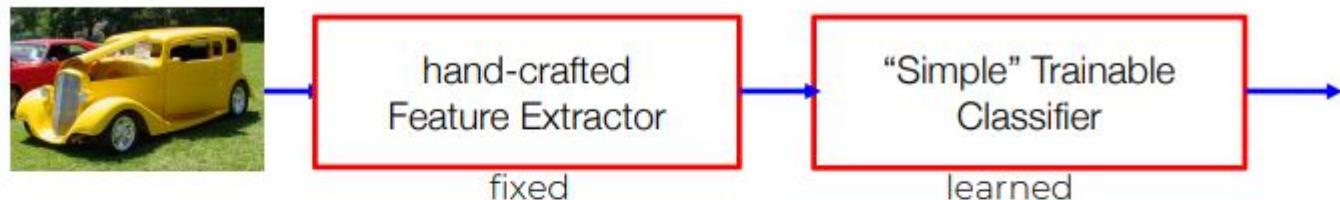


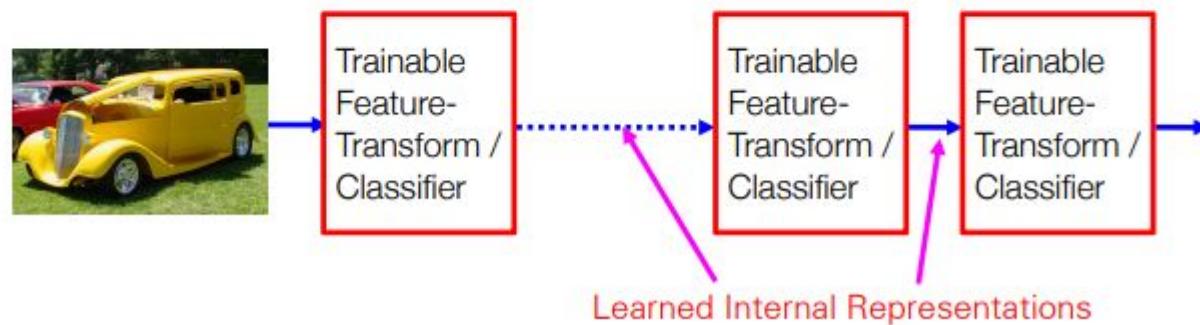
Convolutional Neural Network

So what is deep learning?

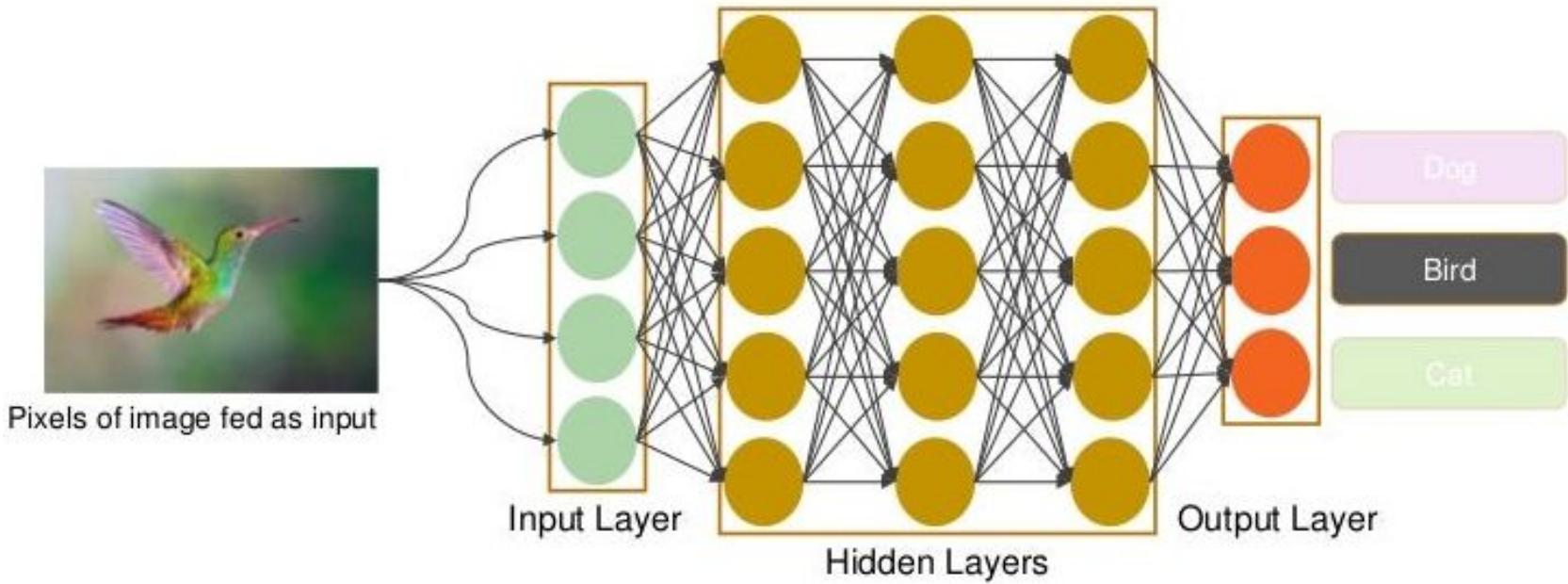
- “Shallow” models



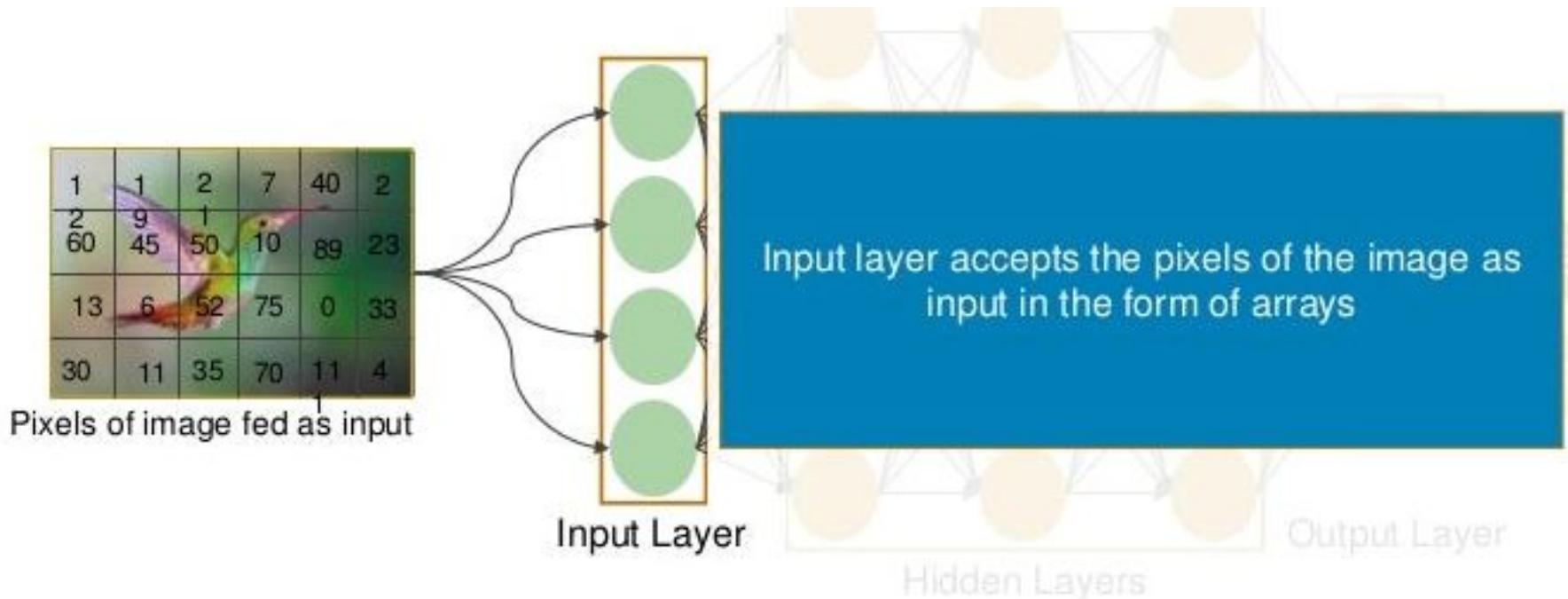
- Deep models



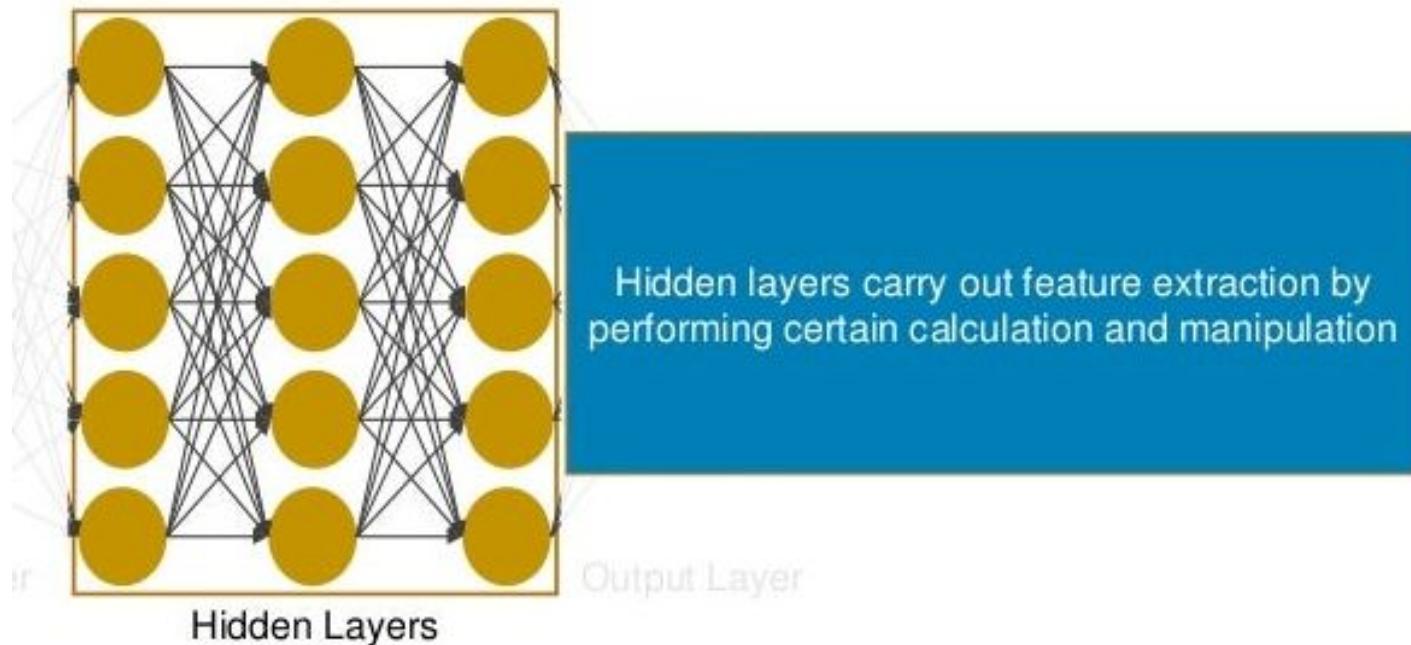
Outline

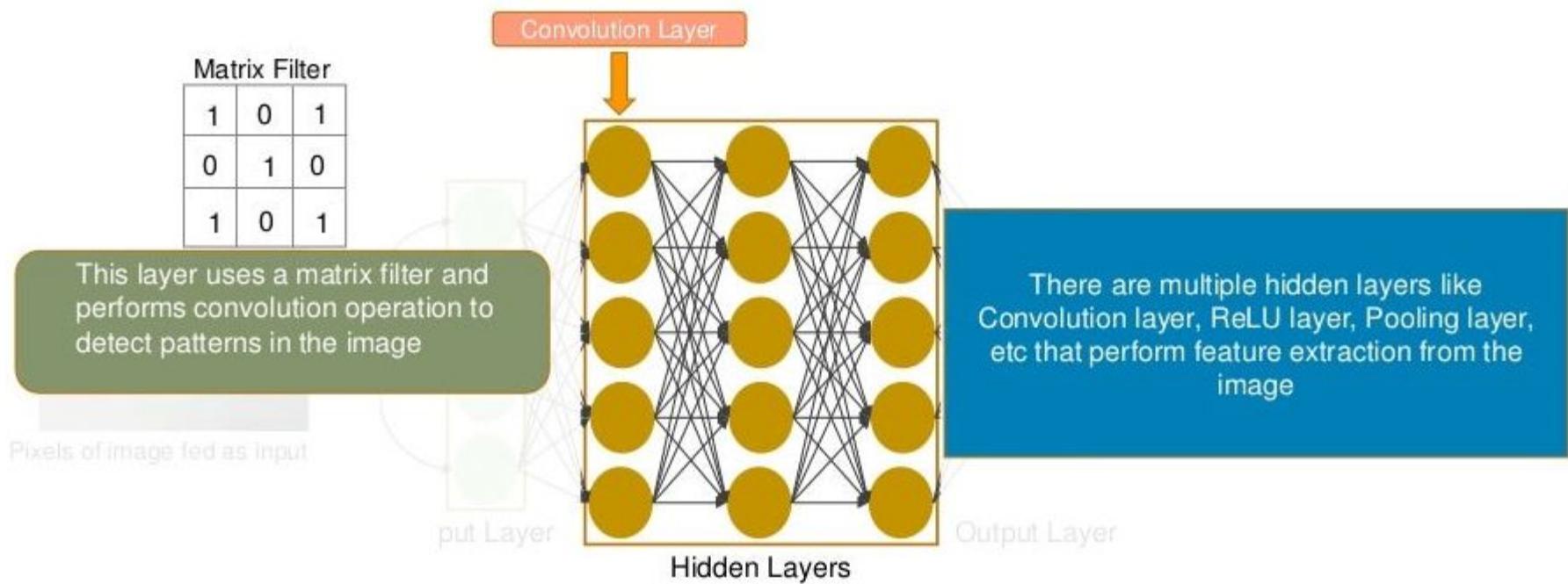


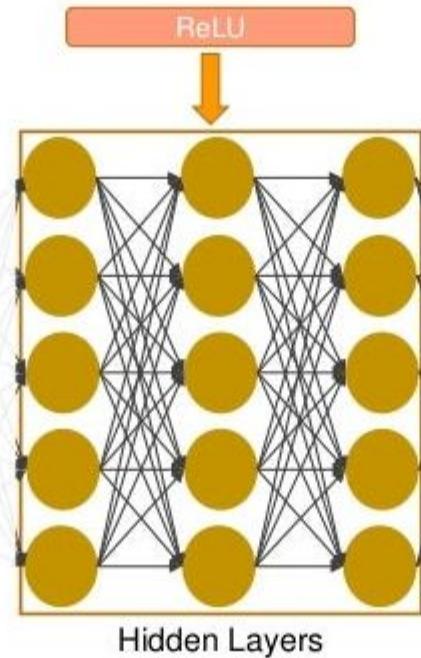
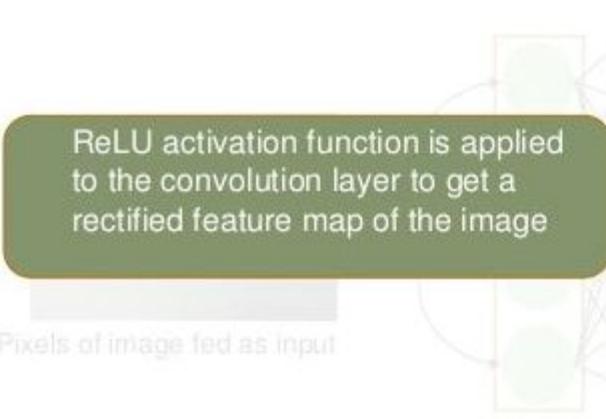
Outline



Outline

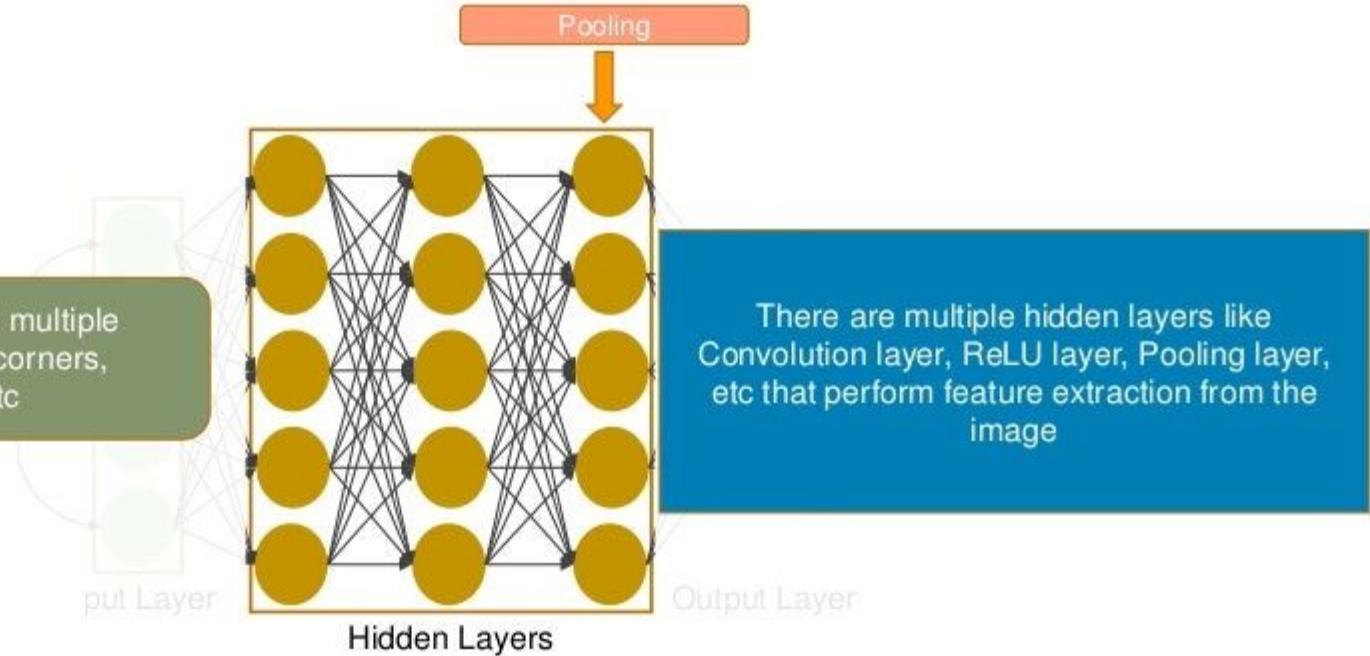






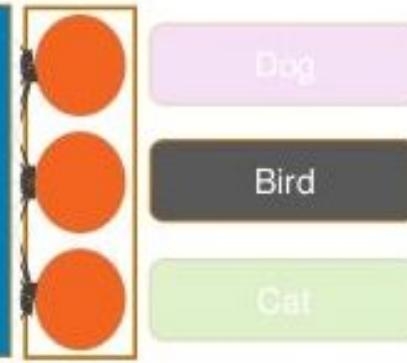
There are multiple hidden layers like Convolution layer, ReLU layer, Pooling layer, etc that perform feature extraction from the image

Pixels of image fed as input



The diagram illustrates a neural network architecture. It starts with an 'Input Layer' represented by a white rectangle containing a small icon of a person. This layer connects to three 'Hidden Layers', each represented by a row of three yellow circles. The final layer is the 'Output Layer', which consists of three colored rectangles (pink, dark grey, and light green) containing the words 'Dog', 'Bird', and 'Cat' respectively. A blue callout box is positioned above the output layer, containing the text: 'Finally there is a fully connected layer that identifies the object in the image'. An orange rectangular box highlights the top circle of the output layer, corresponding to the word 'Dog'.

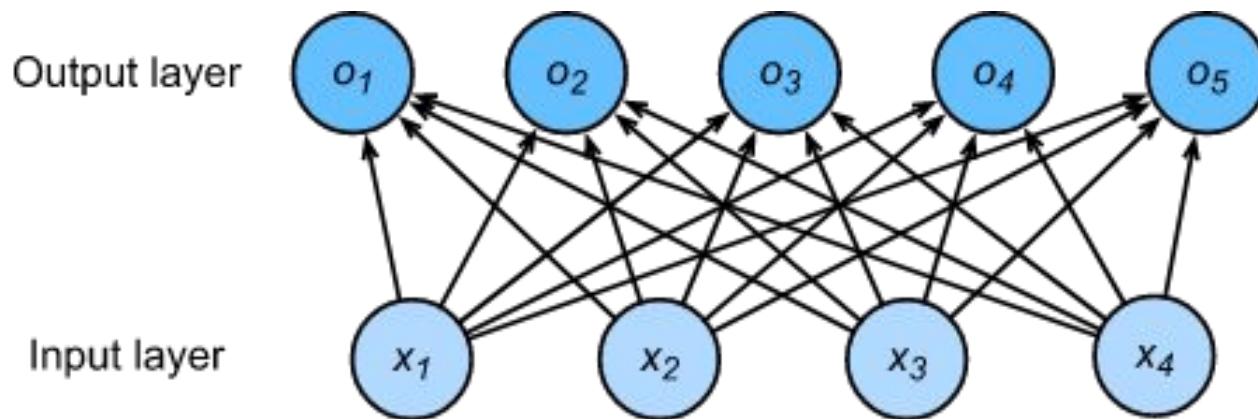
Finally there is a fully connected layer that identifies the object in the image



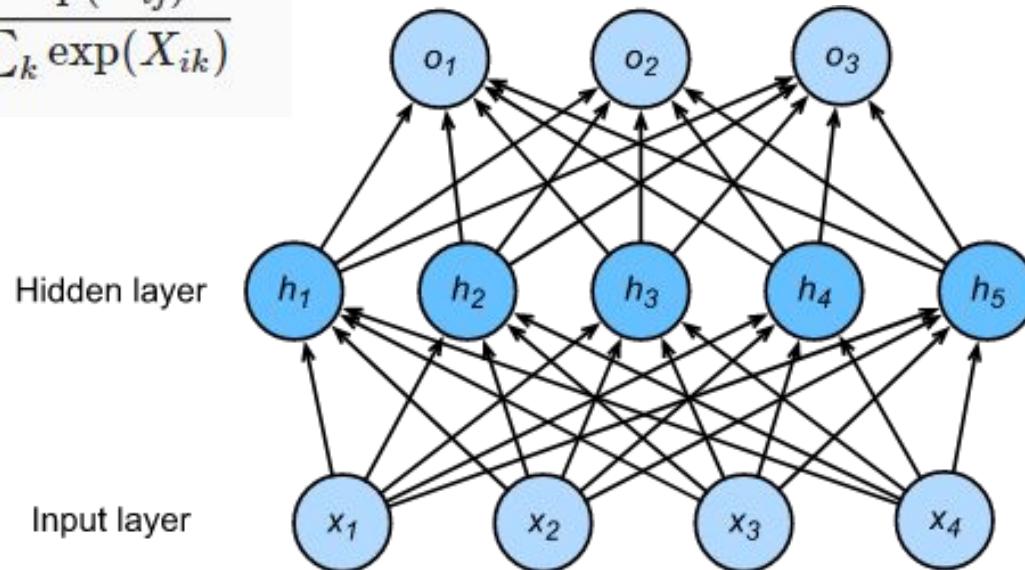
Input Layer

Hidden Layers

Output Layer



$$\text{softmax}(\mathbf{X})_{ij} = \frac{\exp(X_{ij})}{\sum_k \exp(X_{ik})}$$



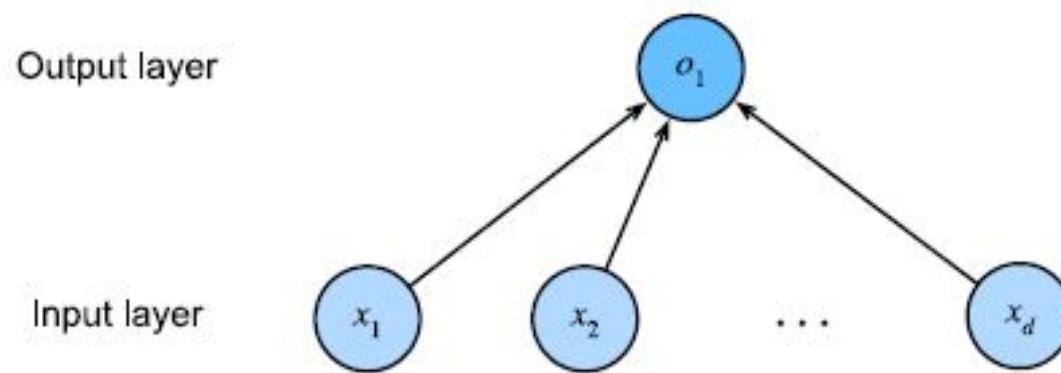


Fig. 3.1.2: Linear regression is a single-layer neural network.

