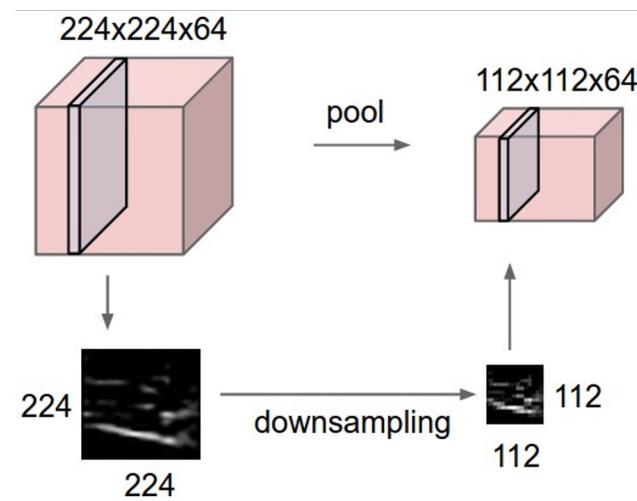
Data Preprocessing
Data augmentation
Transfer Learning
Hyperparameter Selection

Components of CNNs

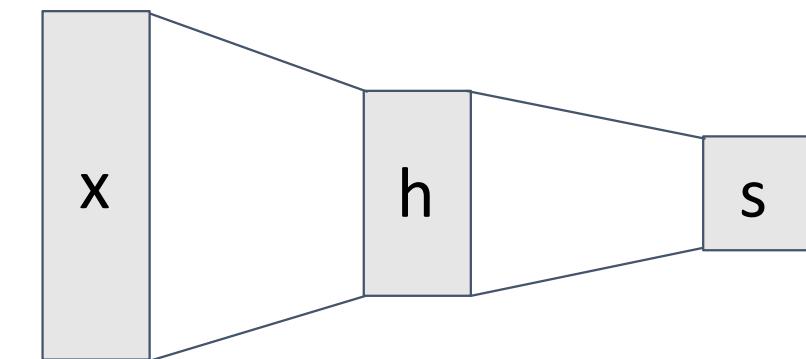
Convolution Layers



Pooling Layers



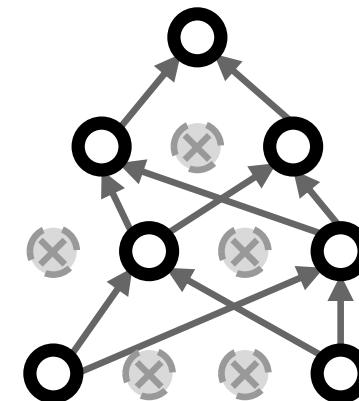
Fully-Connected Layers



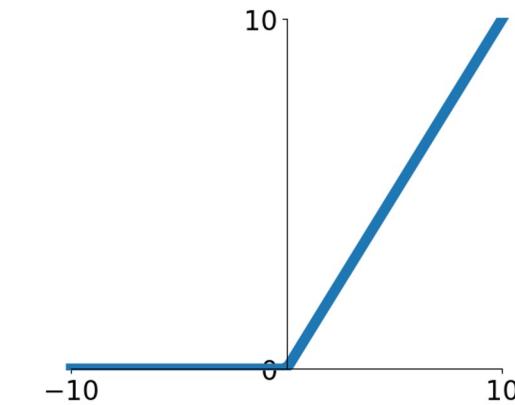
Normalization Layers

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

Dropout (sometimes)

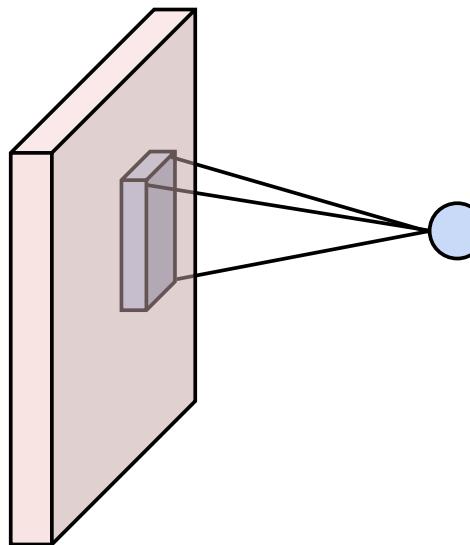


Activation Functions

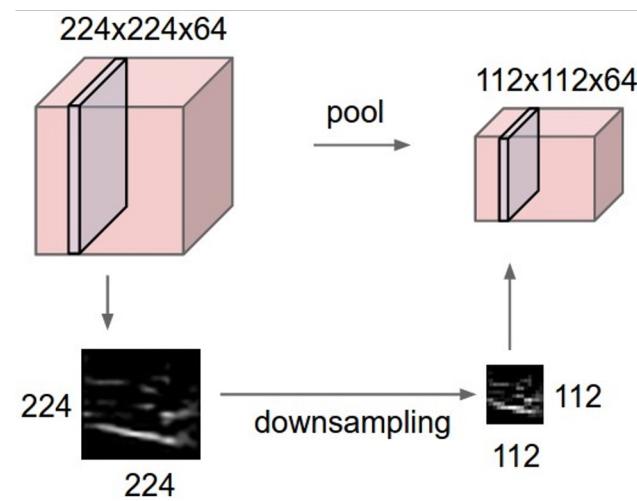


Components of CNNs

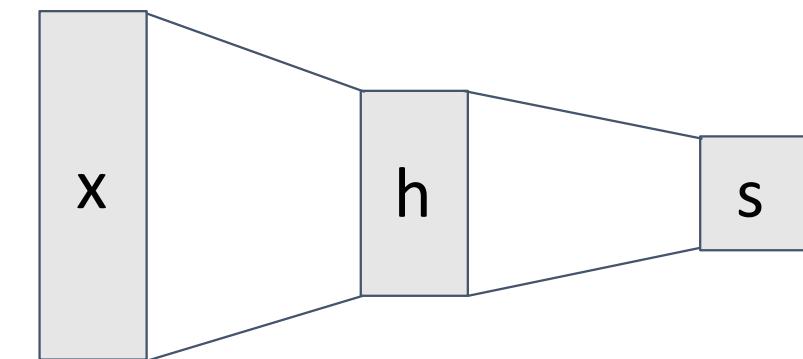
Convolution Layers



Pooling Layers



Fully-Connected Layers

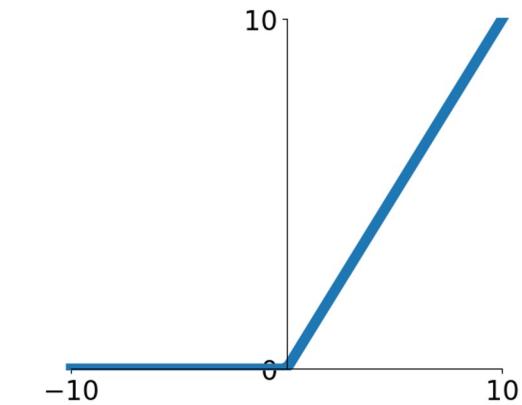


Normalization Layers

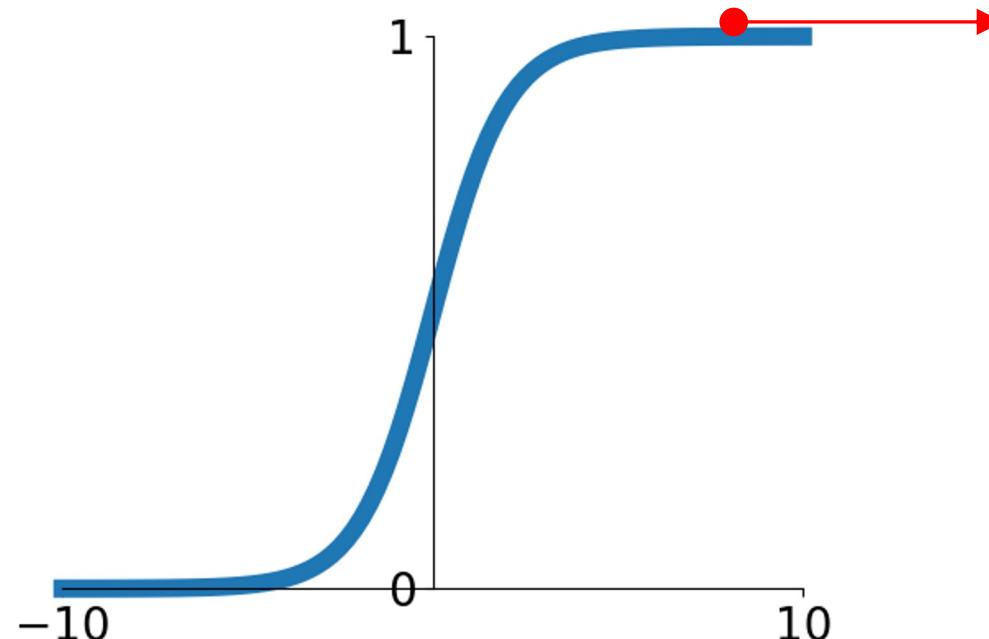
Dropout (sometimes)

Goal: Introduce non-linearities to our model!

Activation Functions



Activation Functions



Sigmoid

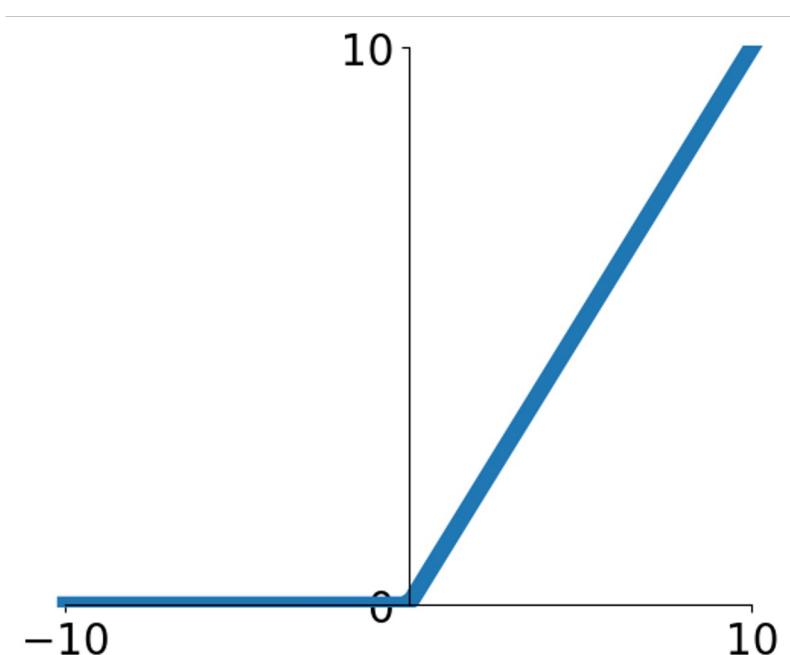
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

Key problem:

Large positive or negative values can “kill” the gradients. Many layers of sigmoids → smaller and smaller gradients in practice

Activation Functions

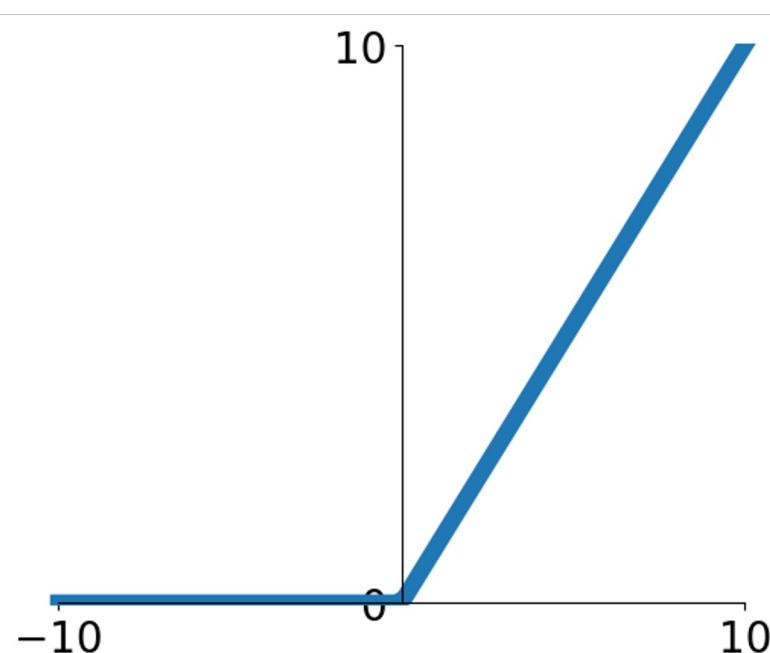


- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid in practice (e.g. 6x)

ReLU
(Rectified Linear Unit)

[Krizhevsky et al., 2012]

Activation Functions



ReLU
(Rectified Linear Unit)

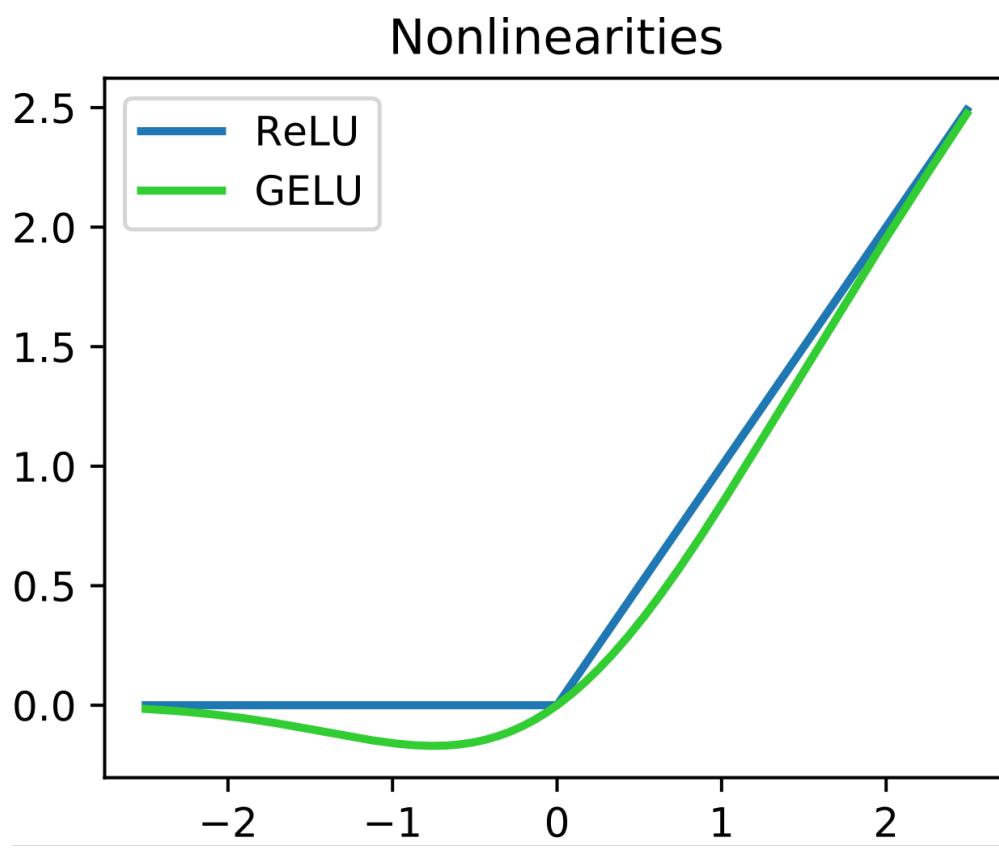
- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid in practice (e.g. 6x)

- Not zero-centered output
- An annoyance:

Dead ReLUs when $x < 0$!

