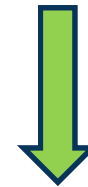


Backwards Pass for Convolution Layer

It is instructive to calculate **the backwards pass** of a convolution layer

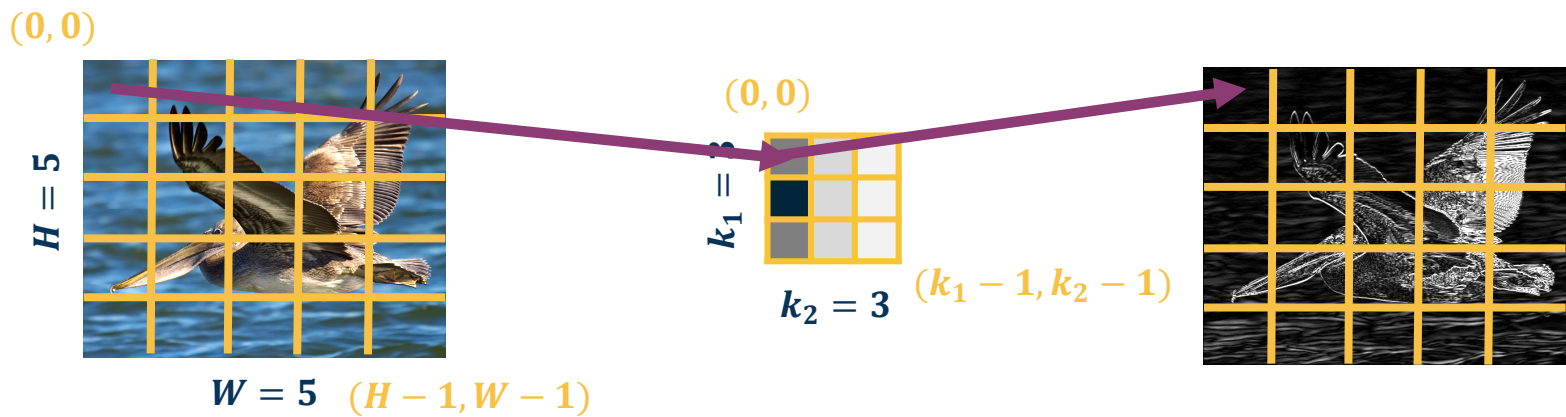
- Similar to fully connected layer, will be **simple vectorized linear algebra operation!**
- We will see a **duality** between cross-correlation and convolution

$$K = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$



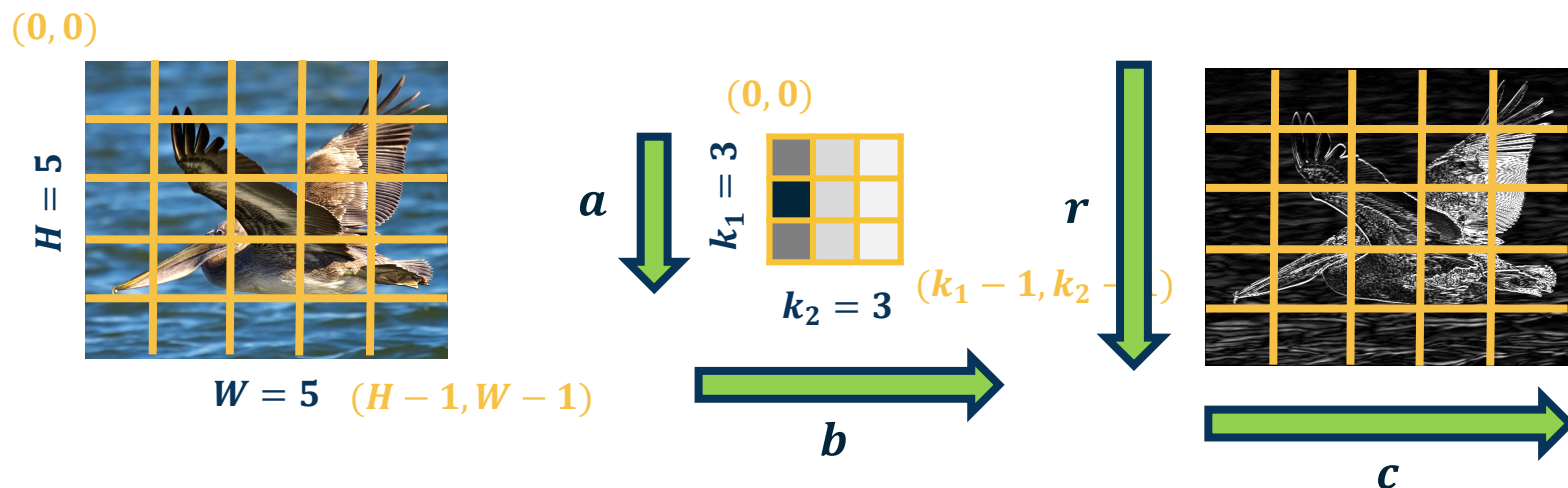
$$K' = \begin{bmatrix} 9 & 8 & 7 \\ 6 & 5 & 4 \\ 3 & 2 & 1 \end{bmatrix}$$

$$y(r, c) = (x * k)(r, c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b) k(a, b)$$



Recap: Cross-Correlation

$$y(r, c) = (x * k)(r, c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b) k(a, b)$$



Some simplification: 1 channel input, 1 kernel (channel output), padding (here 2 pixels on right/bottom) to make output the same size

Iterators

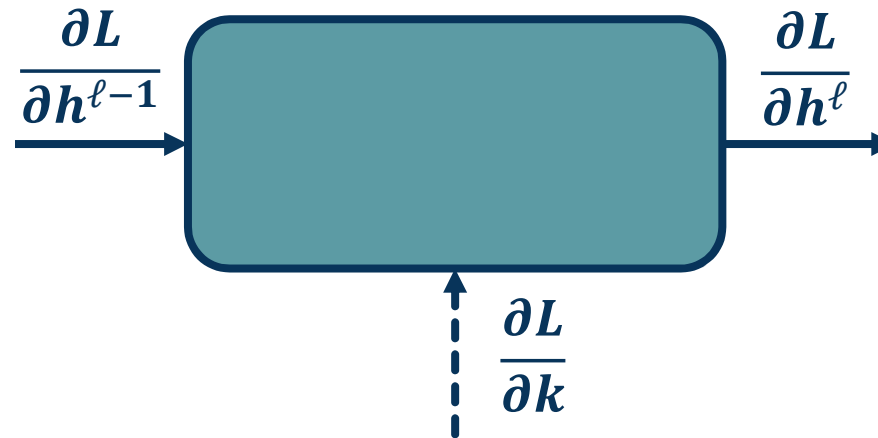
$$y(r, c) = (x * k)(r, c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r+a, c+b) k(a, b)$$

$$|y| = H \times W$$

$$\frac{\partial L}{\partial y} ?$$

Assume size $H \times W$ (add padding, change convention a bit for convenience)

$$\frac{\partial L}{\partial y(r, c)} \text{ to access element}$$



$$\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial h^{\ell-1}}$$

Gradient for passing back

$$\frac{\partial L}{\partial k} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial k}$$

Gradient for weight update

(weights = k, i.e. kernel values)

Backpropagation Chain Rule

Gradient for Convolution Layer

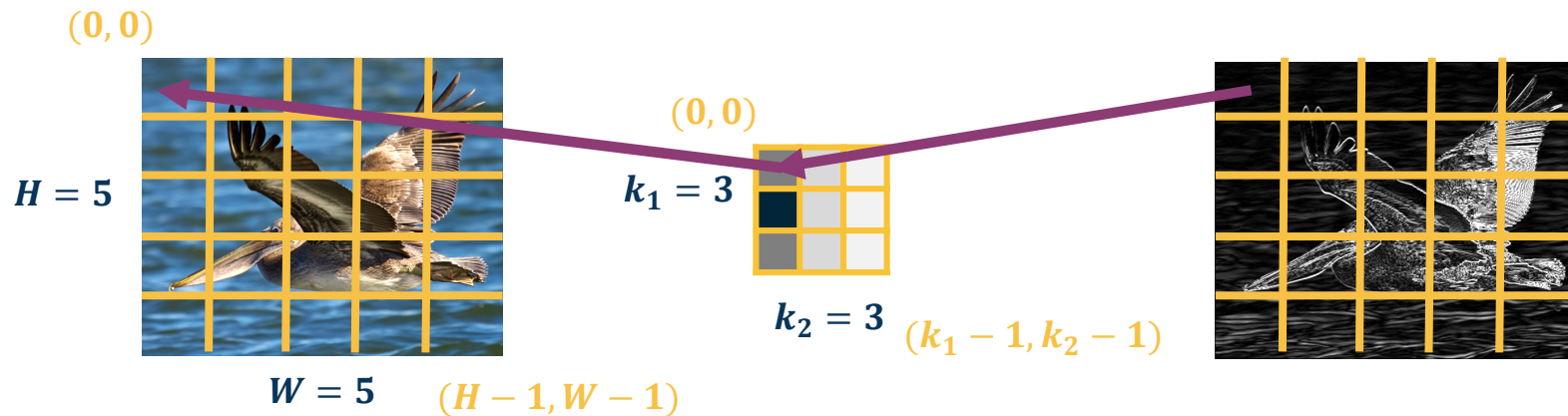
$$\frac{\partial L}{\partial k} = \frac{\partial L}{\partial h^l} \frac{\partial h^l}{\partial k}$$

Gradient for weight update

Calculate one pixel at a time $\frac{\partial L}{\partial k(a, b)}$

What does this weight affect at the output?

Everything!



What a Kernel Pixel Affects at Output

Need to incorporate all upstream gradients:

$$\left\{ \frac{\partial L}{\partial y(0,0)}, \frac{\partial L}{\partial y(0,1)}, \dots, \frac{\partial L}{\partial y(H,W)} \right\}$$

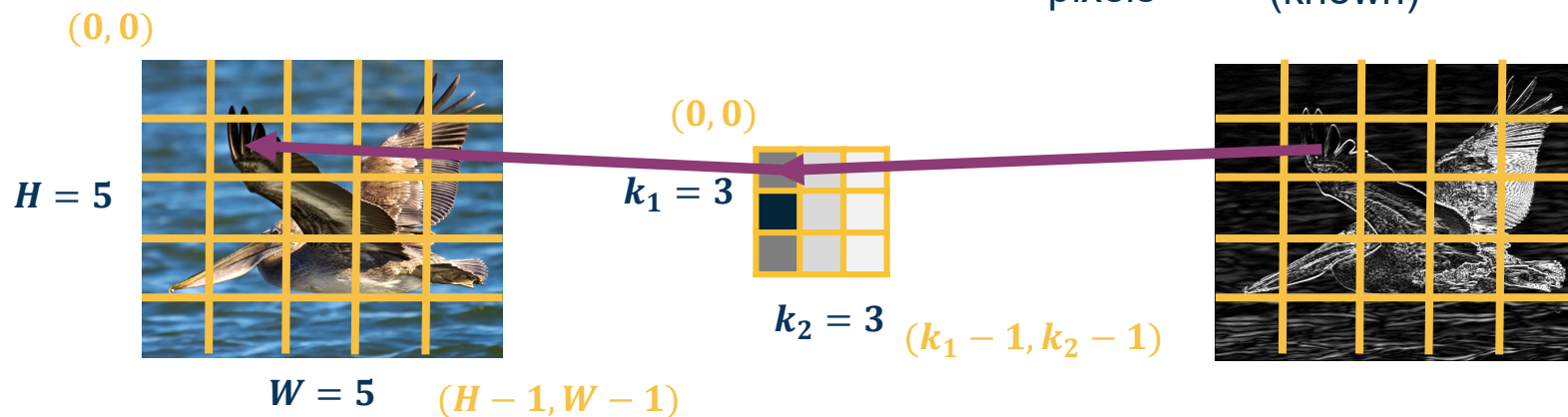
Chain Rule:

$$\frac{\partial L}{\partial k(a', b')} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r, c)} \frac{\partial y(r, c)}{\partial k(a', b')}$$

Sum over all output pixels

Upstream gradient (known)

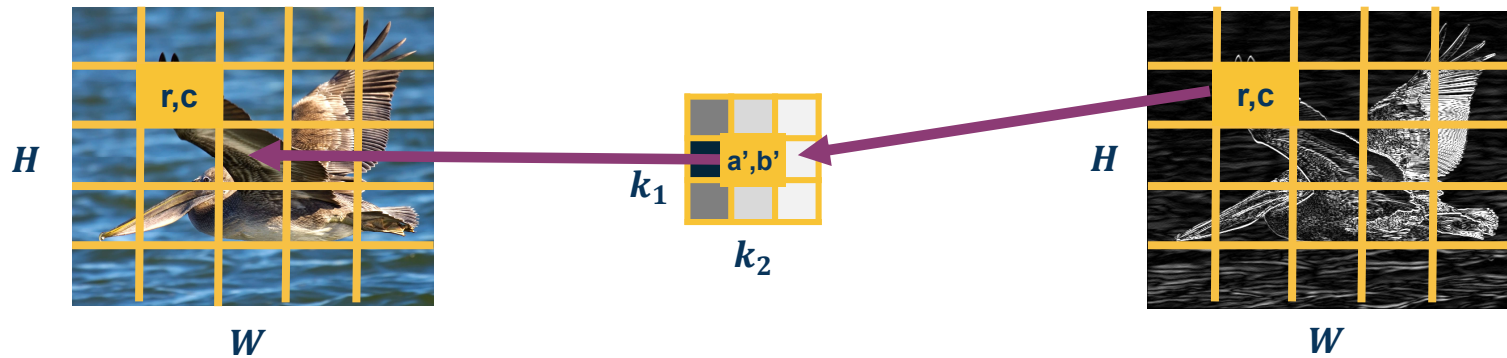
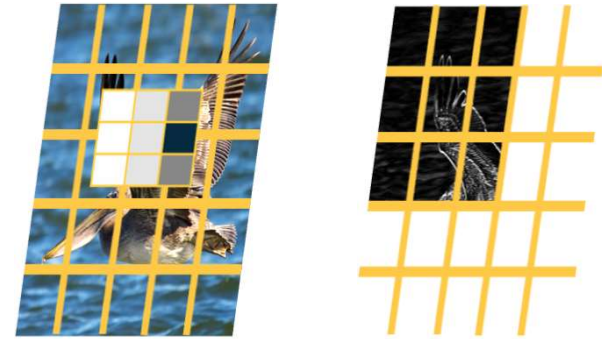
We will compute