From fully connected to convolutions



Classifying Dogs and Cats in Images

- Use a good camera
- RGB image has 36M elements
- The model size of a single hidden layer MLP with a 100 hidden size is 3.6 Billion parameters
- Exceeds the population of dogs and cats on earth
(900M dogs + 600M cats)



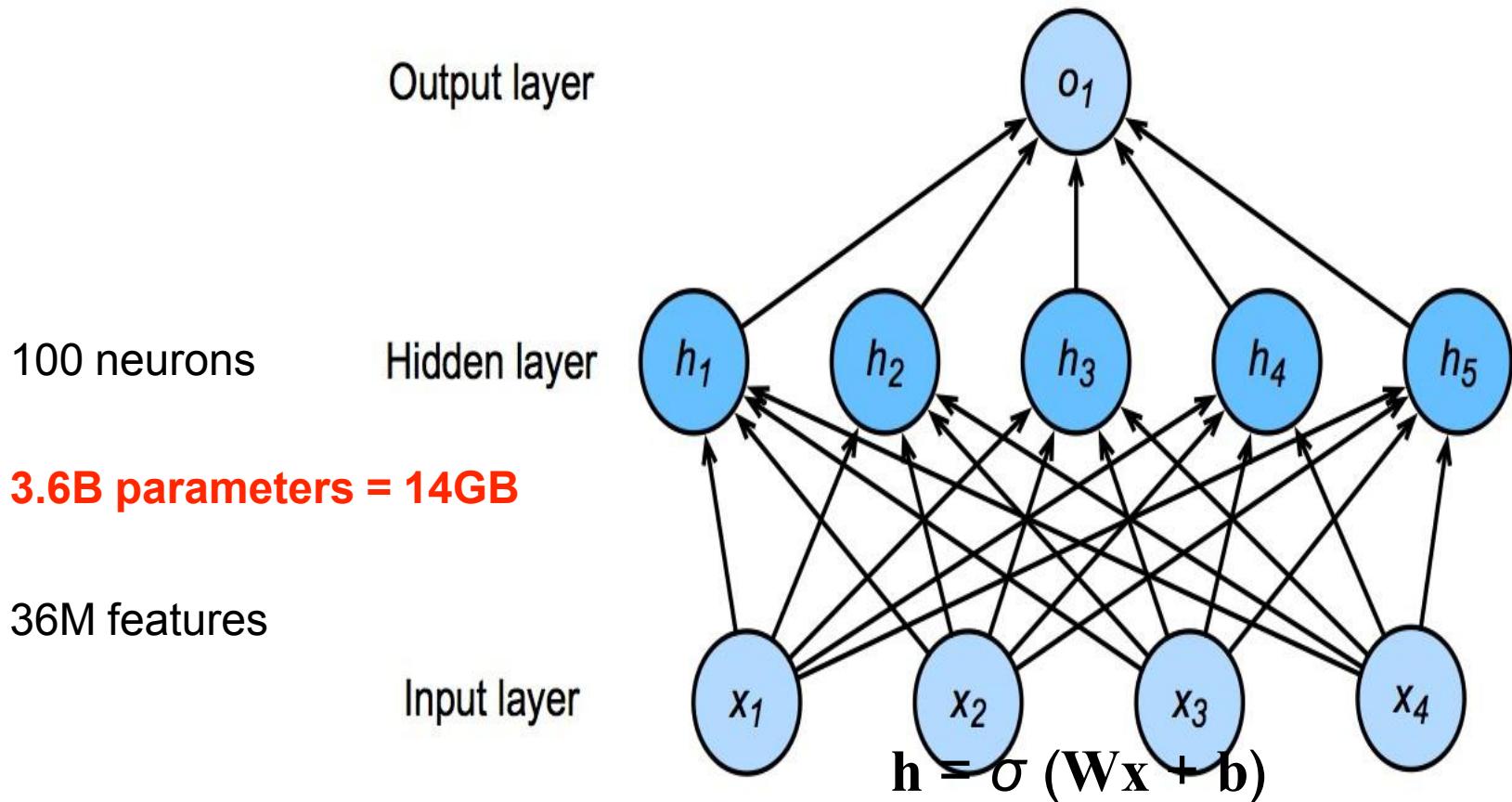
Dual

12MP

wide-angle and
telephoto cameras



Flashback - Network with one hidden layer



Where is
Waldo?



- Translation Invariance
- Locality



Convolution

the *convolution* between two functions, say $f, g : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

$$(f * g)(\mathbf{x}) = \int f(\mathbf{z})g(\mathbf{x} - \mathbf{z})d\mathbf{z}.$$

Discrete

$$(f * g)(i) = \sum_a f(a)g(i - a).$$

Two dimension

$$(f * g)(i, j) = \sum_a \sum_b f(a, b)g(i - a, j - b).$$

Convolutions for Images

The Cross-Correlation Operation

Input	Kernel	Output													
<table border="1" style="border-collapse: collapse; width: 100%;"><tr><td style="padding: 5px;">0</td><td style="padding: 5px;">1</td><td style="padding: 5px;">2</td></tr><tr><td style="padding: 5px;">3</td><td style="padding: 5px;">4</td><td style="padding: 5px;">5</td></tr><tr><td style="padding: 5px;">6</td><td style="padding: 5px;">7</td><td style="padding: 5px;">8</td></tr></table>	0	1	2	3	4	5	6	7	8	$*$	<table border="1" style="border-collapse: collapse; width: 100%;"><tr><td style="padding: 5px;">0</td><td style="padding: 5px;">1</td></tr><tr><td style="padding: 5px;">2</td><td style="padding: 5px;">3</td></tr></table>	0	1	2	3
0	1	2													
3	4	5													
6	7	8													
0	1														
2	3														
	=	<table border="1" style="border-collapse: collapse; width: 100%;"><tr><td style="padding: 5px;">19</td><td style="padding: 5px;">25</td></tr><tr><td style="padding: 5px;">37</td><td style="padding: 5px;">43</td></tr></table>	19	25	37	43									
19	25														
37	43														

Rethinking Dense Layers

- Reshape inputs and output into matrix (width, height)
- Reshape weights into 4-D tensors (h, w) to (h', w')

$$h_{i,j} = \sum w_{i,j,k,l} x_{k,l} = \sum v_{i,j,a,b} x_{i+a, j+b, k, l, a, b}$$

V is re-indexes W such that

$$v_{i,j,a,b} = w_{i,j,i+a,j+b}$$

Idea #1 - Translation Invariance

$$h_{i,j} = \sum_{a,b} v_{i,j,a,b} x_{i+a,j+b}$$

- A shift in x also leads to a shift in h
- v should not depend on (i,j) . Fix via $v_{i,j,a,b} = v_{a,b}$

$$h_{i,j} = \sum_{a,b} v_{a,b} x_{i+a,j+b}$$

That's a 2-D convolution
~~cross-correlation~~

Idea #2 - Locality

$$h_{i,j} = \sum_{a,b} v_{a,b} x_{i+a, j+b}$$

- We shouldn't look very far from $x(i,j)$ in order to assess what's going on at $h(i,j)$

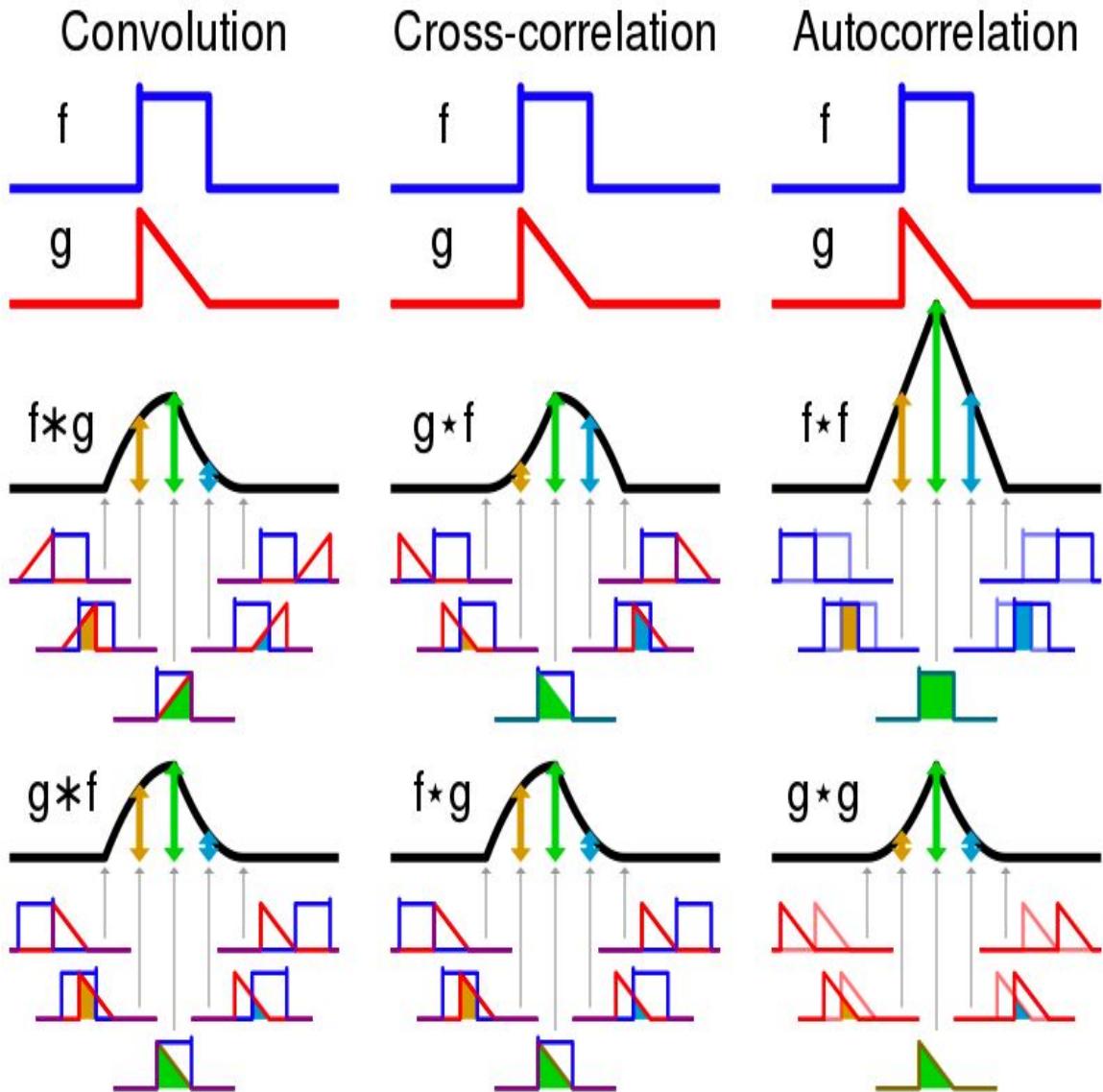
- Outside range $|a|, |b| > \Delta$

parameters vanish $v_{a,b} = 0$

$\Delta \Delta$

$$h_{i,j} = \sum_{b=-\Delta} \sum_{a=-\Delta} v_{a,b} x_{i+a, j+b}$$

Convolution



2-D Cross Correlation

Input

0	1	2
3	4	5
6	7	8

*

Kernel

0	1
2	3

=

Output

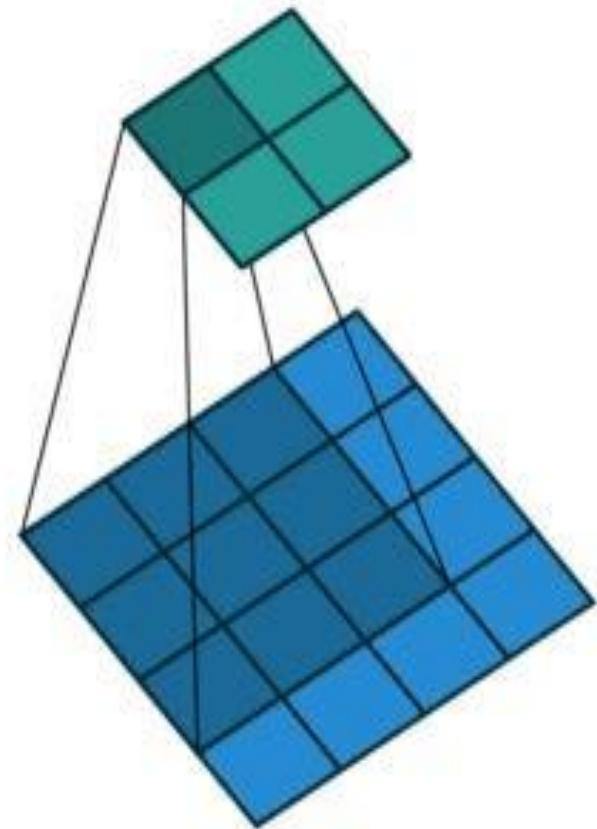
19	25
37	43

$$0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19,$$

$$1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25,$$

$$3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37,$$

$$4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43.$$



2-D Convolution Layer

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 2 & 3 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 19 & 25 \\ \hline 37 & 43 \\ \hline \end{array}$$

- $\mathbf{X} : n_h \times n_w$ input matrix
- $\mathbf{W} : k_h \times k_w$ kernel matrix
- b : scalar bias
- $\mathbf{Y} : (n_w - k_w + 1) \times (n_h - k_h + 1)$ output matrix

$$\mathbf{Y} = \mathbf{X} \star \mathbf{W} + b$$

- \mathbf{W} and b are learnable parameters

Examples

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



Edge Detection



$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$



Sharpen

(wikipedia)

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



Gaussian Blur

Examples



(Rob Fergus)



Feature Map and Receptive Field

The convolutional layer output is sometimes called a **feature map**, as it can be regarded as the learned representations (features) in the spatial dimensions (e.g., width and height) to the subsequent layer.

In CNNs, for any element x of some layer, its ***receptive field*** refers to all the elements (from all the previous layers) that may affect the calculation of x during the forward propagation.

Convolutions



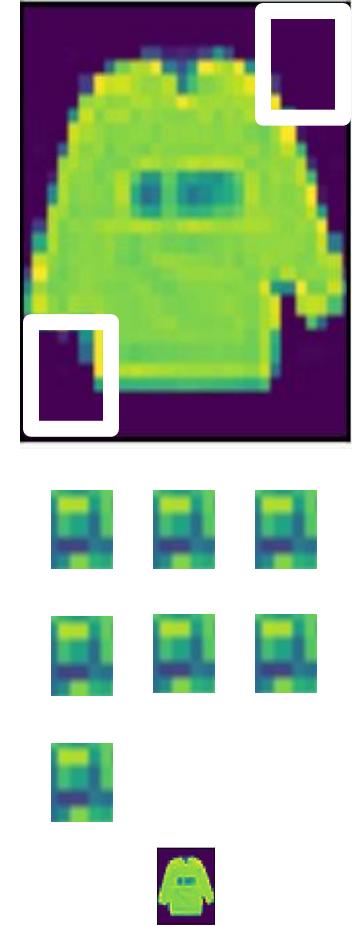


A man in a dark jacket and blue jeans is shown walking across a white crosswalk on a city street. He is captured in four different stages of his walk, from the start of his step to the end, creating a sense of motion. The background shows a busy street with parked cars and buildings, suggesting an urban environment.

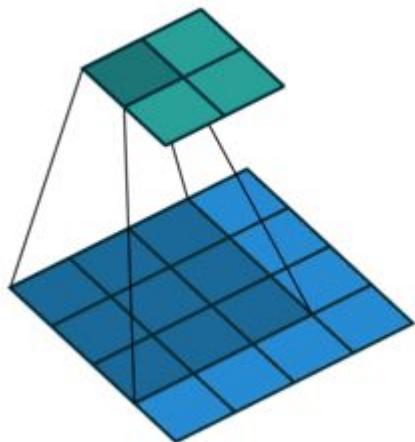
Padding and Stride

Padding

- Given a 32×32 input image
- Apply convolutional layer with 5×5 kernel
 - 28×28 output with 1 layer
 - 4×4 output with 7 layers
- Shape **decreases faster with larger kernels**
 - obliterating any interesting information on the boundaries of the original image
 - Shape reduces from $n_h \times n_w$ to $(n_h - k_h + 1) \times (n_w - k_w + 1)$ using kernel shape is $k_h \times k_w$

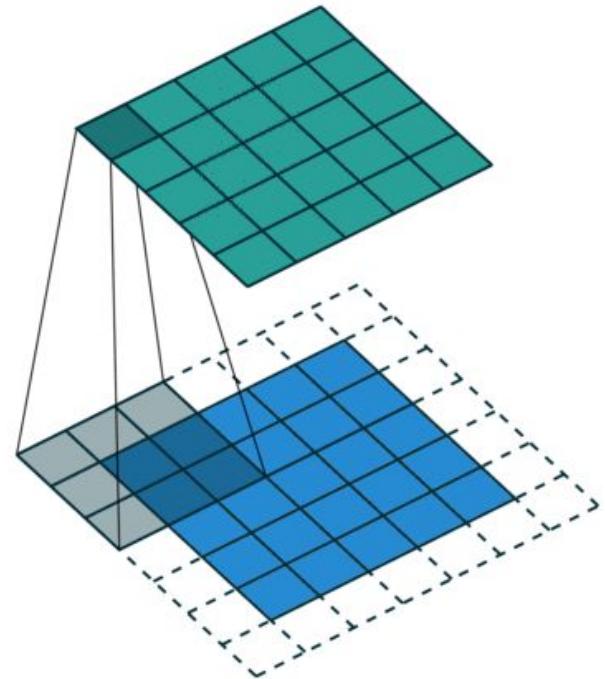


Convolutional Layer



Convolution without padding

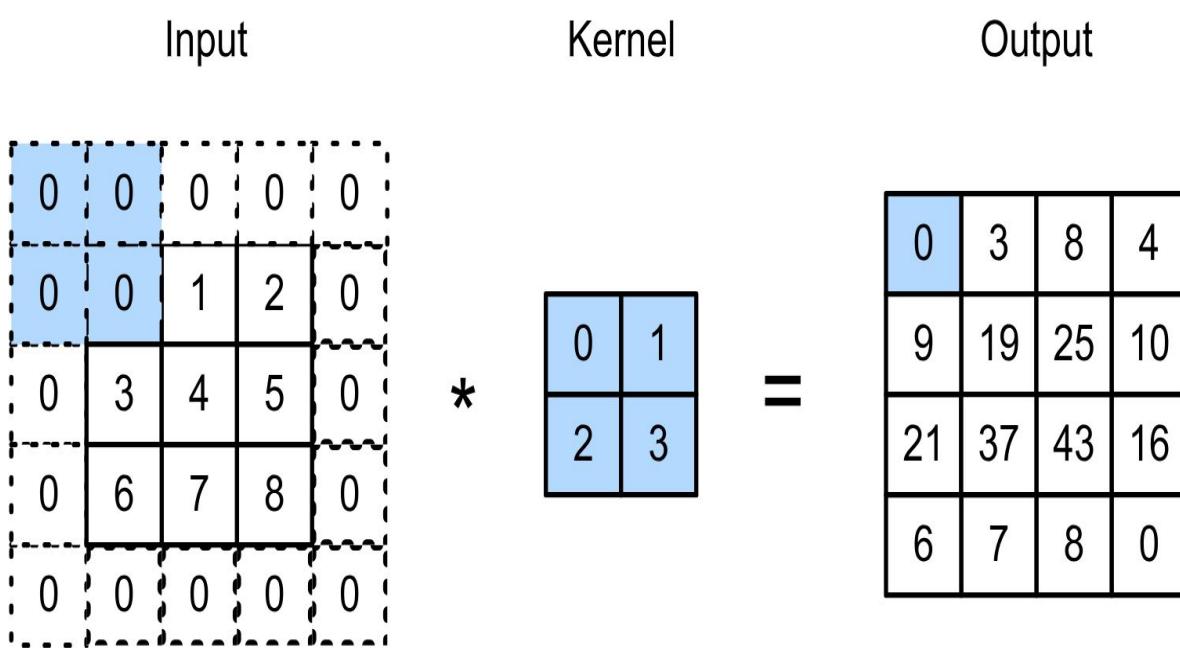
We can also pad the edges of your images with 0-valued pixels as to fully scan the original image and preserve its complete dimensions.



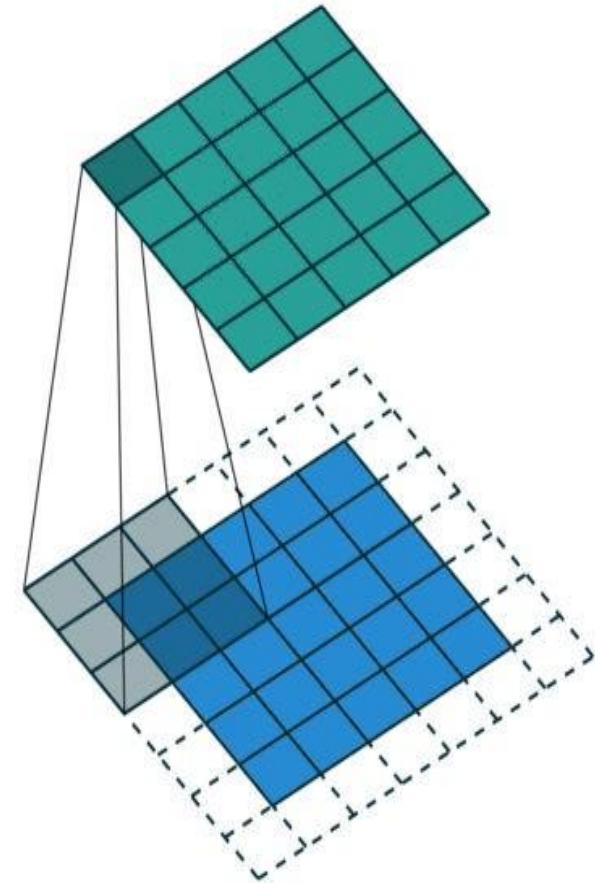
Convolution with padding

Padding

Padding adds rows/columns around input



$$0 \times 0 + 0 \times 1 + 0 \times 2 + 0 \times 3 = 0$$

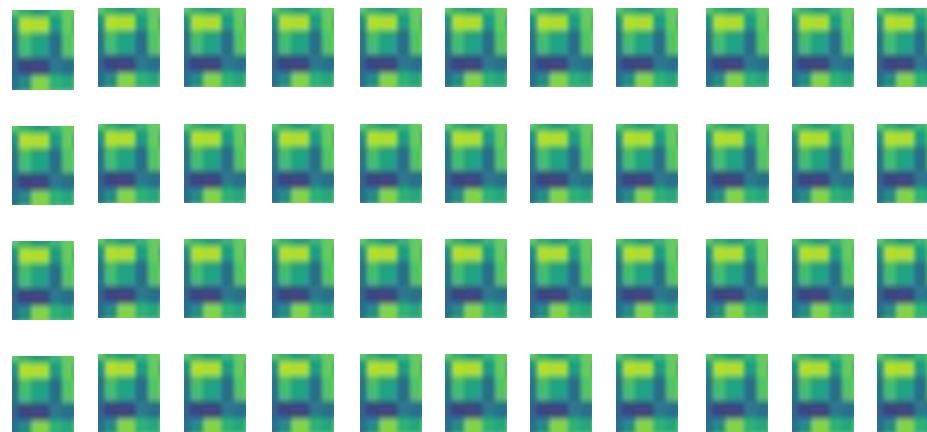


Padding

- Padding p_h rows and p_w columns, output shape will be
$$(n_h - k_h + p_h + 1) \times (n_w - k_w + p_w + 1)$$
- A common choice is $p_h = k_h - 1$ and $p_w = k_w - 1$
 - Odd k_h : pad $\lceil \frac{p_h}{2} \rceil$ on both sides
 - Even k_h : pad $\lceil \frac{p_h}{2} \rceil$ on top, $\lfloor \frac{p_h}{2} \rfloor$ on bottom

Stride

- Padding reduces shape linearly with #layers
 - Given a 224×224 input with a 5×5 kernel, needs 44 layers to reduce the shape to 4×4
 - Requires a large amount of computation



Stride

Sometimes, either for computational efficiency or because we wish to downsample, we move our window more than one element at a time, skipping the intermediate locations.