Activation Functions

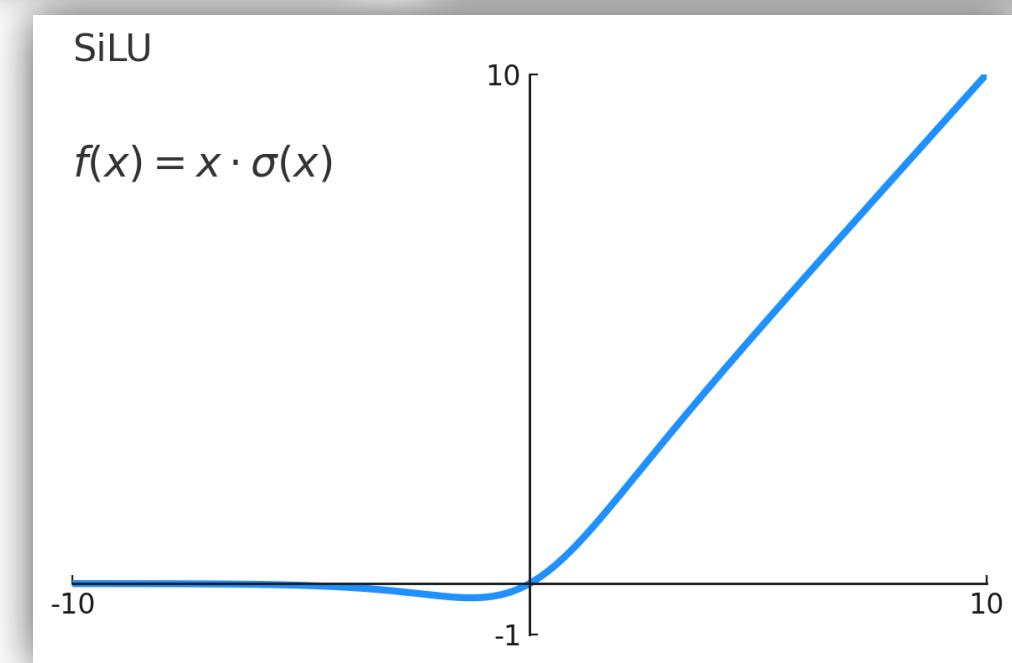
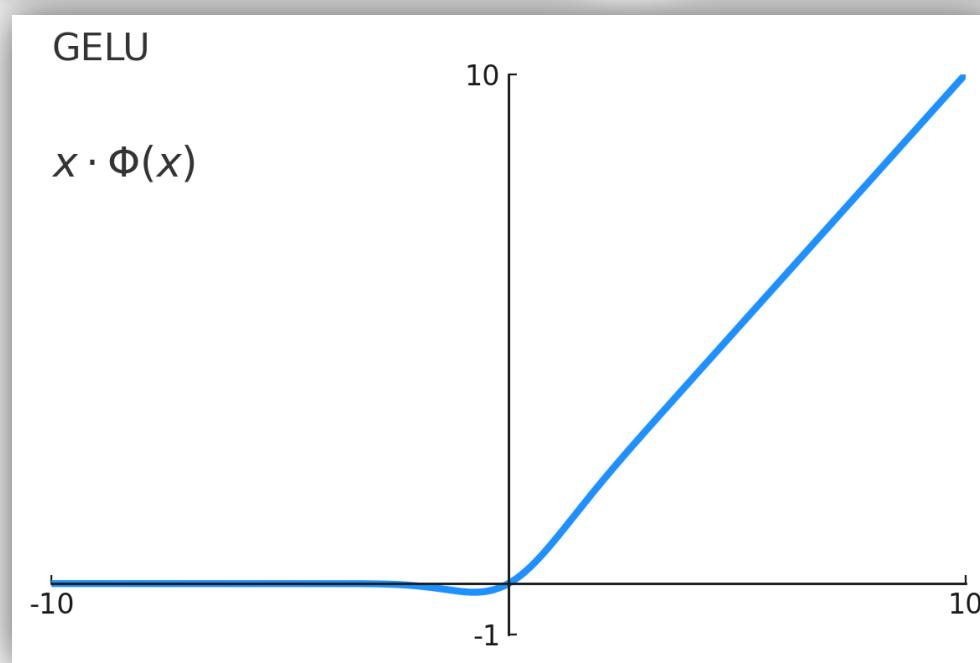
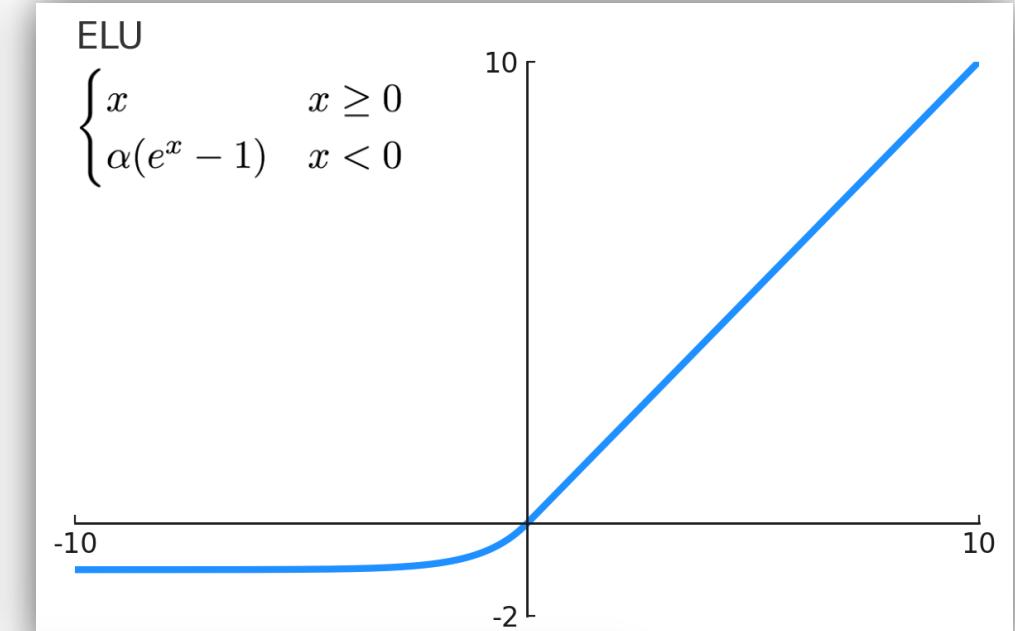
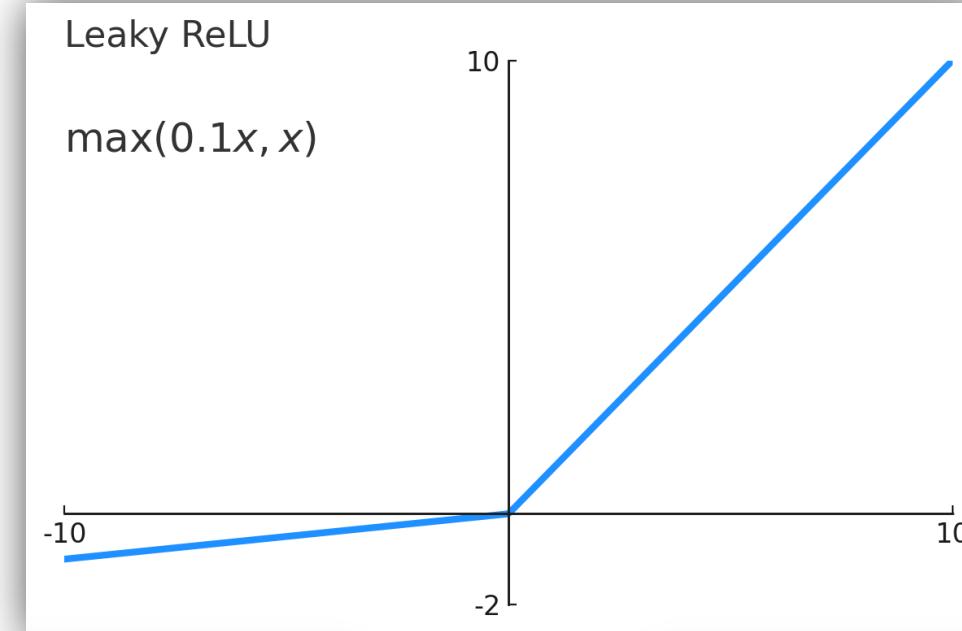
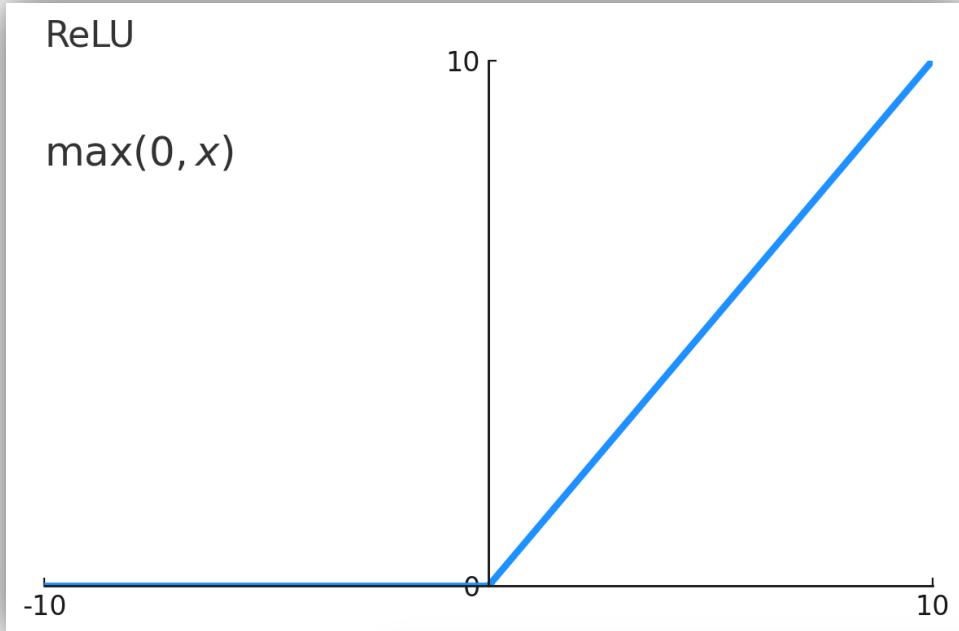
[Hendrycks et al., 2016]



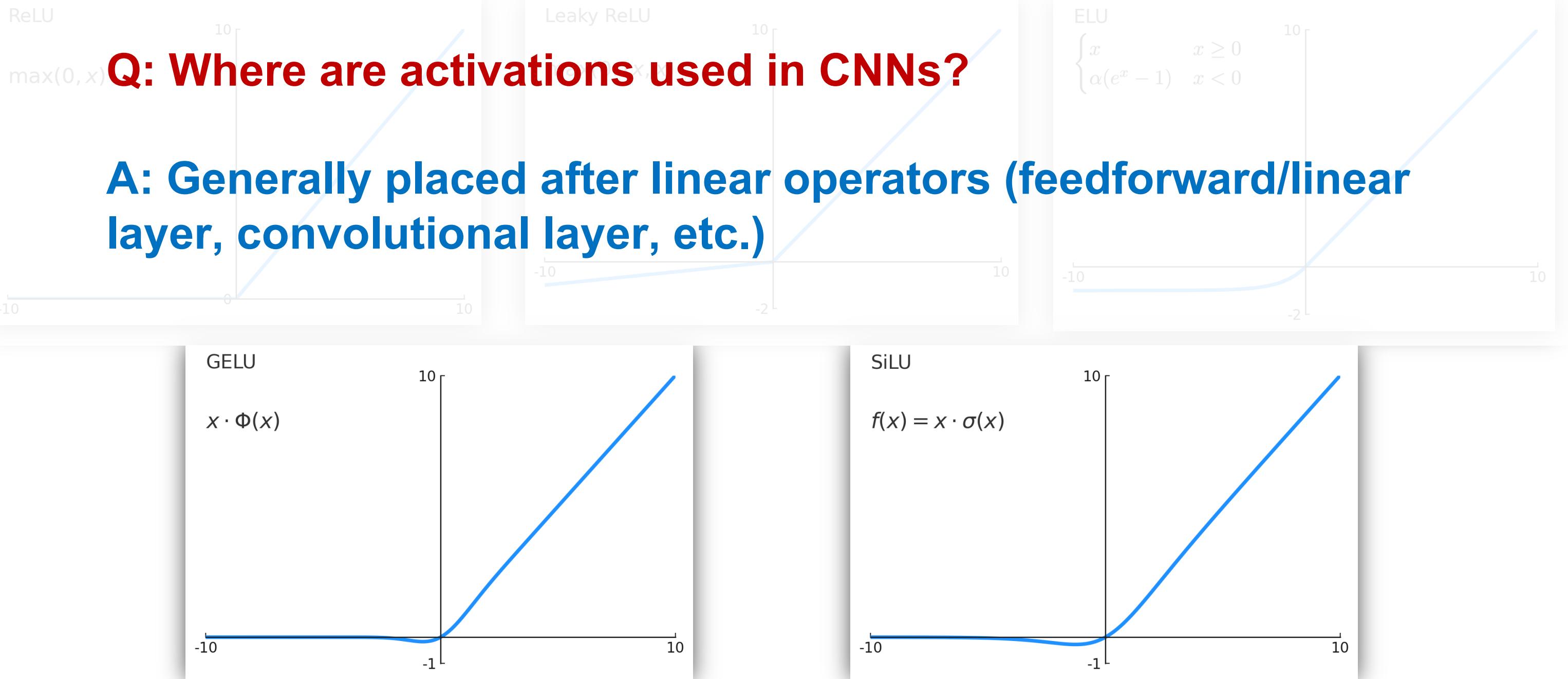
- Computes $f(x) = x^* \Phi(x)$
- Very nice behavior around 0
- Smoothness facilitates training in practice
- Higher computational cost than ReLU
- Large negative values can still have gradient $\rightarrow 0$

GELU
(Gaussian Error
Linear Unit)

Activation Function Zoo



Activation Function Zoo



Lecture Overview – Two Broad Sets of Topics

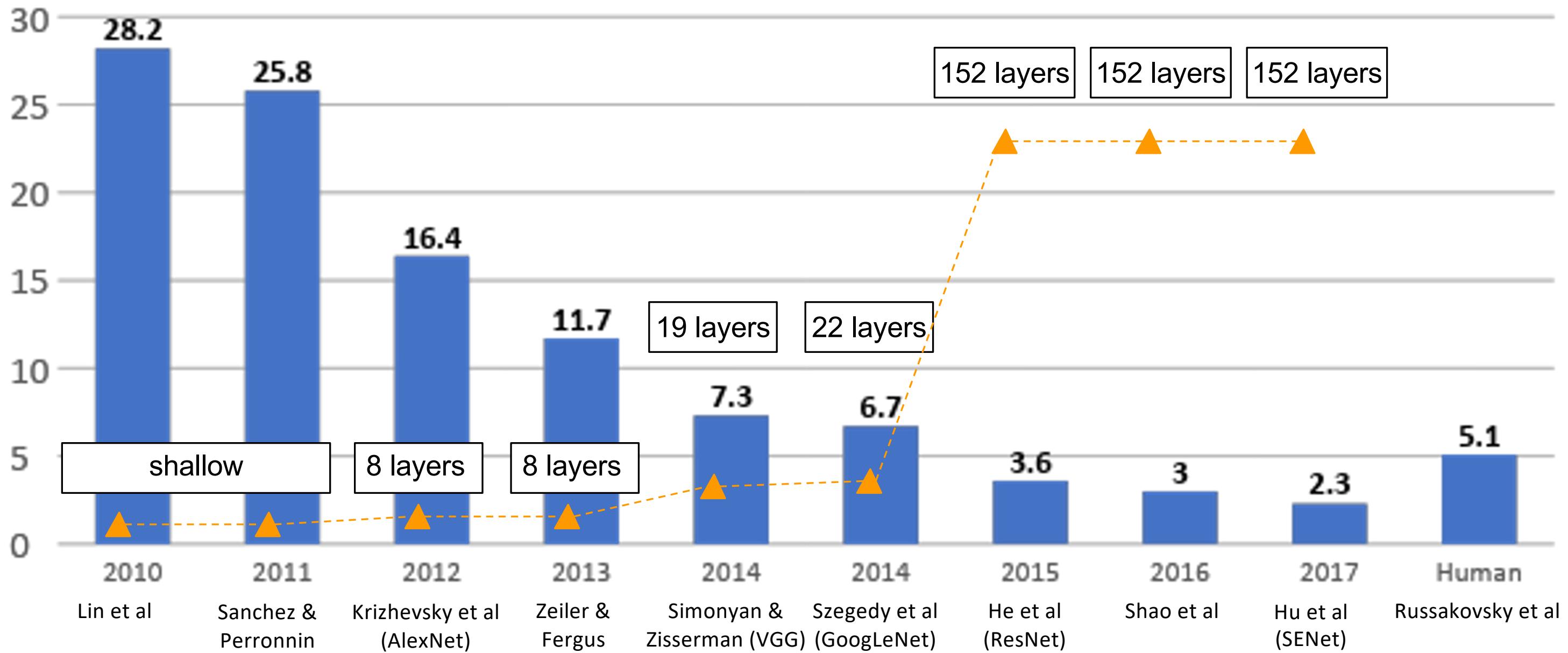
How to build CNNs?