Chain Rule over all Output Pixels

$$\frac{\partial y(r, c)}{\partial k(a', b')} = x(r + a', c + b')$$

$$\frac{\partial L}{\partial k(a', b')} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r, c)} x(r + a', c + b')$$



Chain Rule over all Output Pixels

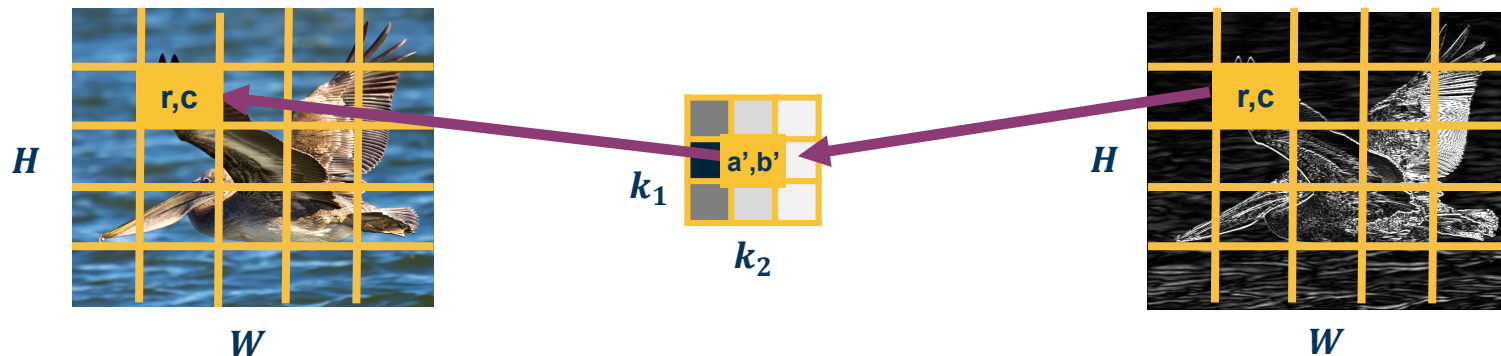
$$\frac{\partial y(r, c)}{\partial k(a', b')} = x(r + a', c + b')$$

$$\frac{\partial L}{\partial k(a', b')} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r, c)} x(r + a', c + b')$$

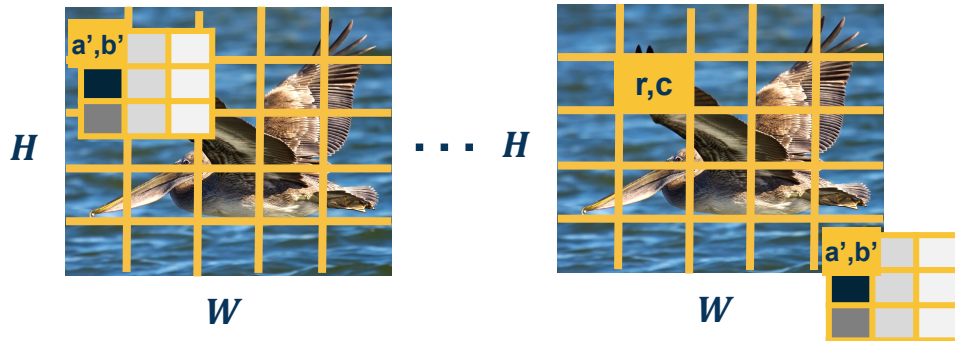
Does this look familiar?

Cross-correlation
between upstream
gradient and input!

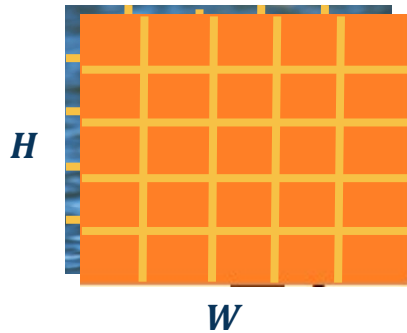
(until $k_1 \times k_2$ output)



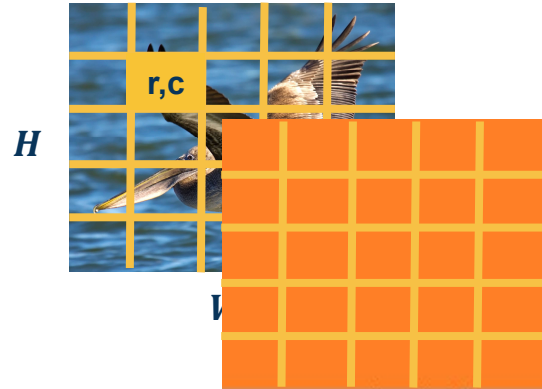
Forward Pass



Backward Pass $k(0,0)$



Backward Pass $k(2,2)$



Does this look familiar?

Cross-correlation
between upstream
gradient and input!
(until $k_1 \times k_2$ output)



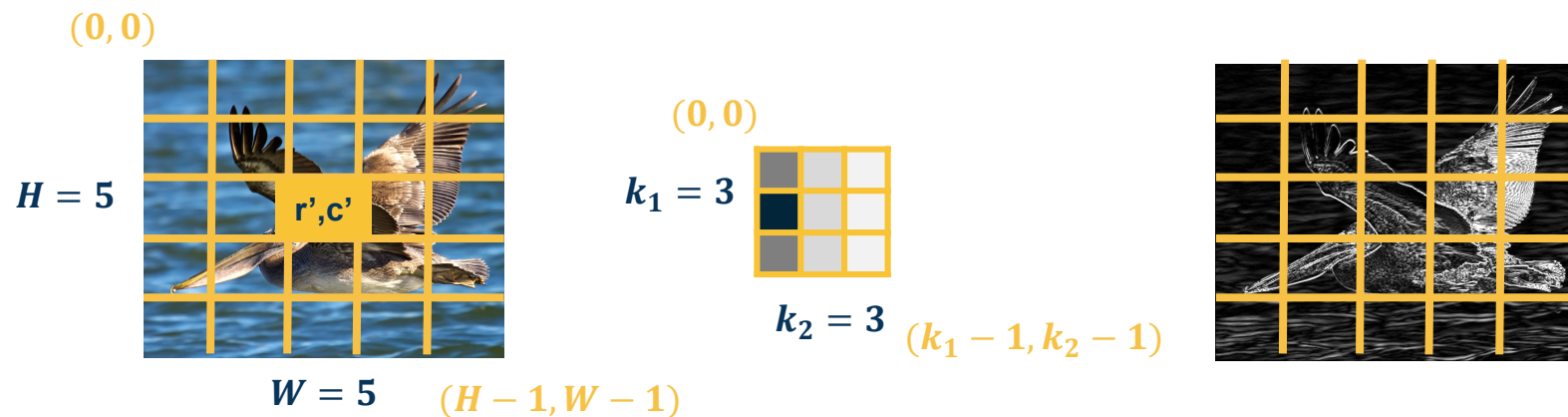
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

Gradient for input (to pass to prior layer)

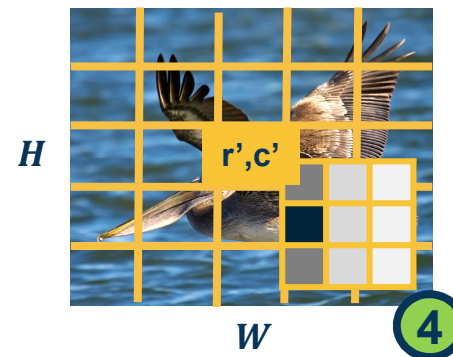
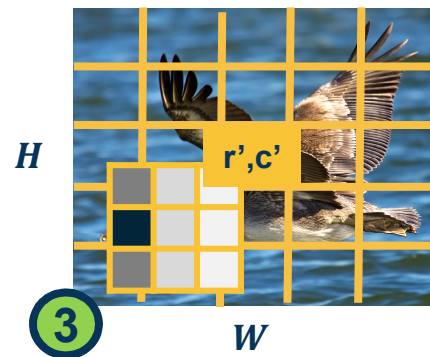
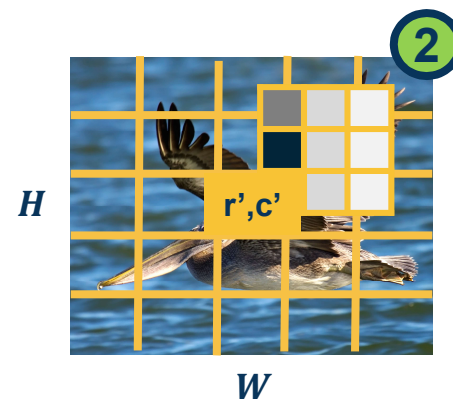
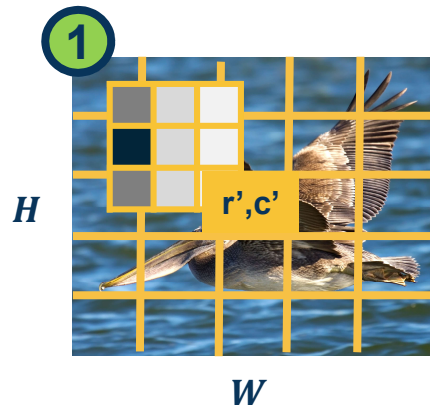
Calculate one pixel at a time $\frac{\partial L}{\partial x(r', c')}$

What does this input pixel affect at the output?

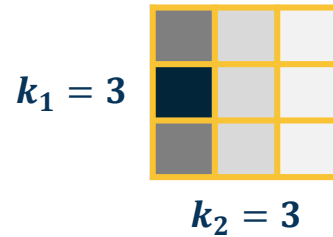
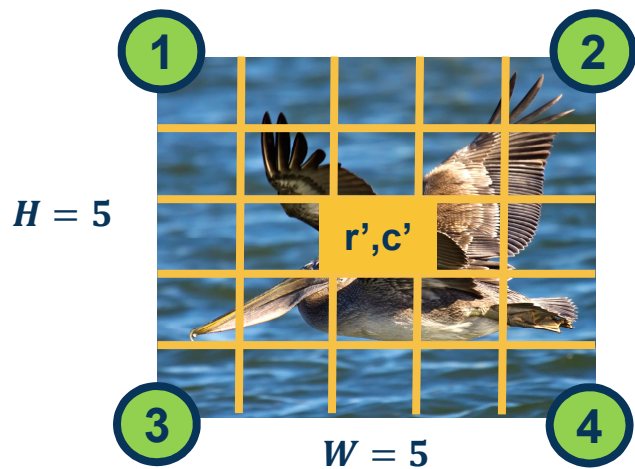
Neighborhood around it (where part of the kernel touches it)



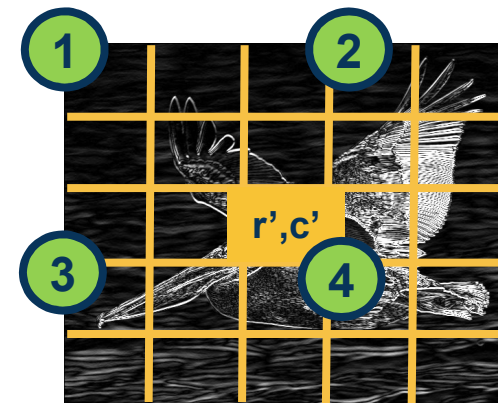
What an Input Pixel Affects at Output



Extents of Kernel Touching the Pixel



$$(r' - k_1 + 1, \\ c' - k_2 + 1)$$



This is where the corresponding locations are for the **output**

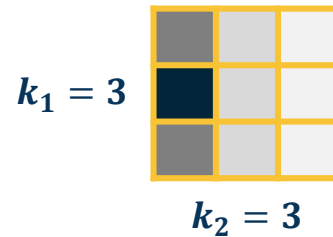
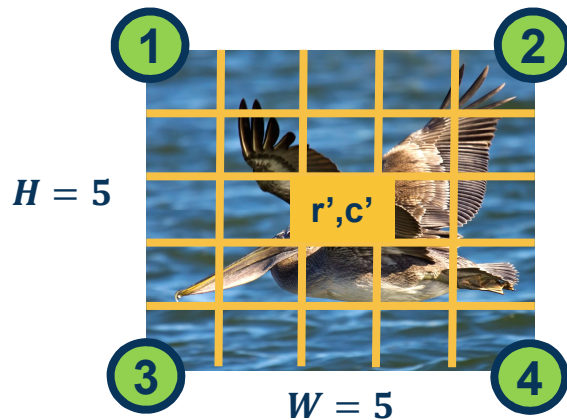
Extents at the Output

Chain rule for affected pixels (sum gradients):