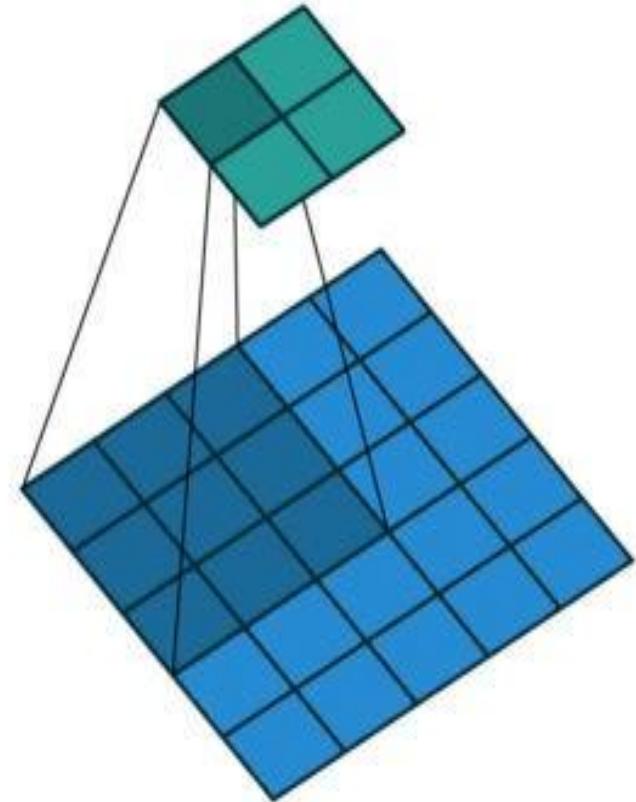
- Stride is the #rows/#columns per slide

Strides of 3 and 2 for height and width

Input	Kernel	Output
$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 3 & 4 & 5 & 0 \\ 0 & 6 & 7 & 8 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$	$*\quad \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}$	$= \begin{bmatrix} 0 & 8 \\ 6 & 8 \end{bmatrix}$

$$0 \times 0 + 0 \times 1 + 1 \times 2 + 2 \times 3 = 8$$

$$0 \times 0 + 6 \times 1 + 0 \times 2 + 0 \times 3 = 6$$



Stride

- Given stride $\frac{s}{h}$ for the height and stride $\frac{s}{w}$ for the width, the output shape is

$$\lfloor (n_h - k_h + p_h + s_h)/s_h \rfloor \times \lfloor (n_w - k_w + p_w + s_w)/s_w \rfloor$$

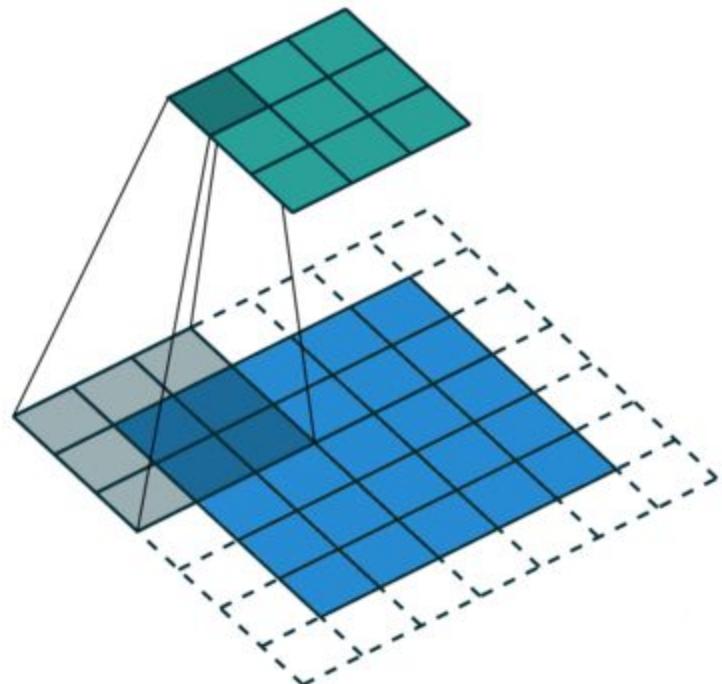
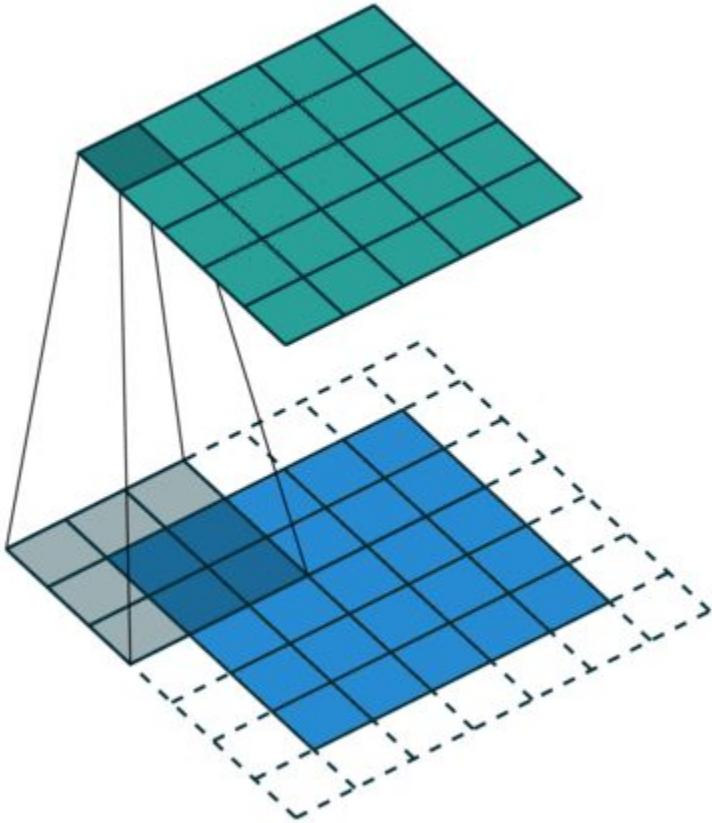
- With $p_h = k_h - 1$ and $p_w = k_w - 1$

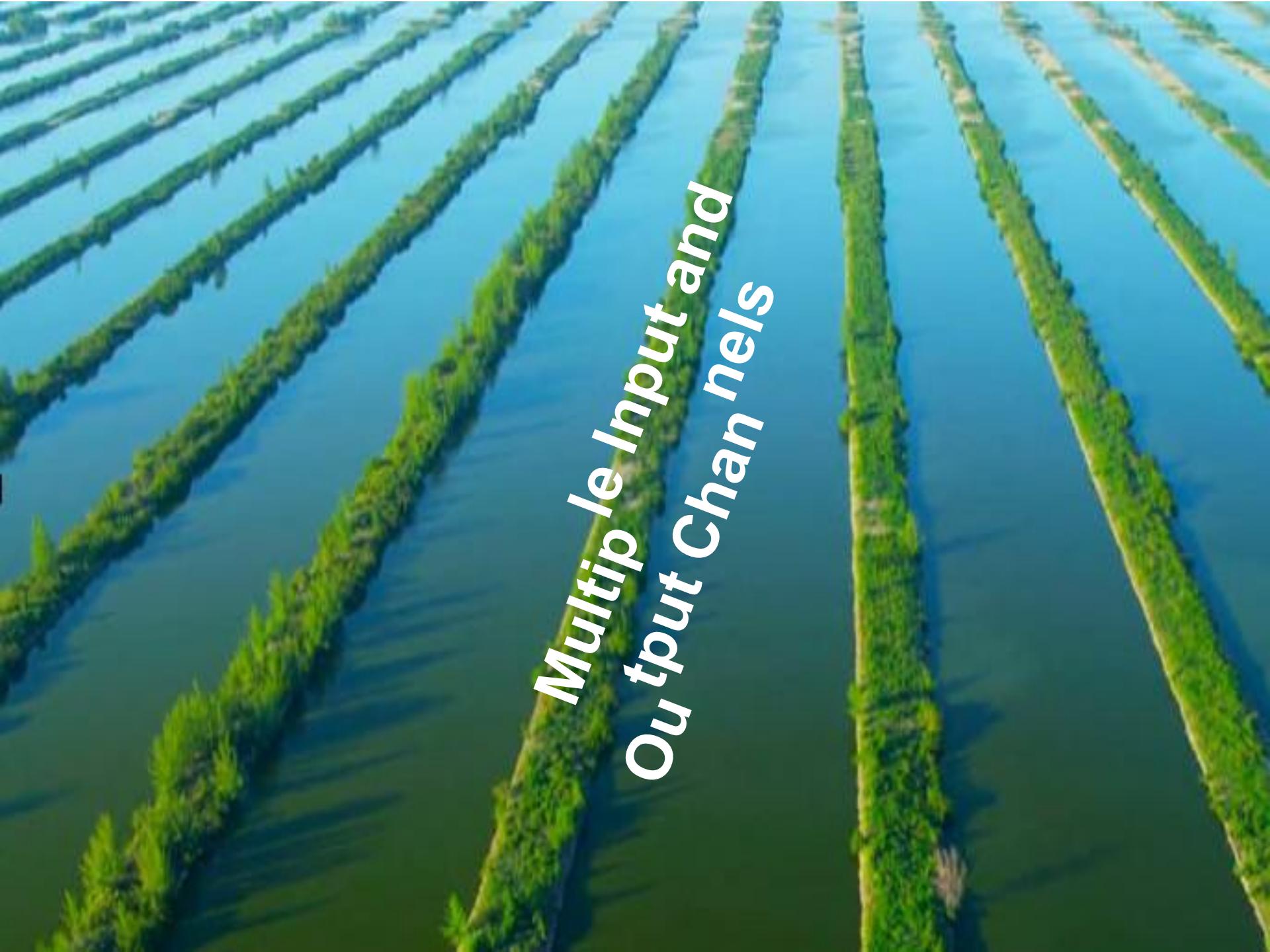
$$\lfloor (n_h + s_h - 1)/s_h \rfloor \times \lfloor (n_w + s_w - 1)/s_w \rfloor$$

- If input height/width are divisible by strides

$$(n_h/s_h) \times (n_w/s_w)$$

Stride



An aerial photograph of a river system, likely a delta or a network of canals, showing multiple parallel waterways flowing through a green landscape. The water is a deep blue-green color. The text is overlaid on the right side of the image.

*Multiple Input and
Output Channels*

Multiple Input Channels

- Color image may have three RGB channels
- Converting to grayscale loses information



Multiple Input Channels

- Color image may have three RGB channels
- Converting to grayscale loses information

f



R (Red)



G (green)

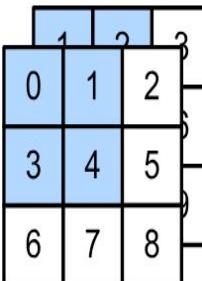
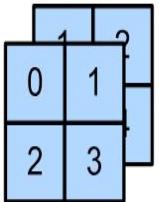
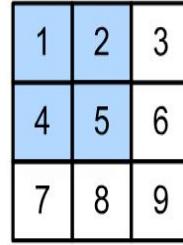
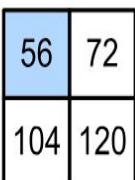
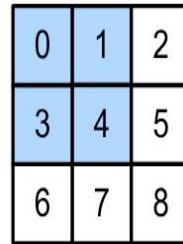
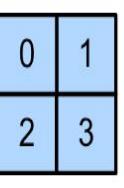


B (blue)



Multiple Input Channels

- Have a kernel **for each channel**, and then **sum results over channels**

Input	Kernel	Input	Kernel	Output
		$*$		$=$
			$*$	
			$+$	$=$
			$*$	

$(1 \times 1 + 2 \times 2 + 4 \times 3 + 5 \times 4)$
 $+(0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3)$
 $= 56$

Multiple Input Channels

- $\mathbf{X} : c_i \times n_h \times n_w$ input
- $\mathbf{W} : c_i \times k_h \times k_w$ kernel
- $\mathbf{Y} : \underset{h}{m} \times \underset{w}{m}$ output

ci

$$\mathbf{Y} = \sum_{i=0} \mathbf{X}_{i,:,:} \star \mathbf{W}_{i,:,:}$$

Multiple Output Channels

- No matter how many inputs channels, so far we always get **single output channel**
 - We can **have multiple 3-D kernels**, each **one generates a output channel**
 - Input $\mathbf{X} : c_i \times n_h \times n_w$
 - Kernel $\mathbf{W} : c_o \times c_i \times k_h \times k_w$
 - Output $\mathbf{Y} : c^o \times m^h \times m^w$
- $$\mathbf{Y}_{i,:,:} = \mathbf{X}_{i,:,:} \star \mathbf{W}_{i,:,:}$$
- for $i = 1, \dots, c_o$

Multiple Input/Output Channels

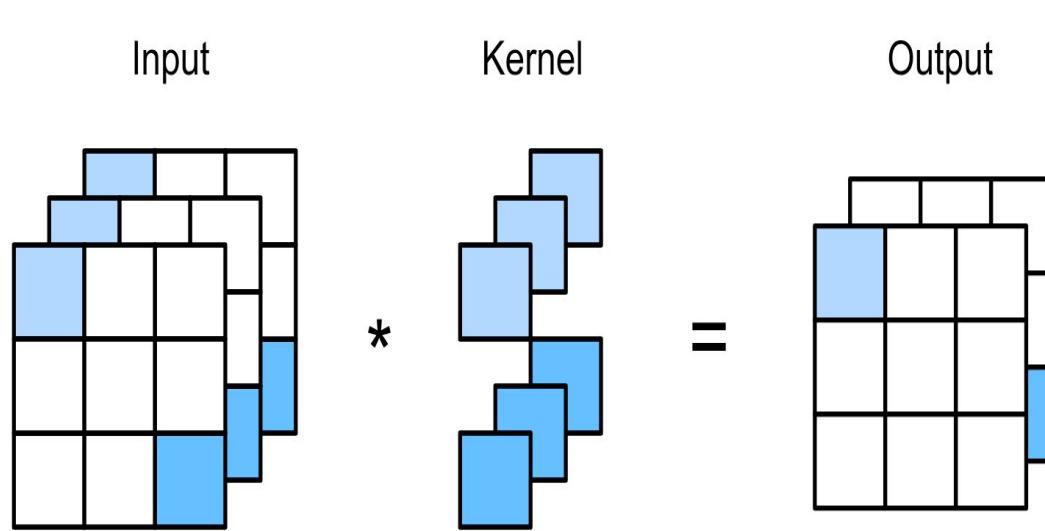
- Each output channel may recognize a particular pattern



- Input channels kernels recognize and combines patterns in inputs

1 x 1 Convolutional Layer

$k_h = k_w = 1$ is a popular choice. It doesn't recognize spatial patterns, but fuse channels.



Equal to a dense layer with $n_h \times n_w \times c_i$ input and $c_o \times c_i$ weight.

2-D Convolution Layer Summary

- Input $\mathbf{X} : c_i \times n_h \times n_w$
 - Kernel $\mathbf{W} : c_o \times c_i \times k_h \times k_w$
 - Bias $\mathbf{B} : c_o \times c_i$
 - Output $\mathbf{Y} : c_o \times m_h \times m_w$
- $$\mathbf{Y} = \mathbf{X} \star \mathbf{W} + \mathbf{B}$$
- Complexity (number of floating point operations FLOP)
 $c_i = c_o = 100$
 $k_h = h_w = 5$
 $m_h = m_w = 64$
 $O(c_i c_o k_h k_w m_h m_w)$ 1GFLOP
 - 10 layers, 1M examples: 10PF
(CPU: 0.15 TF = 18h, GPU: 12 TF = 14min)



A close-up photograph of a swimming pool's floor, showing blue square tiles in a grid pattern. Light reflects off the water surface, creating a diagonal striped effect across the tiles. The perspective is slightly angled, looking down into the pool.

Pooling Layer

Pooling

- Convolution is sensitive to position
 - Detect vertical edges

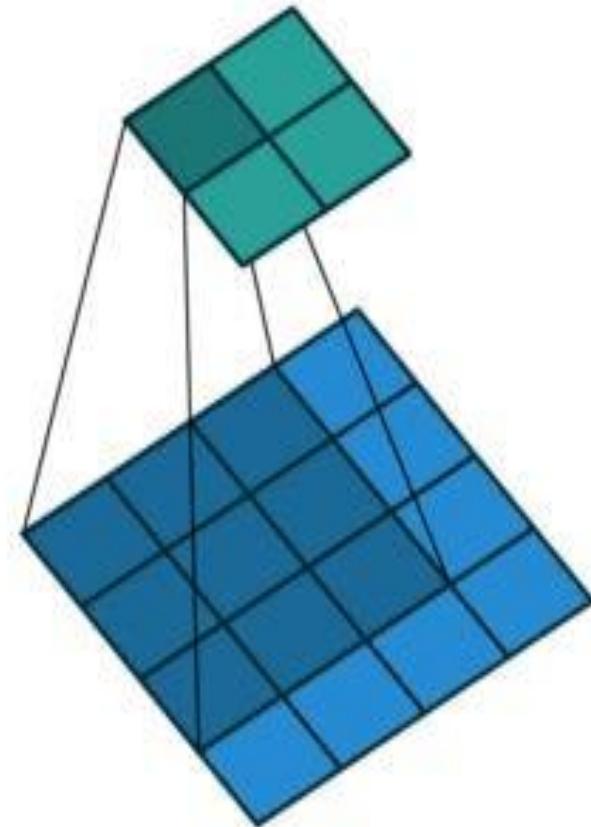
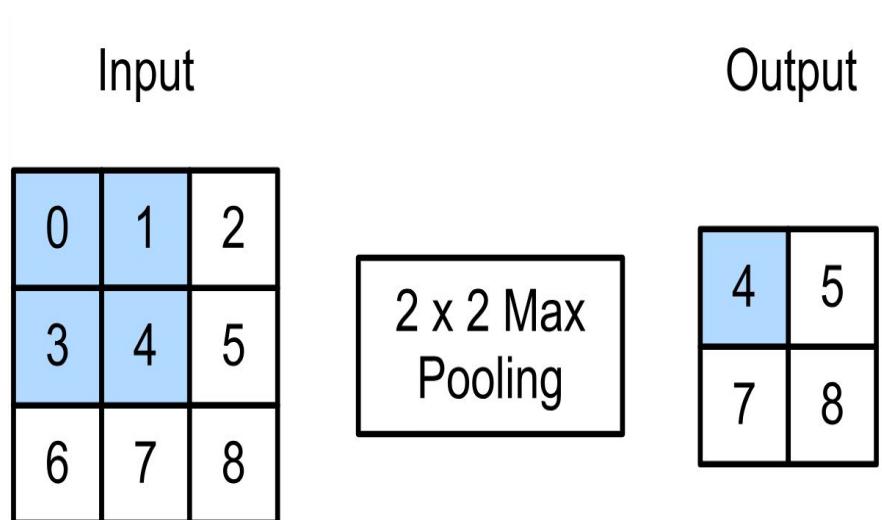
0 output with
1 pixel shift

$$\begin{array}{c} \times \\ \text{X} \end{array} \quad \begin{bmatrix} [1. & 1. & 0. & 0. & 0.] \\ [1. & 1. & 0. & 0. & 0.] \\ [1. & 1. & 0. & 0. & 0.] \\ [1. & 1. & 0. & 0. & 0.] \end{bmatrix} \quad \begin{bmatrix} Y \\ \text{Y} \end{bmatrix} \quad \begin{bmatrix} [0. & 1. & 0. & 0.] \\ [0. & 1. & 0. & 0.] \\ [0. & 1. & 0. & 0.] \\ [0. & 1. & 0. & 0.] \end{bmatrix}$$

- We need some degree of invariance to translation
 - Lighting, object positions, scales, appearance vary among images

2-D Max Pooling

- Returns the maximal value in the sliding window



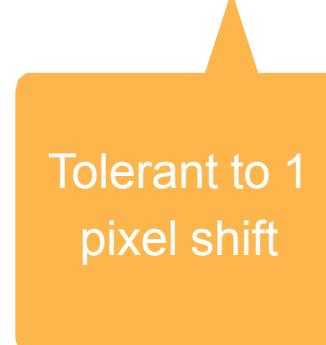
$$\max(0, 1, 3, 4) = 4$$

2-D Max Pooling

- Returns the maximal value in the sliding window

Vertical edge detection Conv output 2 x 2 max pooling

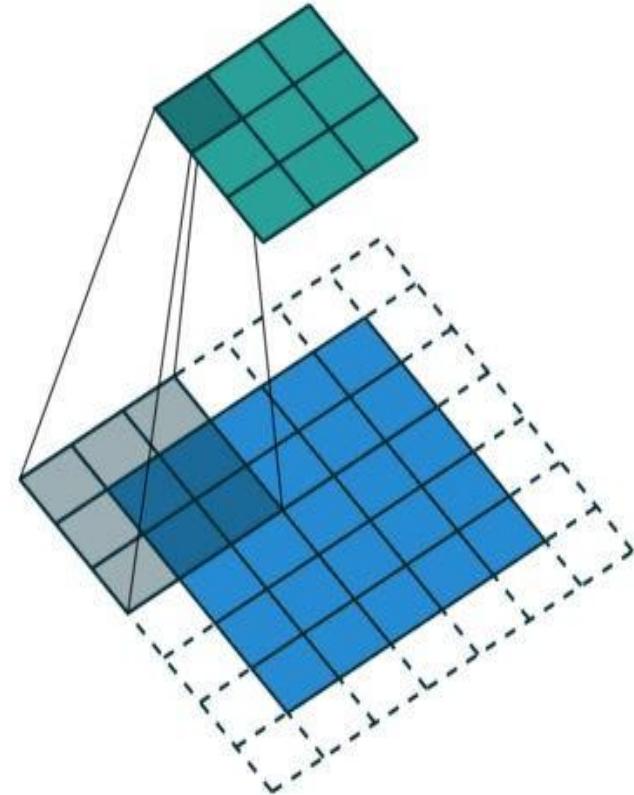
[[1. 1. 0. 0. 0.	[[0. 1. 0. 0.	[[1. 1. 1. 0.
[1. 1. 0. 0. 0.	[0. 1. 0. 0.	[1. 1. 1. 0.
[1. 1. 0. 0. 0.	[0. 1. 0. 0.	[1. 1. 1. 0.
[1. 1. 0. 0. 0.	[0. 1. 0. 0.	[1. 1. 1. 0.