Layers in CNNs
Activation Functions
CNN Architectures
Weight Initialization

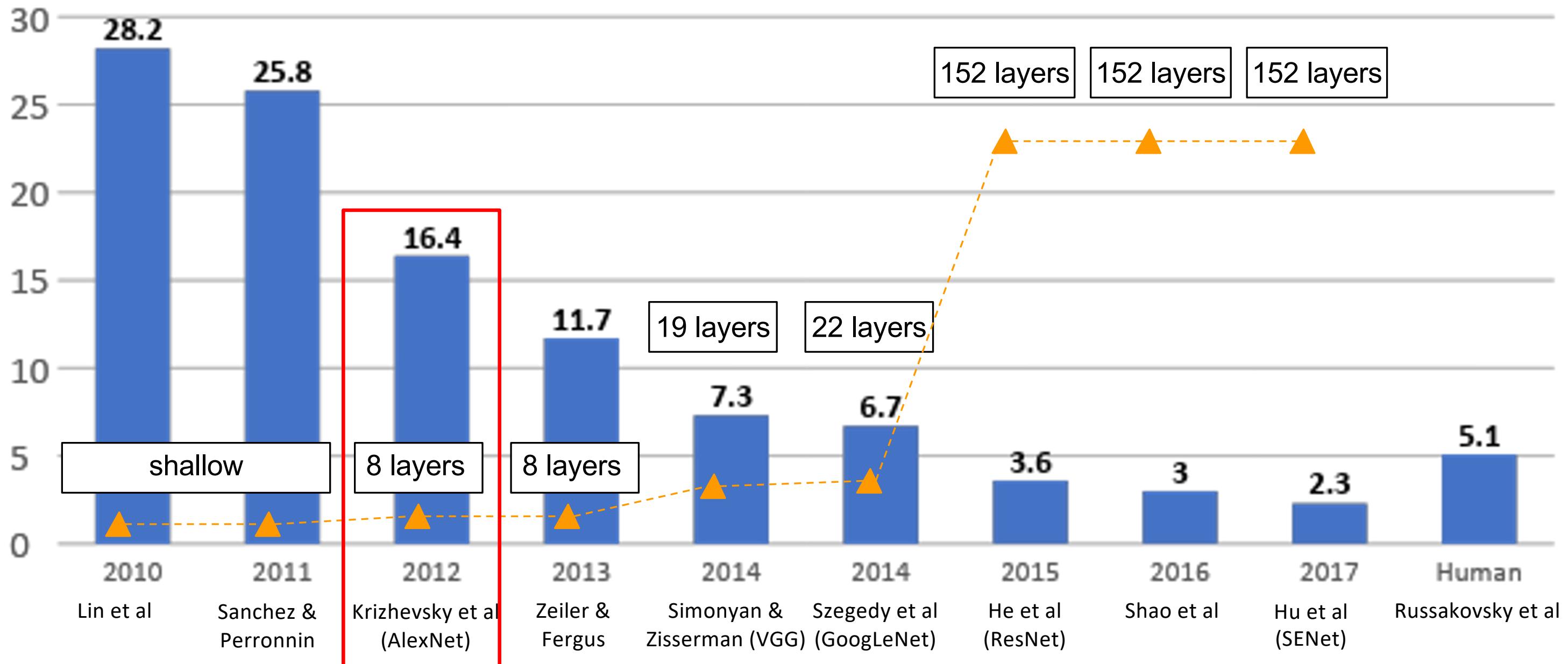
How to train CNNs?

Data Preprocessing
Data augmentation
Transfer Learning
Hyperparameter Selection

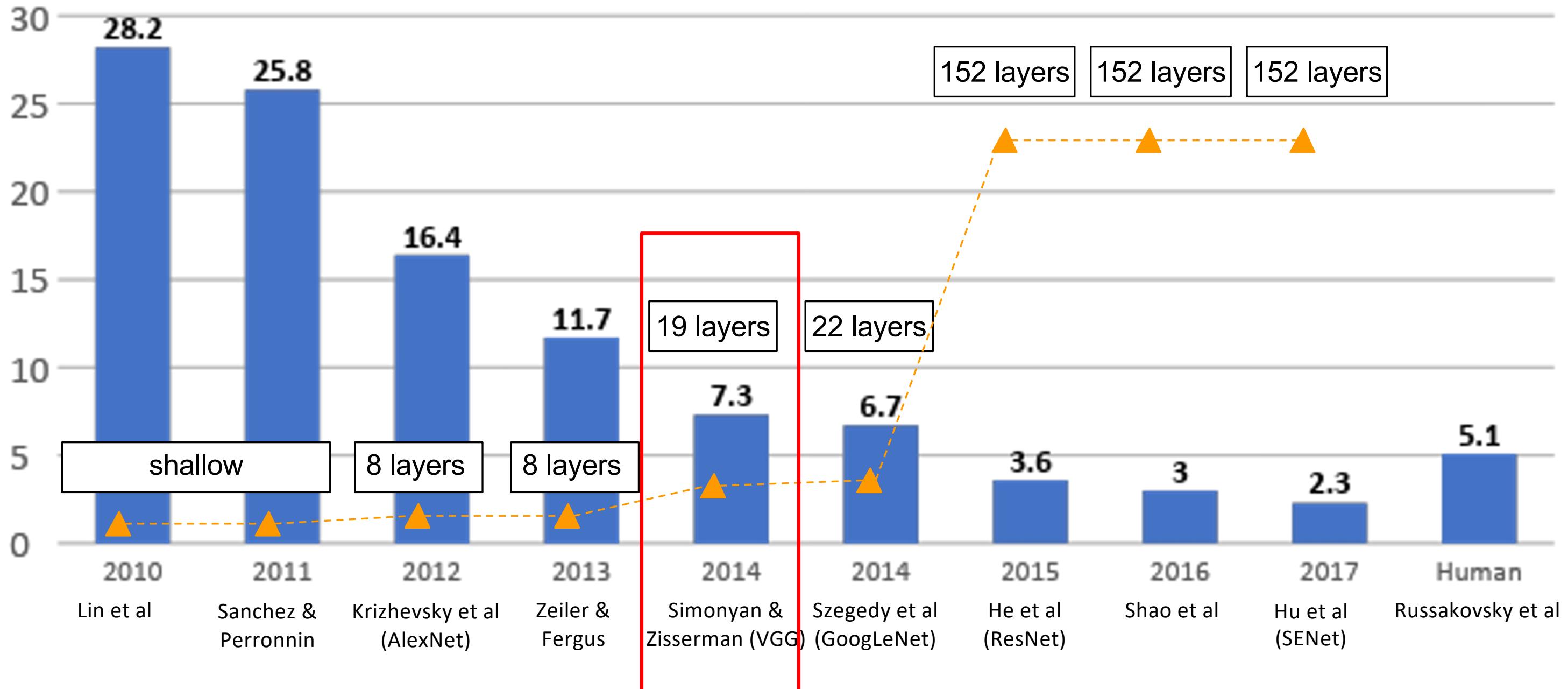
ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Small filters, Deeper networks

8 layers (AlexNet)

-> 16 - 19 layers (VGG16Net)

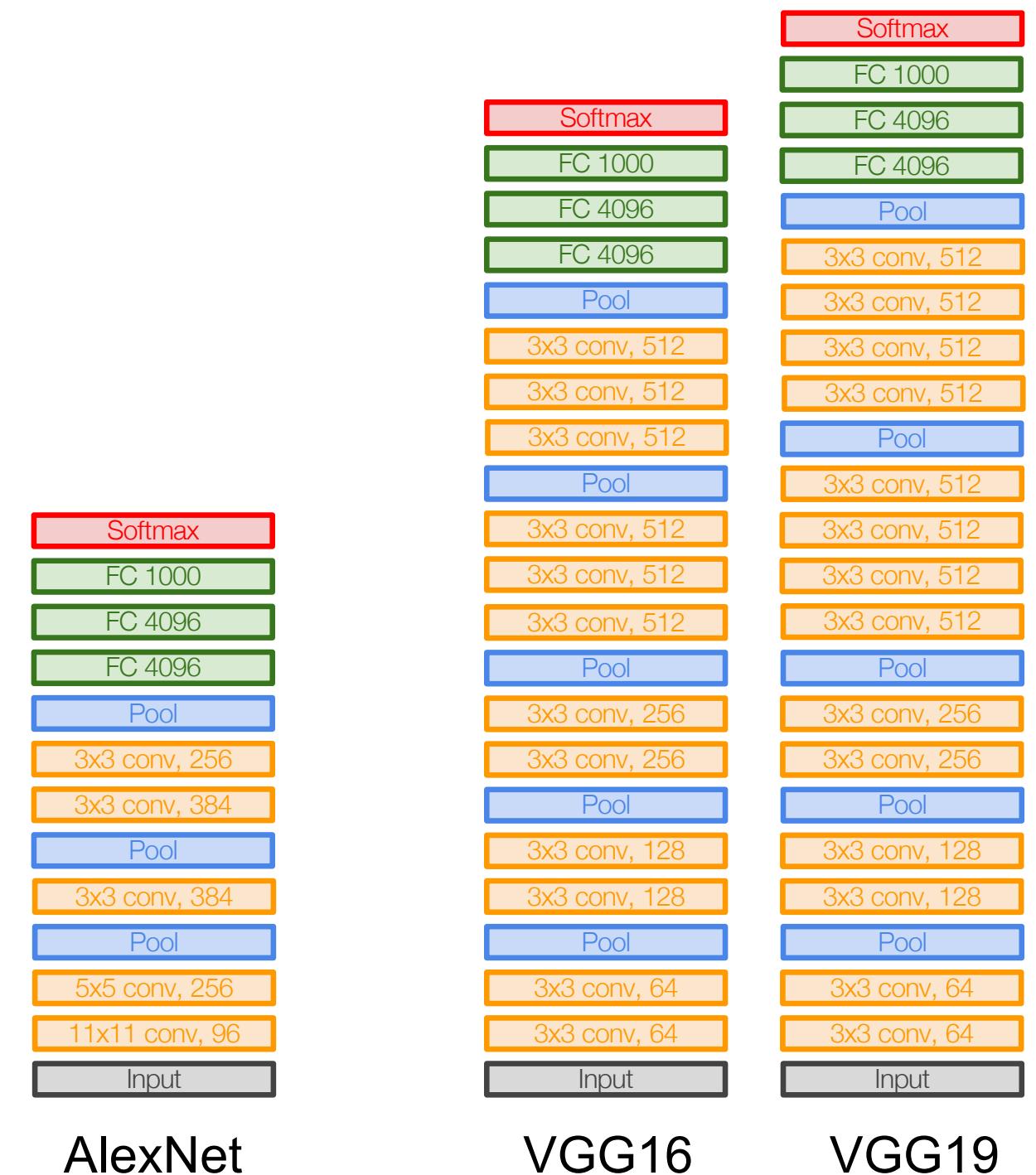
Only 3x3 CONV stride 1, pad 1

and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13

(ZFNet)

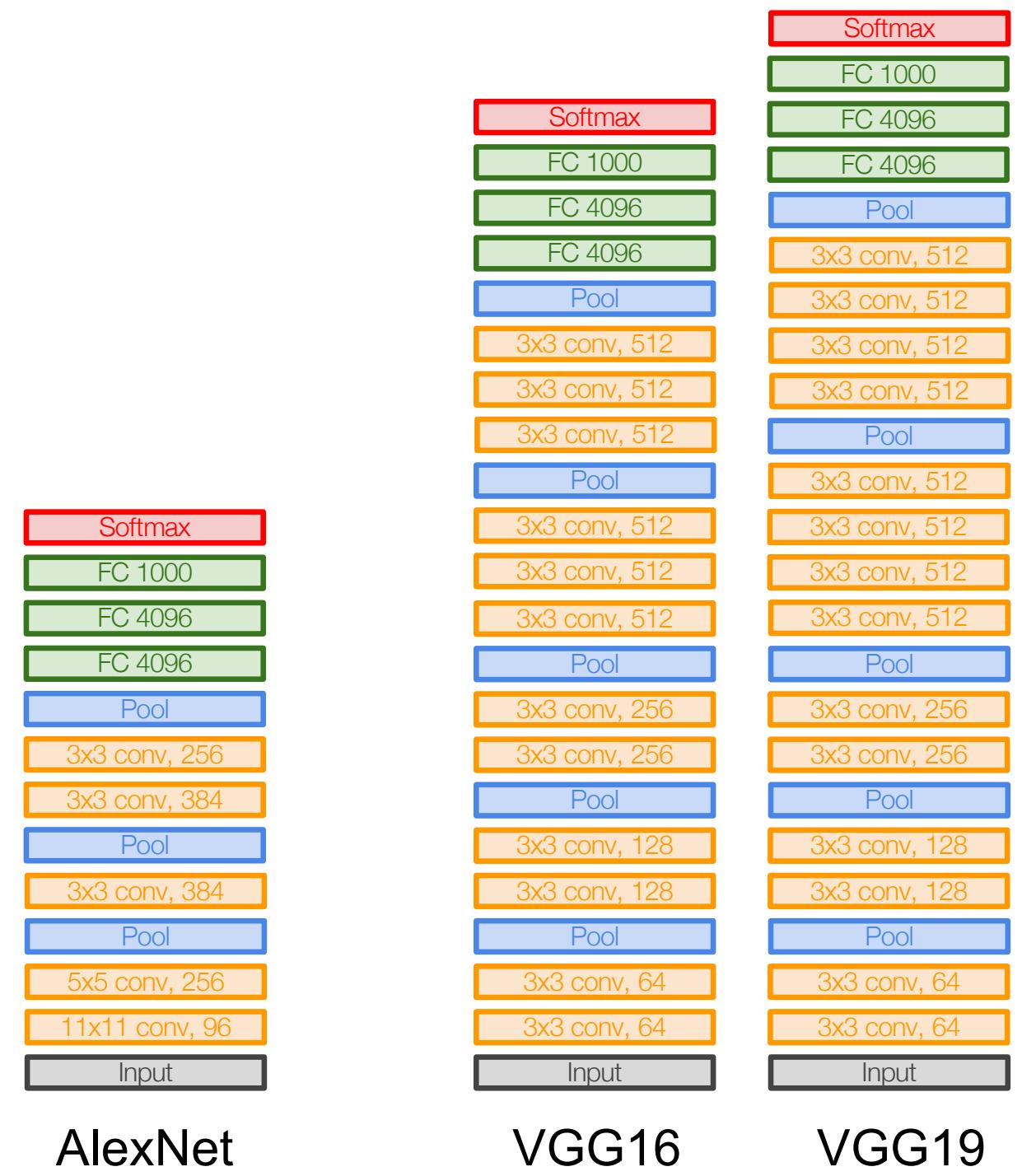
-> 7.3% top 5 error in ILSVRC'14



Case Study: VGGNet

[Simonyan and Zisserman, 2014]

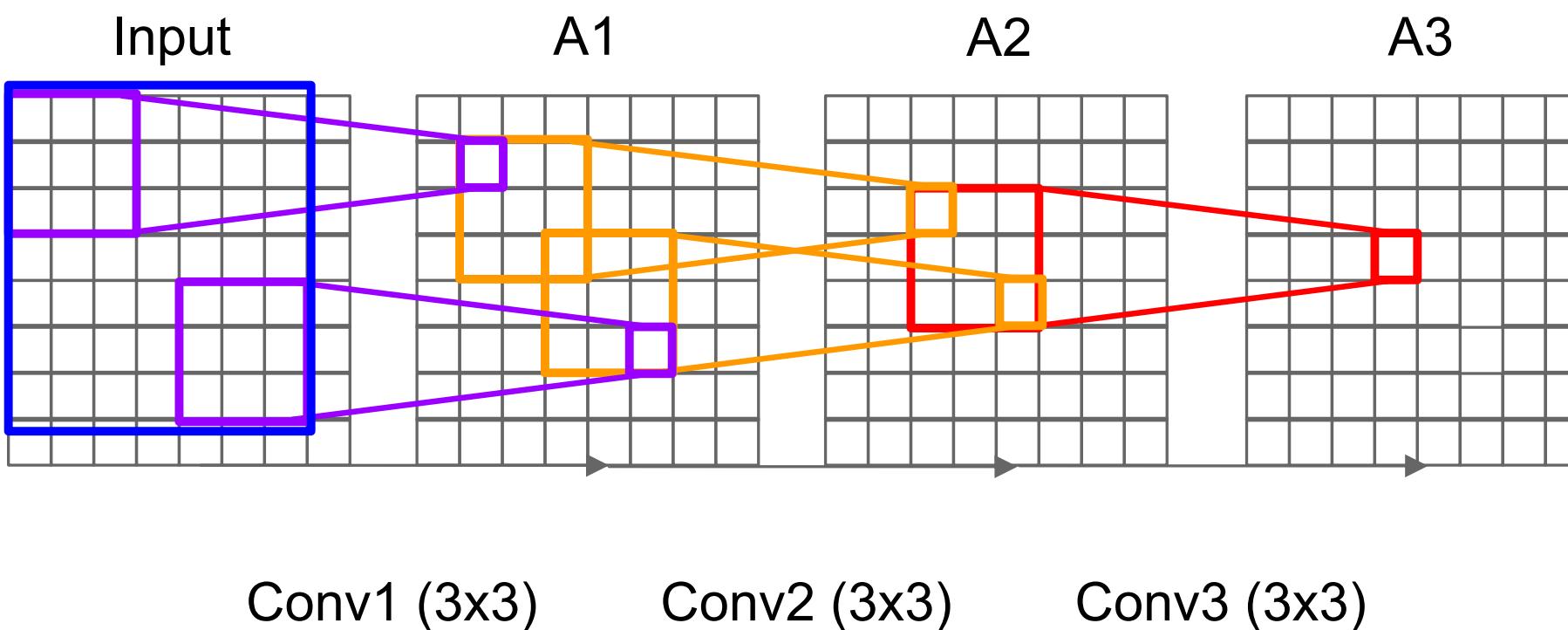
Q: Why use smaller filters? (3x3 conv)



Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: What is the effective receptive field
of three 3x3 conv (stride 1) layers?