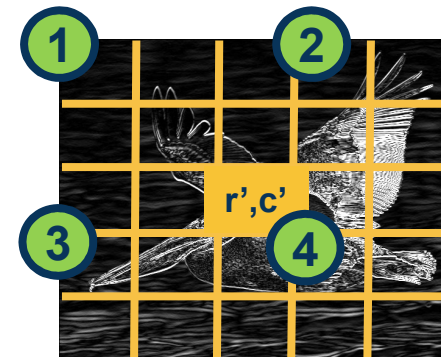
$$\frac{\partial L}{\partial x(r', c')} = \sum_{\text{Pixels } p} \frac{\partial L}{\partial y(p)} \frac{\partial y(p)}{\partial x(r', c')}$$

$$\frac{\partial L}{\partial x(r', c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r' - a, c' - b)} \frac{\partial y(r' - a, c' - b)}{\partial x(r', c')}$$

Let's derive it analytically this time (as opposed to visually)



$(r' - k_1 + 1, c' - k_2 + 1)$



Summing Gradient Contributions

Definition of cross-correlation (use a', b' to distinguish from prior variables):

$$y(\mathbf{r}', \mathbf{c}') = (\mathbf{x} * \mathbf{k})(\mathbf{r}', \mathbf{c}') = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x(\mathbf{r}' + \mathbf{a}', \mathbf{c}' + \mathbf{b}') k(\mathbf{a}', \mathbf{b}')$$

Plug in what we actually wanted :

$$y(\mathbf{r}' - \mathbf{a}, \mathbf{c}' - \mathbf{b}) = (\mathbf{x} * \mathbf{k})(\mathbf{r}', \mathbf{c}') = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x(\mathbf{r}' - \mathbf{a} + \mathbf{a}', \mathbf{c}' - \mathbf{b} + \mathbf{b}') k(\mathbf{a}', \mathbf{b}')$$

What is $\frac{\partial y(\mathbf{r}' - \mathbf{a}, \mathbf{c}' - \mathbf{b})}{\partial x(\mathbf{r}', \mathbf{c}')} = \mathbf{k}(\mathbf{a}, \mathbf{b})$

(we want term with $x(\mathbf{r}', \mathbf{c}')$ in it;
this happens when $\mathbf{a}' = \mathbf{a}$ and $\mathbf{b}' = \mathbf{b}$)

Calculating the Gradient



Plugging in to earlier equation:

$$\begin{aligned}\frac{\partial L}{\partial x(r', c')} &= \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r' - a, c' - b)} \frac{\partial y(r' - a, c' - b)}{\partial x(r', c')} \\ &= \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r' - a, c' - b)} k(a, b)\end{aligned}$$

Again, all operations can be implemented via matrix multiplications (same as FC layer)!

Does this look familiar?

Convolution between upstream gradient and kernel!

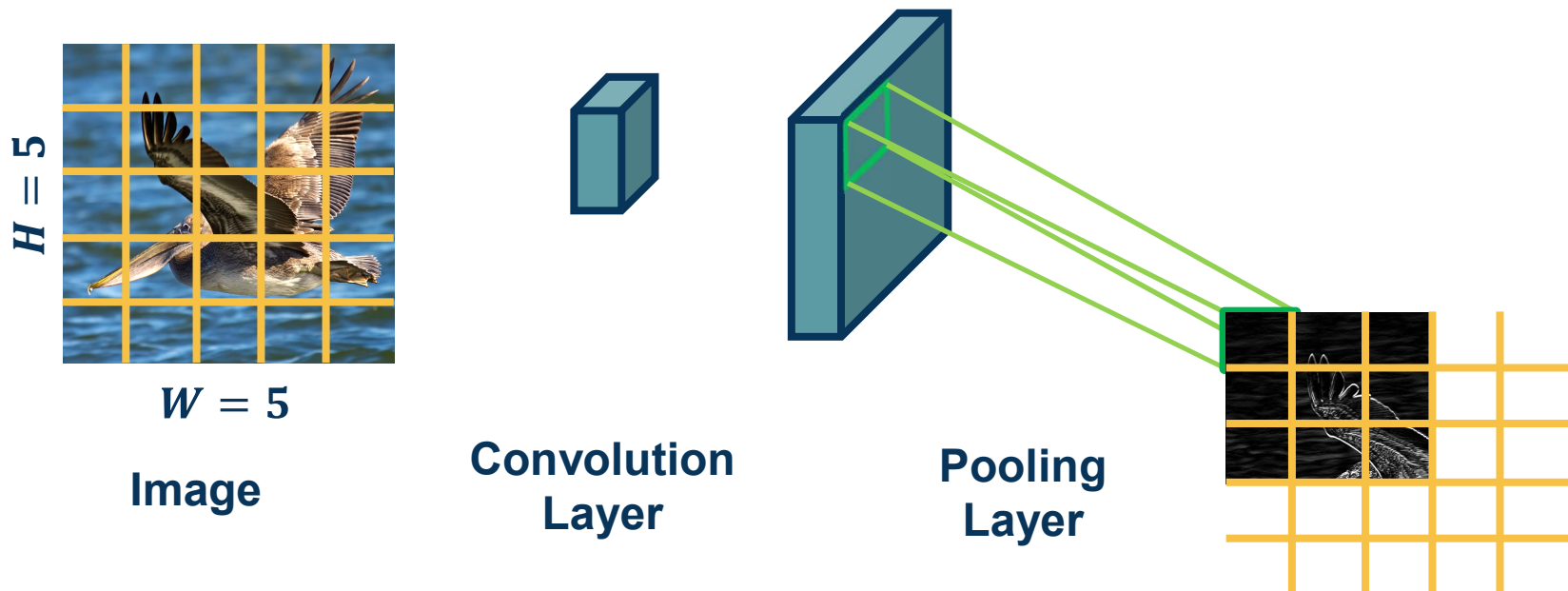
(can implement by flipping kernel and cross-correlation)