Padding, Stride, and Multiple Channels

- Pooling layers have similar padding and stride as convolutional layers
- No learnable parameters
- Apply pooling for each input channel to obtain the corresponding output channel

#output channels = #input channels



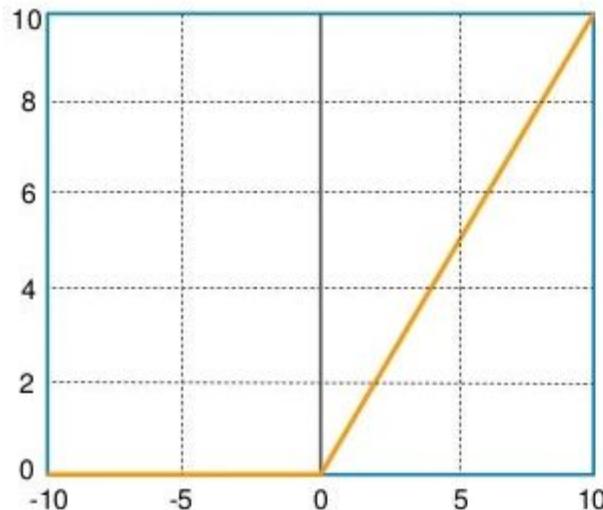
Average Pooling

- Max pooling: the strongest pattern signal in a window
- Average pooling: replace max with mean in max pooling
 - The average signal strength in a window

Max pooling Average pooling



ReLU



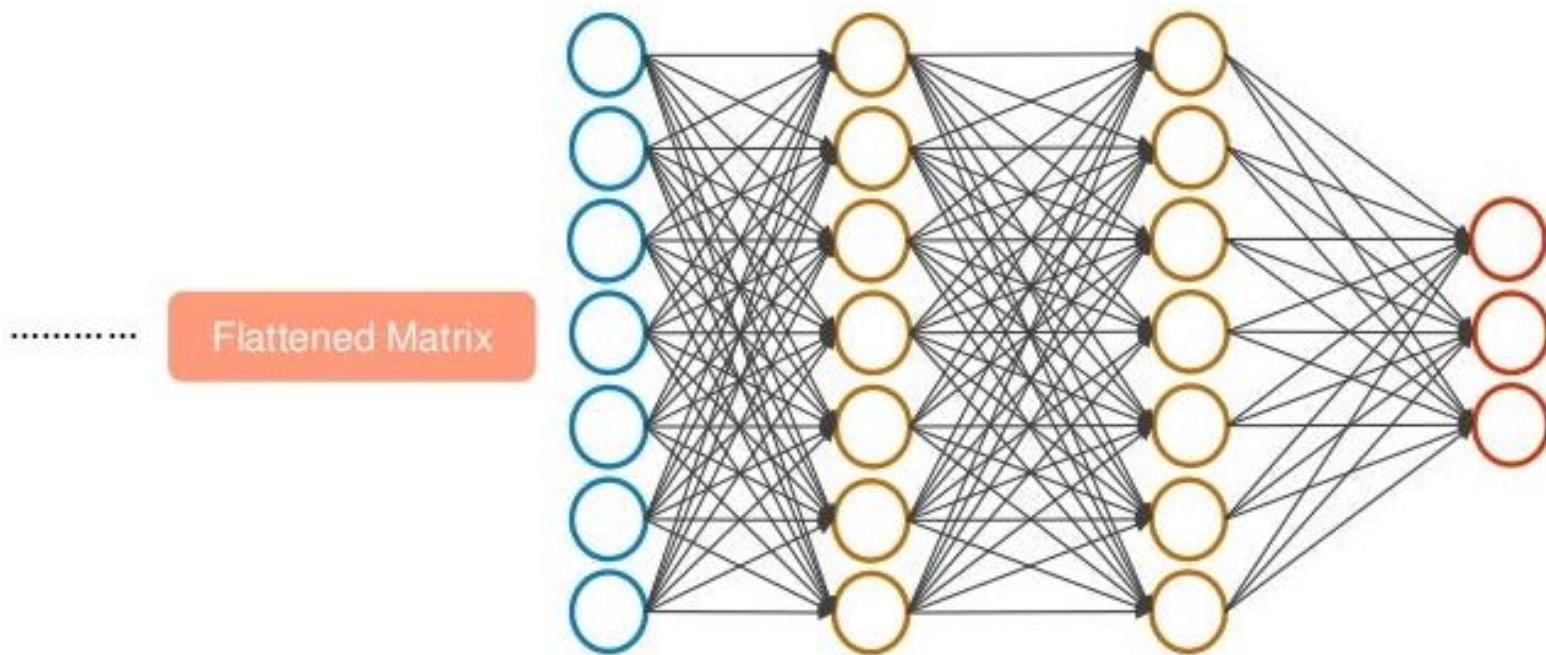
$$R(z) = \max(0, z)$$

Sets all negative pixels to 0
Performs element wise operation

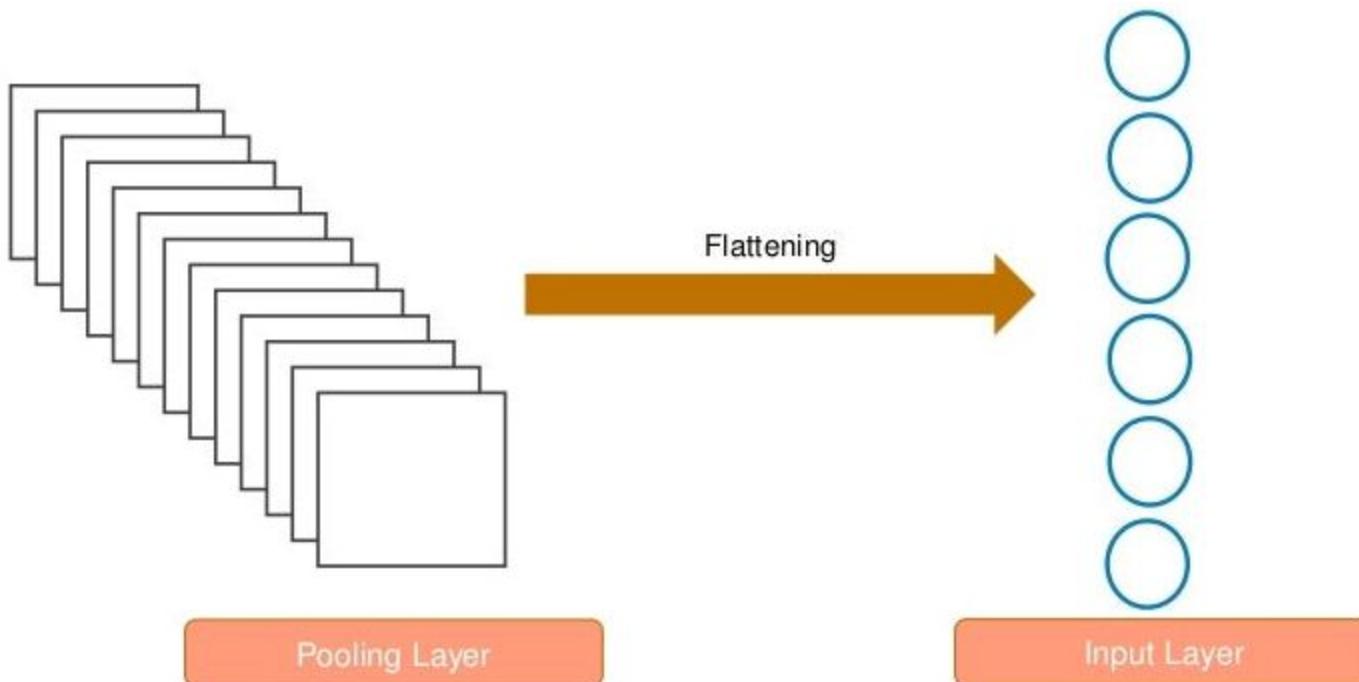
ReLU

Introduces non-linearity to the network
The output is a rectified feature map

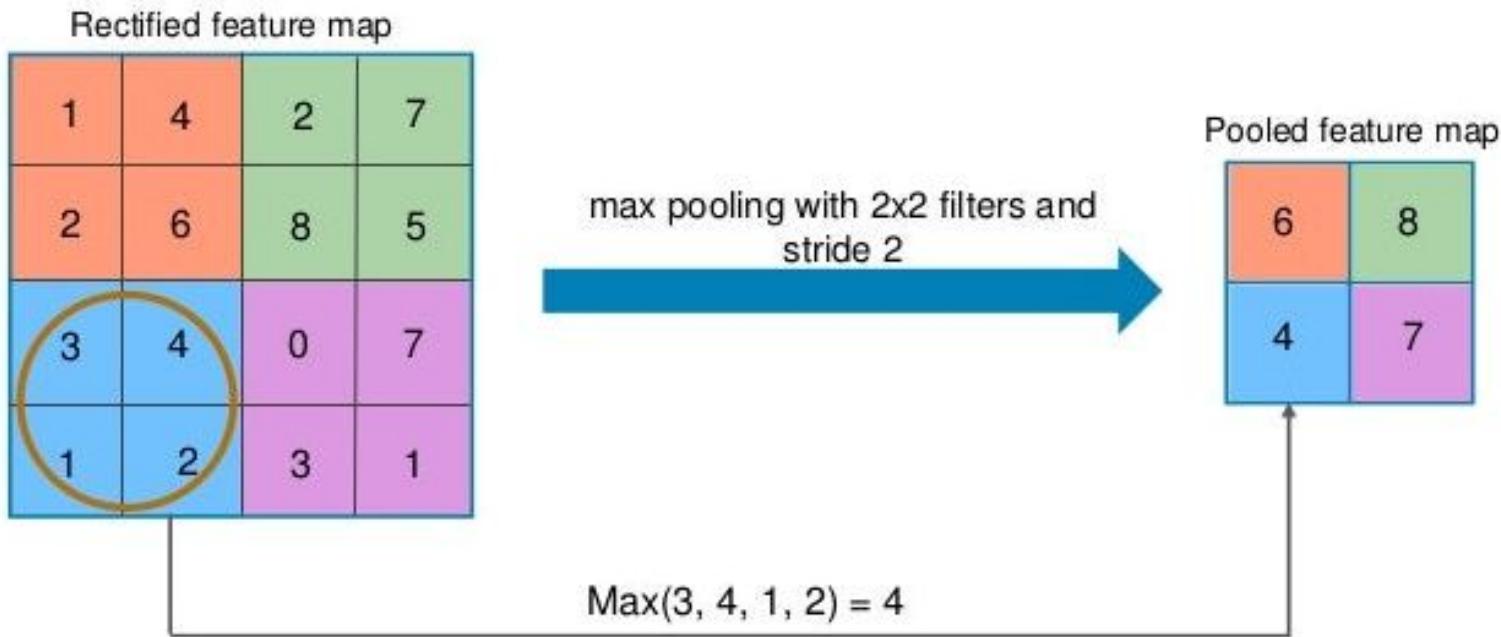
Fully connected layer



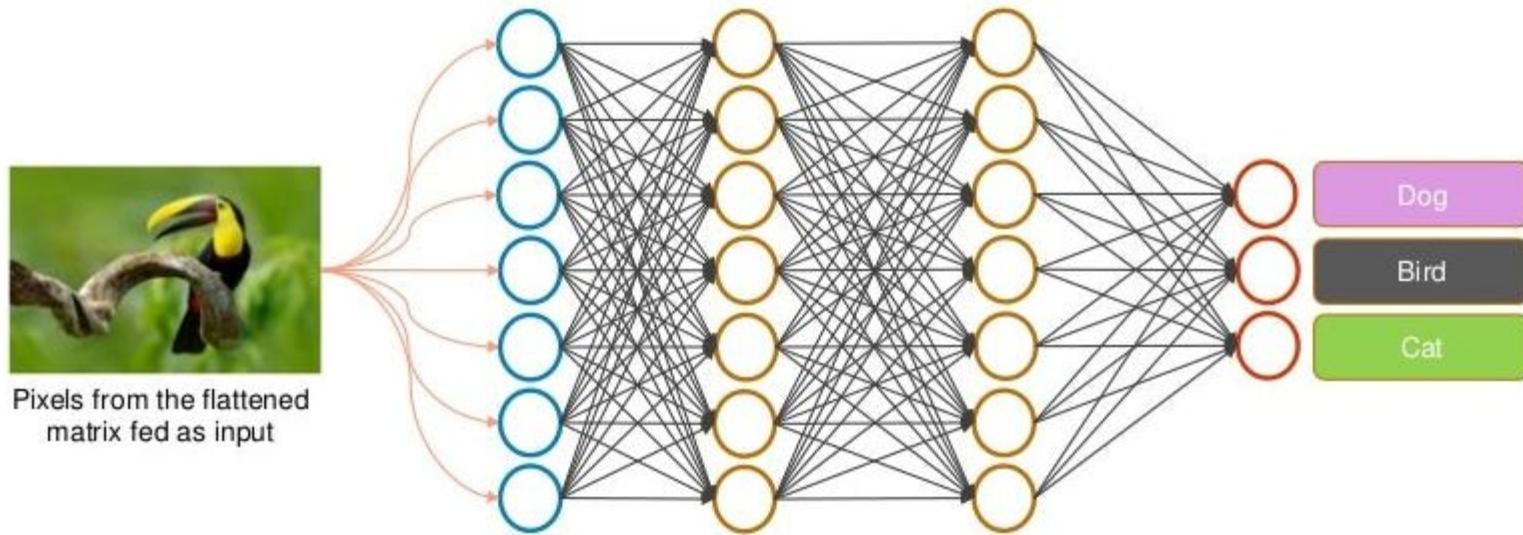
Flattening



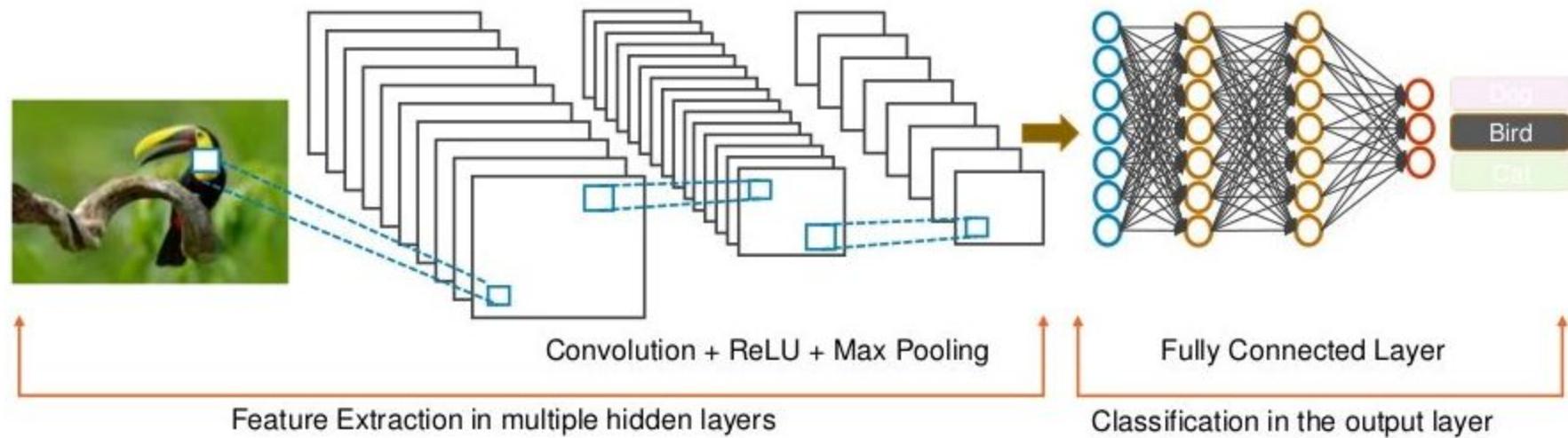
Max pooling

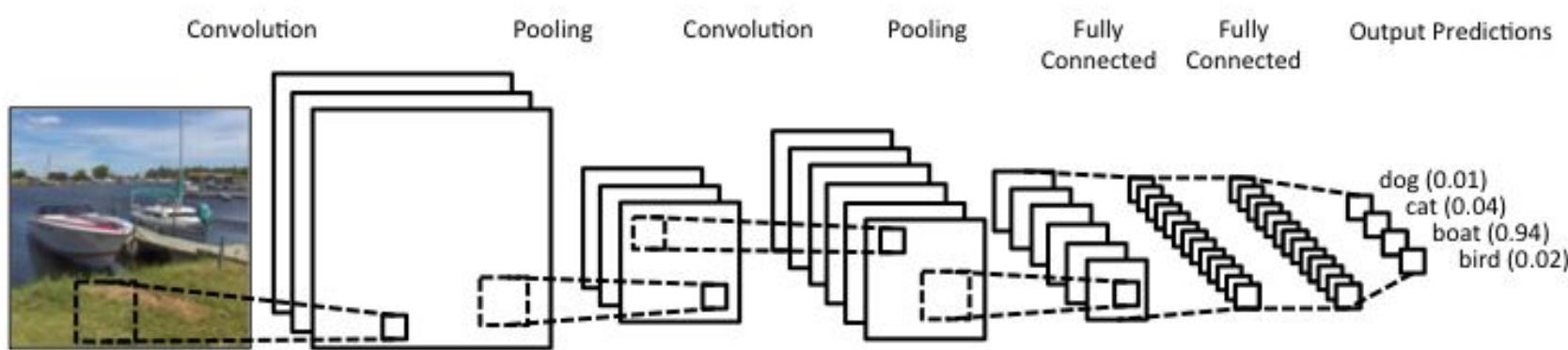


Fully Connected Layer

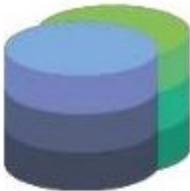


Final CNN Model

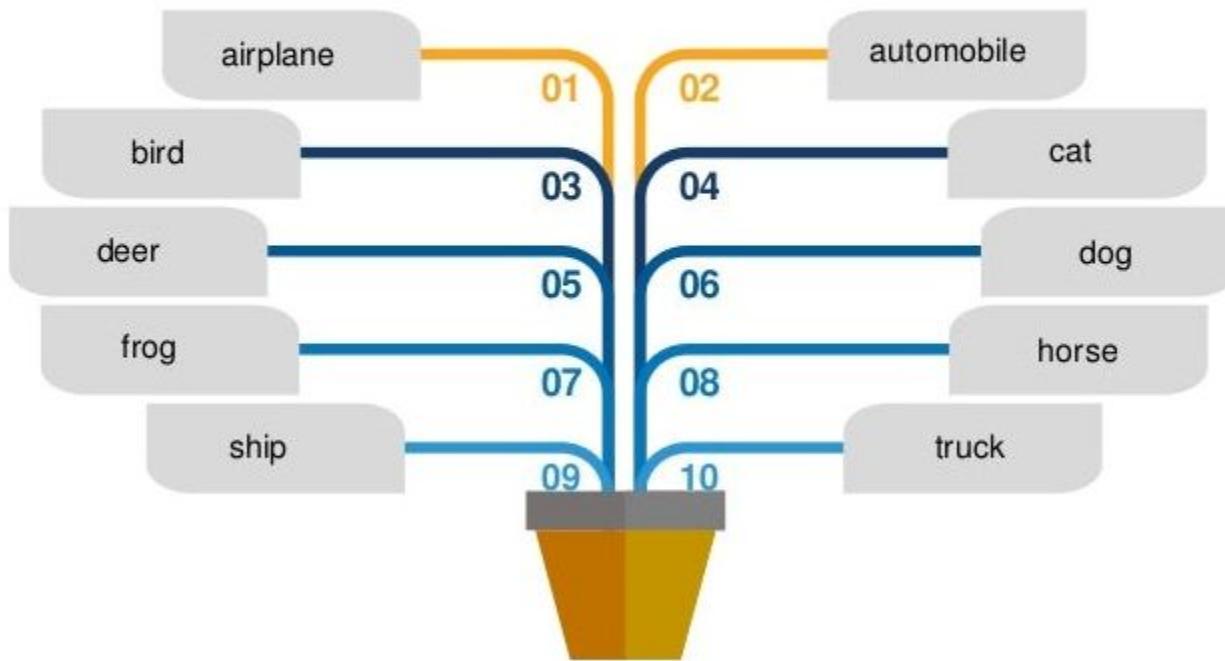




Dataset



We will be using CIFAR-10 data set (from Canadian Institute For Advanced Research) for classifying images across 10 categories



airplane



automobile



bird



cat



deer



dog



frog



horse



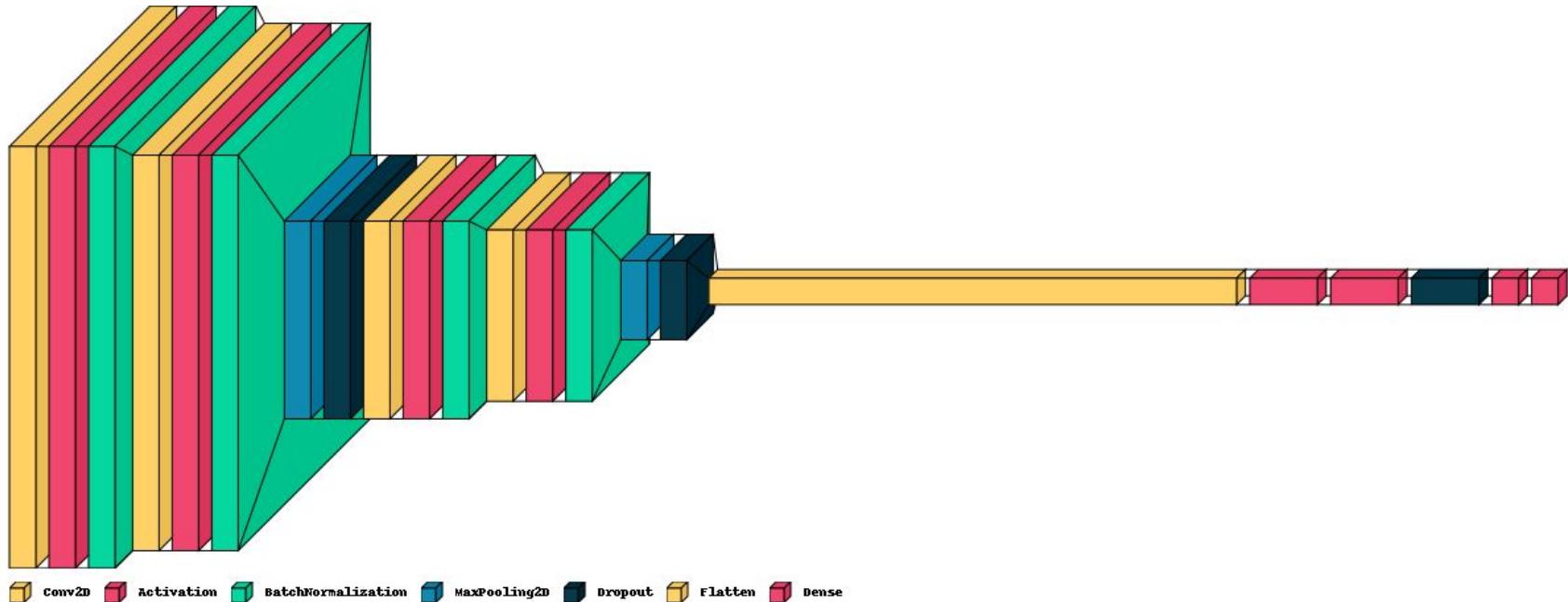
ship



truck



Model



Model parameters

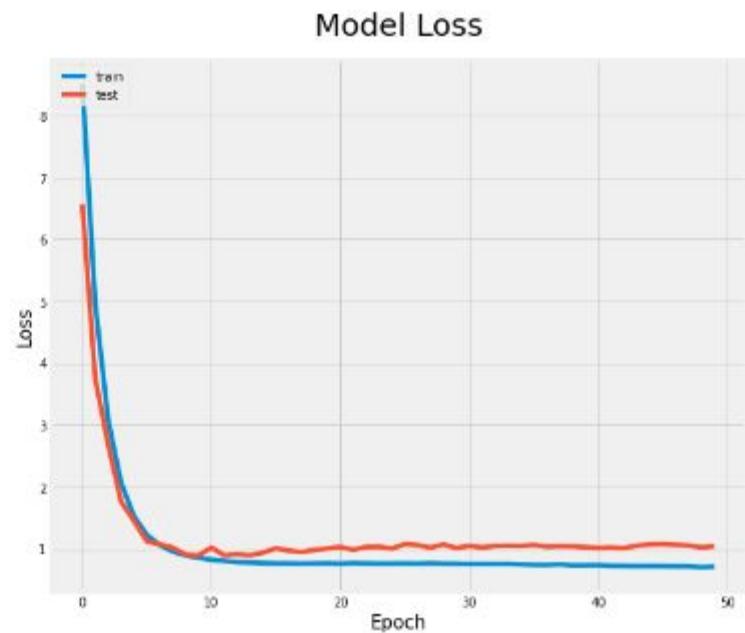
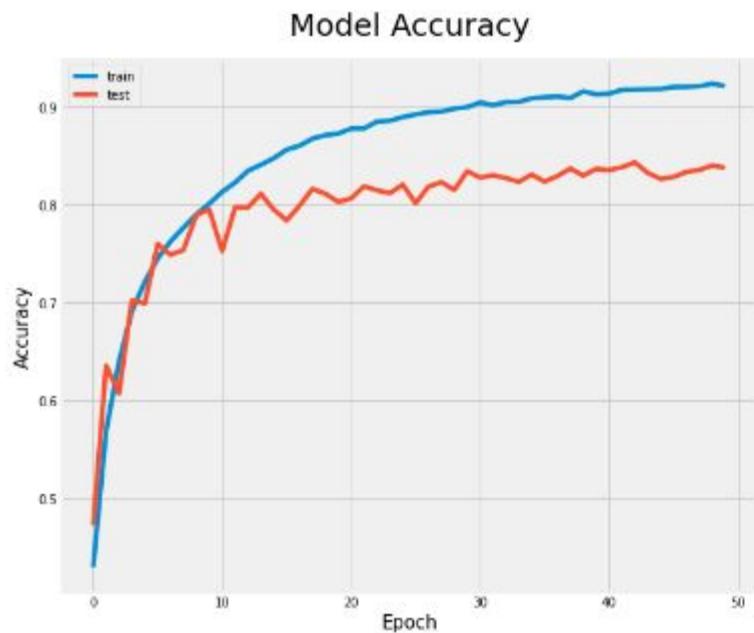
```
Model: "sequential"
```

Layer (type)	Output Shape	Param #
<hr/>		
conv2d (Conv2D)	(None, 32, 32, 64)	1792
<hr/>		
activation (Activation)	(None, 32, 32, 64)	0
<hr/>		
batch_normalization (BatchNo	(None, 32, 32, 64)	256
<hr/>		
conv2d_1 (Conv2D)	(None, 30, 30, 64)	36928
<hr/>		
<hr/>		
activation_4 (Activation)	(None, 512)	0
<hr/>		
dropout_2 (Dropout)	(None, 512)	0
<hr/>		
dense_1 (Dense)	(None, 10)	5130
<hr/>		
activation_5 (Activation)	(None, 10)	0
<hr/>		
Total params: 2,626,634		
Trainable params: 2,625,866		
Non-trainable params: 768		
<hr/>		
None		

Training

```
2022-06-20 09:25:04.301225: I tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:185] None of the MLIR Optimization Passes are enabled (registered 2)
Epoch 1/50
12500/12500 [=====] - 513s 41ms/step - loss: 4.2169 - accuracy: 0.3639 - val_loss: 4.6308
- val_accuracy: 0.1687
Epoch 2/50
12500/12500 [=====] - 499s 40ms/step - loss: 2.0688 - accuracy: 0.5234 - val_loss: 2.8501
- val_accuracy: 0.2733
Epoch 3/50
9207/12500 [=====>.....] - ETA: 2:02 - loss: 1.8332 - accuracy: 0.5930
```

Accuracy and Loss



Prediction

Predictions of CIFAR-10 Data



Batch Normalization

The regularization techniques help to improve a model and allows it to converge faster.