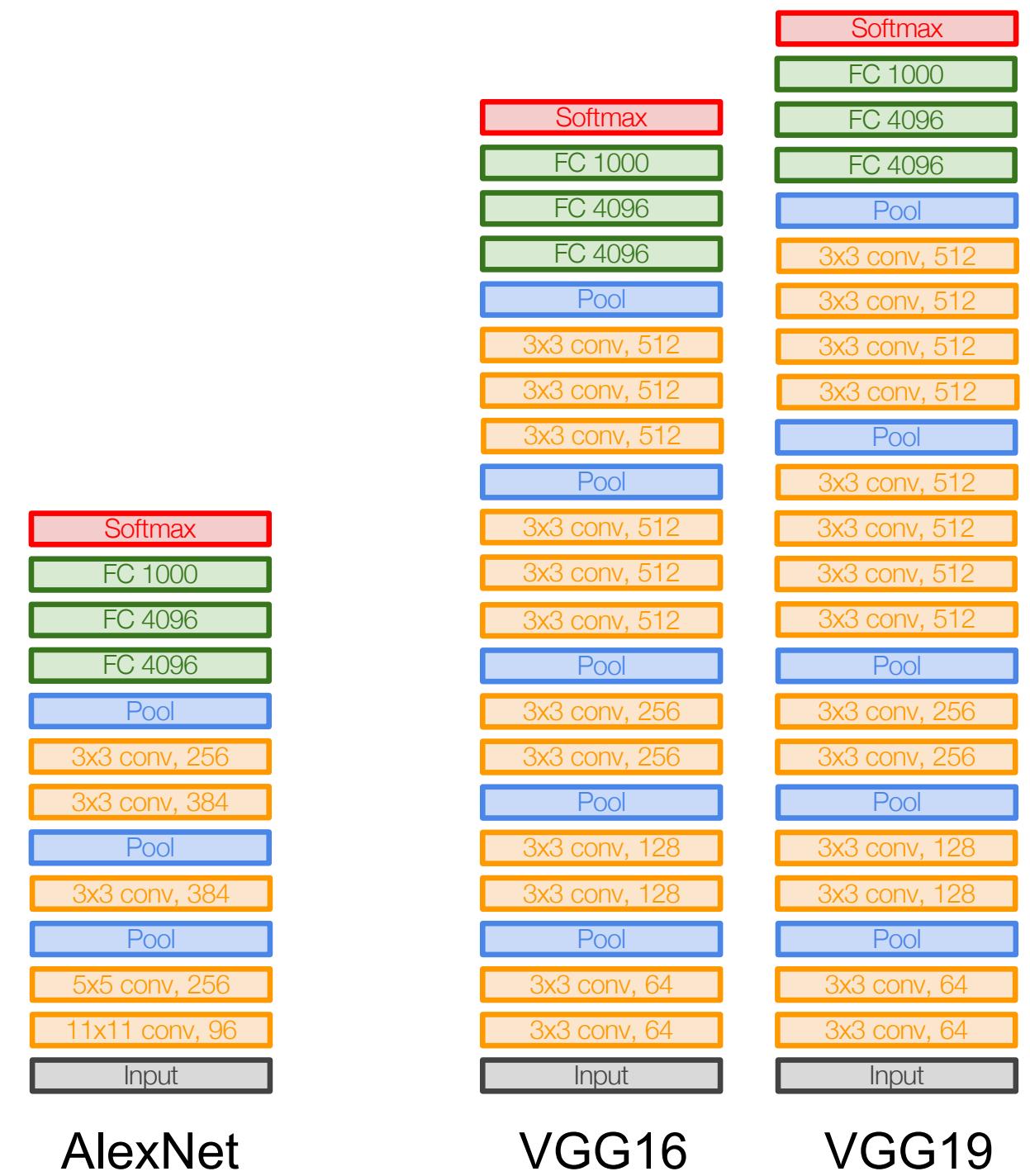
Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

[7x7]



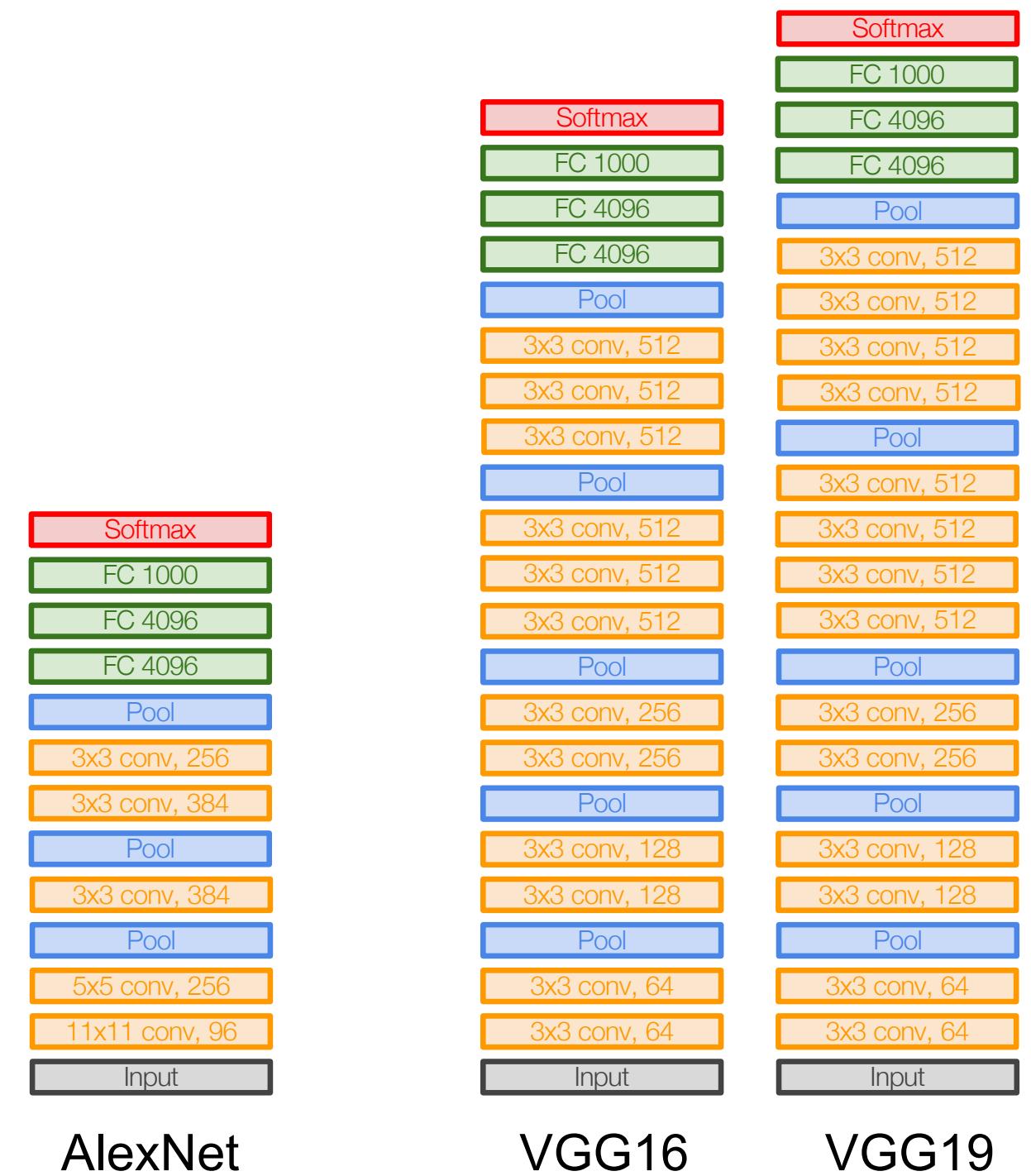
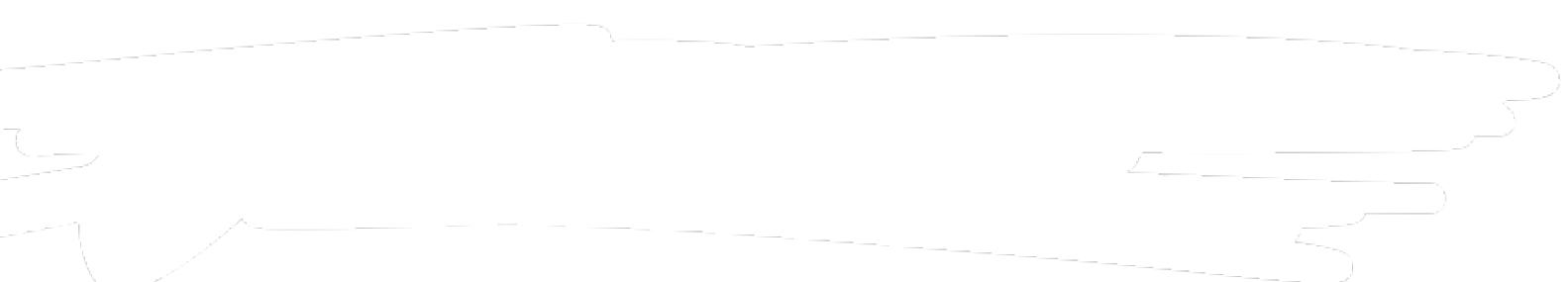
Case Study: VGGNet

[Simonyan and Zisserman, 2014]

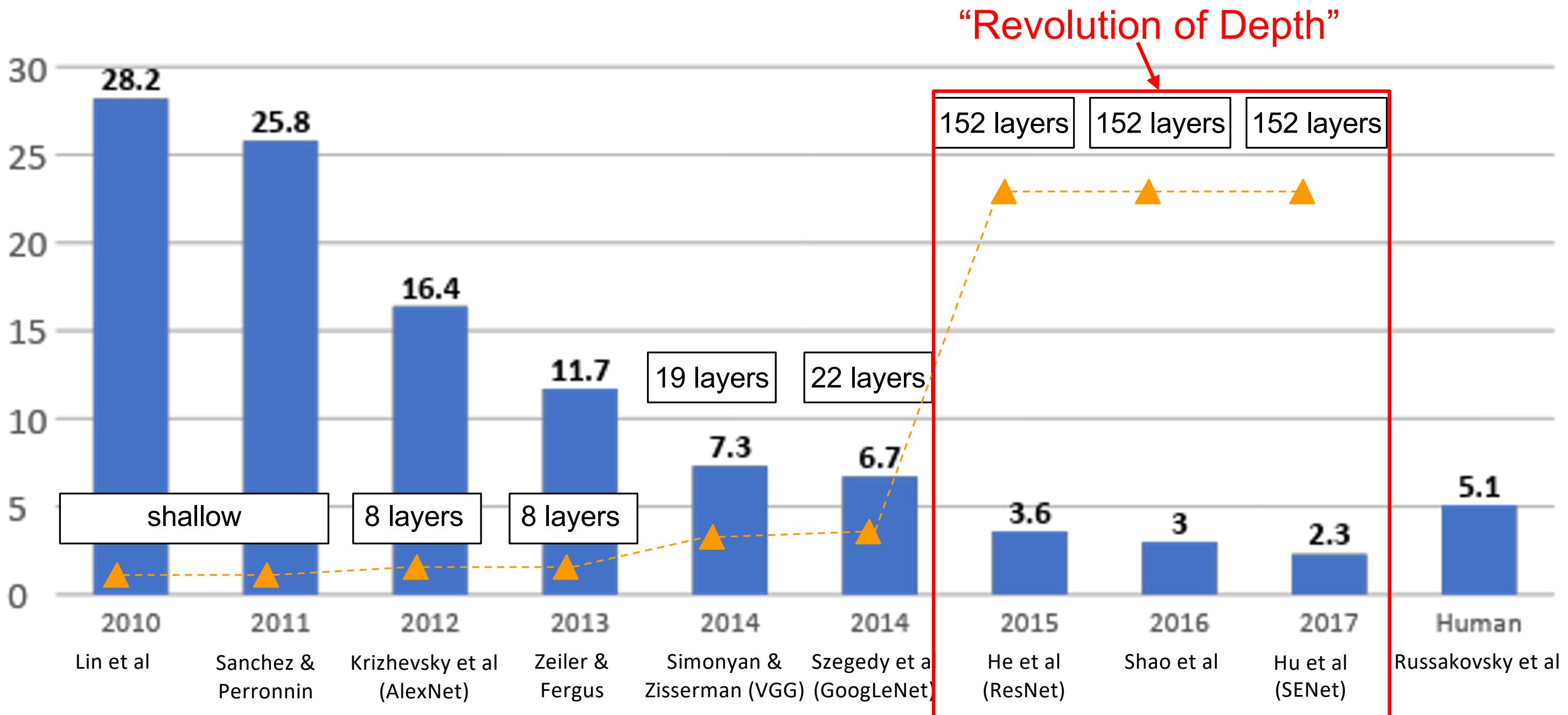
Q: Why use smaller filters? (3x3 conv)

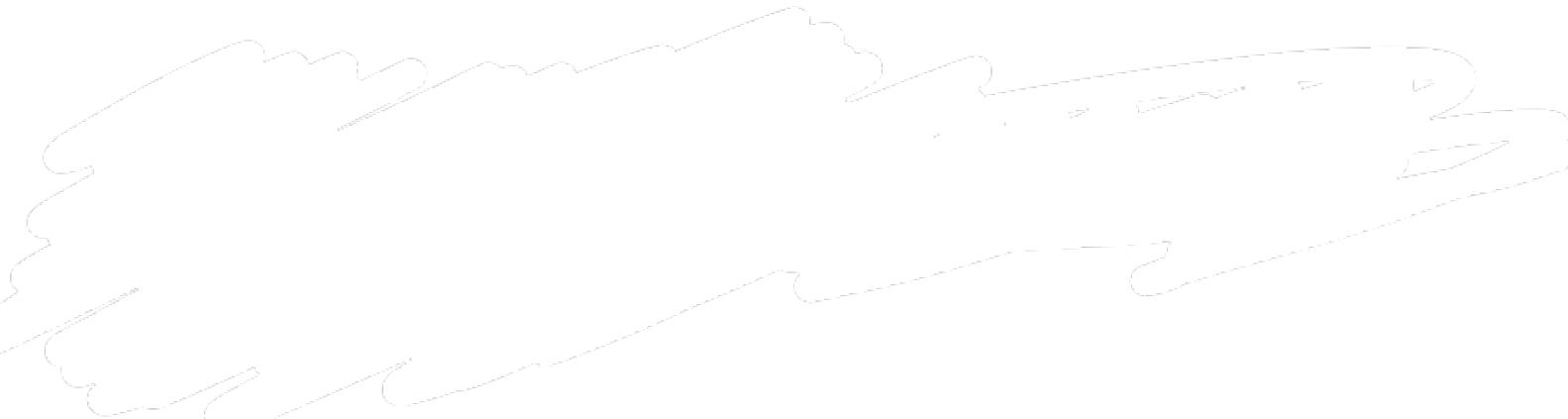
Stack of three 3x3 conv (stride 1) layers
has same **effective receptive field** as
one 7x7 conv layer

But deeper, more non-linearities

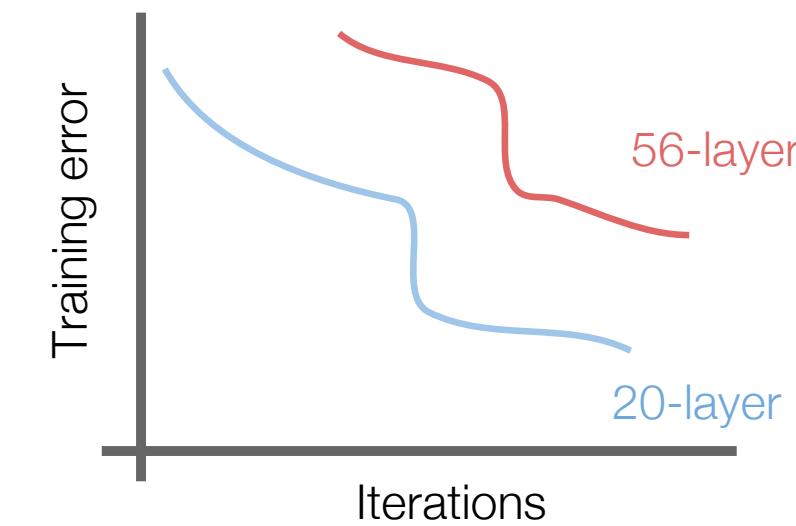
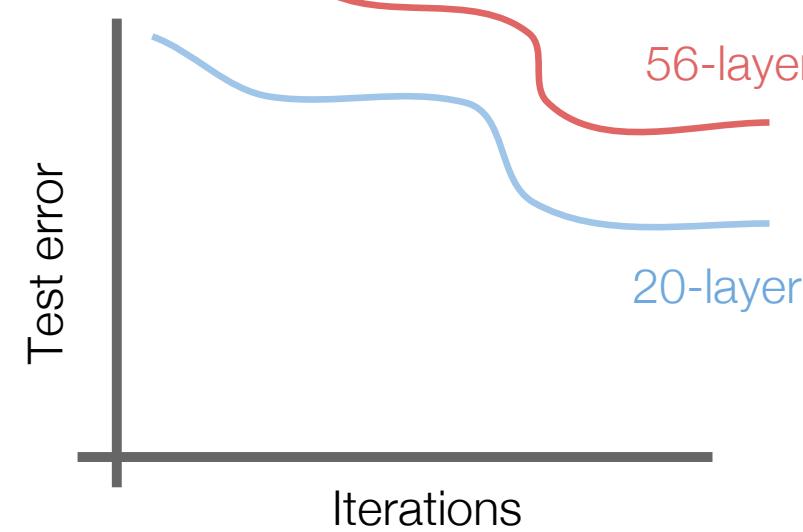


ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners





What happens when we continue stacking deeper layers on a “plain” convolutional neural network?