Backwards is Convolution



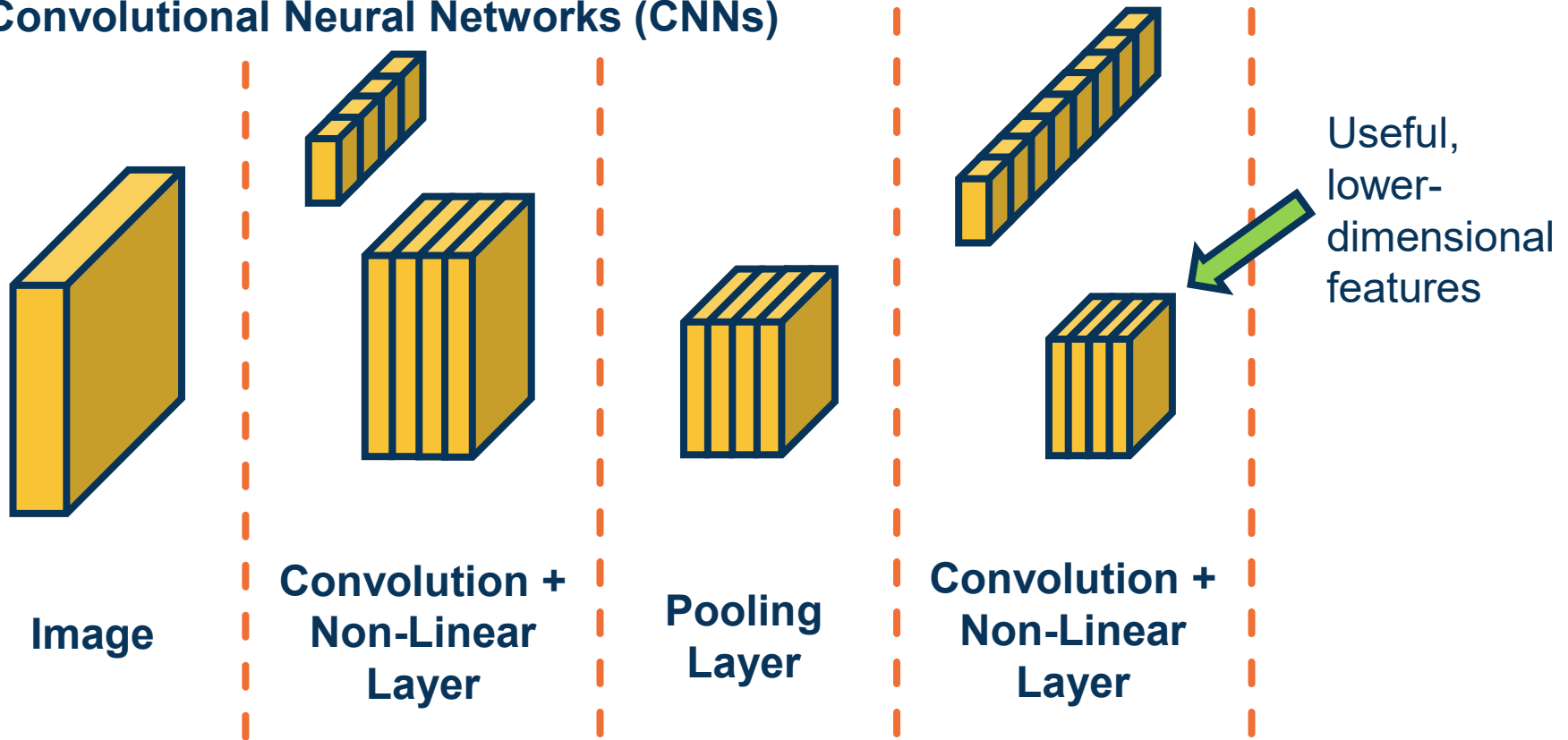
Simple Convolutional Neural Networks

Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer

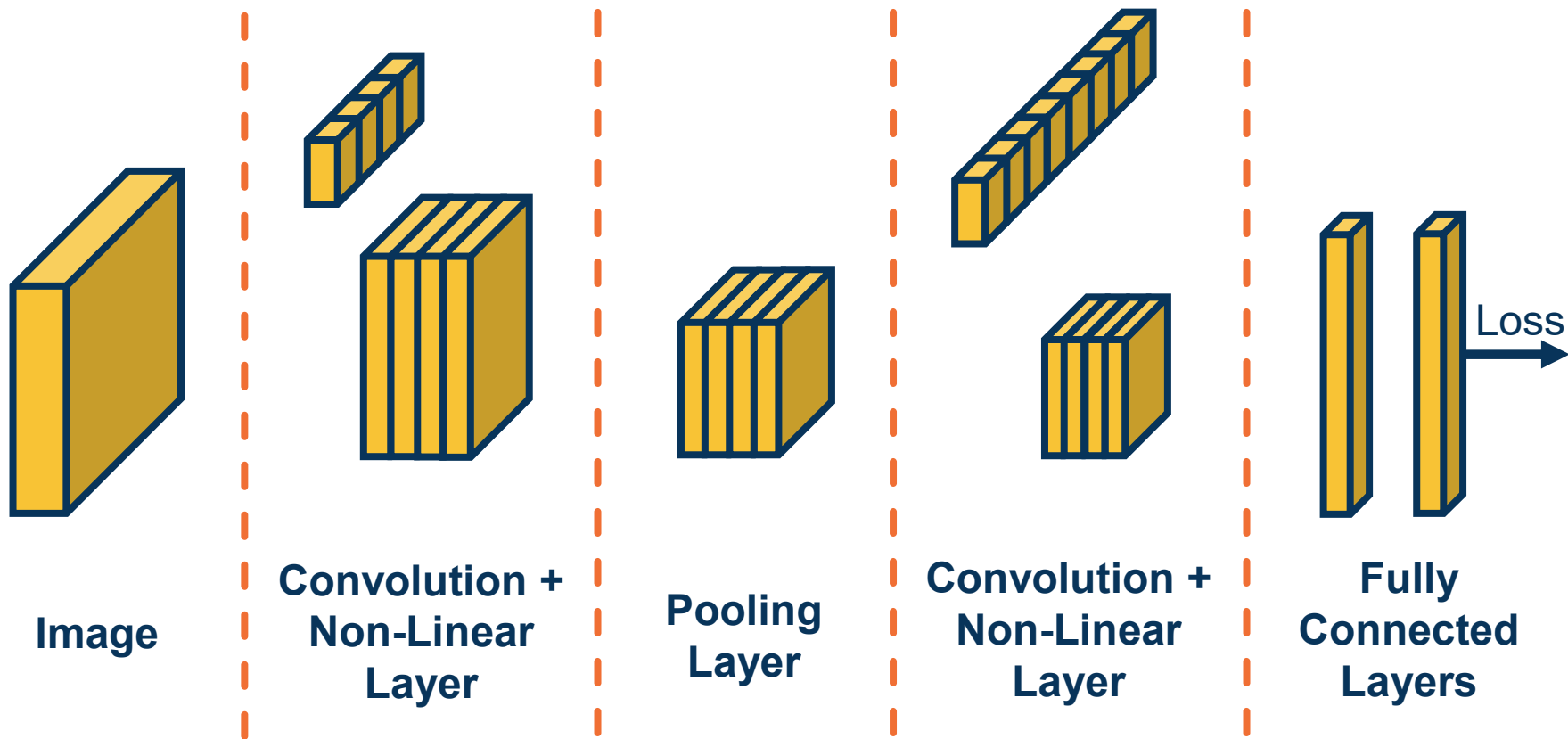


Combining Convolution & Pooling Layers

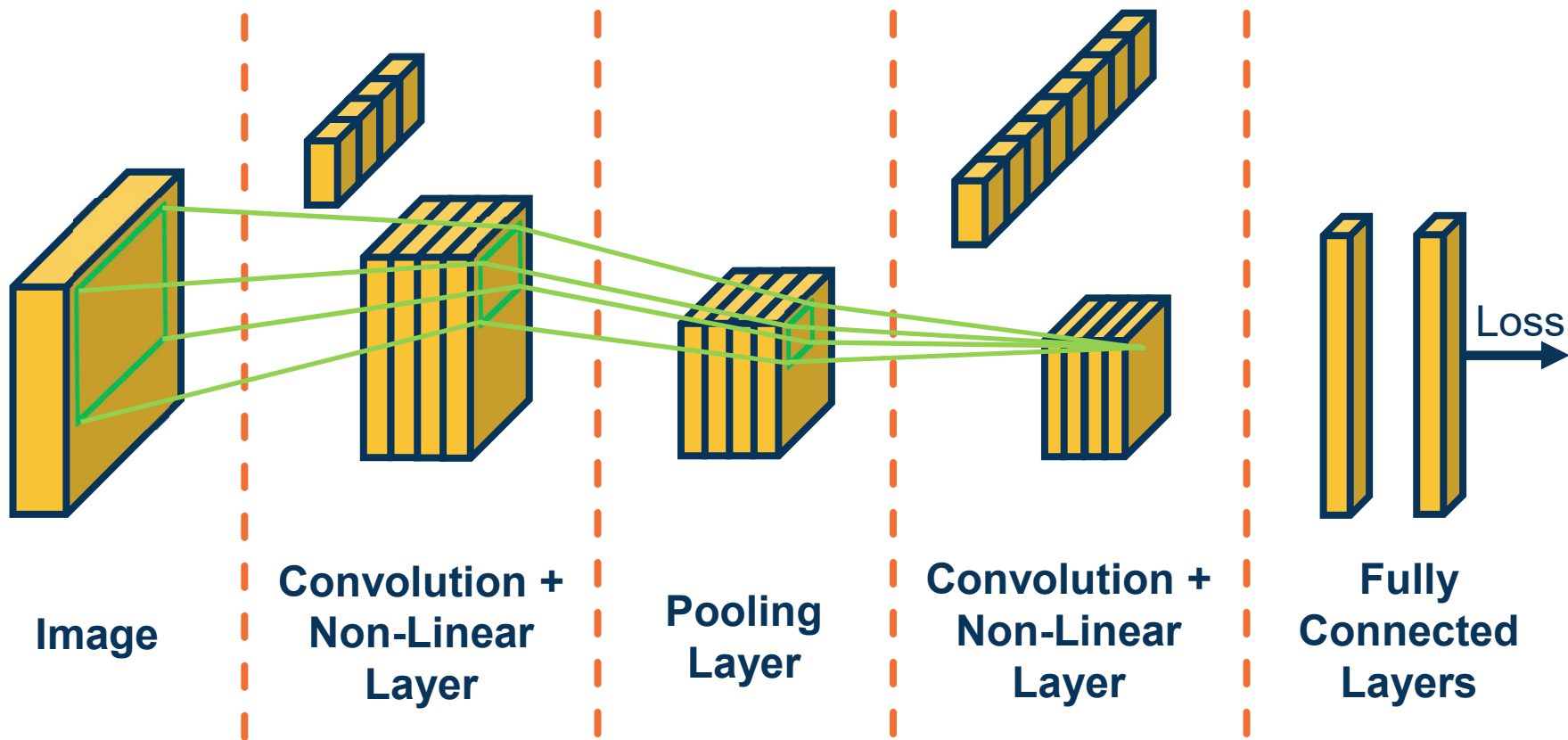
Convolutional Neural Networks (CNNs)



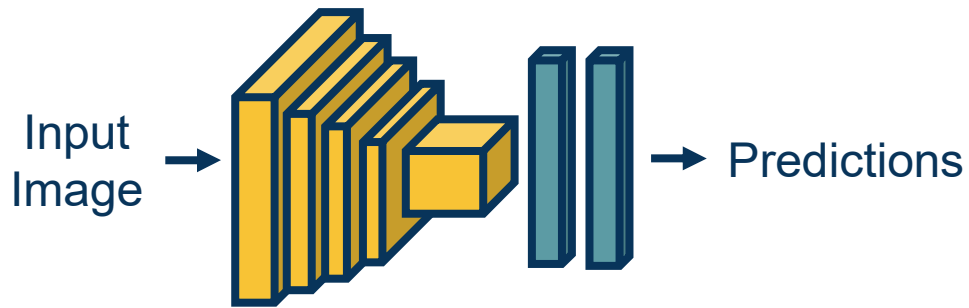
Alternating Convolution and Pooling



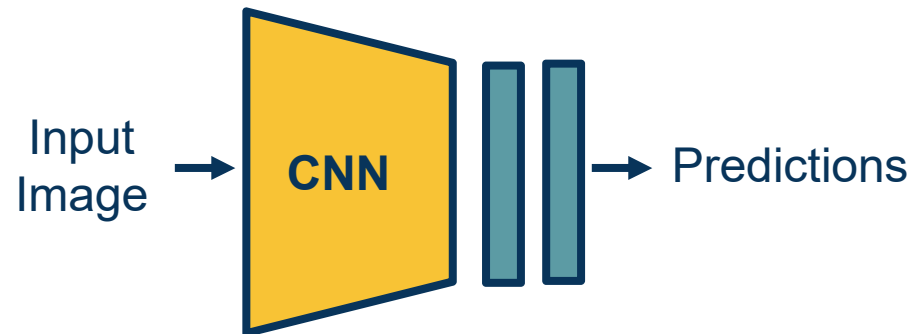
Adding a Fully Connected Layer



Receptive Fields



**Convolutional Neural
Networks**



Typical Depiction of CNNs

These architectures have existed **since 1980s**

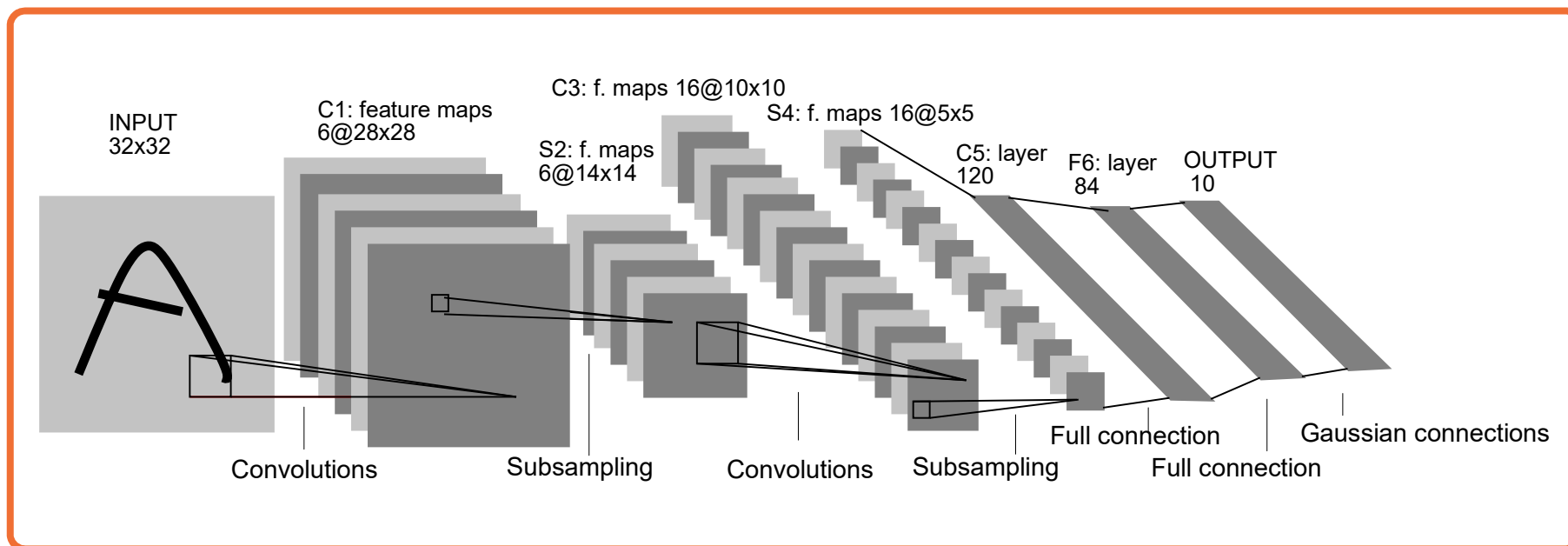


Image Credit: Yann LeCun, Kevin Murphy

LeNet Architecture



Handwriting Recognition

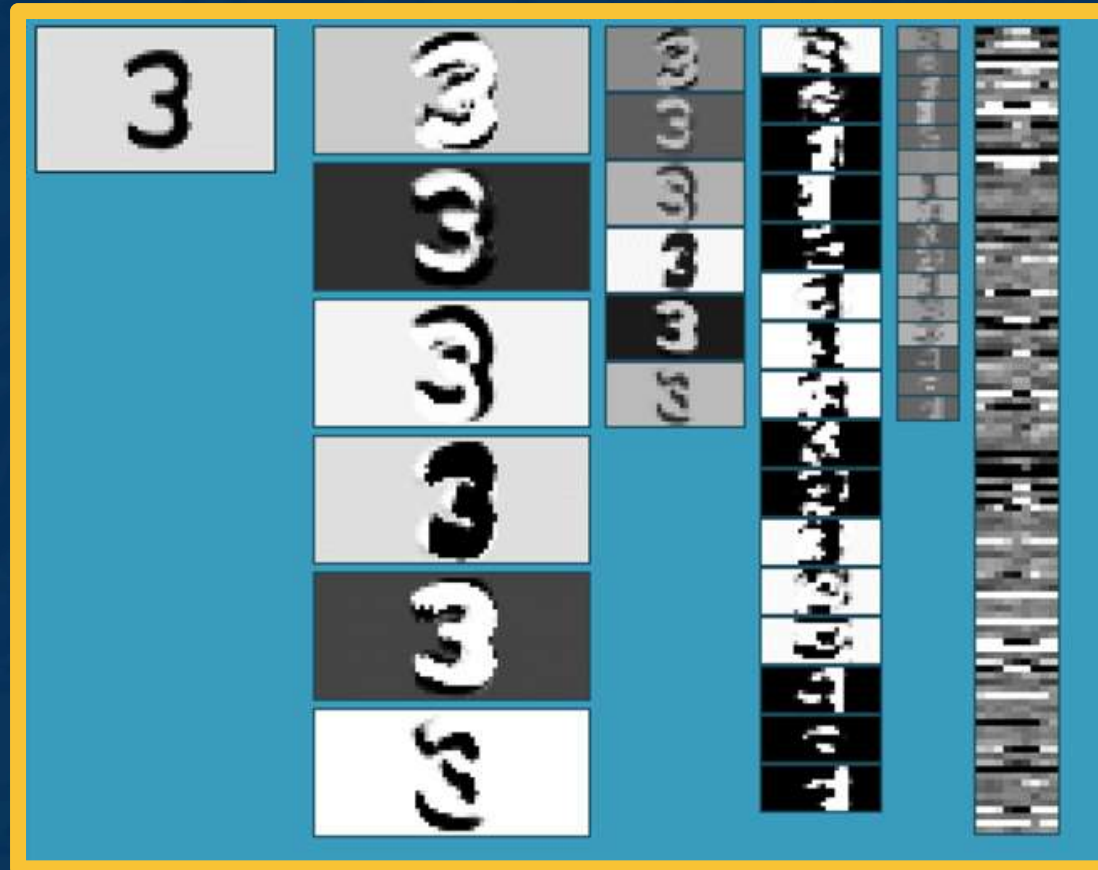


Image Credit:
Yann LeCun



Translation Equivariance (Conv Layers) & Invariance (Output)