We have several regularization tools at our end, some of them

- early stopping,
- dropout,
- weight initialization techniques,
- batch normalization.

The regularization helps in preventing the over-fitting of the model and the learning process becomes more efficient.

Regularization

Sometimes situation occurs when the model is performing very well on the training data but is unable to predict the test data accurately. The reason is your model is **over fitting**.

The solution to such a problem is **regularization**.

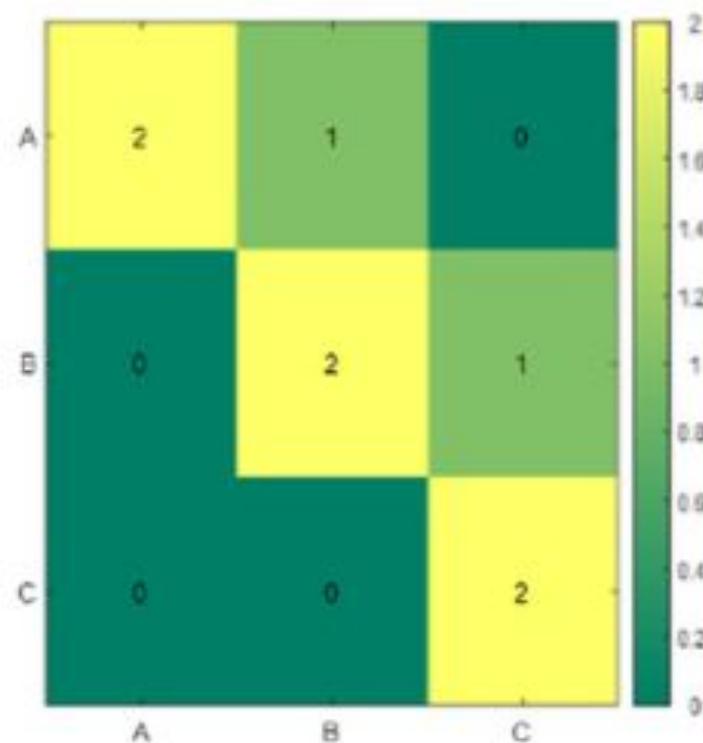
Pooling Notebook

```
cm =
```

```
 2  1  0  
 0  2  1  
 0  0  2
```

```
grp =
```

```
A  
B  
C
```



Confusion Matrix

		Actual Value	
		Yes (1)	No (0)
Predicted Value	Yes (1)	TP	FP
	No (0)	FN	TN

TP= True Positive

FP= False Positive

FN= False Negative

TN= True Negative

Type I Error

		Actual Value	
		Yes (1)	No (0)
Predicted Value	Yes (1)	500 (TP)	100 (FP)
	No (0)	200 (FN)	200 (TN)

Type I Error = False Positives

Type II Error

		Actual Value	
		Yes (1)	No (0)
Predicted Value	Yes (1)	500 (TP)	100 (FP)
	No (0)	200 (FN)	200 (TN)

Type II Error = False Negatives

Accuracy

		Actual Value	
		Yes (1)	No (0)
Predicted Value	Yes (1)	500 (TP)	100 (FP)
	No (0)	200 (FN)	200 (TN)

$$\begin{aligned} \text{Accuracy} &= (TP+TN)/(TP+FP+FN+TN) \\ &= (500+200)/(500+100+200+200) \\ &= 70\% \end{aligned}$$

Accuracy = (TP+TN)/Total customers

Receiver Operating Characteristics curve

- The values described are used to calculate different measurements of the quality of the test.
- The first one is **sensitivity, SE**, which is the probability of having a **positive test** among the patients who have a **positive diagnosis**.

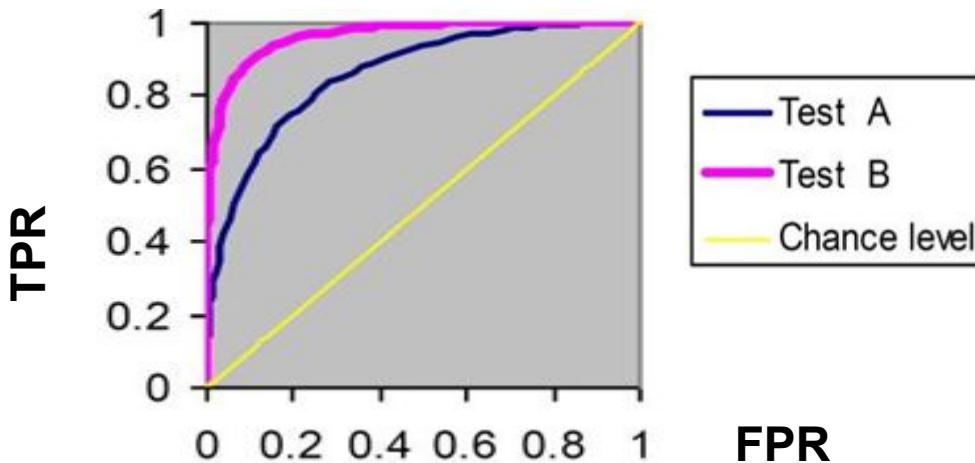
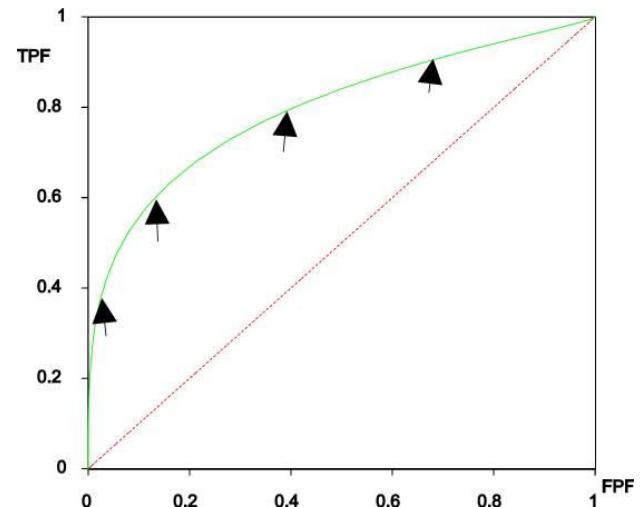
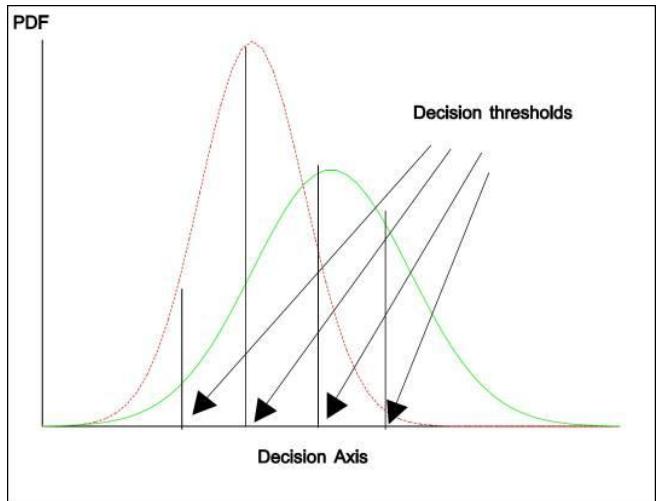
$$SE = TP / (TP + FN) = TP / P.$$

- **Specificity, SP**, is the probability of having a **negative test** among the patients who have a **negative diagnosis**.

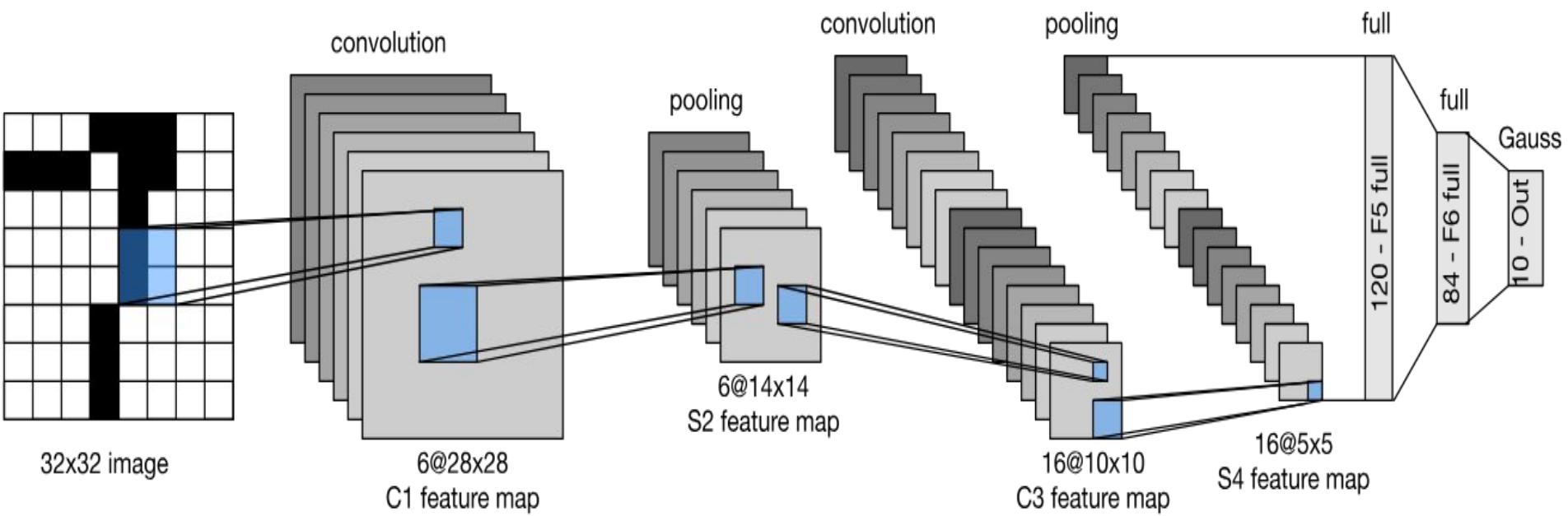
$$SP = TN / (FP + TN) = TN / P'.$$

- **Efficiency** is defined as , $EFF = TP + TN$.

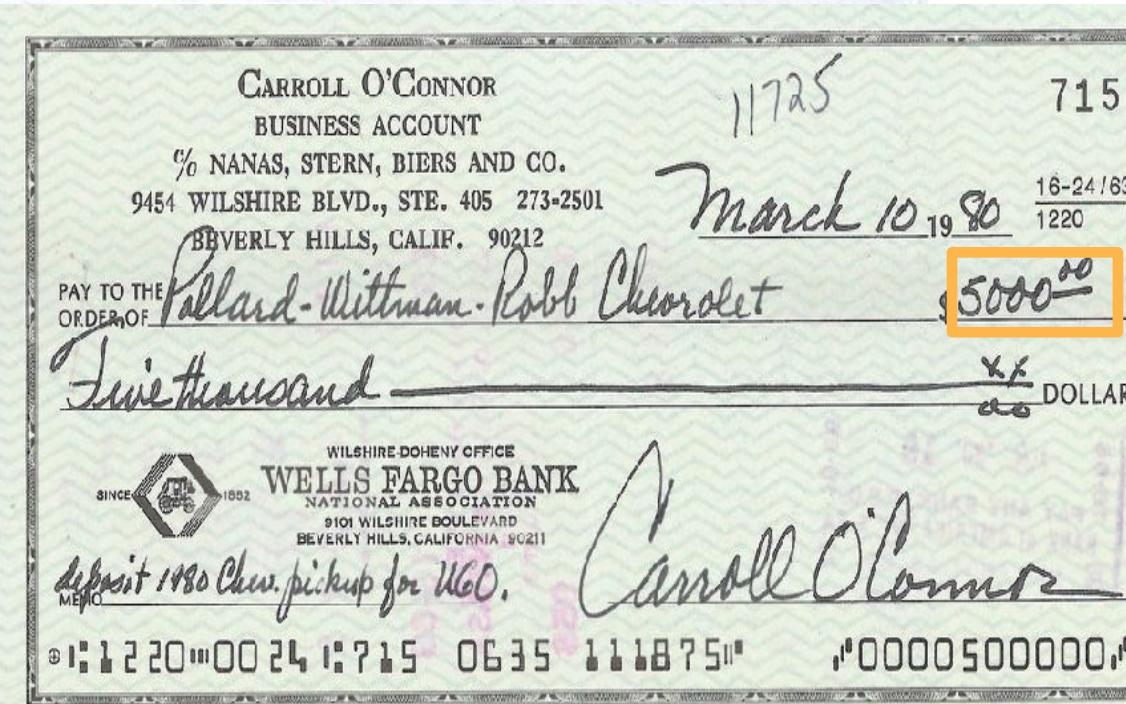
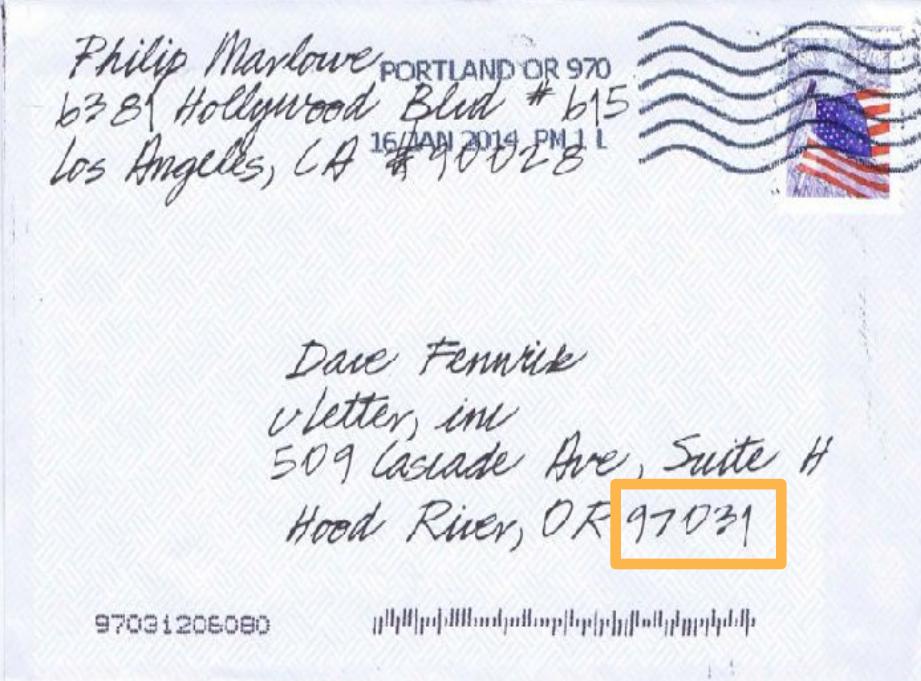
Receiver Operating Characteristics curve



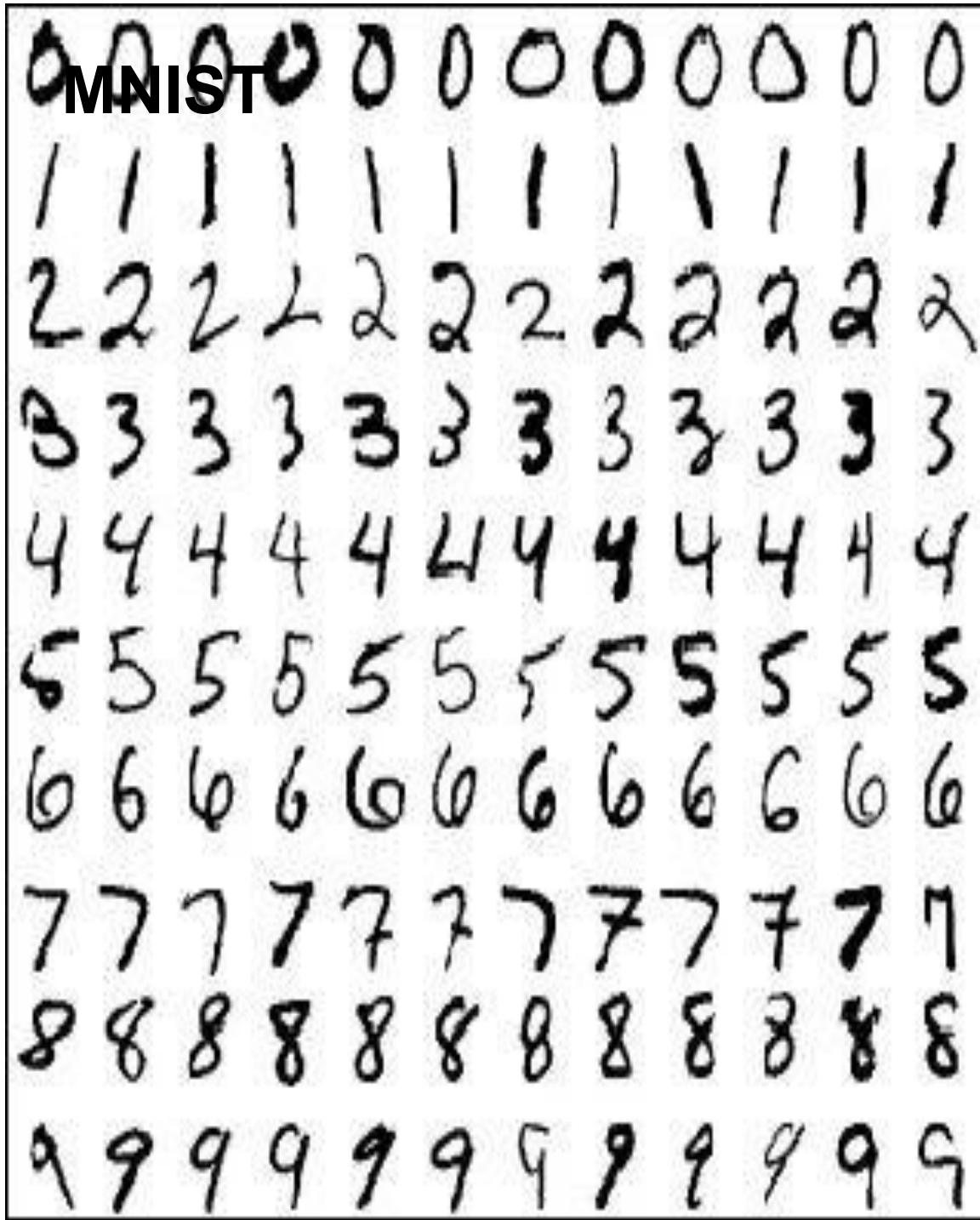
LeNet



Handwritten Digit Recognition



- Centered and scaled
- 50,000 training data
- 10,000 test data
- 28 x 28 images
- 10 classes