56-layer model performs worse on both test and training error
-> The deeper model performs worse, but it's not caused by overfitting!

Case Study: ResNet

[He et al., 2015]

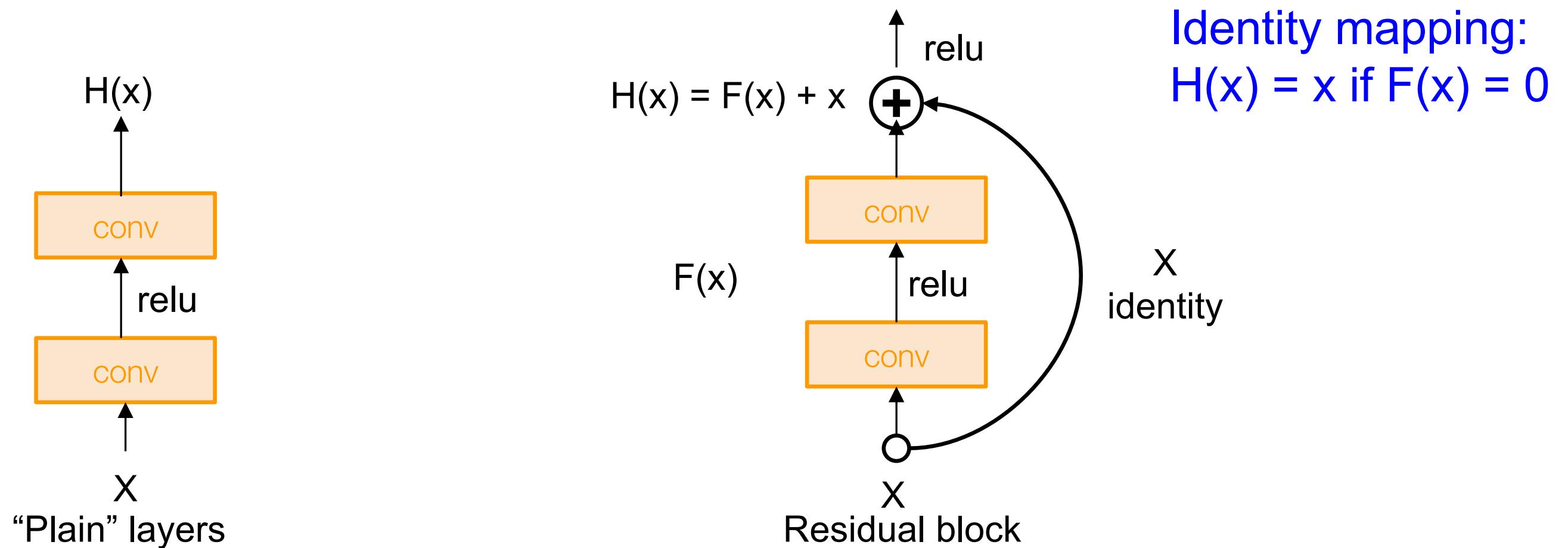
Fact: Deep models have more representation power
(more parameters) than shallower models.

Hypothesis: the problem is an *optimization* problem,
deeper models are harder to optimize

Case Study: ResNet

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

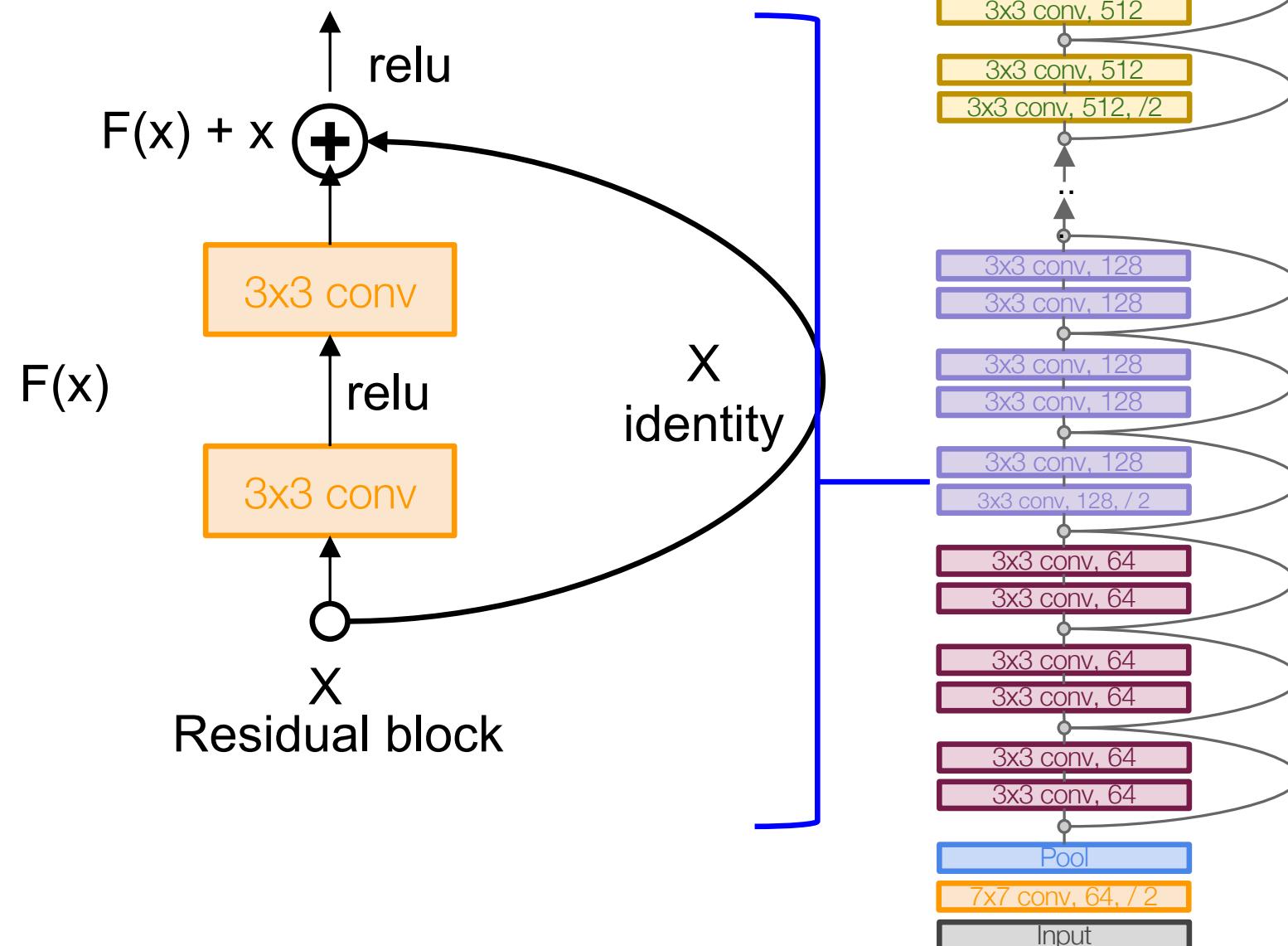


Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers

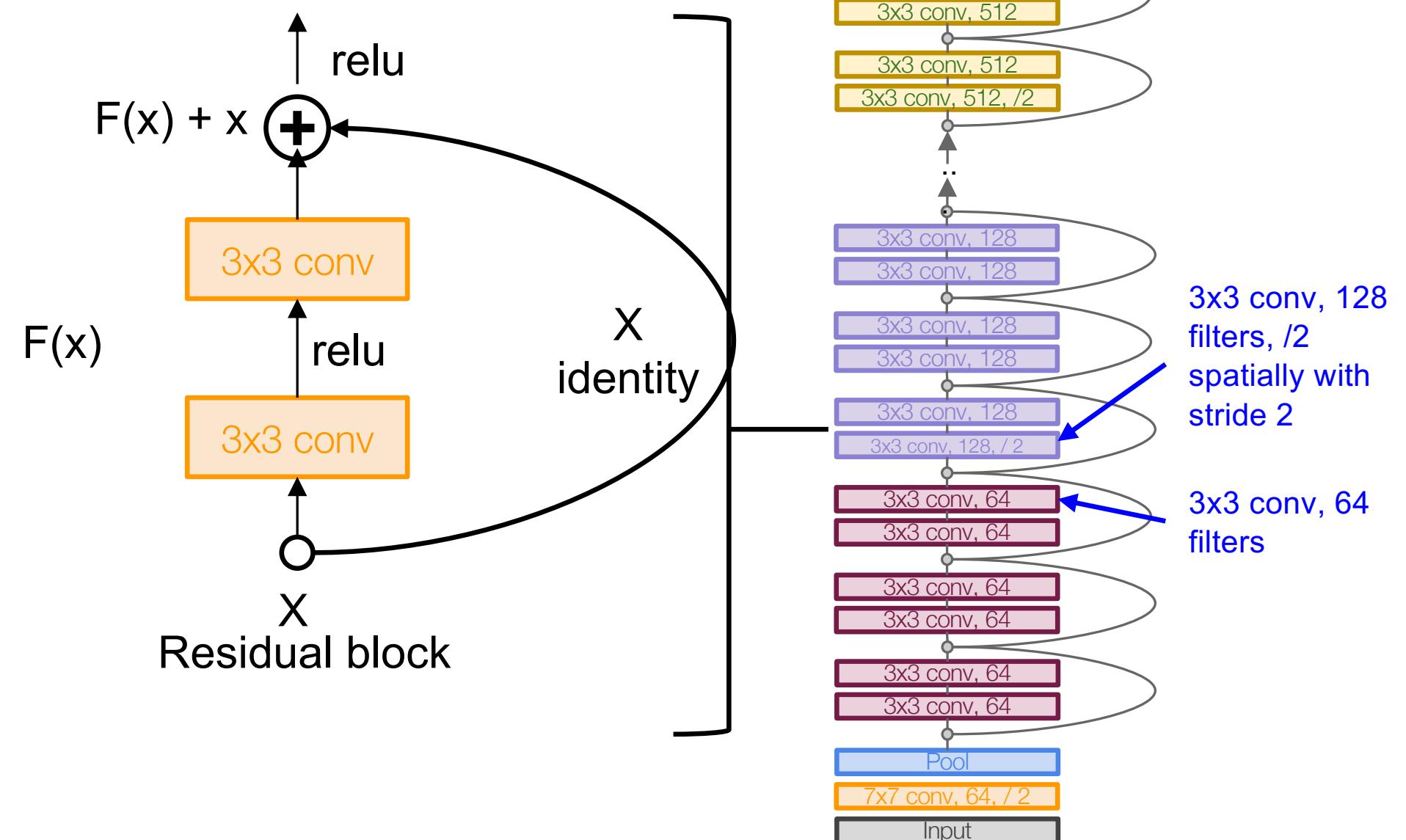


Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
Reduce the activation volume by half.

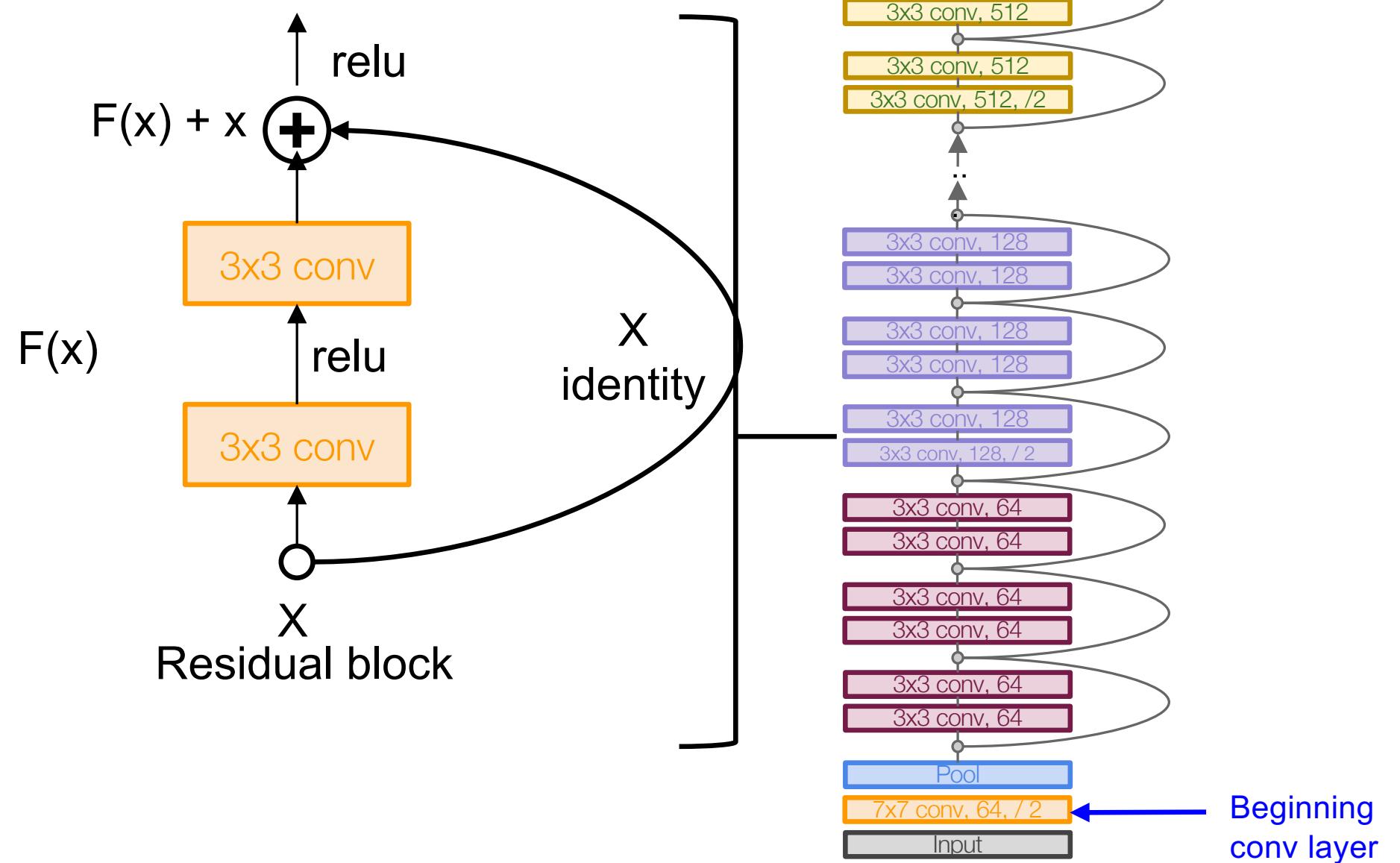


Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

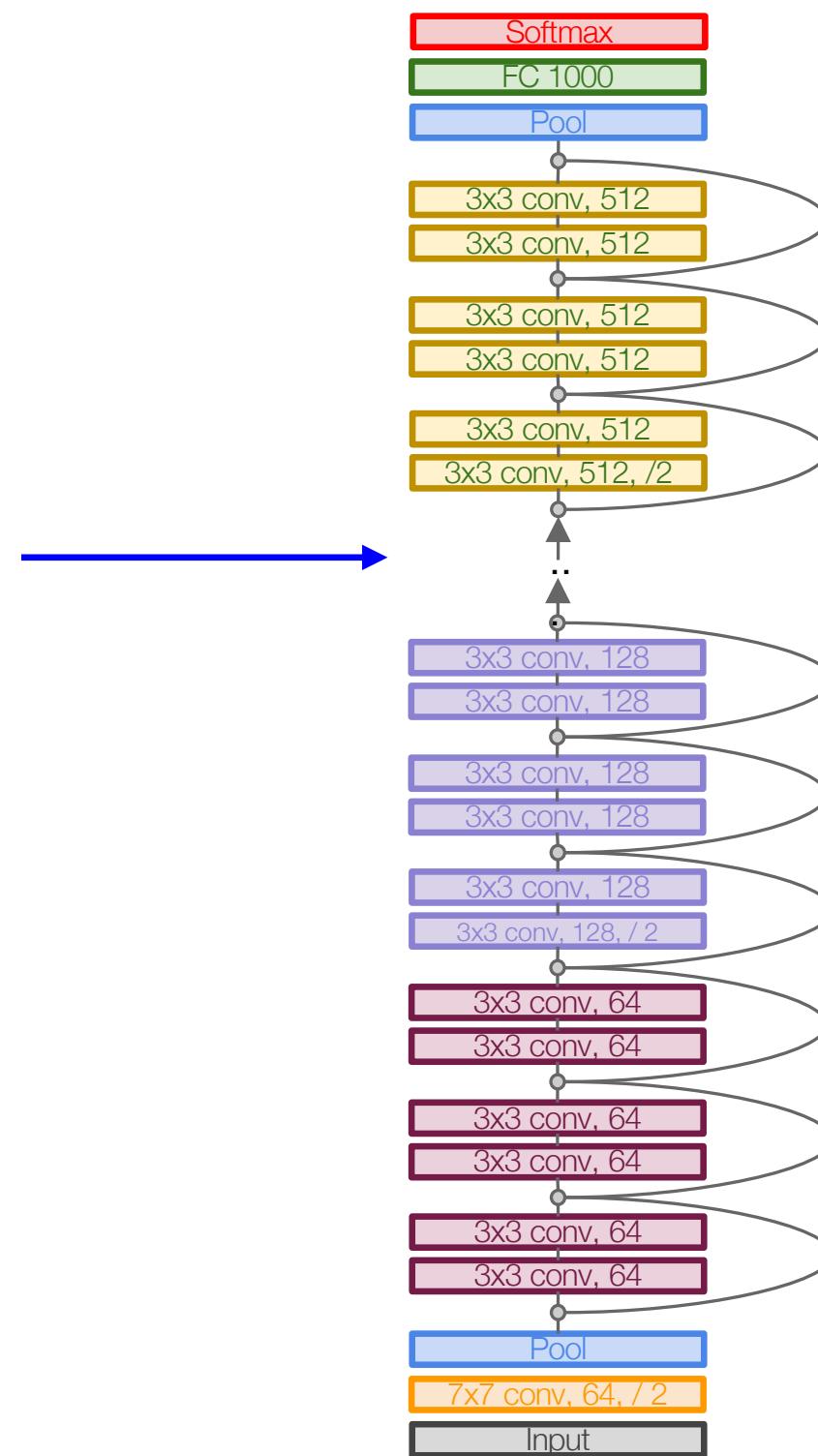
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning (stem)



Case Study: ResNet

[He et al., 2015]

Total depths of 18, 34, 50, 101, or 152 layers for ImageNet



Lecture Overview – Two Broad Sets of Topics

How to build CNNs?