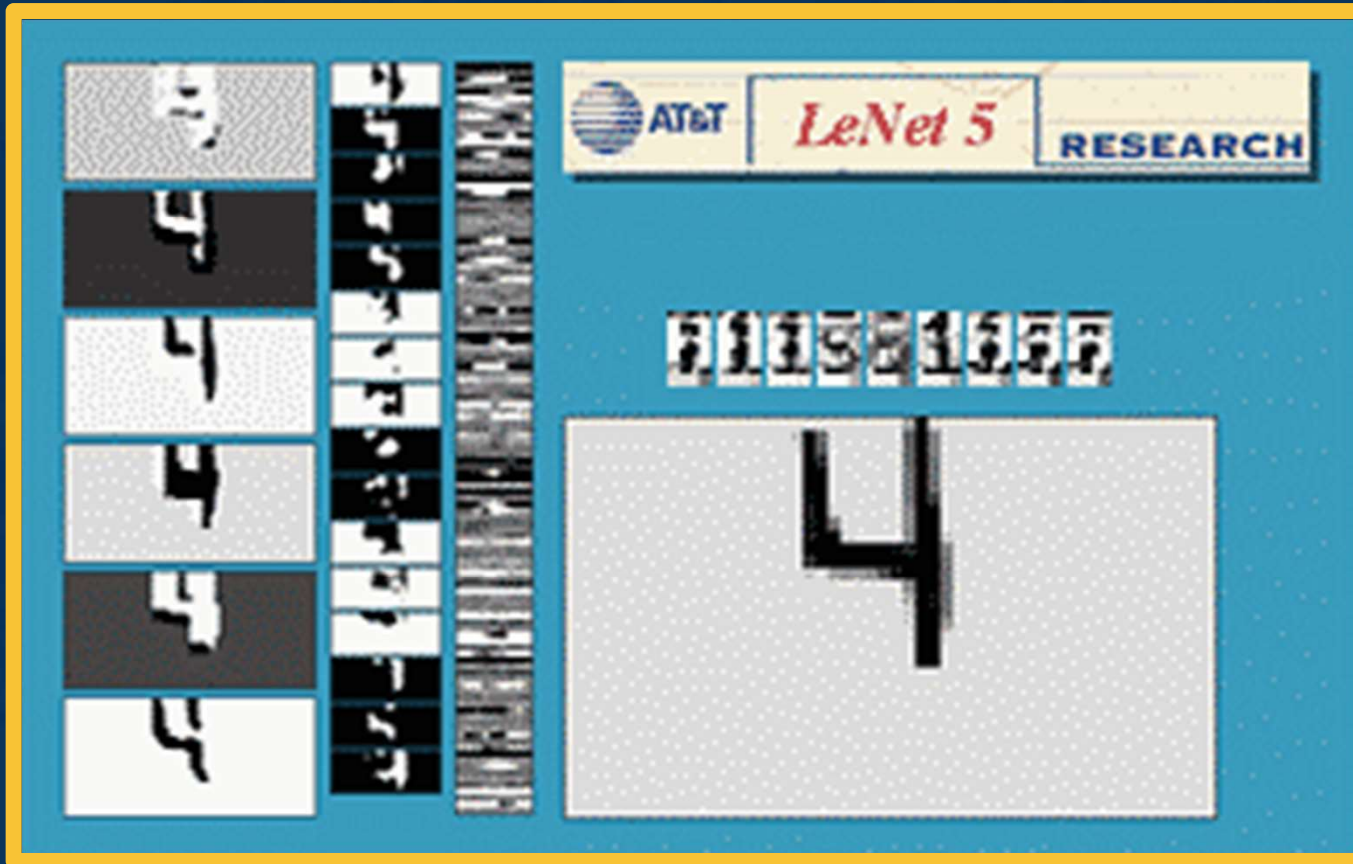


Image Credit:
Yann LeCun



(Some) Rotation Invariance

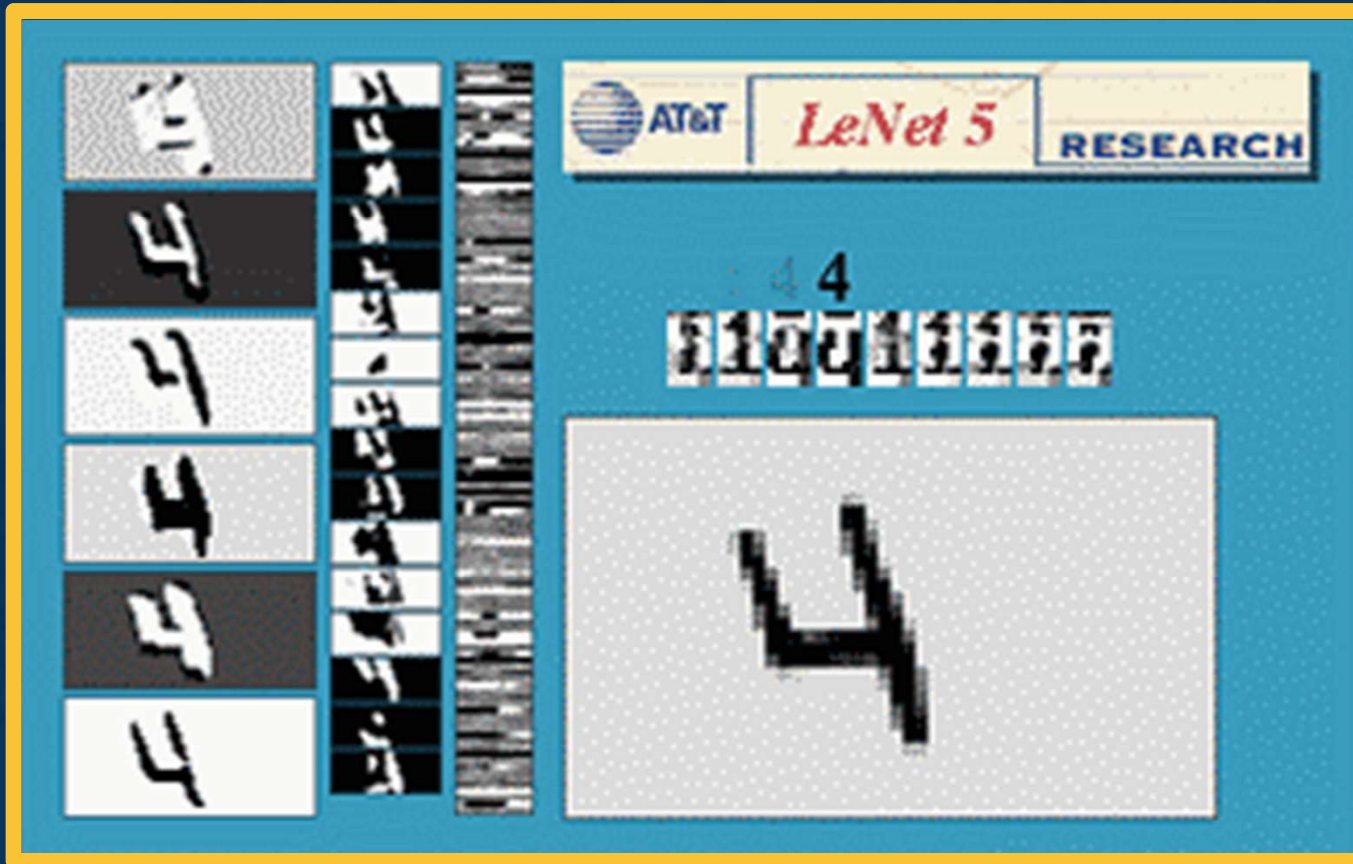


Image Credit:
Yann LeCun



(Some) Scale Invariance

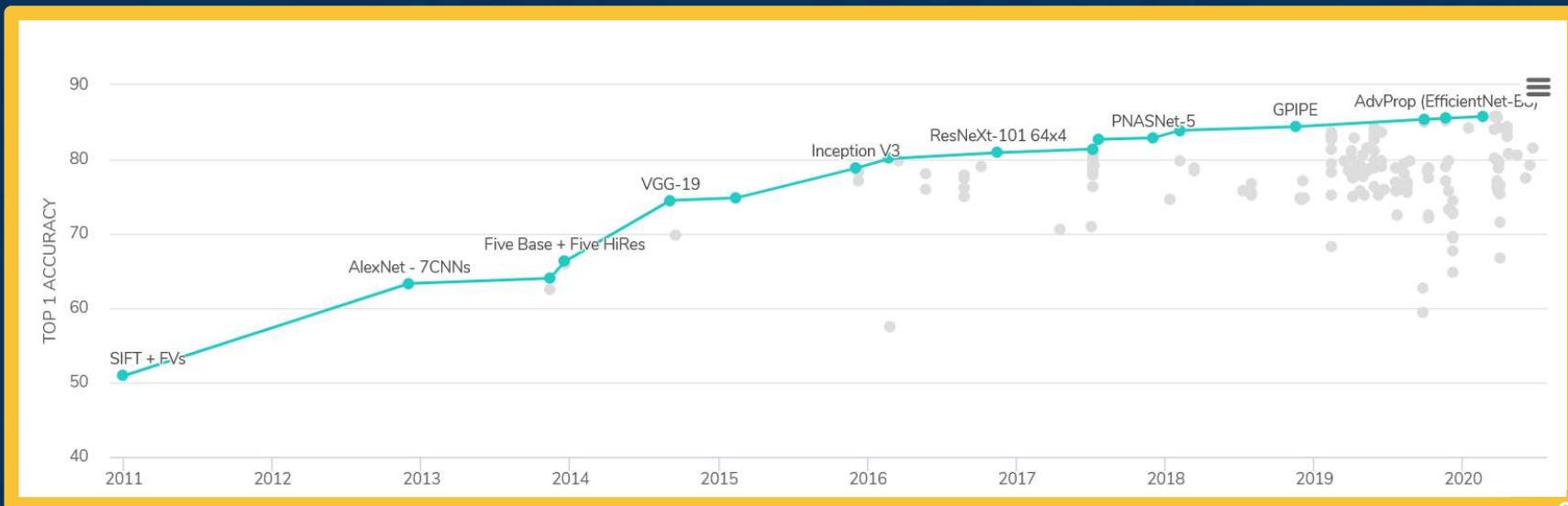


Image Credit:
Yann LeCun



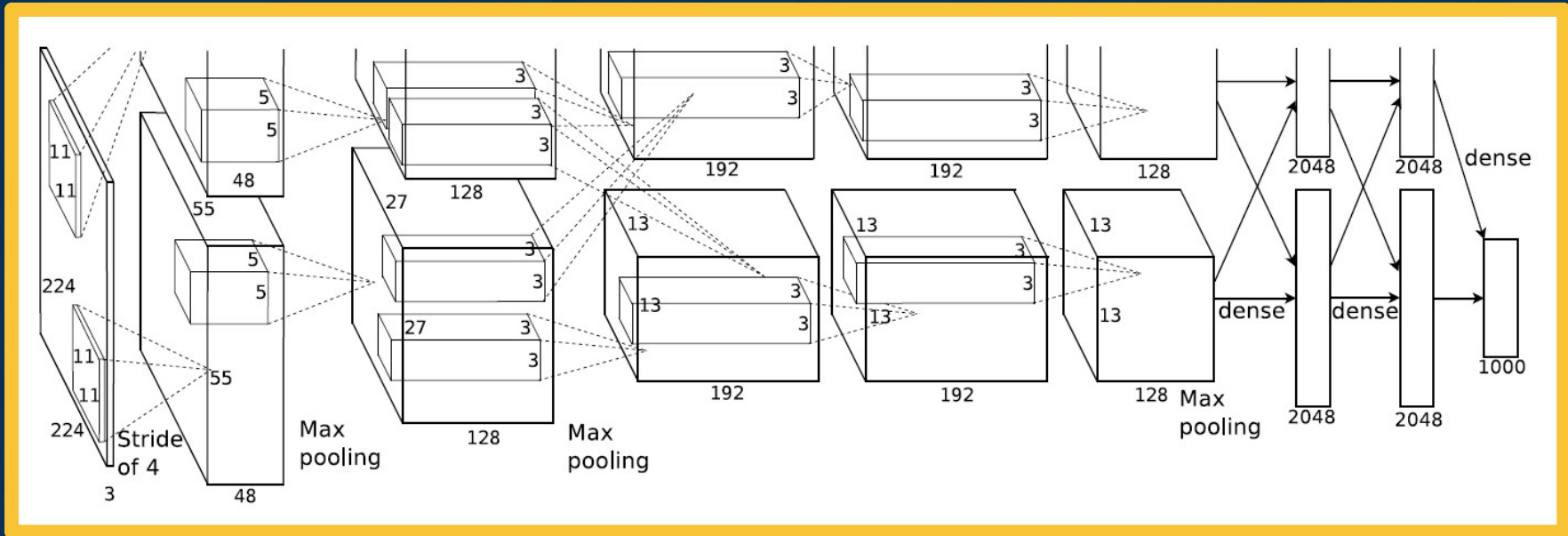
Advanced Convolutional Networks

The Importance of Benchmarks



From: <https://paperswithcode.com>

AlexNet - Architecture



From: Krizhevsky et al., *ImageNet Classification with Deep Convolutional Neural Networks*, 2012.

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

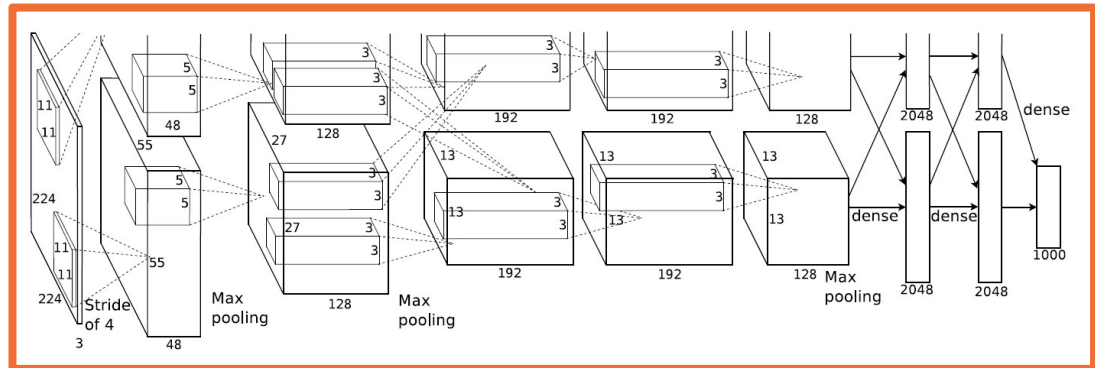
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



Key aspects:

- ReLU instead of sigmoid or tanh
- Specialized normalization layers
- PCA-based data augmentation
- Dropout
- Ensembling

From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

AlexNet – Layers and Key Aspects

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases)
 CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728
 CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864
 POOL2: [112x112x64] memory: 112*112*64=800K params: 0
 CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728
 CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456
 POOL2: [56x56x128] memory: 56*56*128=400K params: 0
 CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912
 CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824
 CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824
 POOL2: [28x28x256] memory: 28*28*256=200K params: 0
 CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648
 CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296
 CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296
 POOL2: [14x14x512] memory: 14*14*512=100K params: 0
 CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
 CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
 CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
 POOL2: [7x7x512] memory: 7*7*512=25K params: 0
 FC: [1x1x4096] memory: 4096 params: 7*7*512*4096 = 102,760,448
 FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216
 FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

From: Simonyan & Zimmerman, Very Deep Convolutional Networks for Large-Scale Image Recognition
 From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

VGG



```

INPUT: [224x224x3]    memory: 224*224*3=150K  params: 0    (not counting biases)
CONV3-64: [224x224x64] memory: 224*224*64=3.2M  params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64] memory: 224*224*64=3.2M  params: (3*3*64)*64 = 36,864
POOL2: [112x112x64]  memory: 112*112*64=800K  params: 0
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  params: (3*3*128)*128 = 147,456
POOL2: [56x56x128]  memory: 56*56*128=400K  params: 0
CONV3-256: [56x56x256] memory: 56*56*256=800K  params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256] memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256] memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
POOL2: [28x28x256]  memory: 28*28*256=200K  params: 0
CONV3-512: [28x28x512] memory: 28*28*512=400K  params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512] memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512] memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]  memory: 14*14*512=100K  params: 0
CONV3-512: [14x14x512] memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]   memory: 7*7*512=25K    params: 0
FC: [1x1x4096]     memory: 4096          params: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]     memory: 4096          params: 4096*4096 = 16,777,216
FC: [1x1x1000]     memory: 1000         params: 4096*1000 = 4,096,000