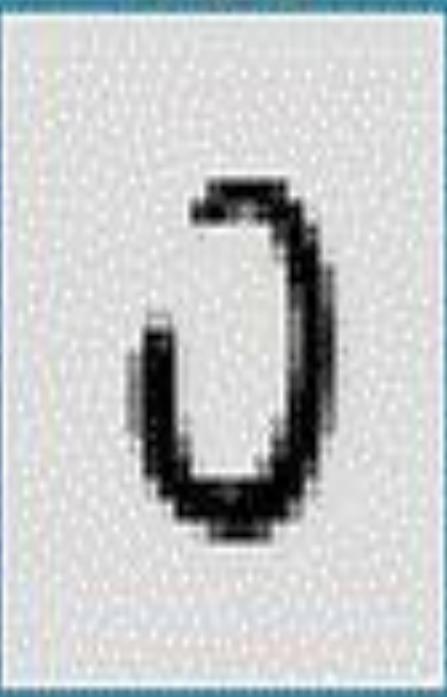




AT&T *LeNet 5* RESEARCH

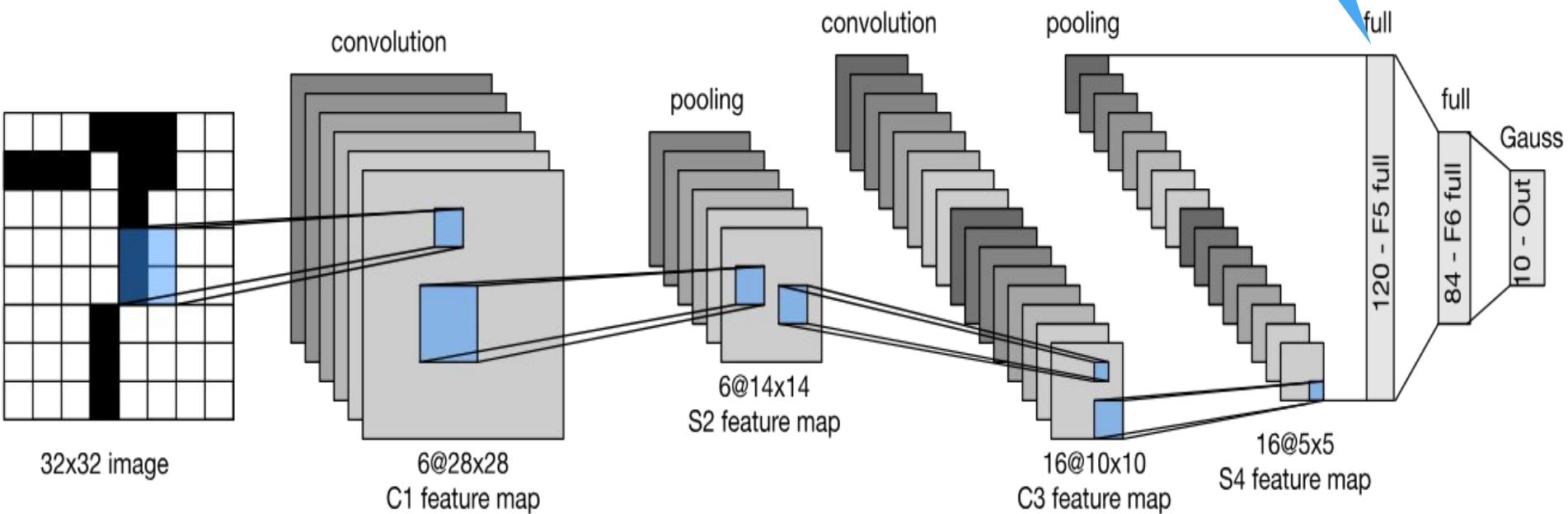
answer: 0

0
103



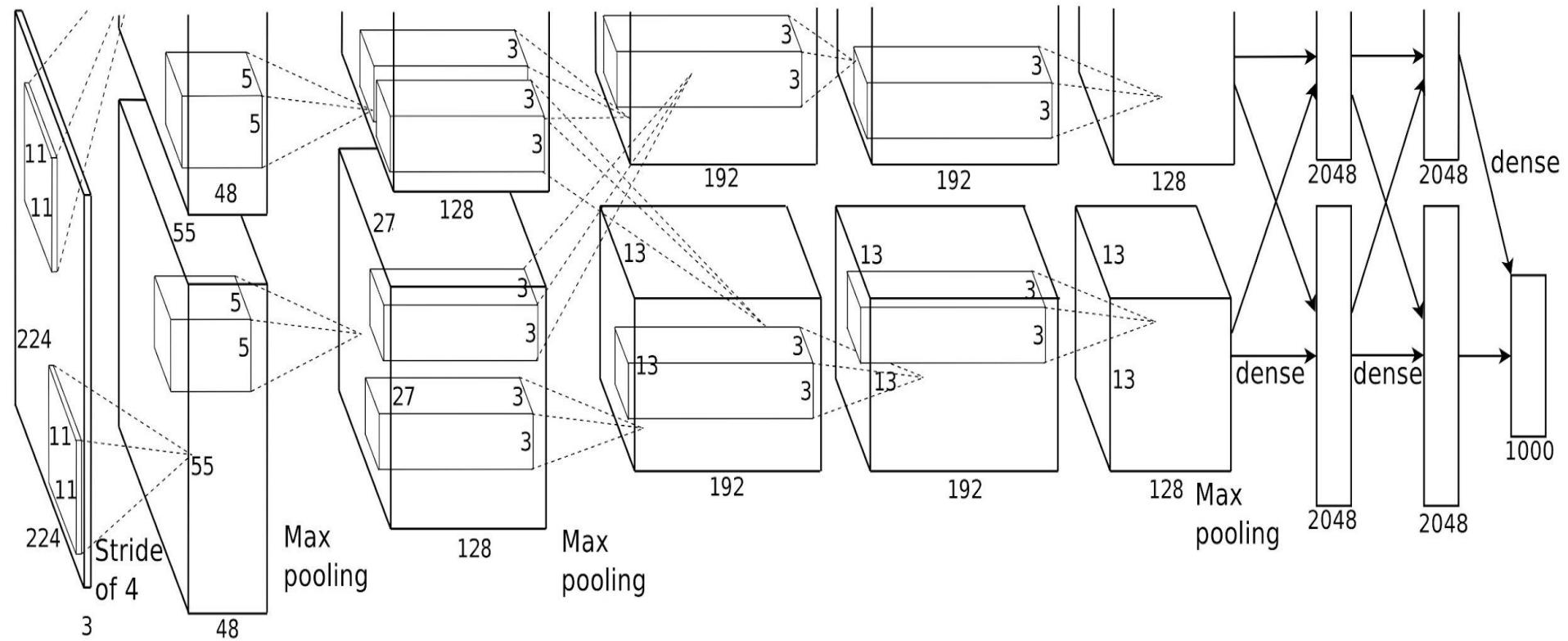
Y. LeCun, L.
Bottou, Y. Bengio,
P. Haffner, 1998
Gradient-based
learning applied to
document
recognition

Expensive if we
have many
outputs



LeNet Notebook

AlexNet

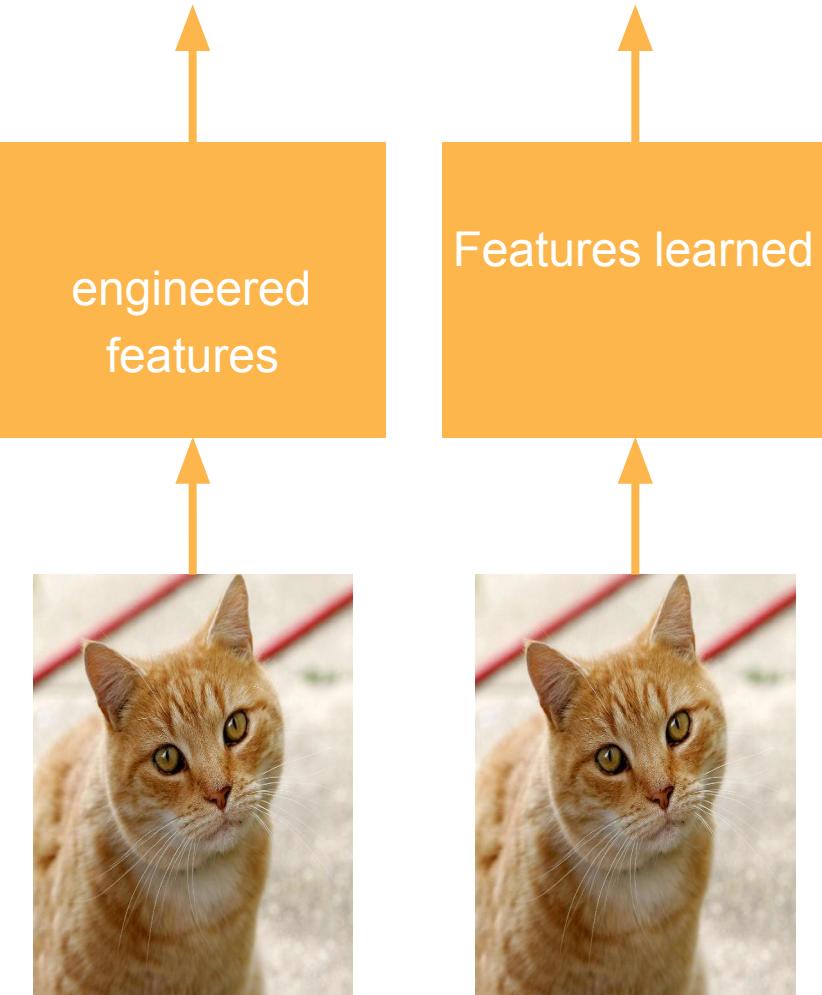


AlexNet

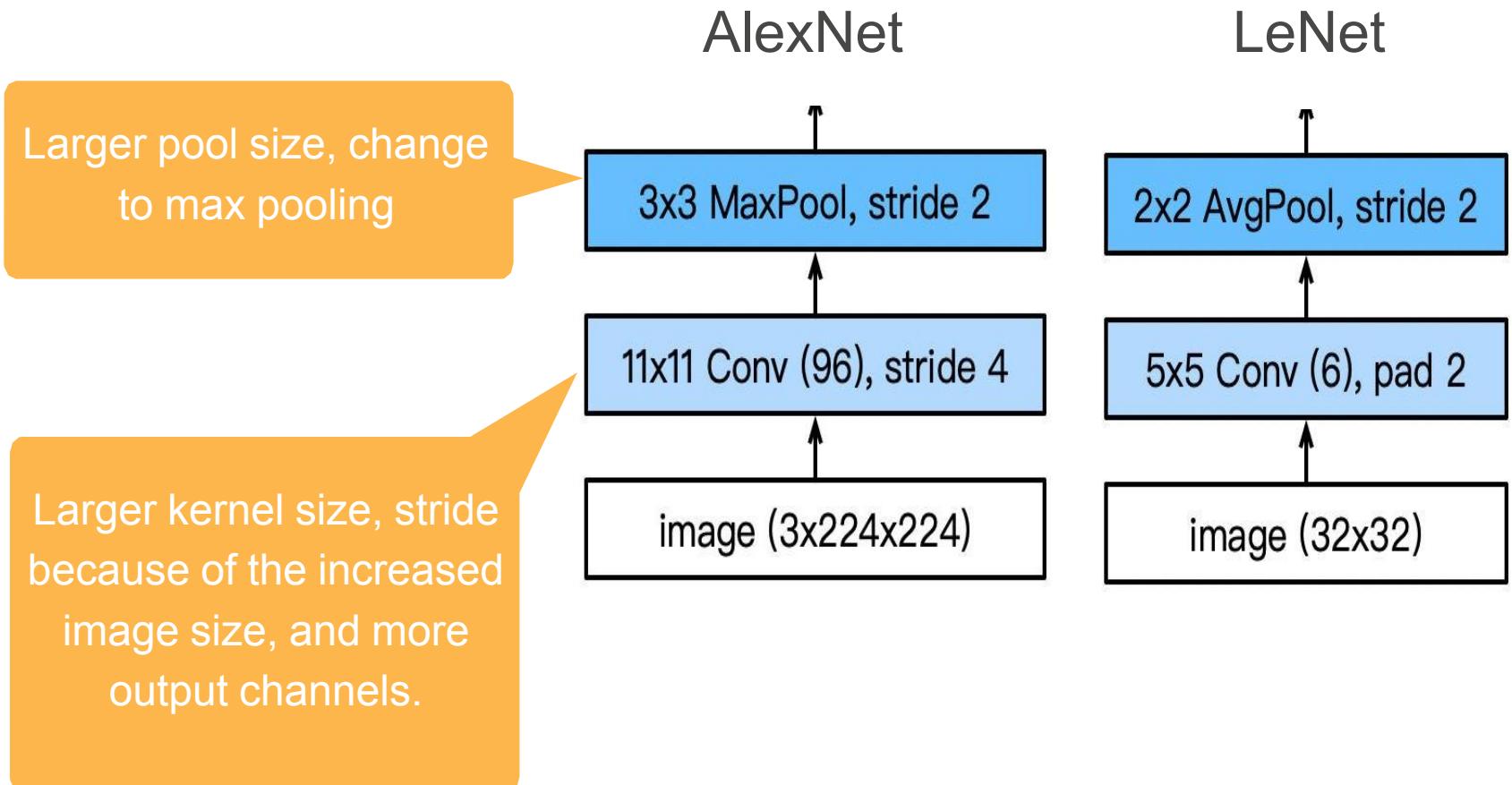
SVM

Softmax
regression

- AlexNet won ImageNet competition in 2012
- Key modifications
 - Dropout (regularization)
 - ReLu (training)
 - MaxPooling
- Paradigm shift for computer vision

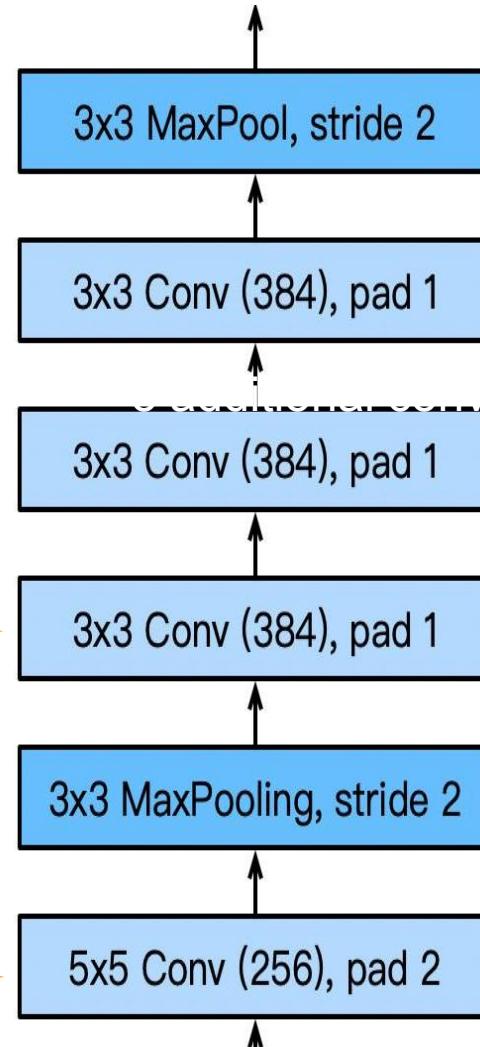


AlexNet Architecture

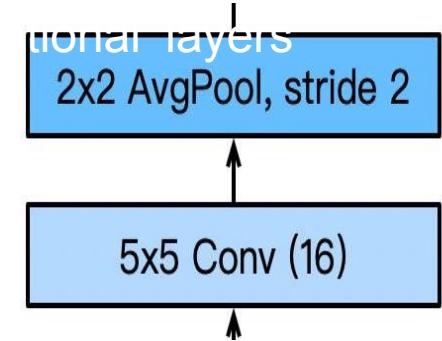


AlexNet Architecture

AlexNet

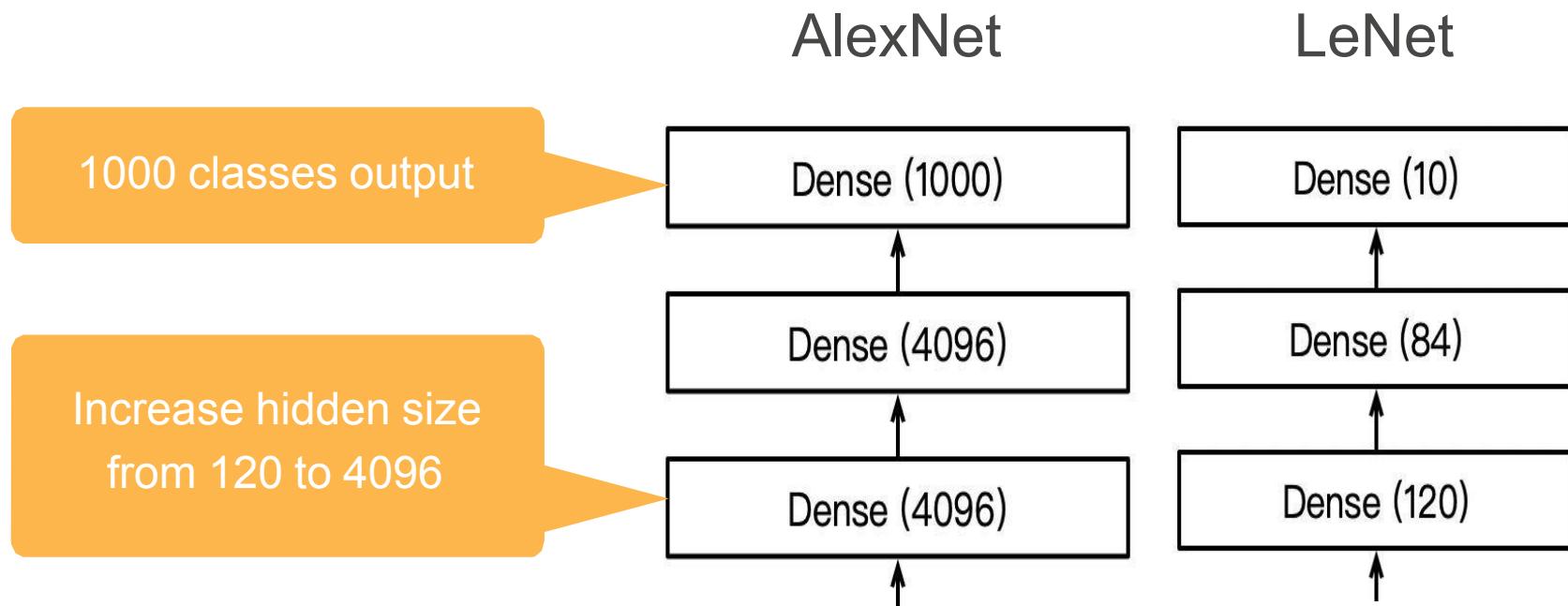


LeNet



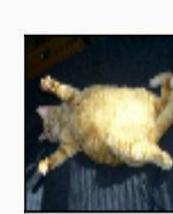
More output channels.

AlexNet Architecture



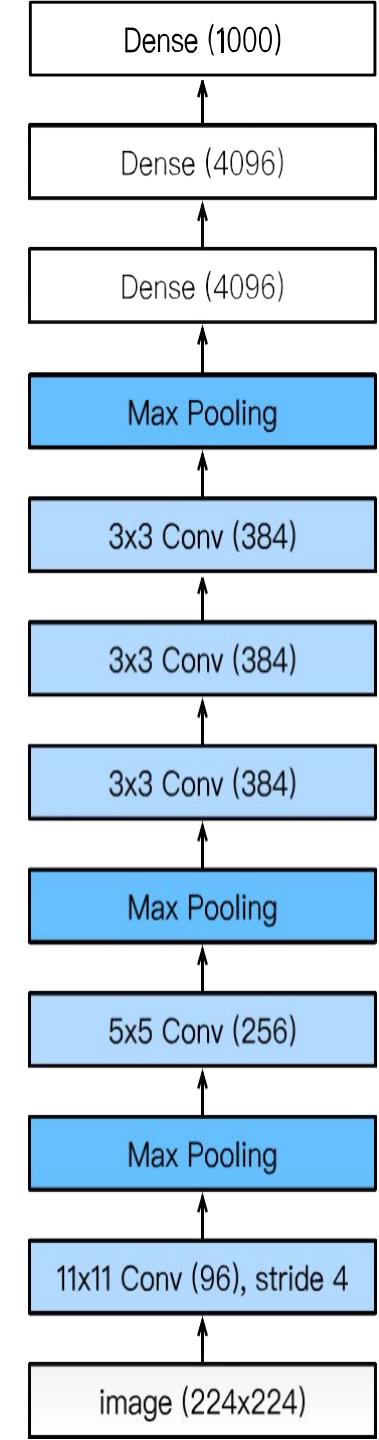
More Tricks

- Change activation function from sigmoid to ReLu
(no more vanishing gradient)
- Add a dropout layer after two hidden dense layers
(better robustness / regularization)
- Data augmentation



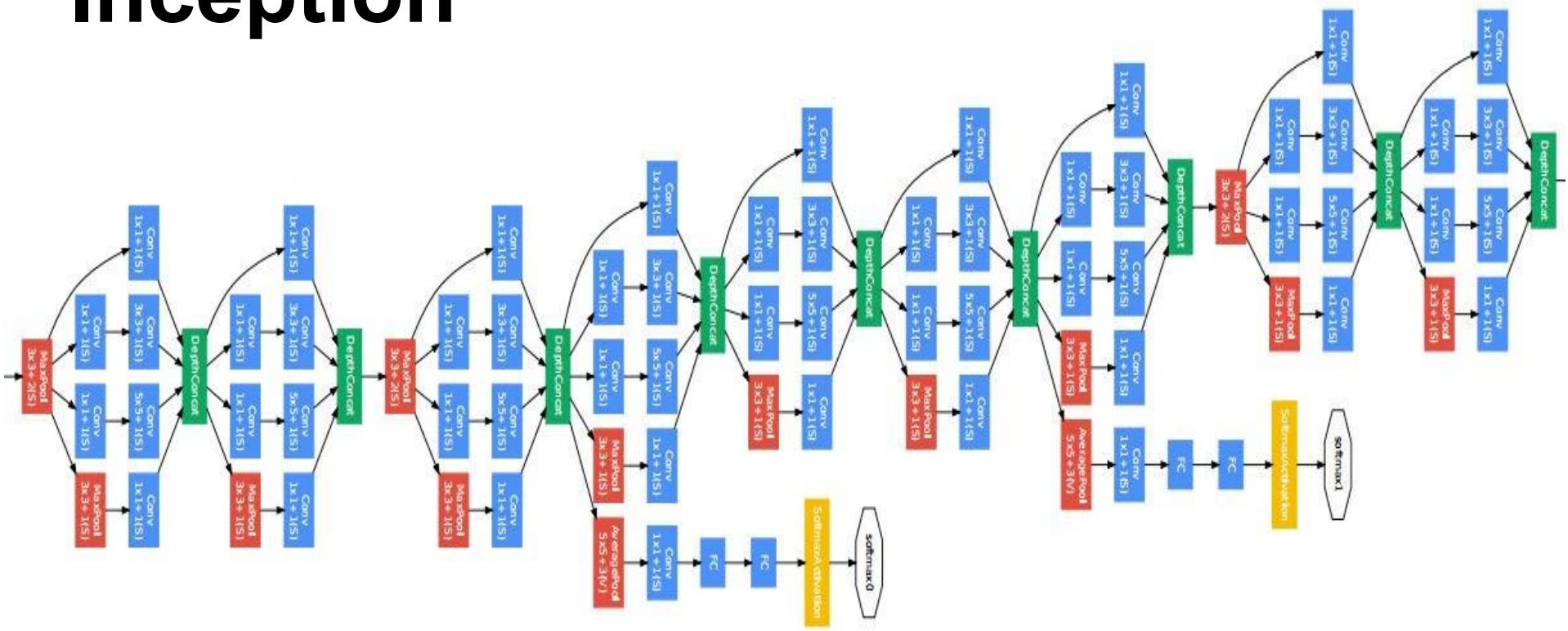
Complexity

	#parameters		FLOP	
	AlexNet	LeNet	AlexNet	LeNet
Conv1	35K	150	101M	1.2M
Conv2	614K	2.4K	415M	2.4M
Conv3-5	3M		445M	
Dense1	26M	0.48M	26M	0.48M
Dense2	16M	0.1M	16M	0.1M
Total	46M	0.6M	1G	4M
Increase	11x	1x	250x	1x



AlexNet Notebook

Inception



Picking the best convolution ...