Layers in CNNs
Activation Functions
CNN Architectures
Weight Initialization

How to train CNNs?

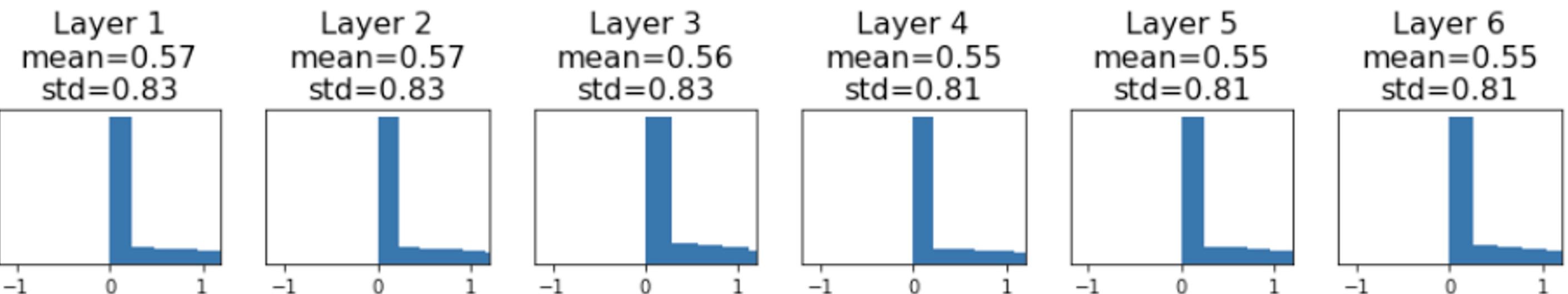
Data Preprocessing
Data augmentation
Transfer Learning
Hyperparameter Selection

How to initialize weights in neural network layers?

One solution: Kaiming / MSRA Initialization

```
dims = [4096] * 7    ReLU correction: std = sqrt(2 / Din)
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(2/Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

“Just right”: Activations are nicely scaled for all layers!



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

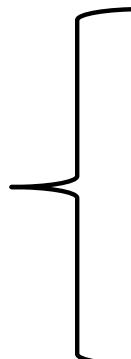
Lecture Overview – Two Broad Sets of Topics

How to build CNNs?



- Layers in CNNs
- Activation Functions
- CNN Architectures
- Weight Initialization

How to train CNNs?



- Data Preprocessing
- Data augmentation
- Transfer Learning
- Hyperparameter Selection

TLDR for Image Normalization: center and scale for each channel

- Subtract per-channel mean and Divide by per-channel std (almost all modern models)
(stats along each channel = 3 numbers)
- Requires pre-computing means and std for each pixel channel (given your dataset)