```

Most memory usage in convolution layers

Most parameters in FC layers

From: Simonyan & Zimmerman, *Very Deep Convolutional Networks for Large-Scale Image Recognition*
 From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Parameters and Memory



Key aspects:

Repeated application of:

- 3x3 conv (stride of 1, padding of 1)
- 2x2 max pooling (stride 2)

Very large number of parameters

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

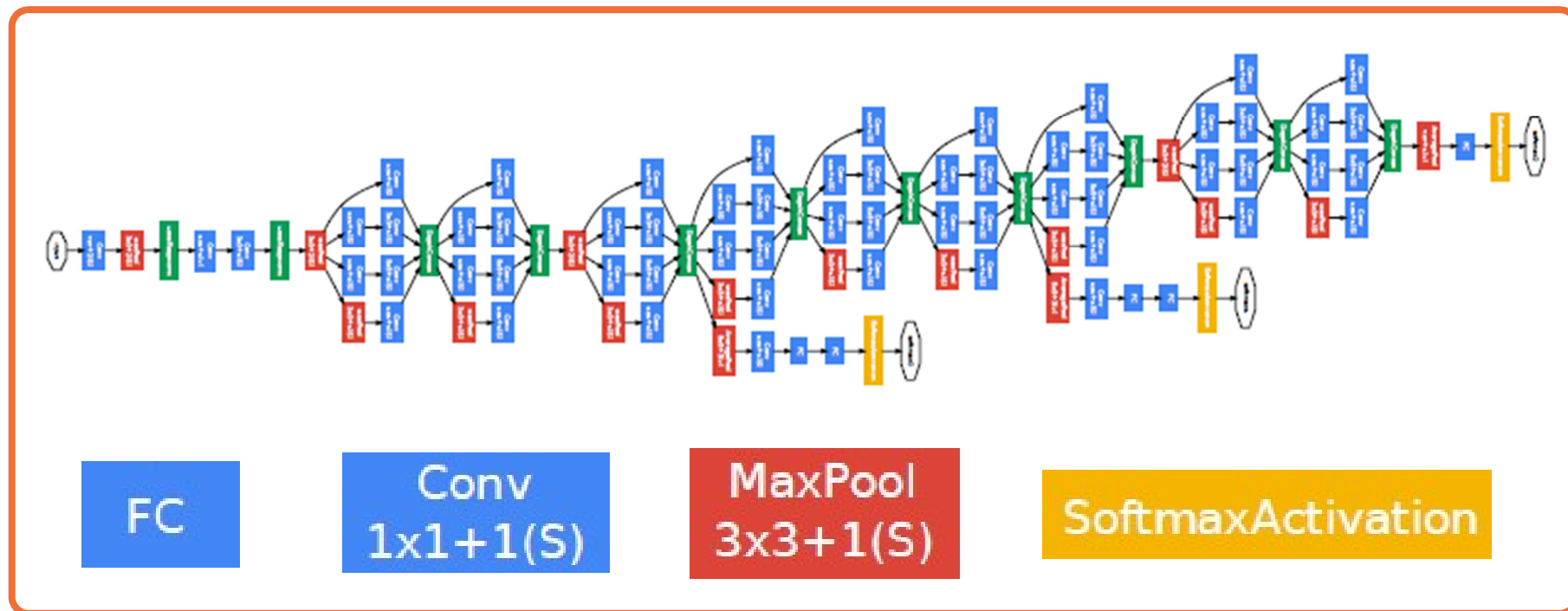
From: Simonyan & Zimmerman, *Very Deep Convolutional Networks for Large-Scale Image Recognition*

From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

VGG – Key Characteristics



But have become **deeper and more complex**

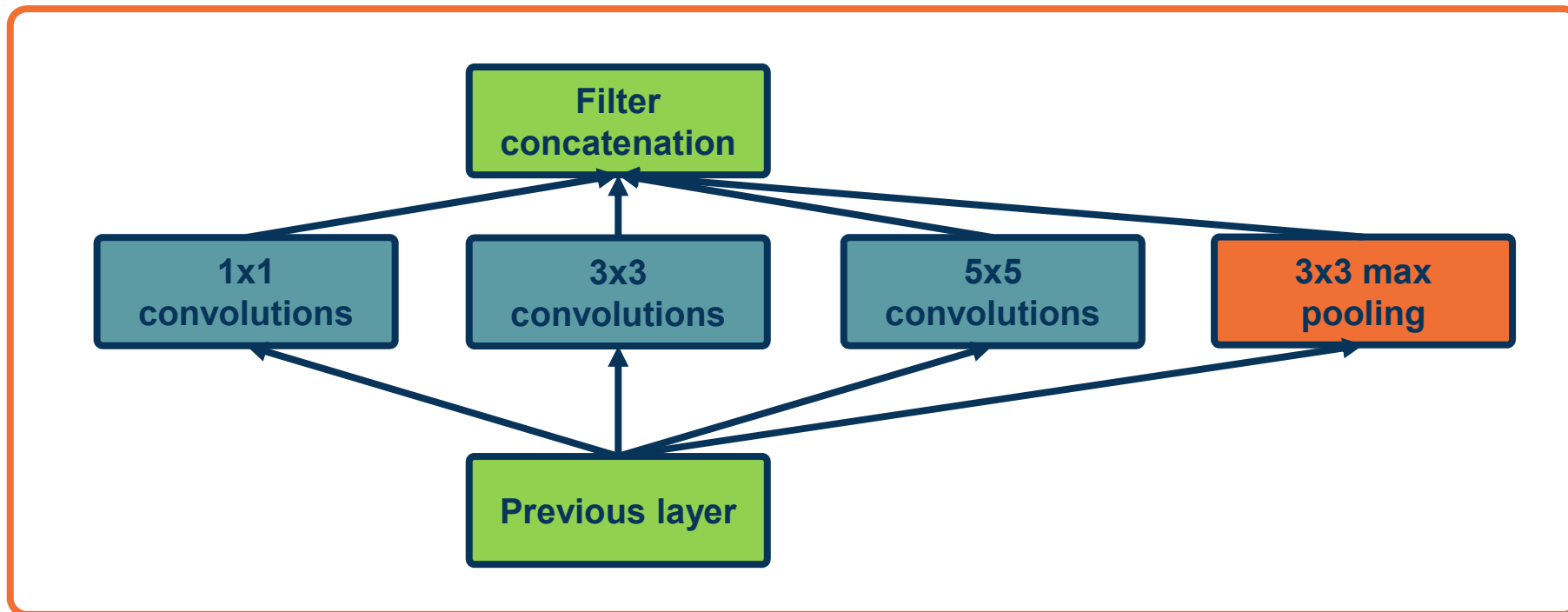


From: Szegedy et al. *Going deeper with convolutions*

Inception Architecture

Georgia
Tech

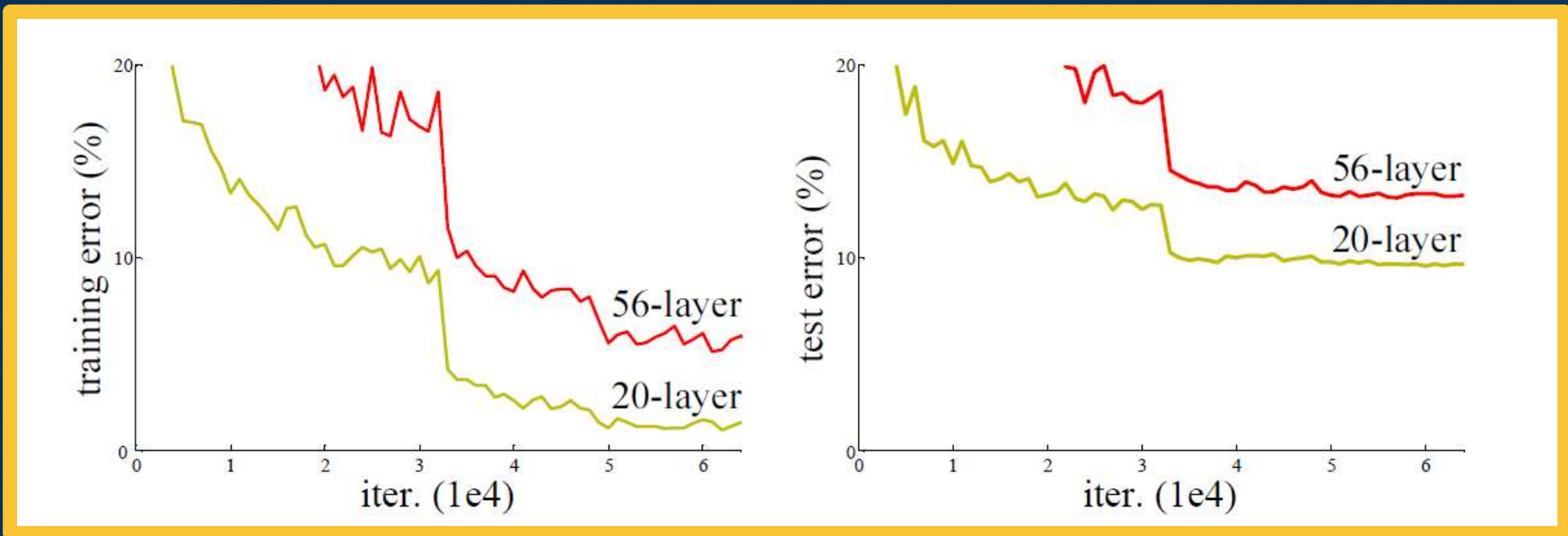
Key idea: Repeated blocks and multi-scale features



From: Szegedy et al. Going deeper with convolutions

Inception Module

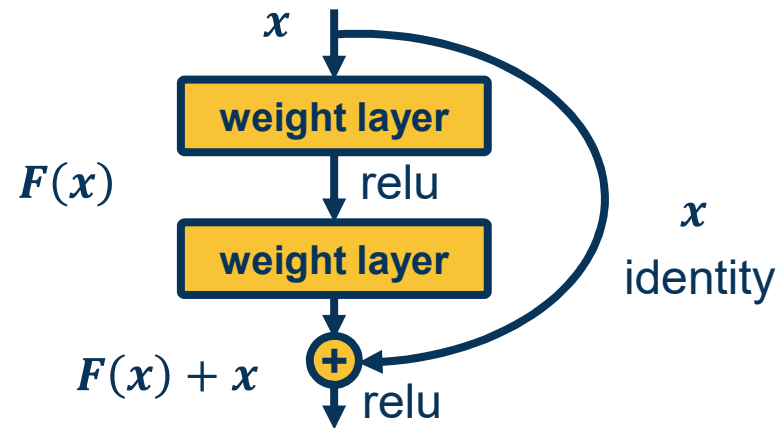
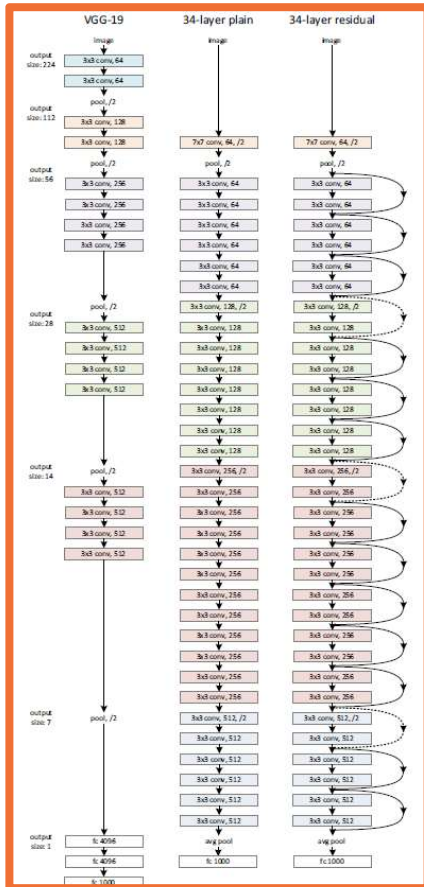
The Challenge of Depth



From: He et al., *Deep Residual Learning for Image Recognition*

Optimizing very deep networks is challenging!





Key idea: Allow information from a layer to propagate to any future layer (forward)

Same is true for gradients!