1x1 3

x5

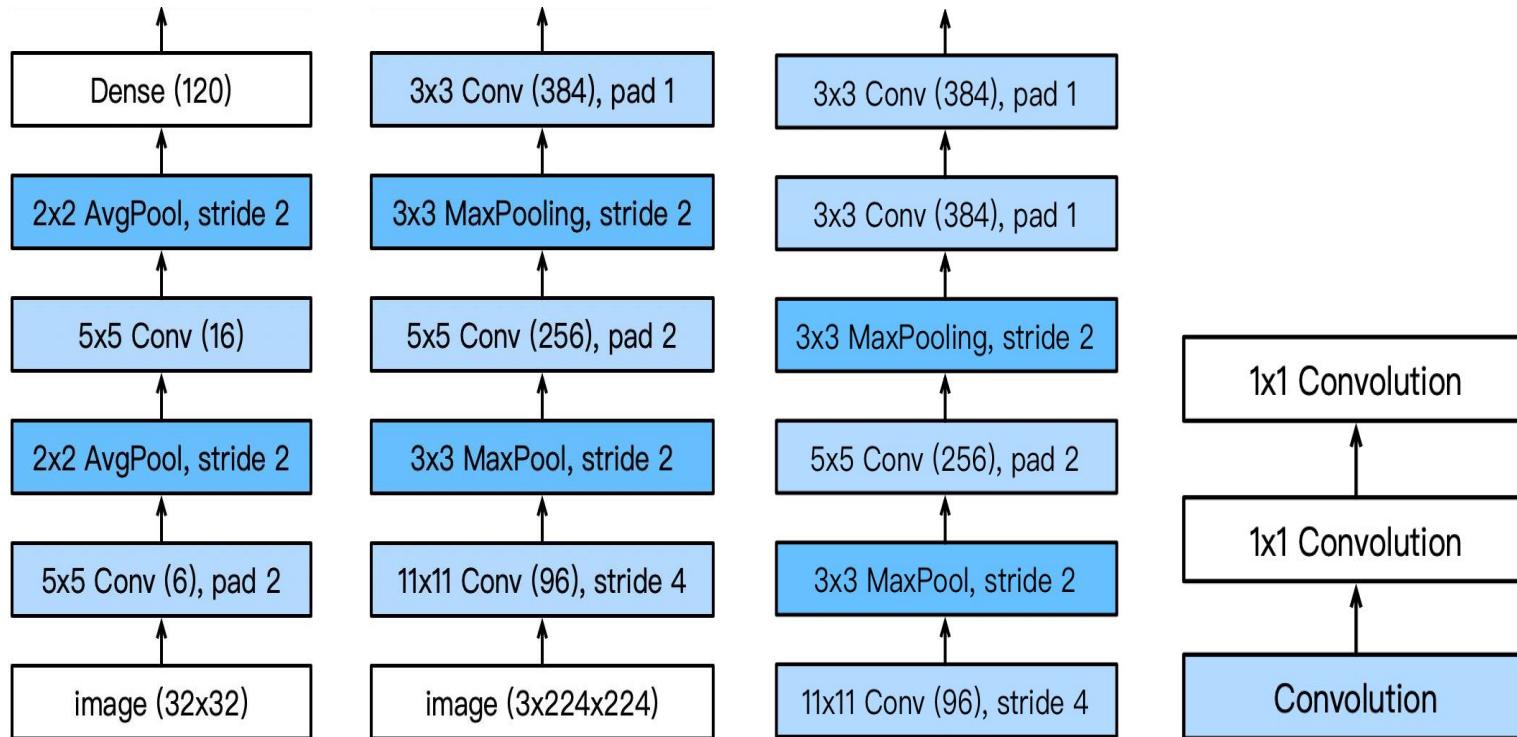
Max pooling

Multiple 1x1

LeNet AlexNet

VGG

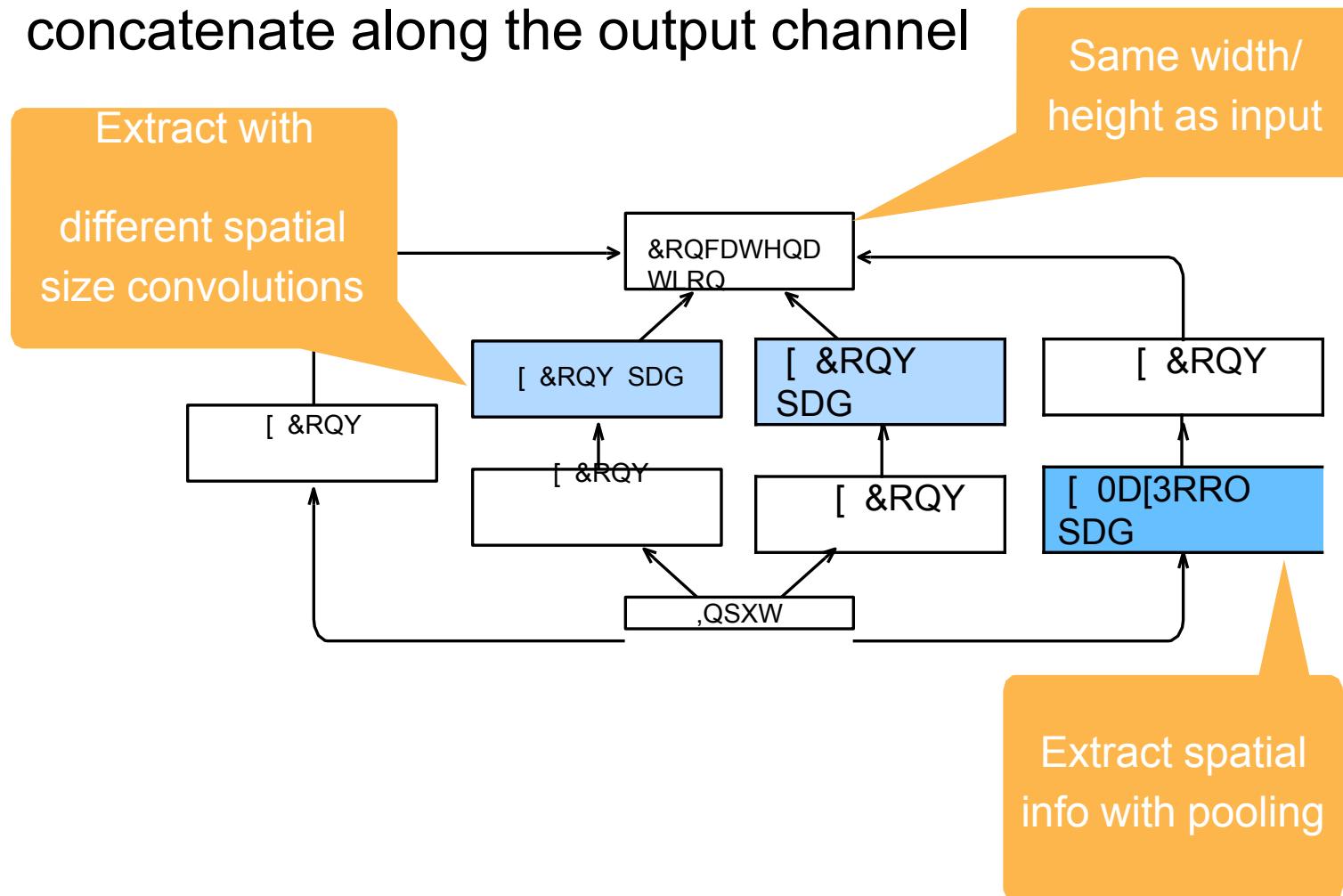
NiN



Why choose? Just pick them all.

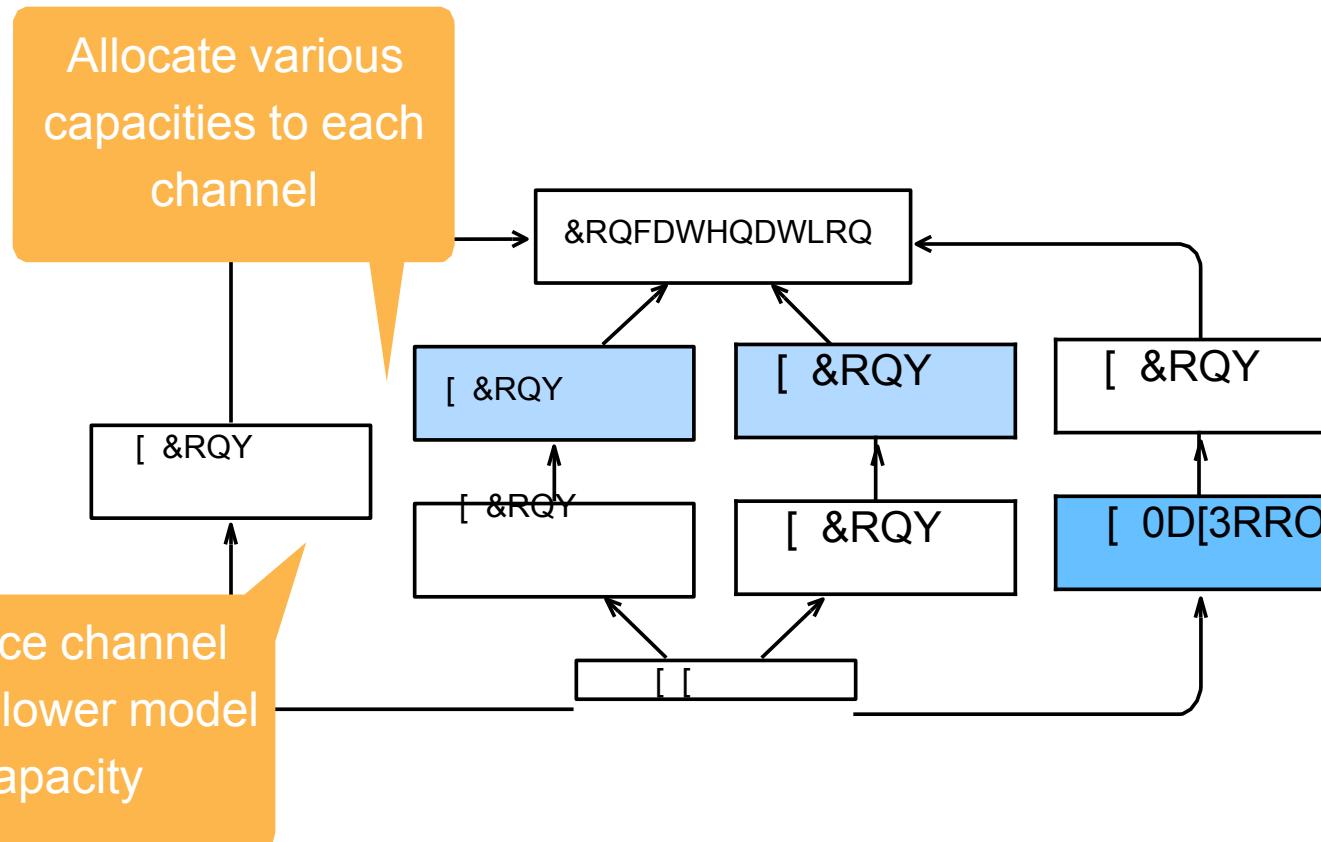
Inception Blocks

4 paths extract information from different aspects, then concatenate along the output channel



Inception Blocks

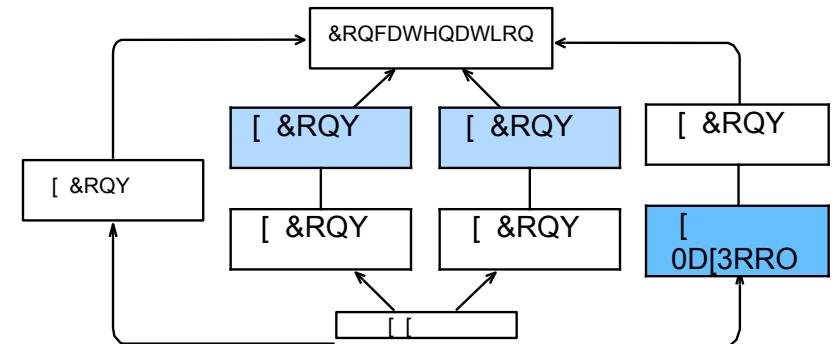
The first inception block with channel sizes specified



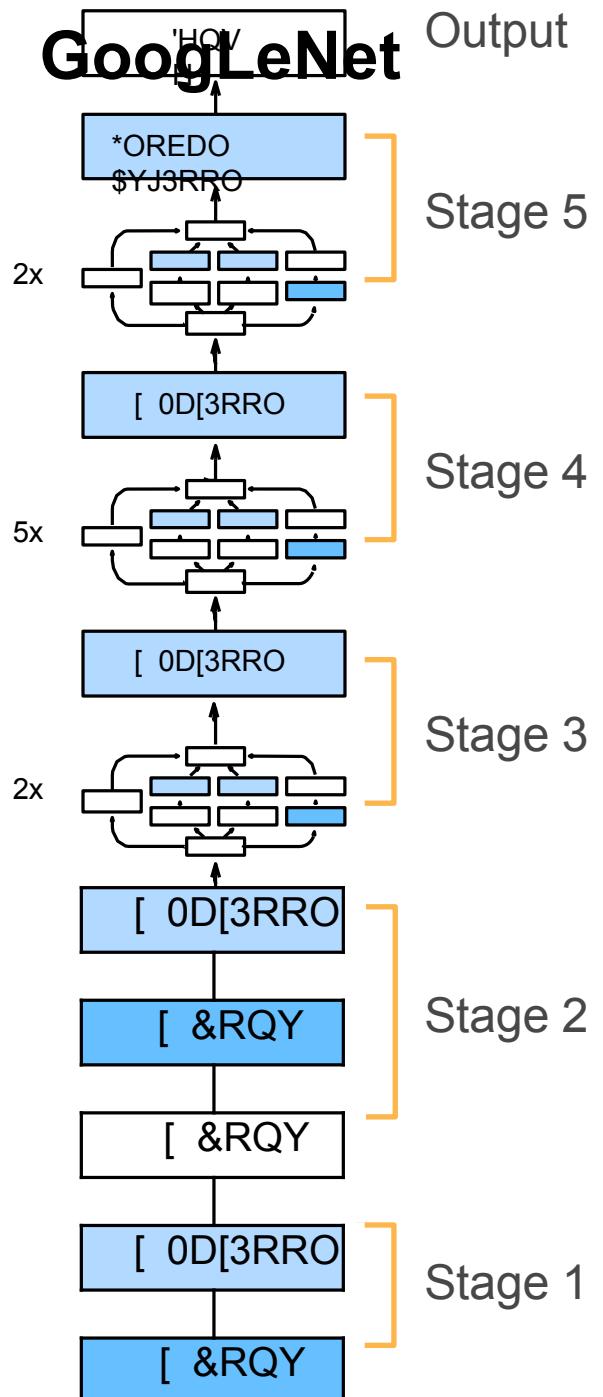
Inception Blocks

- Inception blocks have fewer parameters and less computation complexity than a single 3x3 or 5x5 convolutional layer
- Mix of different functions (powerful function class)
- Memory and compute efficiency (good generalization)

	#parameters	FLOPS
Inception	0.16 M	128 M
3x3 Conv	0.44 M	346 M
5x5 Conv	1.22 M	963 M

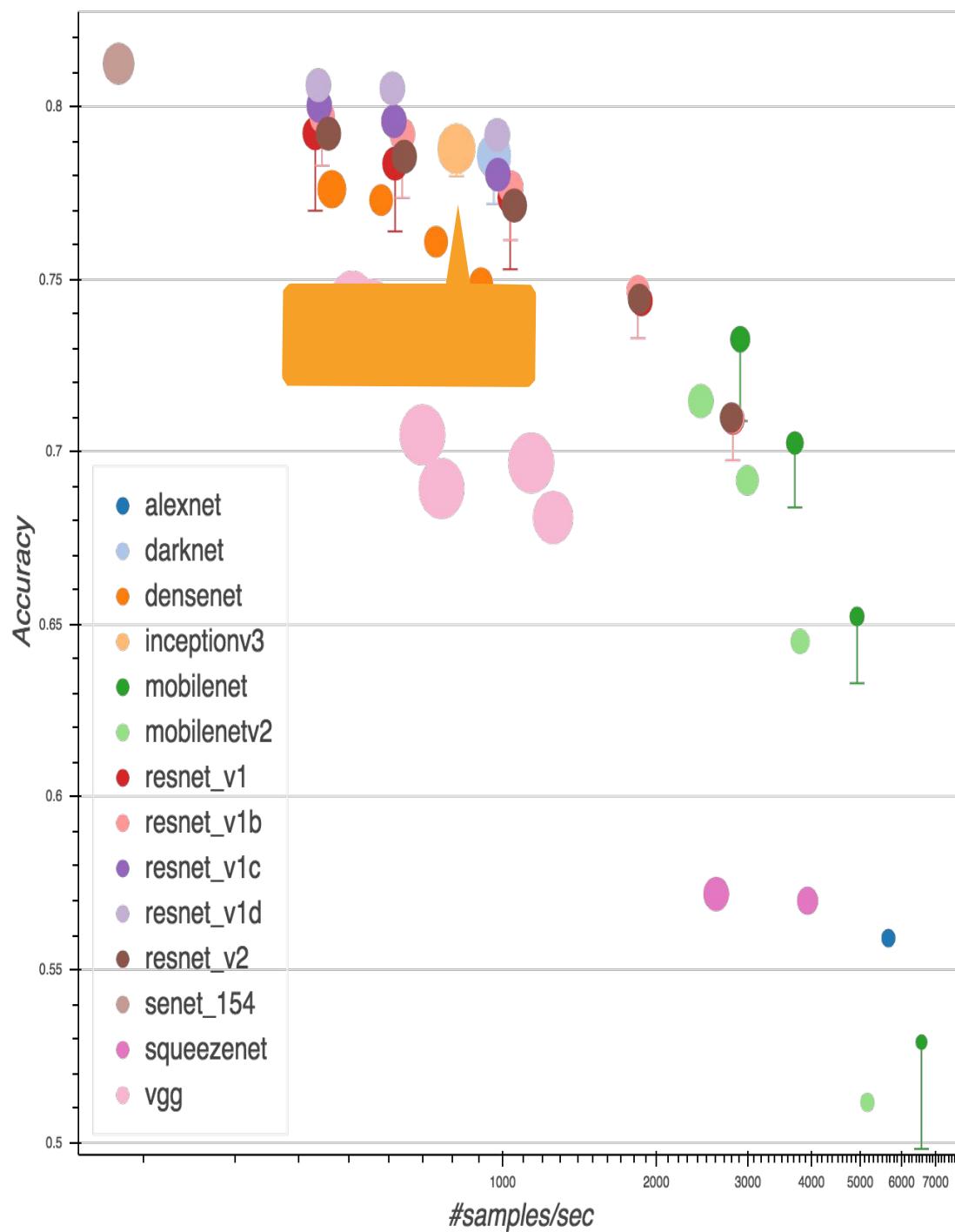


- 5 stages with 9 inception blocks



The many flavors of Inception Networks

- Inception-BN (v2) - Add batch normalization
- Inception-V3 - Modified the inception block
 - Replace 5x5 by multiple 3x3 convolutions
 - Replace 5x5 by 1x7 and 7x1 convolutions
 - Replace 3x3 by 1x3 and 3x1 convolutions
 - Generally deeper stack
- Inception-V4 - Add residual connections (more later)



GluonCV Model Zoo

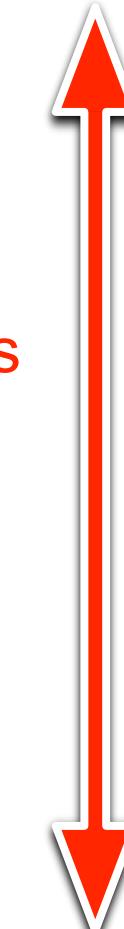
https://gluon-cv.mxnet.io/model_zoo/classification.html



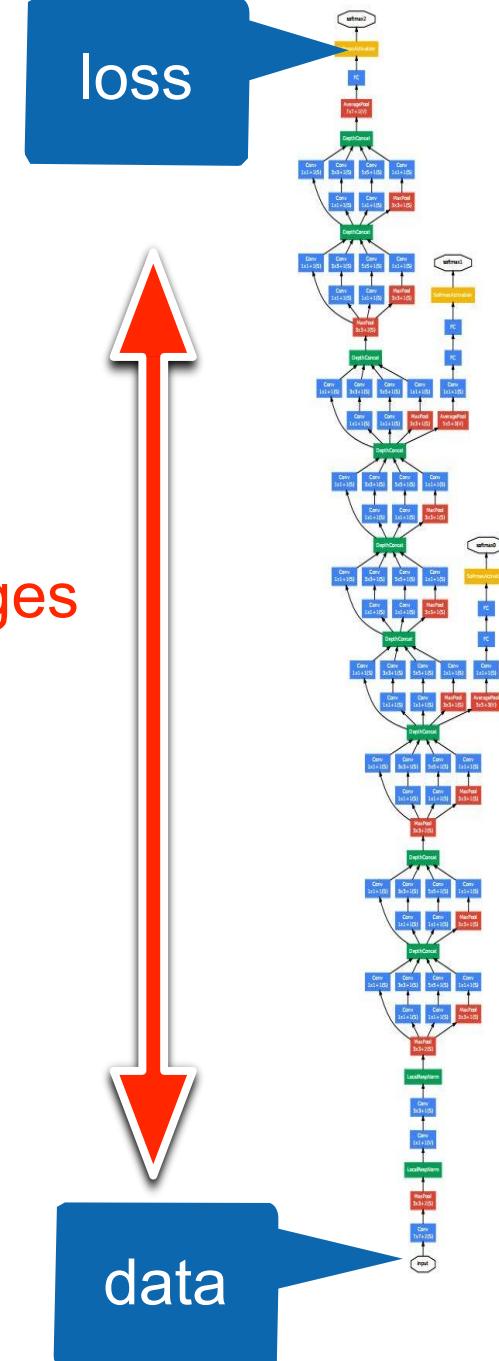
Batch Normalization

- Loss occurs at last layer
 - Last layers learn quickly
- Data is inserted at bottom layer
 - Bottom layers change - **everything** changes
 - Last layers need to relearn many times
 - Slow convergence
- This is like covariate shift

Can we avoid changing last layers while learning first layers?



data



Batch Normalization

- Can we avoid changing last layers while learning first layers?
- Fix mean and variance

$$\mu_B = \frac{1}{|B|} \sum_{i \in B} x_i \text{ and } \sigma_B^2 = \frac{1}{|B|} \sum_{i \in B} (x_i - \mu_B)^2 + \epsilon$$

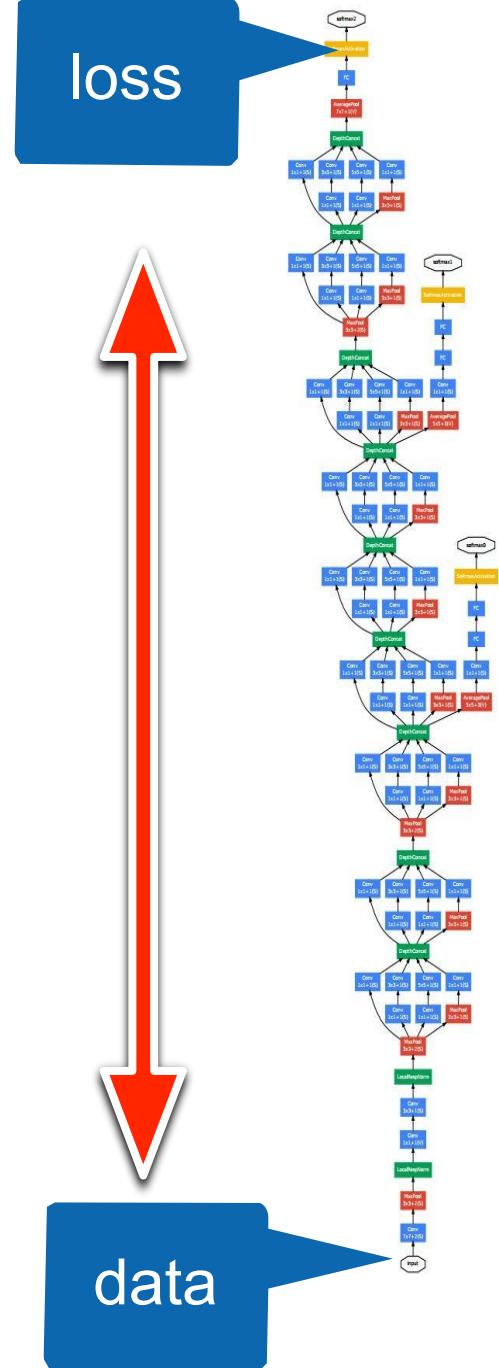
and adjust it separately

$$x_{i+1} = \gamma \frac{x_i - \mu_B}{\sigma_B} + \beta$$

variance

mean

data



This was the original motivation ...