```
norm_pixel[i,j,c] = (pixel[i,j,c] - np.mean(pixel[:, :, c])) / np.std(pixel[:, :, c])
```

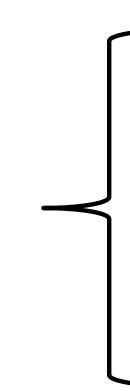
Lecture Overview – Two Broad Sets of Topics

How to build CNNs?



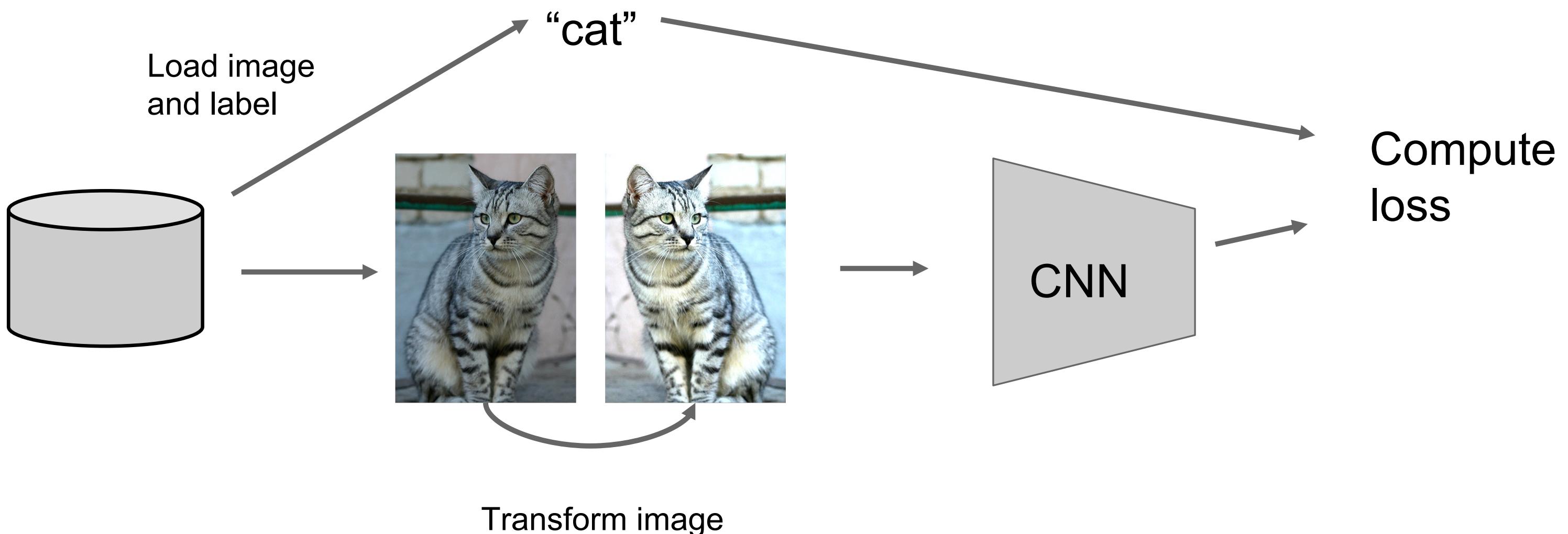
- Layers in CNNs
- Activation Functions
- CNN Architectures
- Weight Initialization

How to train CNNs?



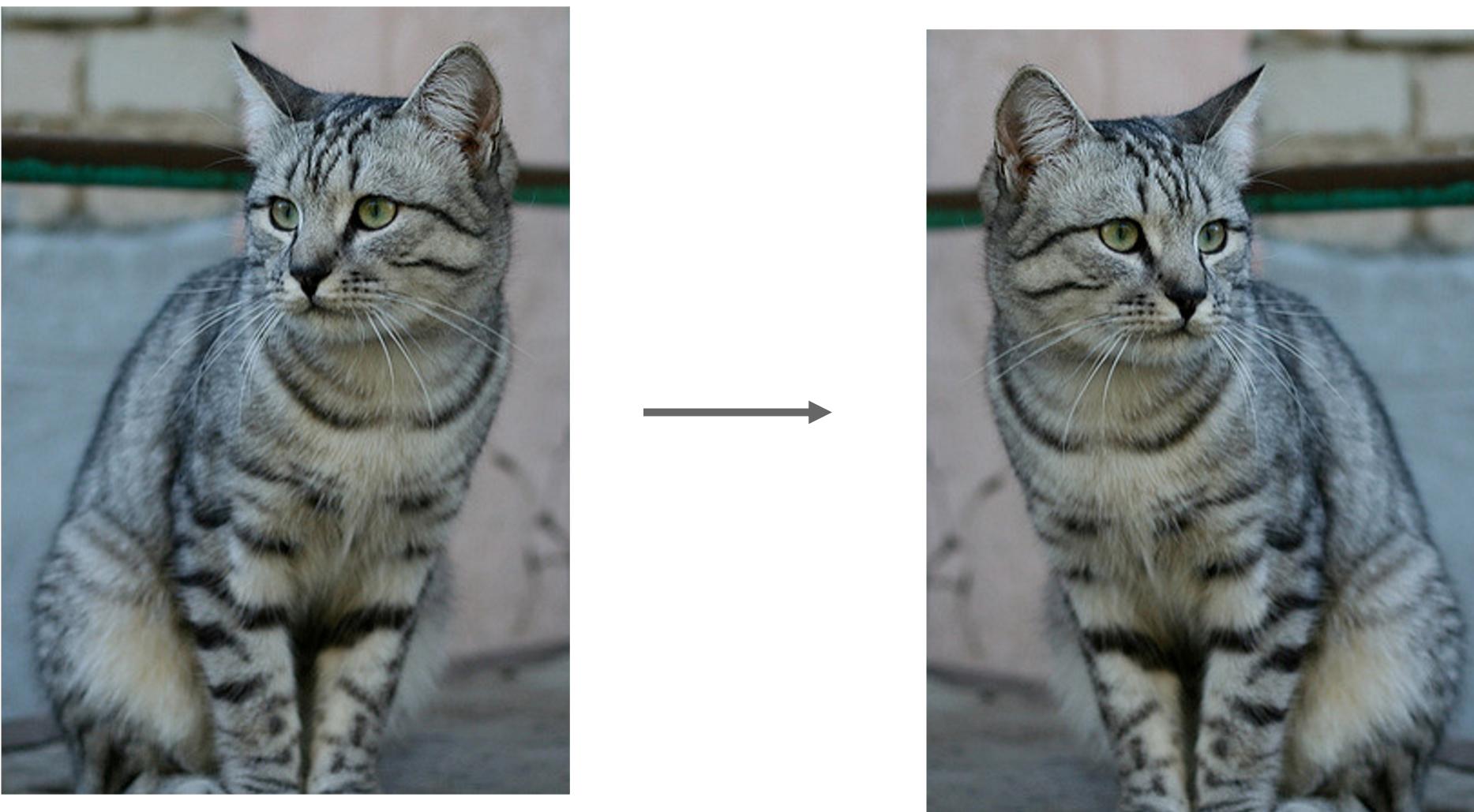
- Data Preprocessing
- Data augmentation**
- Transfer Learning
- Hyperparameter Selection

Regularization: Data Augmentation



Data Augmentation

Horizontal Flips



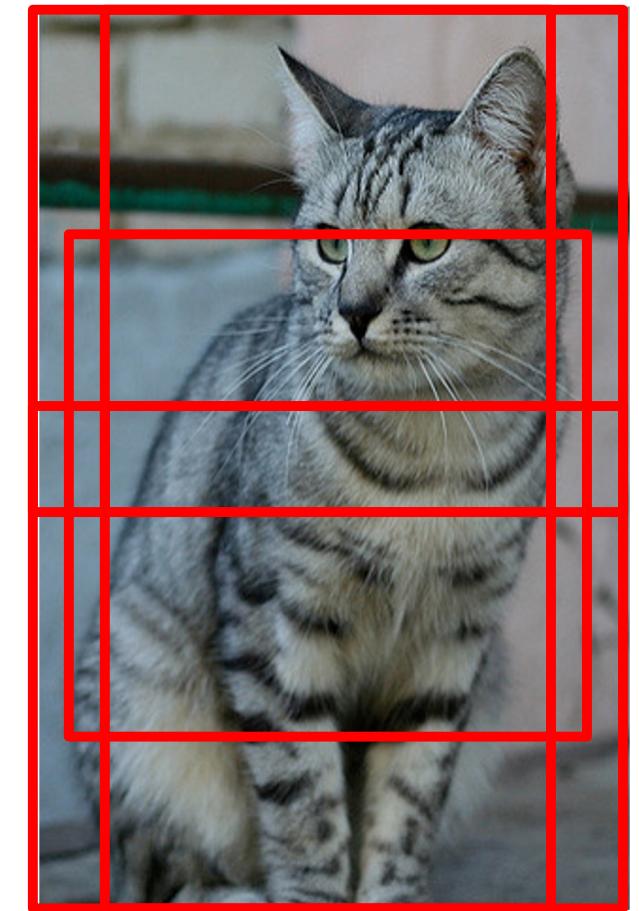
Data Augmentation

Random crops and scales

Training: sample random crops / scales

ResNet:

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224×224 patch



Test Time Augmentation: average a fixed set of crops

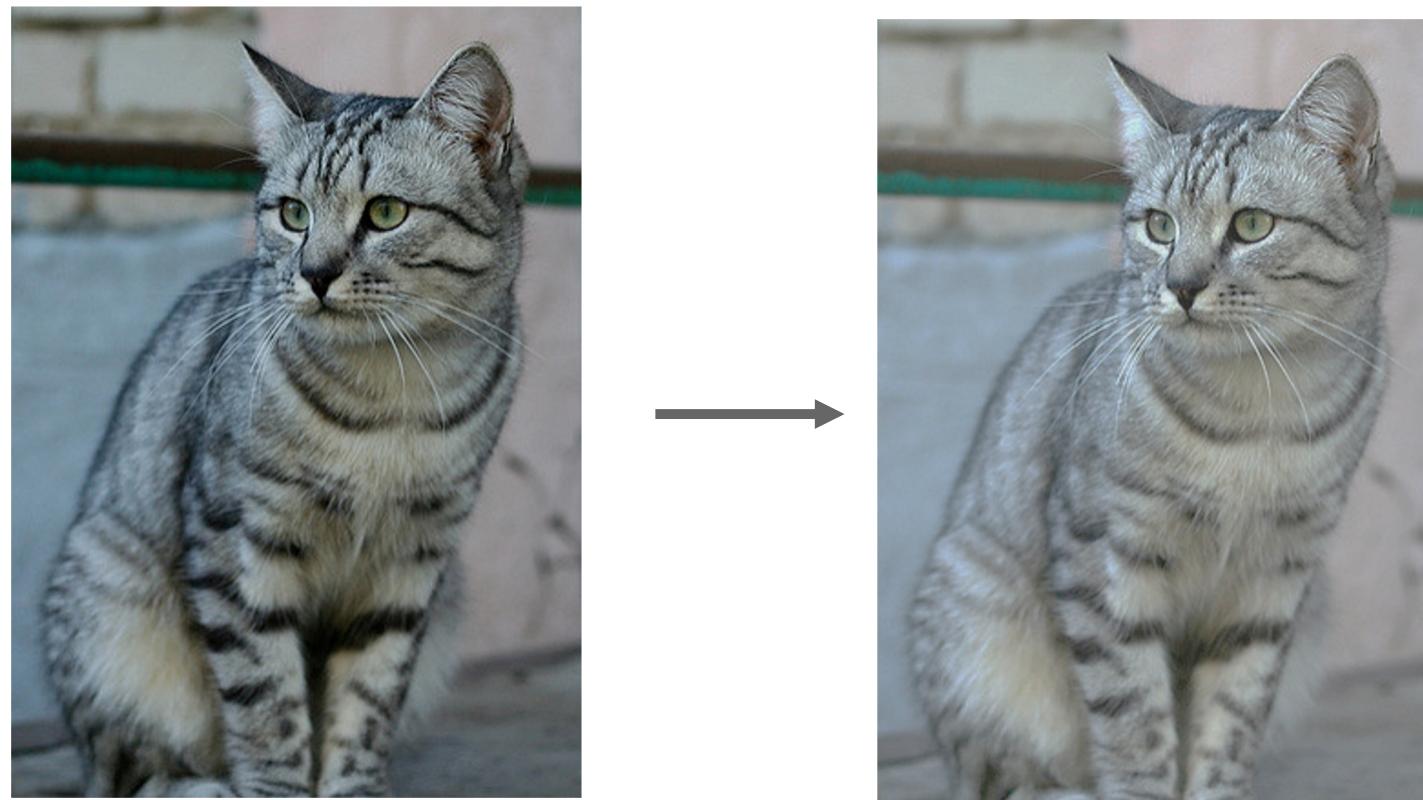
ResNet:

1. Resize image at 5 scales: {224, 256, 384, 480, 640}
2. For each size, use 10 224×224 crops: 4 corners + center, + flips

Data Augmentation

Color Jitter

Simple: Randomize
contrast and brightness



Regularization: Cutout

Training: Set random image regions to zero

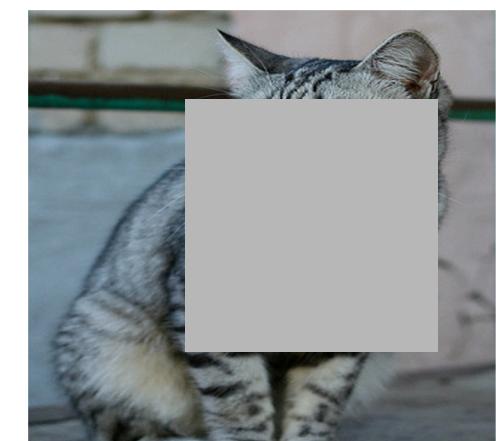
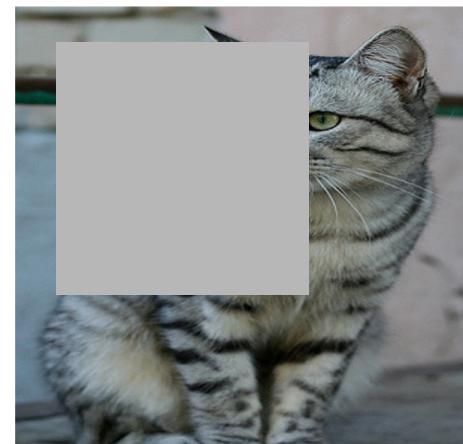
Testing: Use full image

Examples:

Dropout

Data Augmentation

Cutout / Random Crop

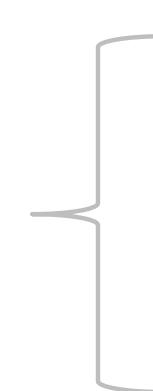


Works very well for small datasets like CIFAR,
less common for large datasets like ImageNet

DeVries and Taylor, "Improved Regularization of
Convolutional Neural Networks with Cutout", arXiv 2017

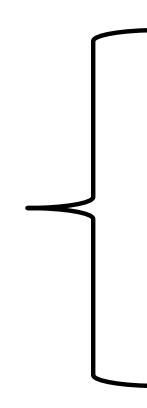
Lecture Overview – Two Broad Sets of Topics

How to build CNNs?



- Layers in CNNs
- Activation Functions
- CNN Architectures
- Weight Initialization

How to train CNNs?



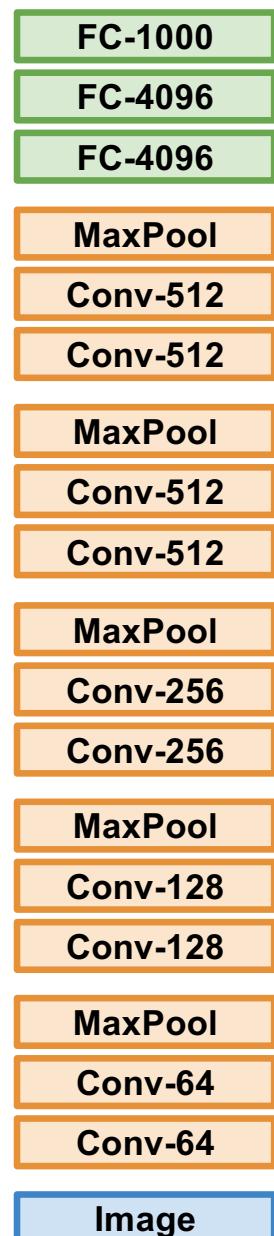
- Data Preprocessing
- Data augmentation
- Transfer Learning**
- Hyperparameter Selection

What if you **don't have** a lot of data? Can
you still train CNNs?

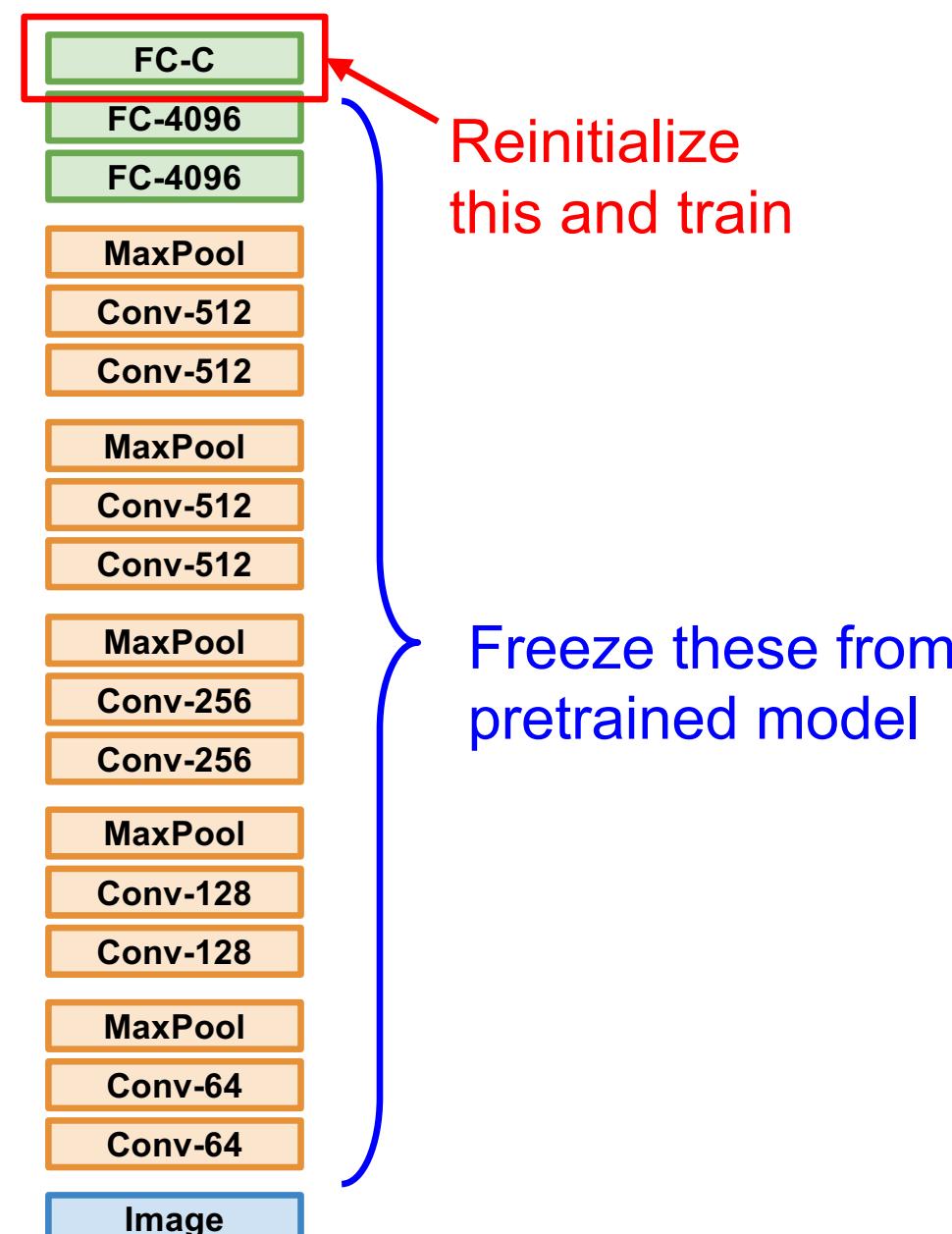
Transfer Learning with CNNs

Donahue et al, “DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition”, ICML 2014
Razavian et al, “CNN Features Off-the-Shelf: An Astounding Baseline for Recognition”, CVPR Workshops 2014

1. Train on Imagenet



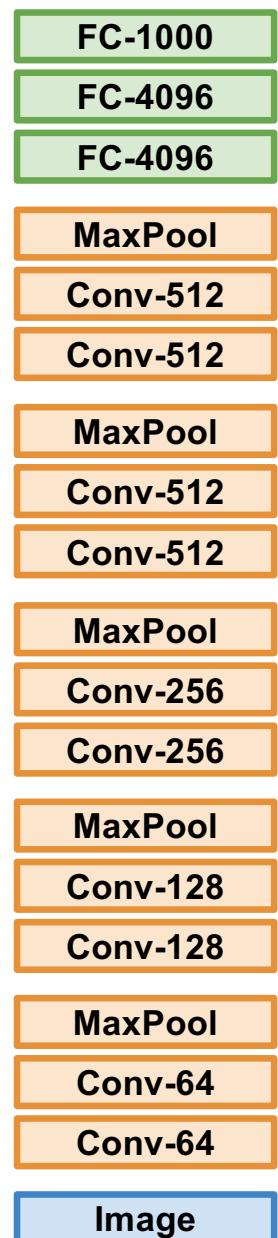
2. Small Dataset (C classes)



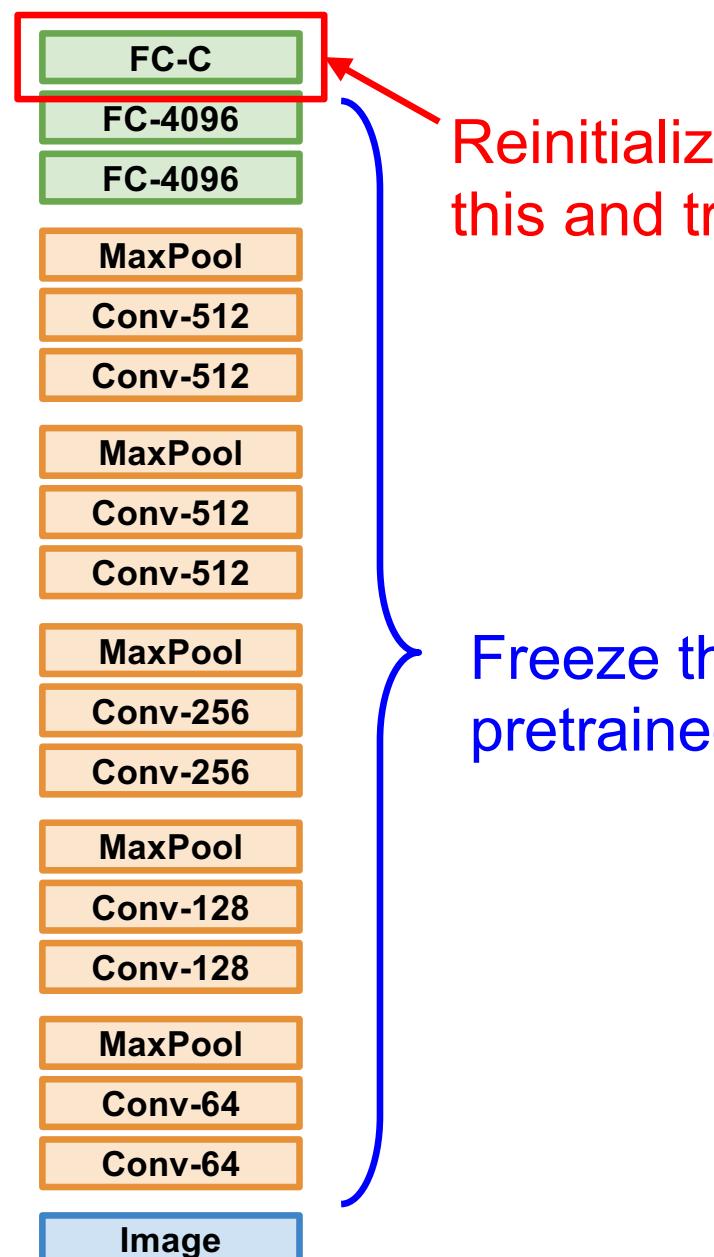
Transfer Learning with CNNs