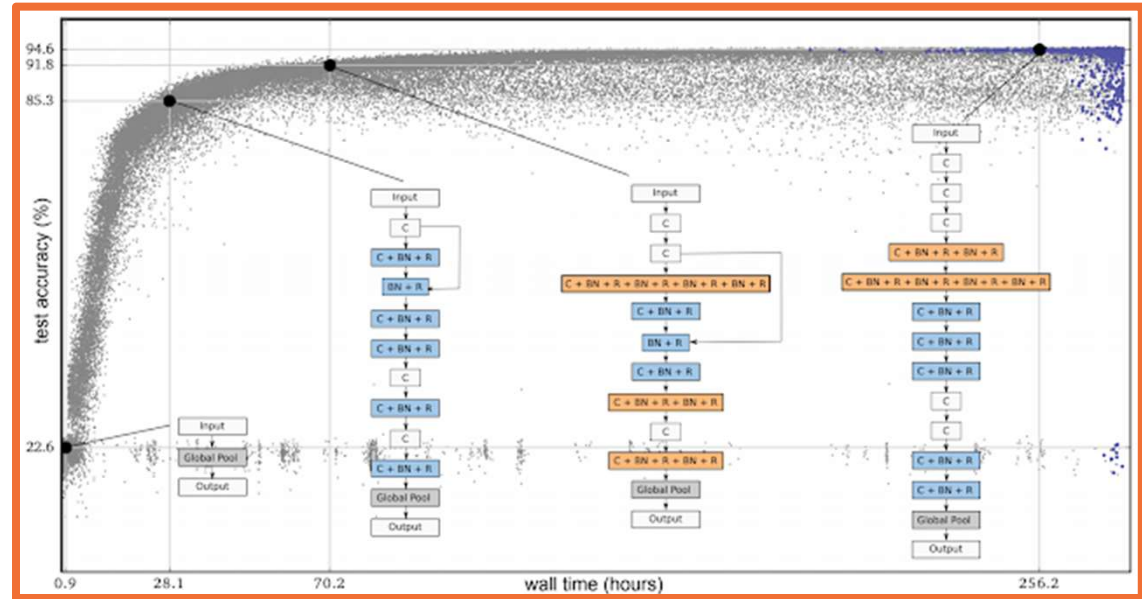
From: He et al., Deep Residual Learning for Image Recognition

Residual Blocks and Skip Connections



Several ways to *learn* architectures:

- Evolutionary learning and reinforcement learning
- Prune over-parameterized networks
- Learning of **repeated blocks** typical

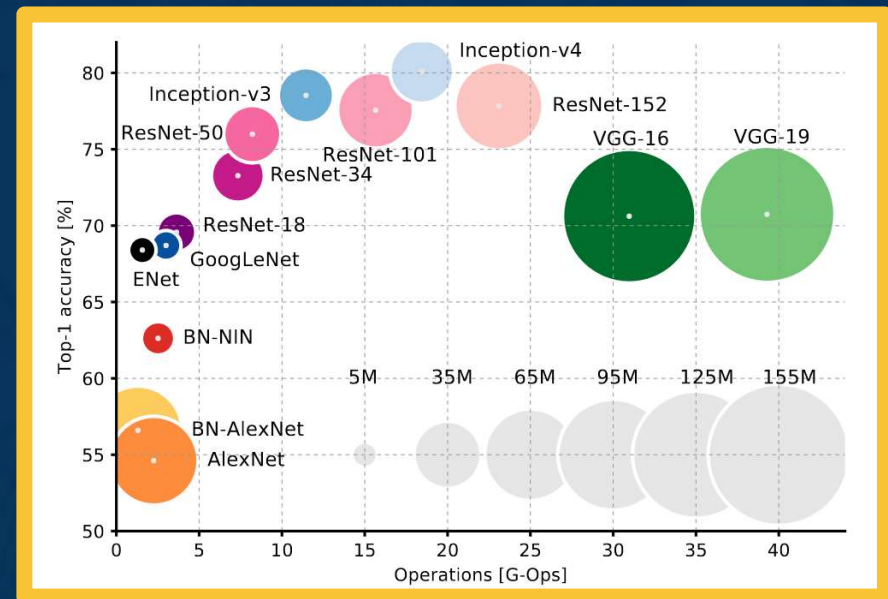
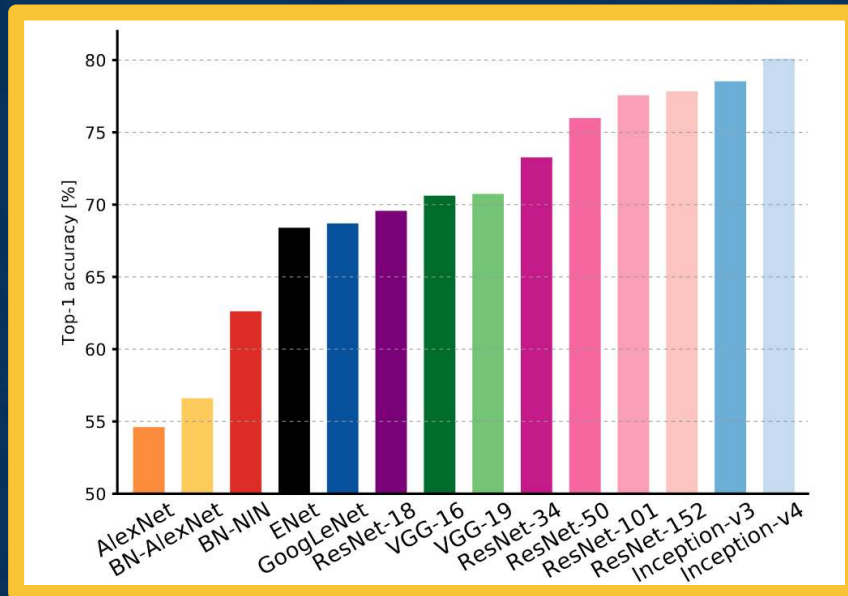


From: <https://ai.googleblog.com/2018/03/using-evolutionary-automl-to-discover.html>

Evolving Architectures and AutoML



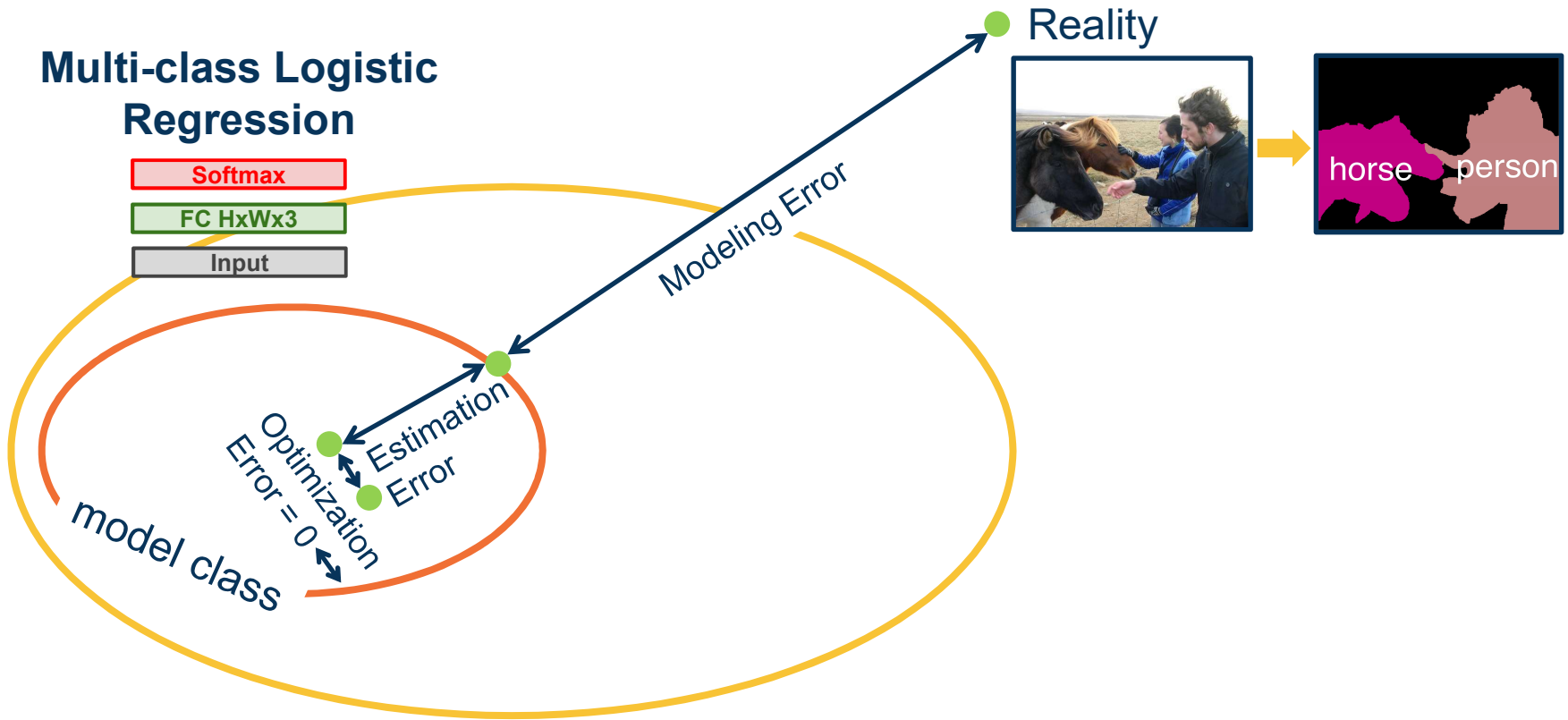
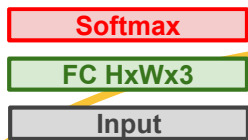
Computational Complexity