Donahue et al, “DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition”, ICML 2014
Razavian et al, “CNN Features Off-the-Shelf: An Astounding Baseline for Recognition”, CVPR Workshops 2014

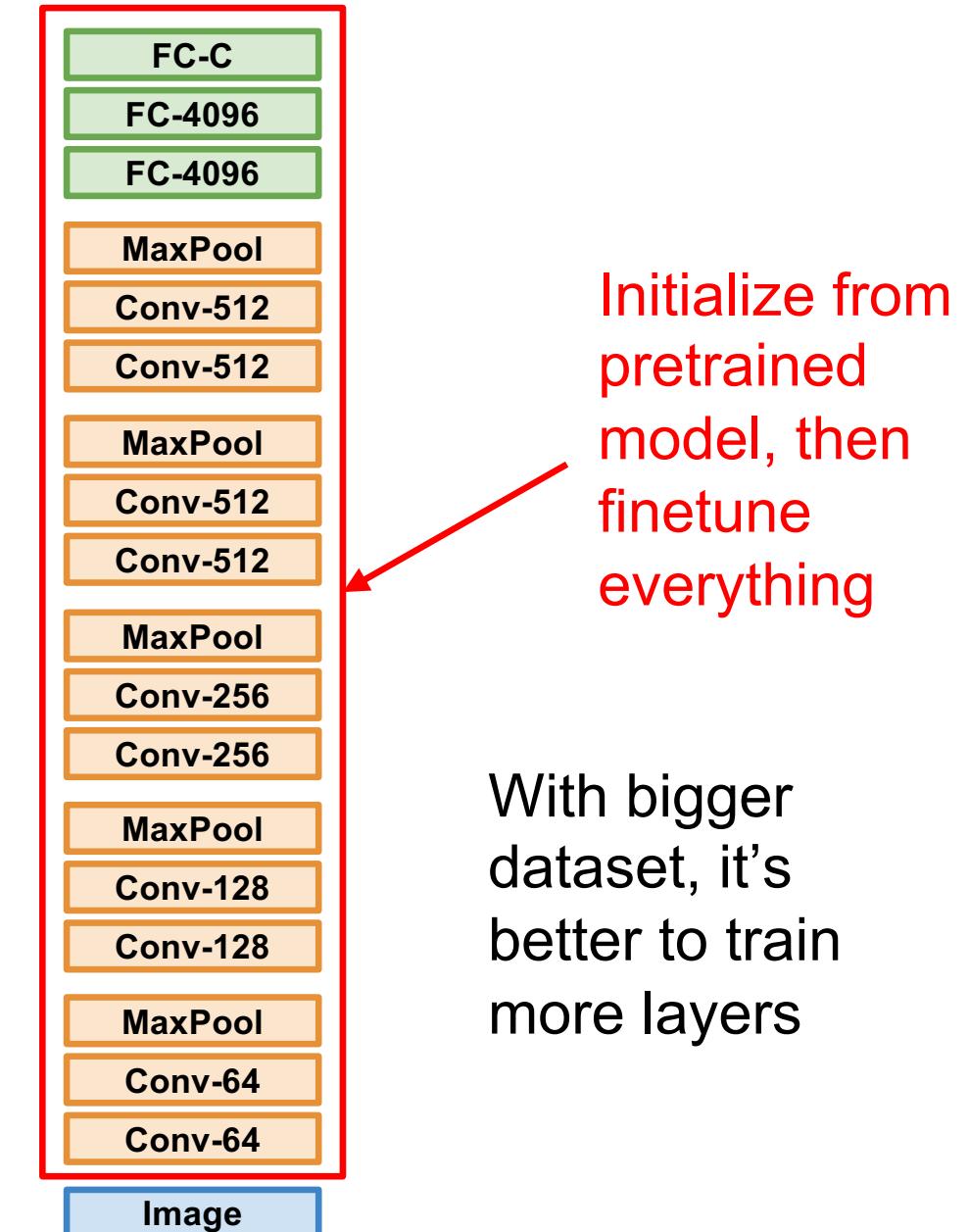
1. Train on Imagenet

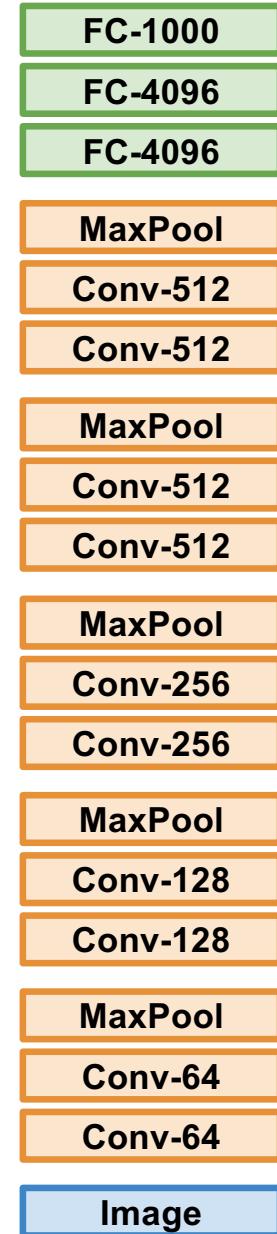


2. Small Dataset (C classes)



3. Bigger dataset





More generic
More specific

	very similar dataset	very different dataset
very little data	Use Linear Classifier on final layer	Try another model or collect more data ☹
quite a lot of data	Finetune all model layers	Either finetune all model layers or train from scratch!

Lecture Overview – Two Broad Sets of Topics

How to build CNNs?

Layers in CNNs
Activation Functions
CNN Architectures
Weight Initialization

How to train CNNs?

Data Preprocessing
Data augmentation
Transfer Learning
Hyperparameter Selection

Choosing Hyperparameters

Step 1: Check initial loss

Step 2: Overfit a small sample

Step 3: Find LR that makes loss go down

Use the architecture from the previous step, use all training data, turn on small weight decay, find a learning rate that makes the loss drop significantly within ~100 iterations

Good learning rates to try: 1e-1, 1e-2, 1e-3, 1e-4, 1e-5

Choosing Hyperparameters

Step 1: Check initial loss

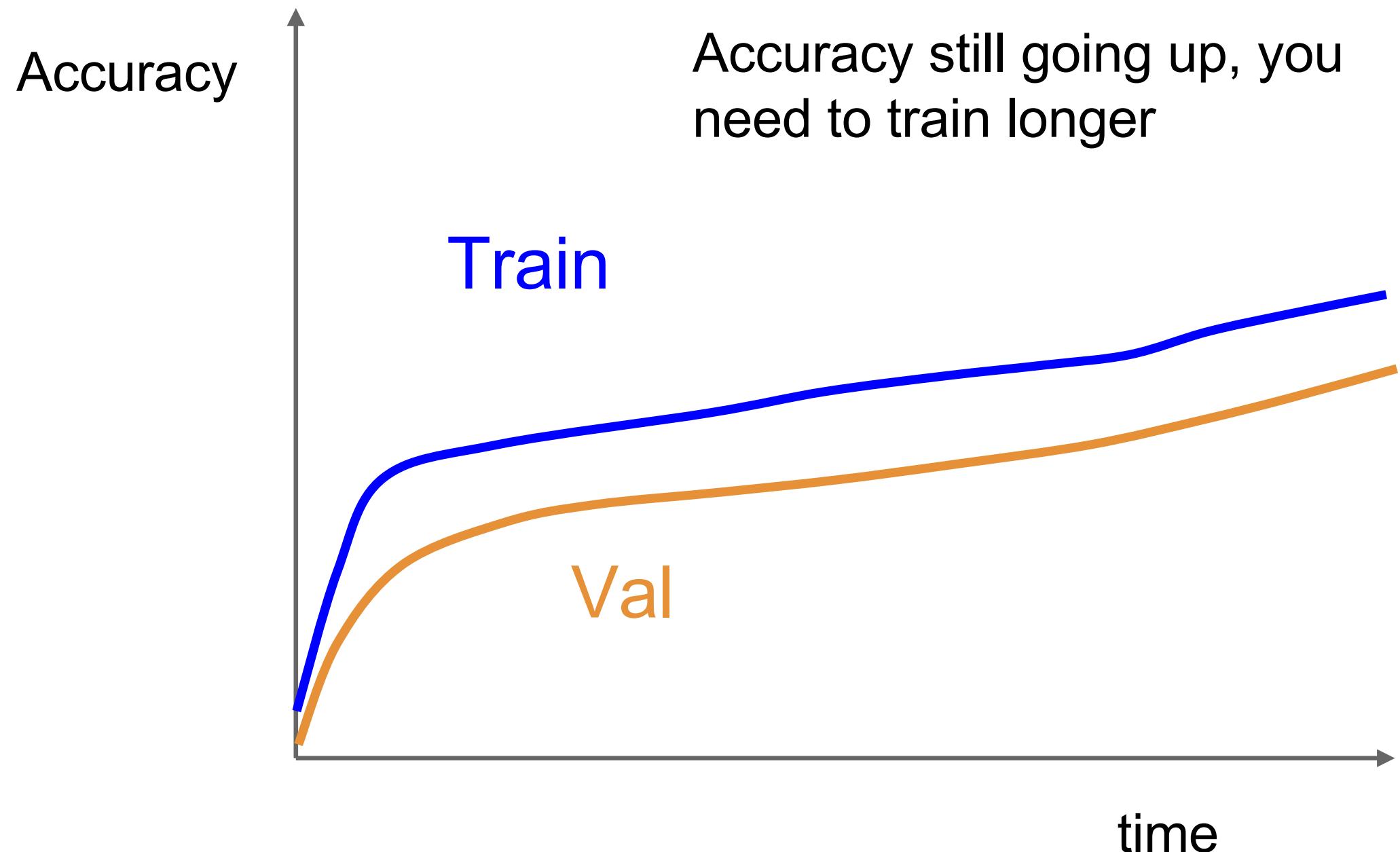
Step 2: Overfit a small sample

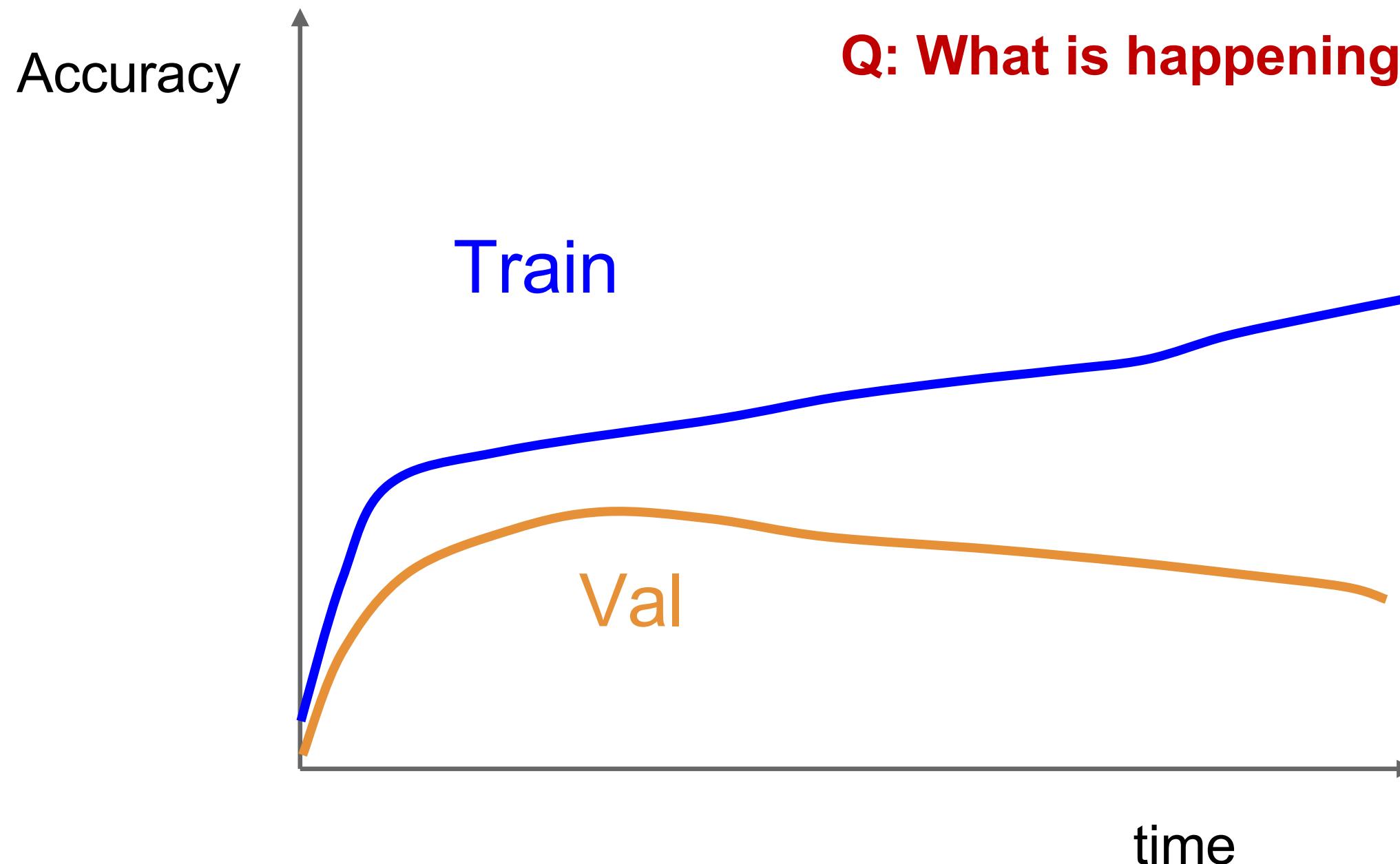
Step 3: Find LR that makes loss go down

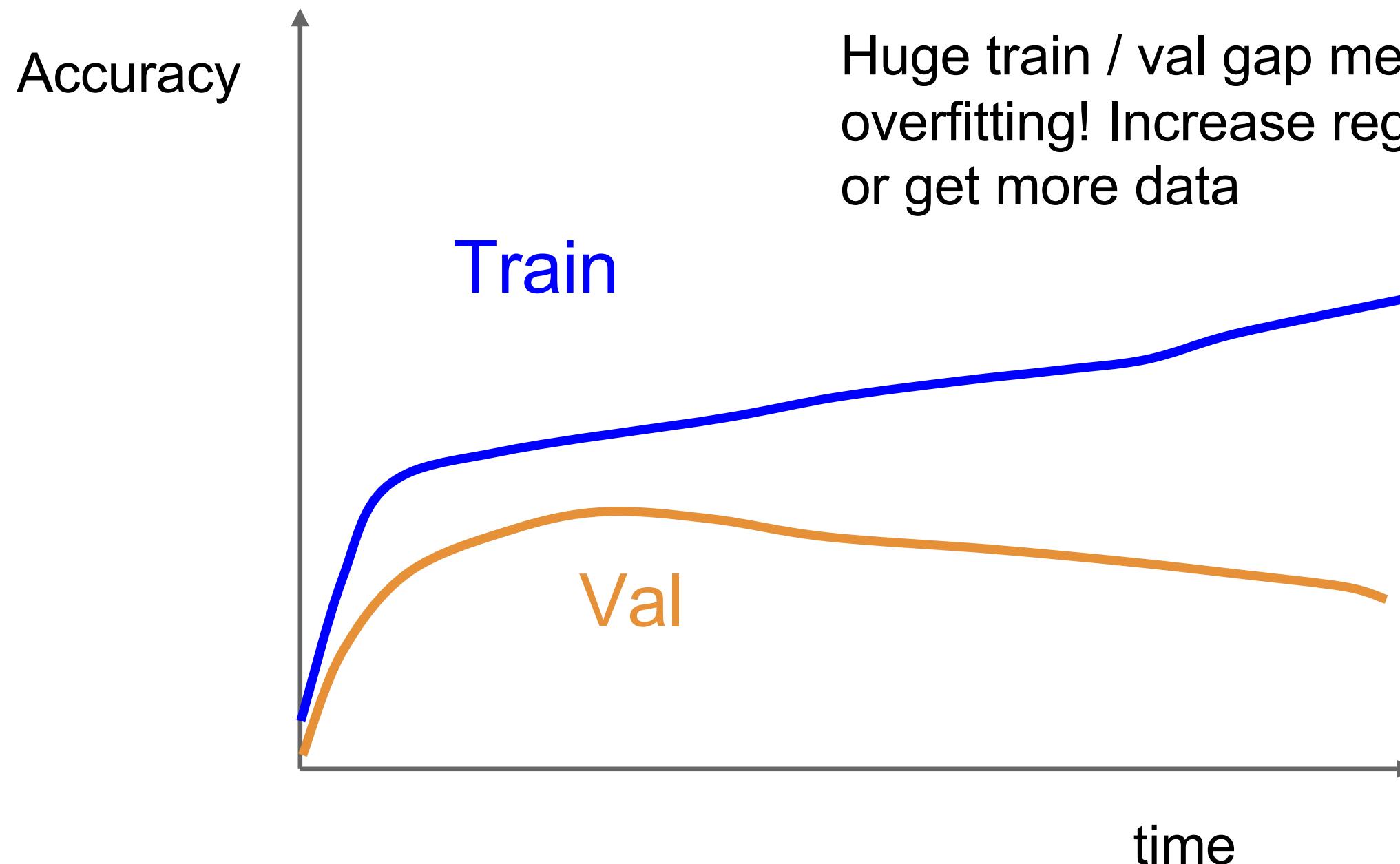
Step 4: Coarse grid of hyperparams, train for ~1-5 epochs

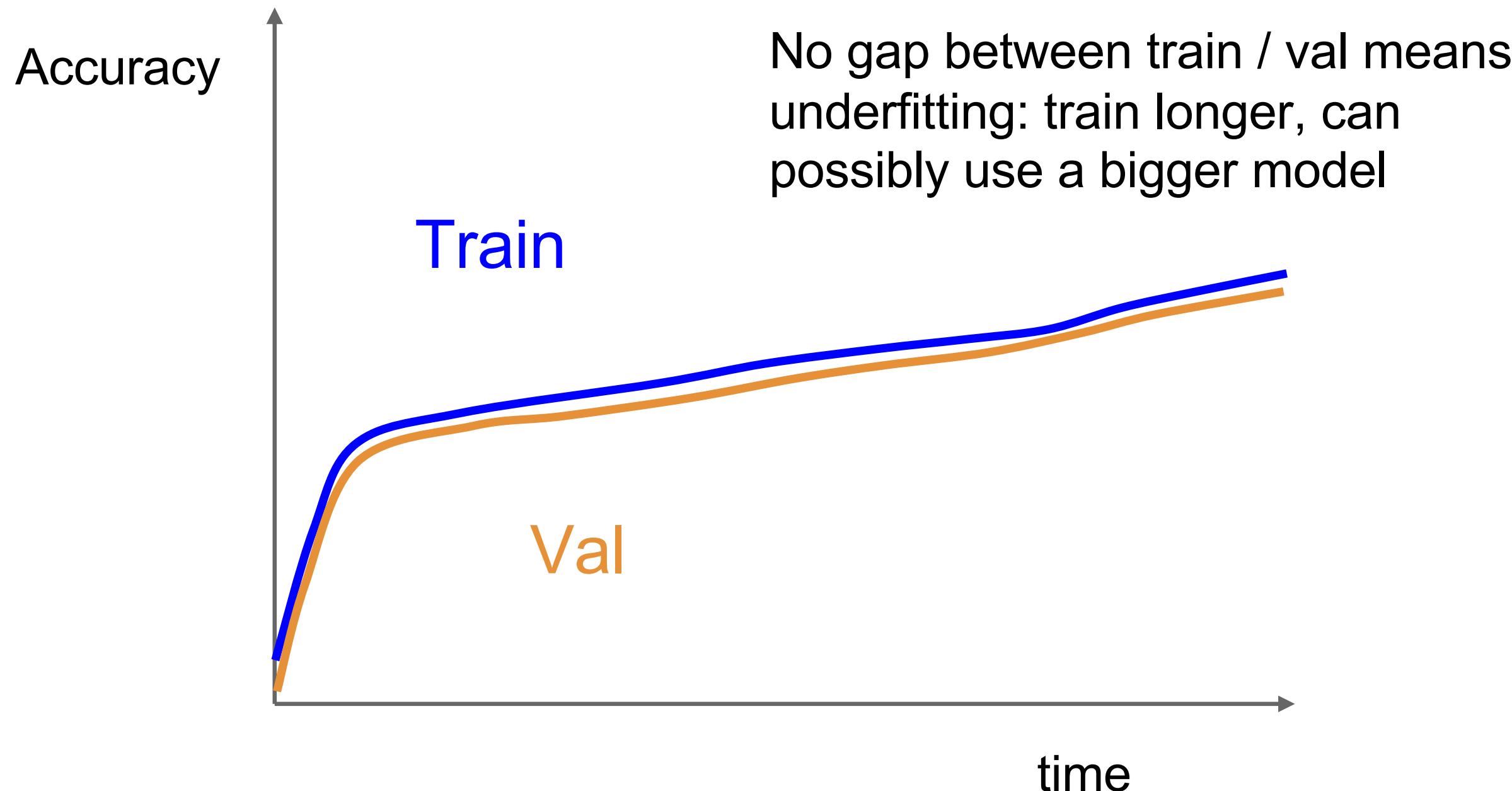
Step 5: Refine grid, train longer

Step 6: Look at loss and accuracy curves (next slides)









Choosing Hyperparameters

Step 1: Check initial loss

Step 2: Overfit a small sample

Step 3: Find LR that makes loss go down

Step 4: Coarse grid, train for ~1-5 epochs

Step 5: Refine grid, train longer

Step 6: Look at loss and accuracy curves

Step 7: **GOTO step 5**

Summary

We reviewed 8 topics at a high level:

1. Layers in CNNs (Conv, FC, Norm, Dropout)
2. Activation Functions in NNs (ReLU, GELU, etc.)
3. CNN Architectures (VGG, ResNets)
4. Weight Initialization (Maintain Activation Distribution)

Summary

We reviewed 8 topics at a high level:

5. Data Preprocessing (subtract mean, divide std)
6. Data augmentation (cropping, jitter)
7. Transfer Learning (train on ImageNet first)
8. Hyperparameter (Checking Losses + Random Search)