From: *An Analysis Of Deep Neural Network Models For Practical Application*

Transfer Learning & Generalization

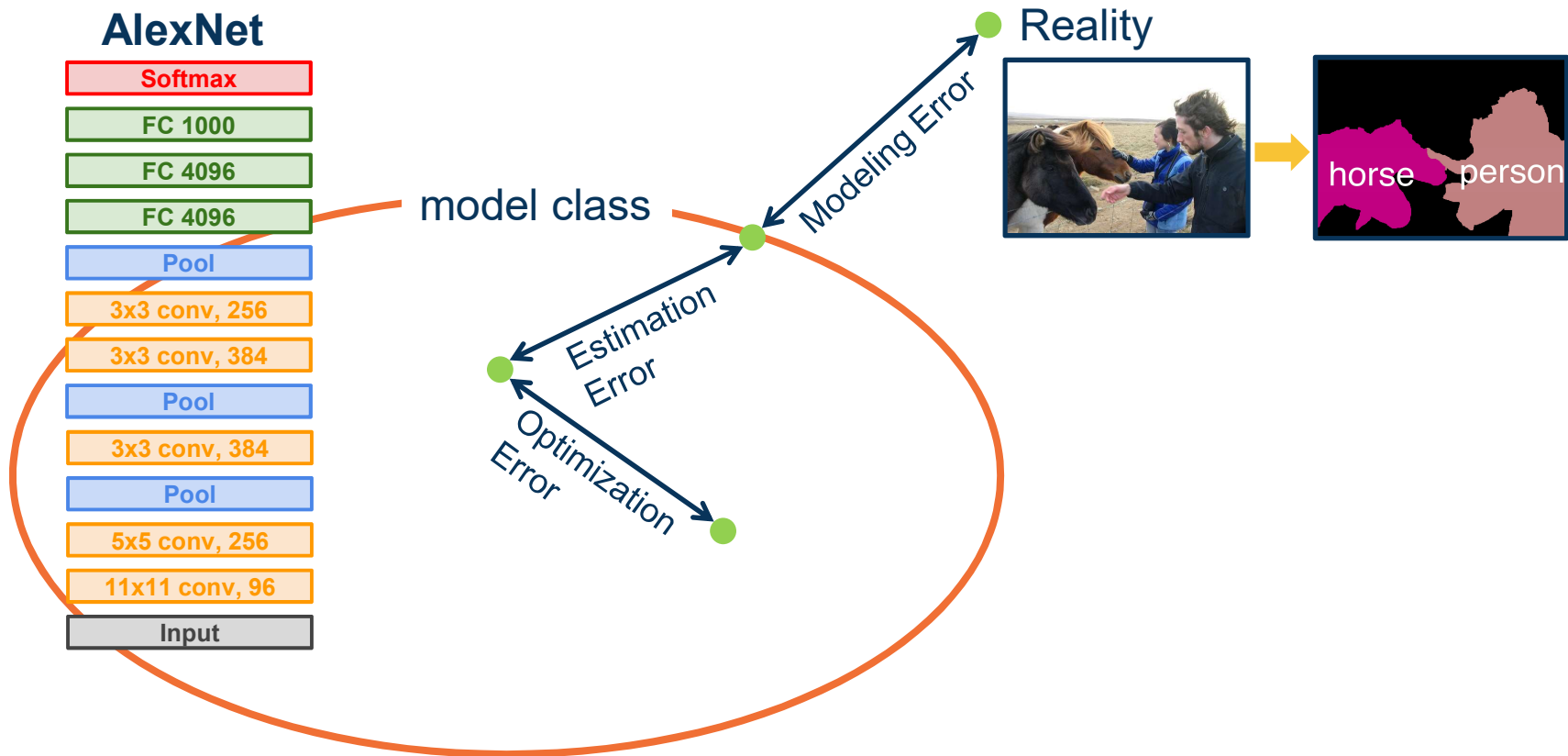
Multi-class Logistic Regression



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generalization

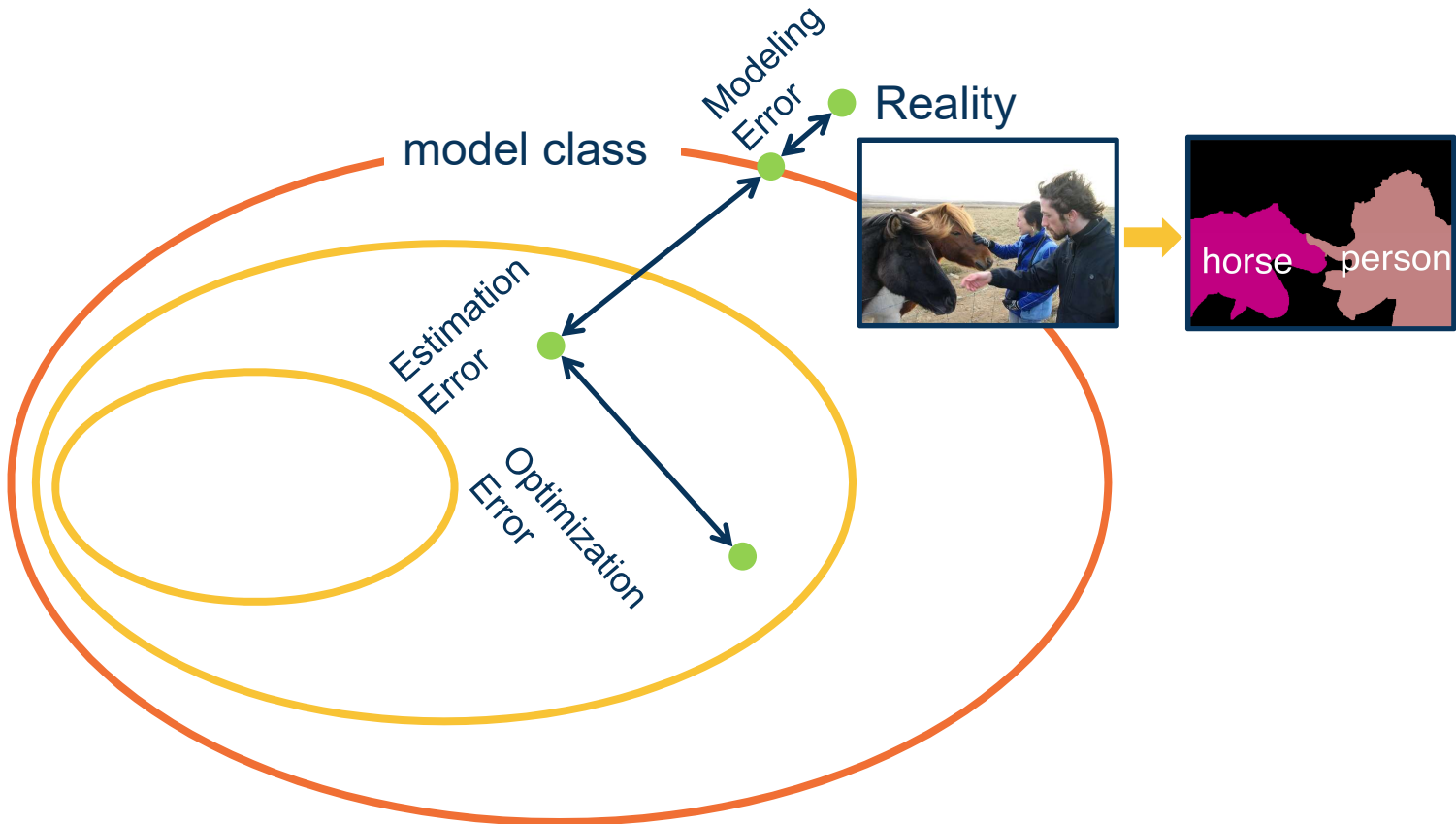




From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generalization

VGG19



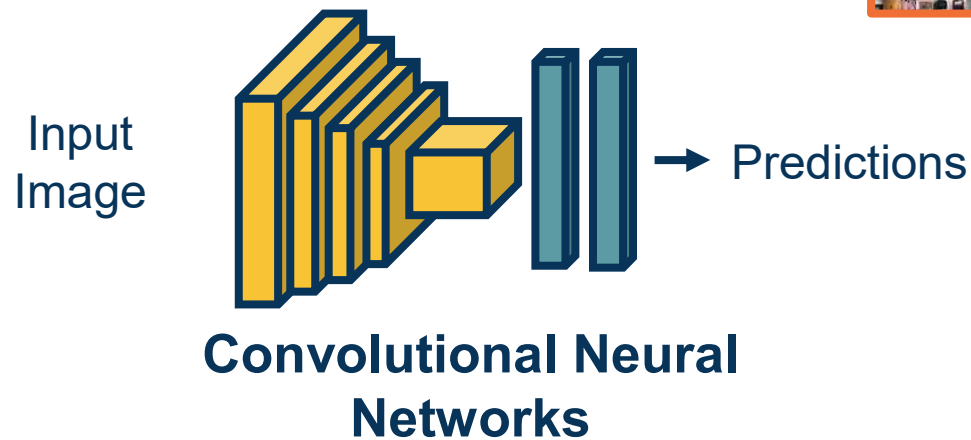
From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generalization



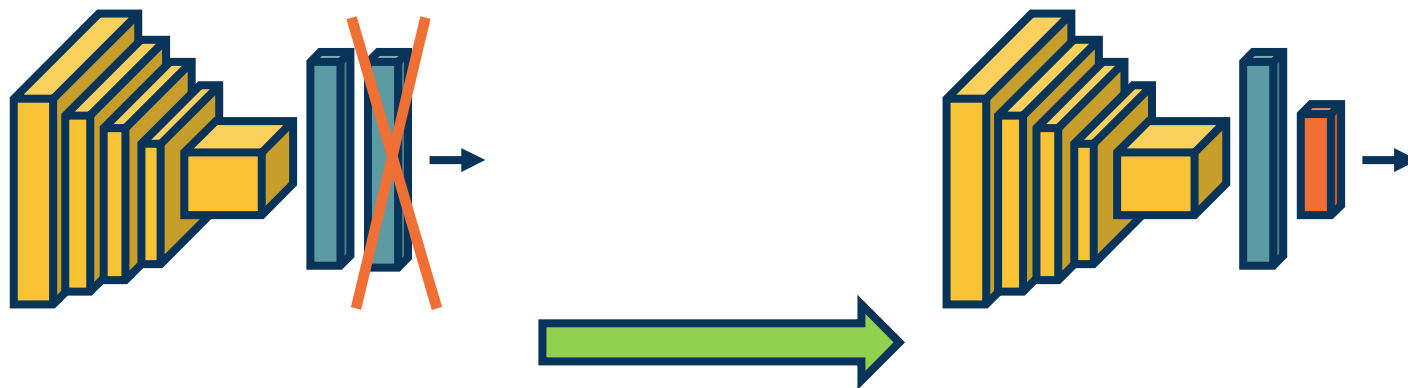
What if we don't have enough data?

Step 1: Train on large-scale dataset



Transfer Learning – Training on Large Dataset

Step 2: Take your custom data and **initialize** the network with weights trained in Step 1

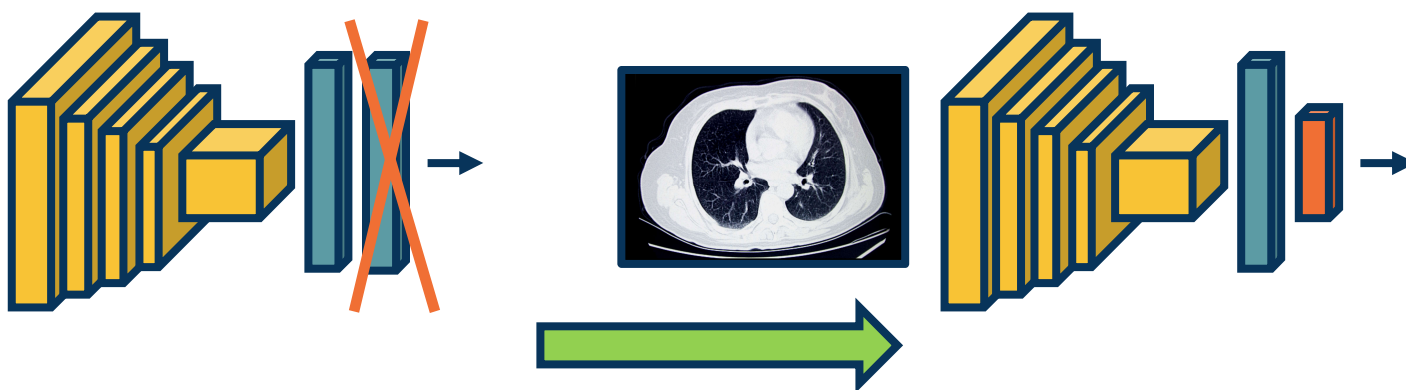


Replace last layer with new fully-connected for output nodes per new category

Initializing with Pre-Trained Network

Step 3: (Continue to) train on new dataset

- **Finetune:** Update all parameters
- **Freeze** feature layer: Update only last layer weights (used when not enough data)

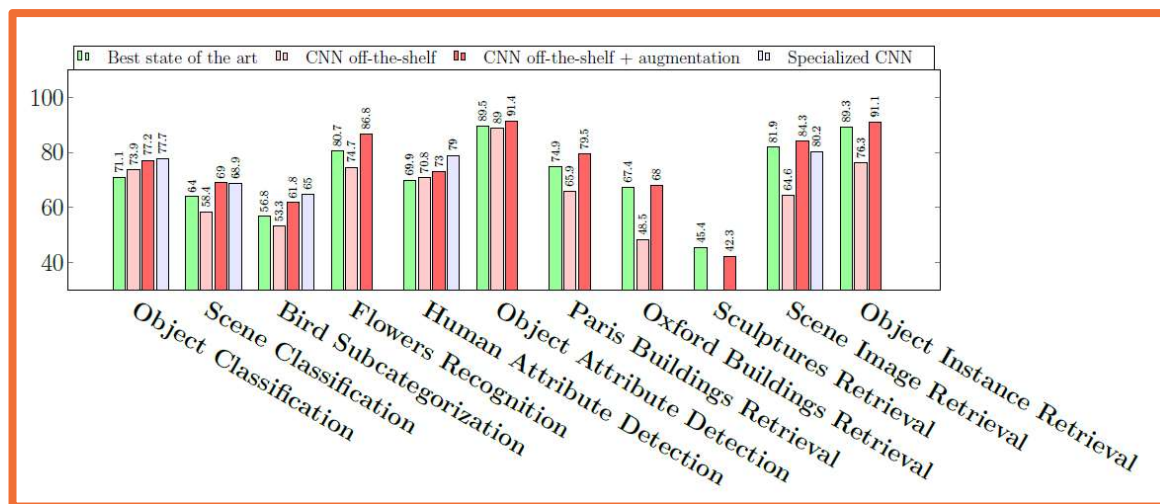


Replace last layer with new fully-connected for output nodes per new category

Finetuning on New Dataset

This works extremely well! It was surprising upon discovery.

- Features learned for 1000 object categories will work well for 1001st!
- Generalizes even across tasks (classification to object detection)



From: Razavian et al., CNN Features off-the-shelf: an Astounding Baseline for Recognition

Surprising Effectiveness of Transfer Learning

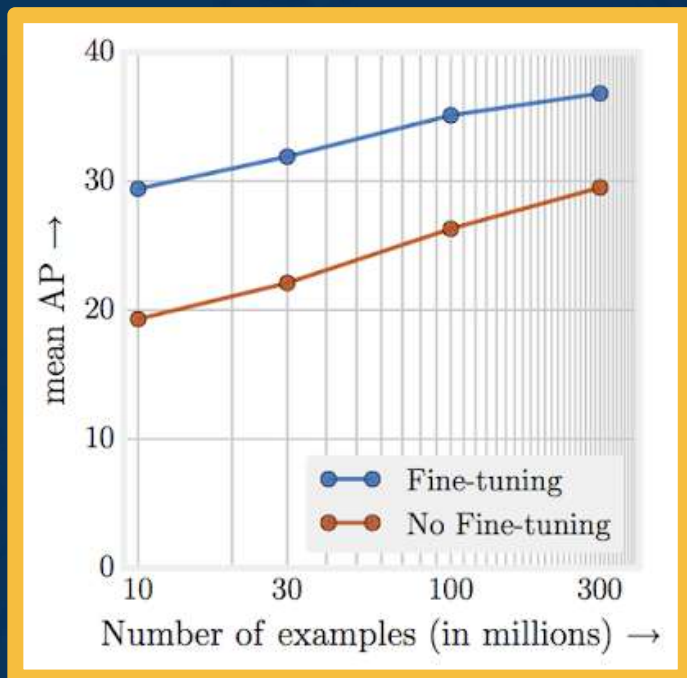
Learning with Less Labels

But it doesn't always work that well!

- ◆ If the **source** dataset you train on is very different from the **target** dataset, transfer learning is not as effective
- ◆ If you have enough data for the target domain, it just results in faster convergence
 - ◆ See He et al., “Rethinking ImageNet Pre-training”



Effectiveness of More Data



From: *Revisiting the Unreasonable Effectiveness of Data*
<https://ai.googleblog.com/2017/07/revisiting-unreasonable-effectiveness.html>

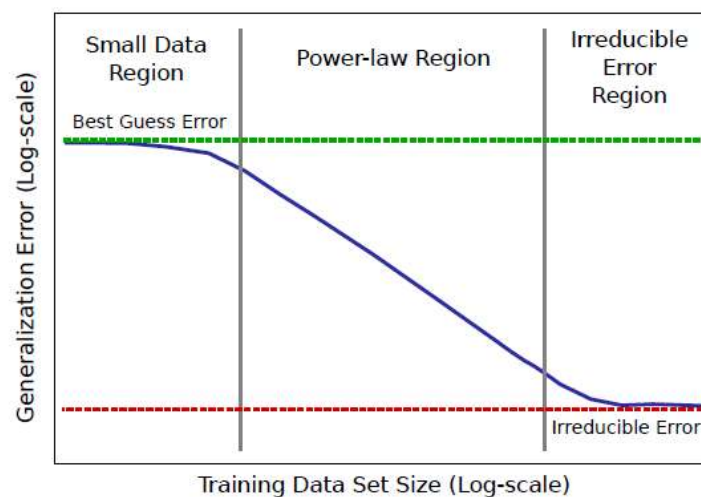


Figure 6: Sketch of power-law learning curves

From: *Hestness et al., Deep Learning Scaling Is Predictable*

There is a large number of different low-labeled settings in DL research

Setting	Source	Target	Shift Type
Semi-supervised	Single labeled	Single unlabeled	None
Domain Adaptation	Single labeled	Single unlabeled	Non-semantic
Domain Generalization	Multiple labeled	Unknown	Non-semantic
Cross-Category Transfer	Single labeled	Single unlabeled	Semantic
Few-Shot Learning	Single labeled	Single few-labeled	Semantic
Un/Self-Supervised	Single unlabeled	Many labeled	Both/Task

Non-Semantic Shift



Semantic Shift