What Batch Norms really do

- Doesn't really reduce covariate shift (Lipton et al., 2018)
- Regularization by noise injection

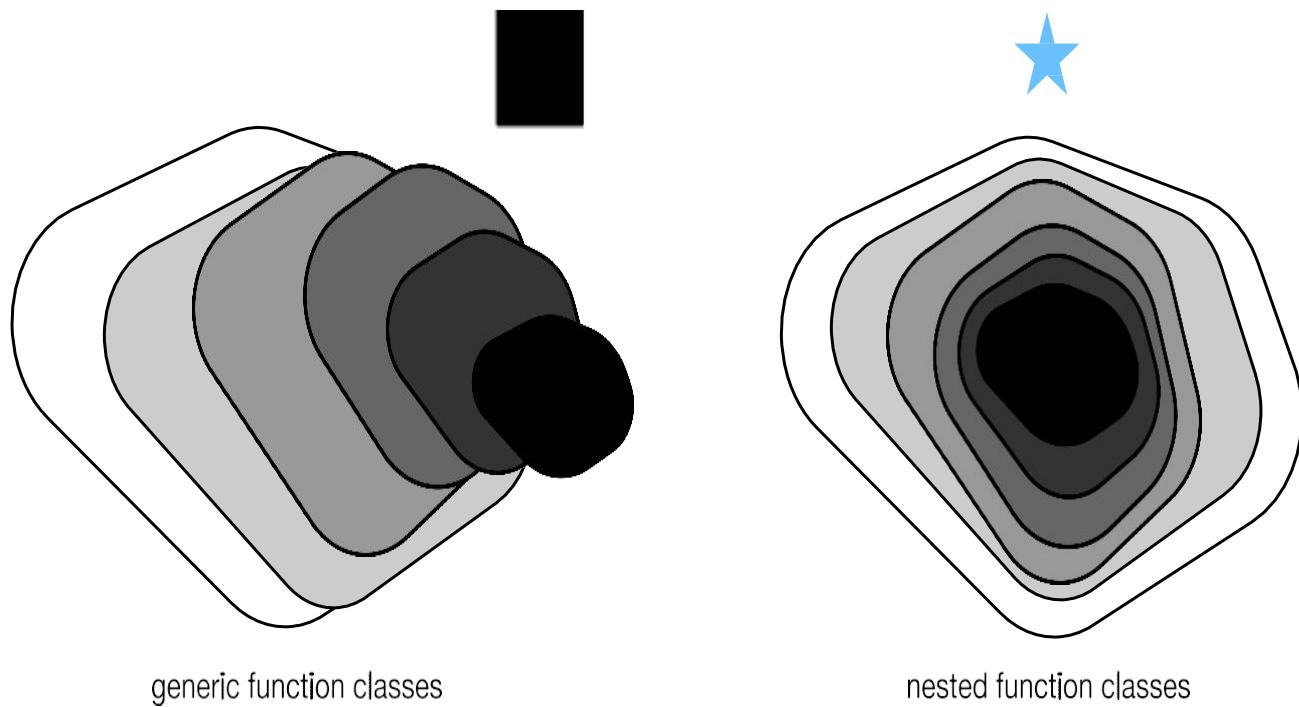
$$x_{i+1} = \gamma \frac{x_i - \hat{\mu}_B}{\hat{\sigma}_B} + \beta$$

- Random shift per minibatch
- Random scale per minibatch
- No need to mix with dropout (both are capacity control)
- Ideal minibatch size of 64 to 256

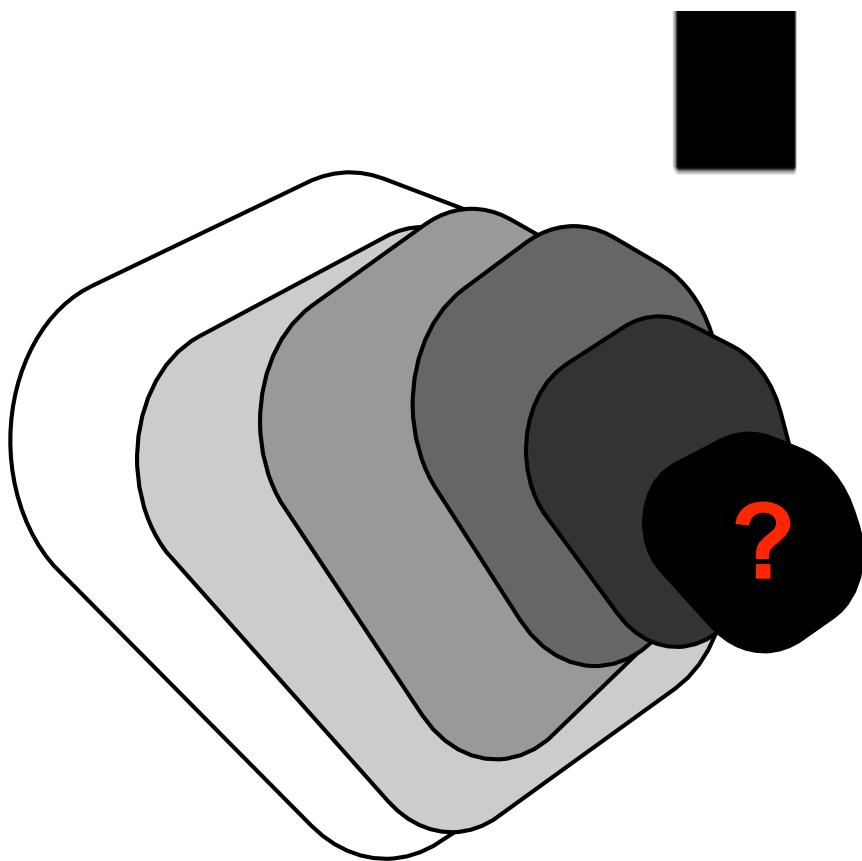
ERmanp
dieroicm
al

ERmanp
dieroicma
l

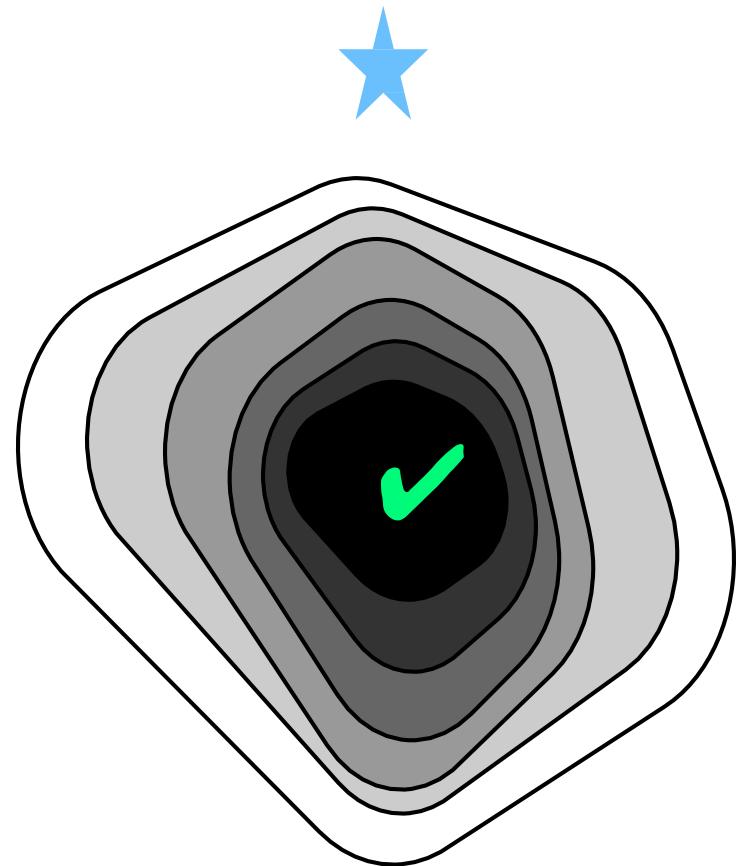
Residual Networks



Does adding layers improve accuracy?



generic function classes

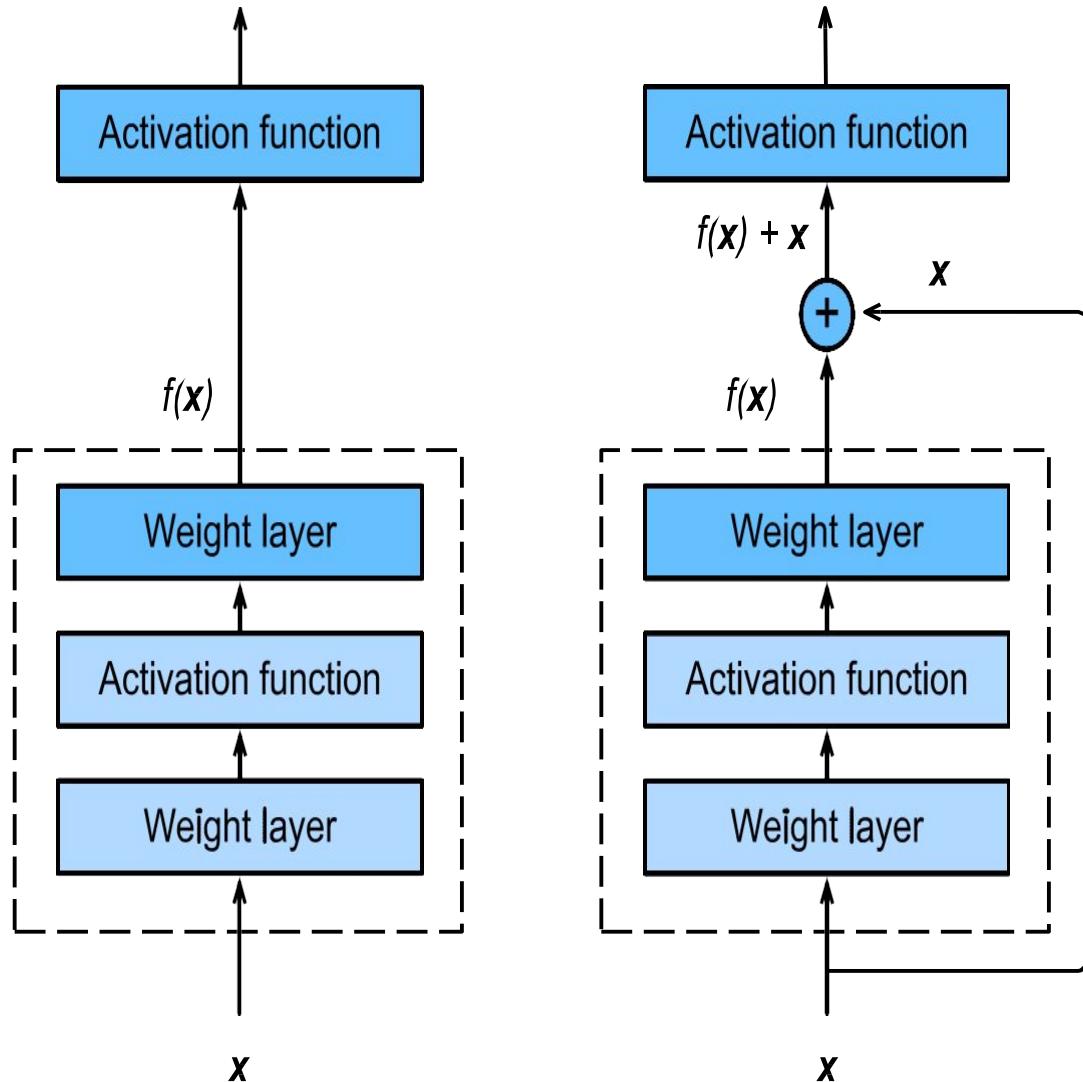


nested function classes

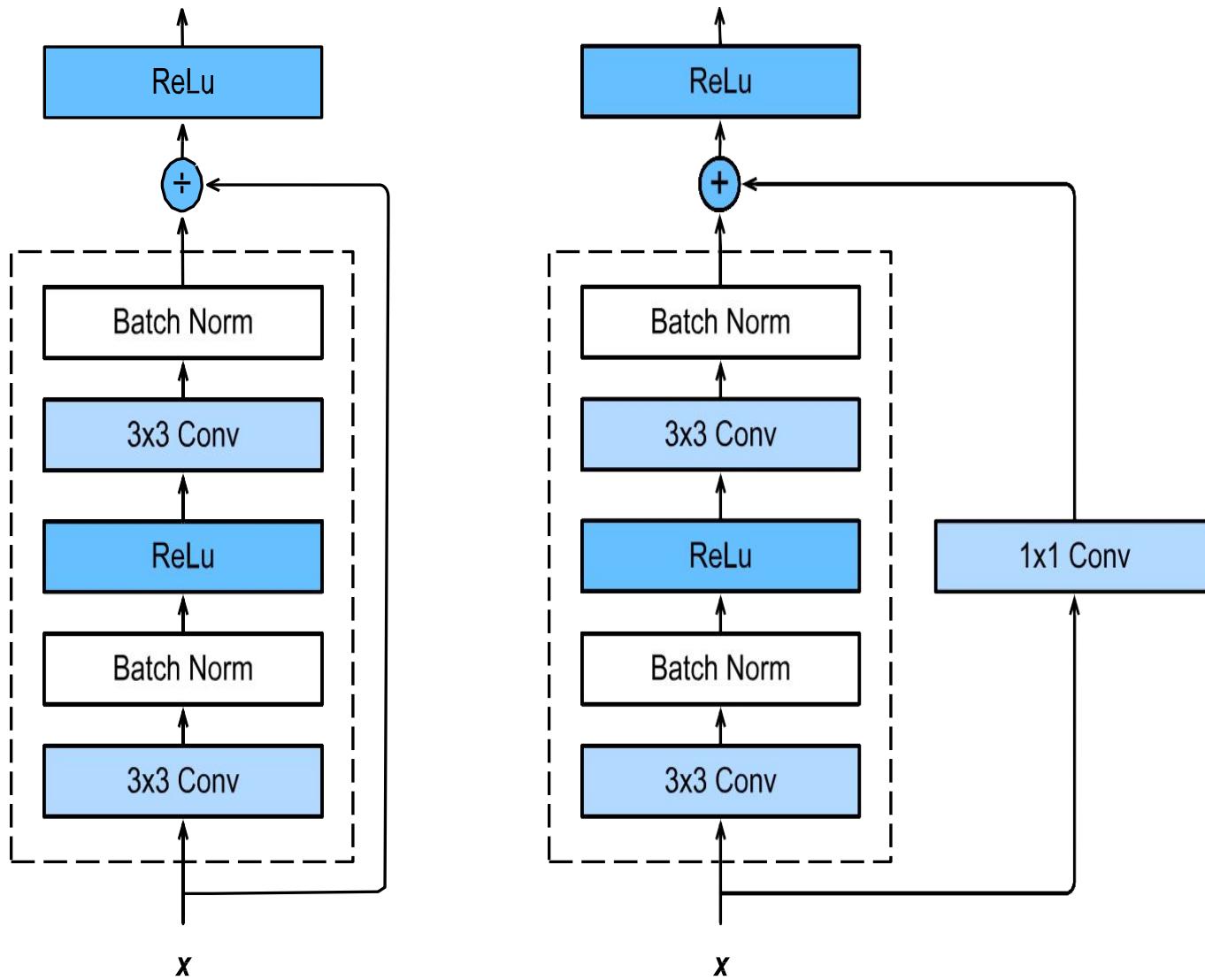
Residual Networks

- Adding a layer **changes** function class
- We want to **add to** the function class
- ‘Taylor expansion’ style parametrization

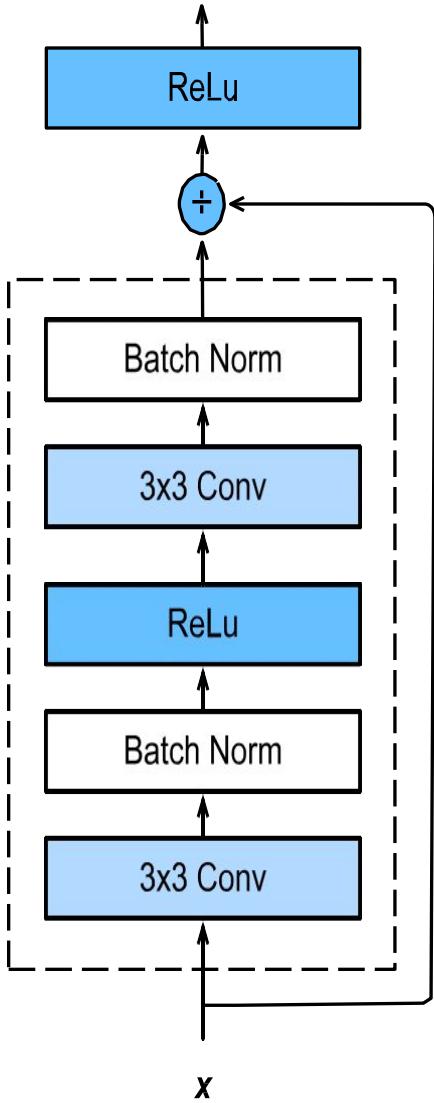
$$f(x) = x + g(x)$$



ResNet Block in detail

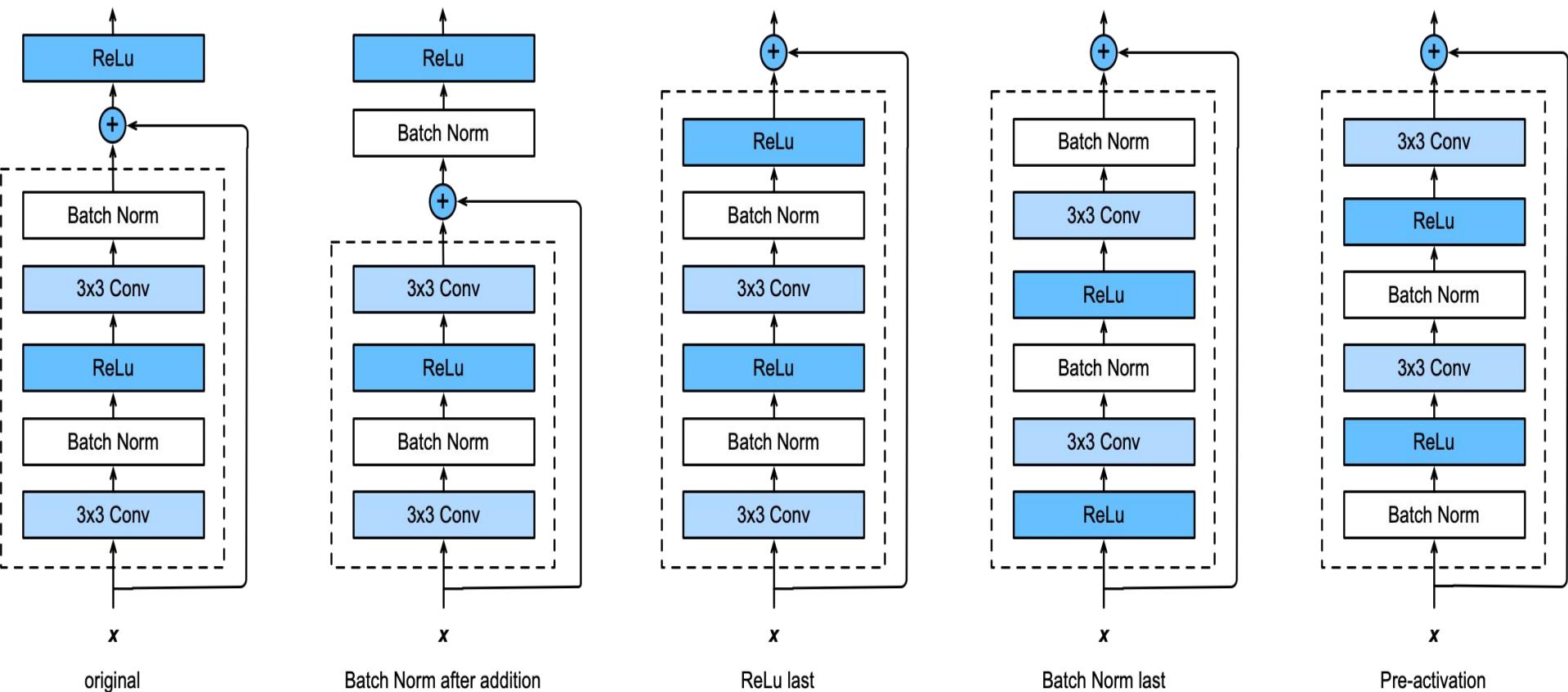


In code



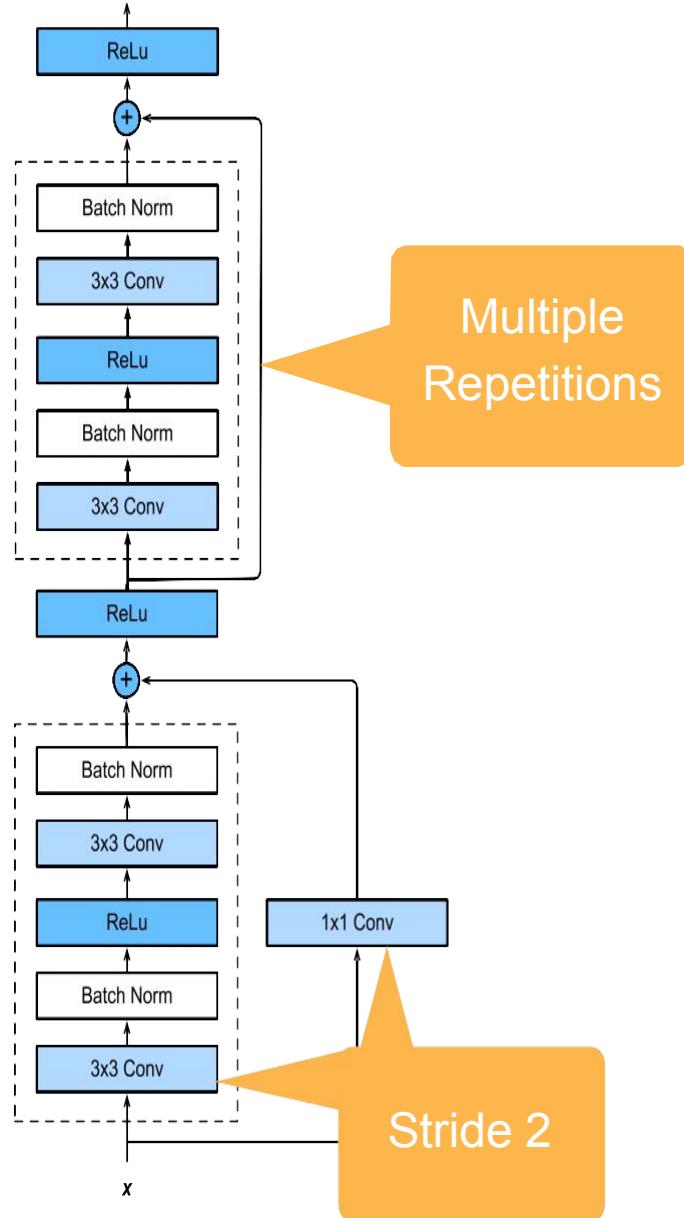
```
def forward(self, X):
    Y = npx.relu(self.bn1(self.conv1(X)))
    Y = self.bn2(self.conv2(Y))
    if self.conv3:
        X = self.conv3(X)
    return npx.relu(Y + X)
```

The many flavors of ResNet blocks



Try every permutation

ResNet Module



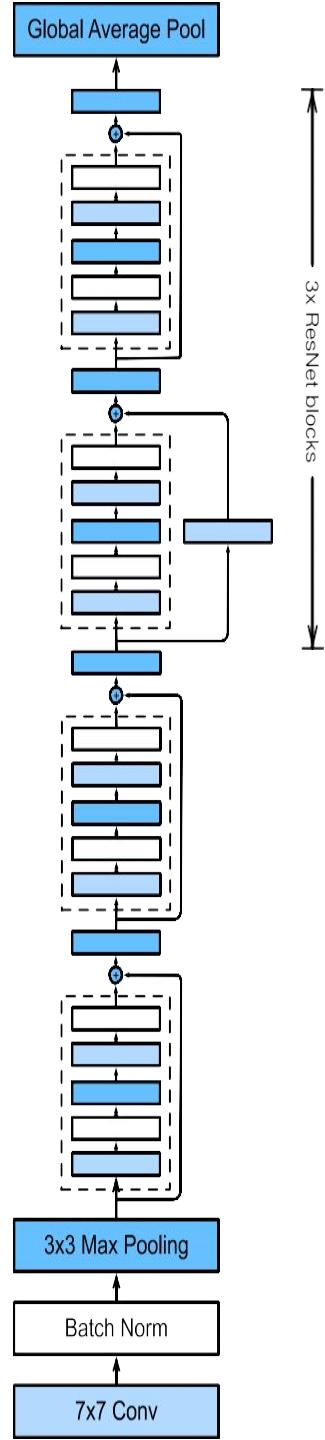
- Downsample per module (stride=2)
- Enforce some nontrivial nonlinearity per module (via 1x1 convolution)
- Stack up in blocks

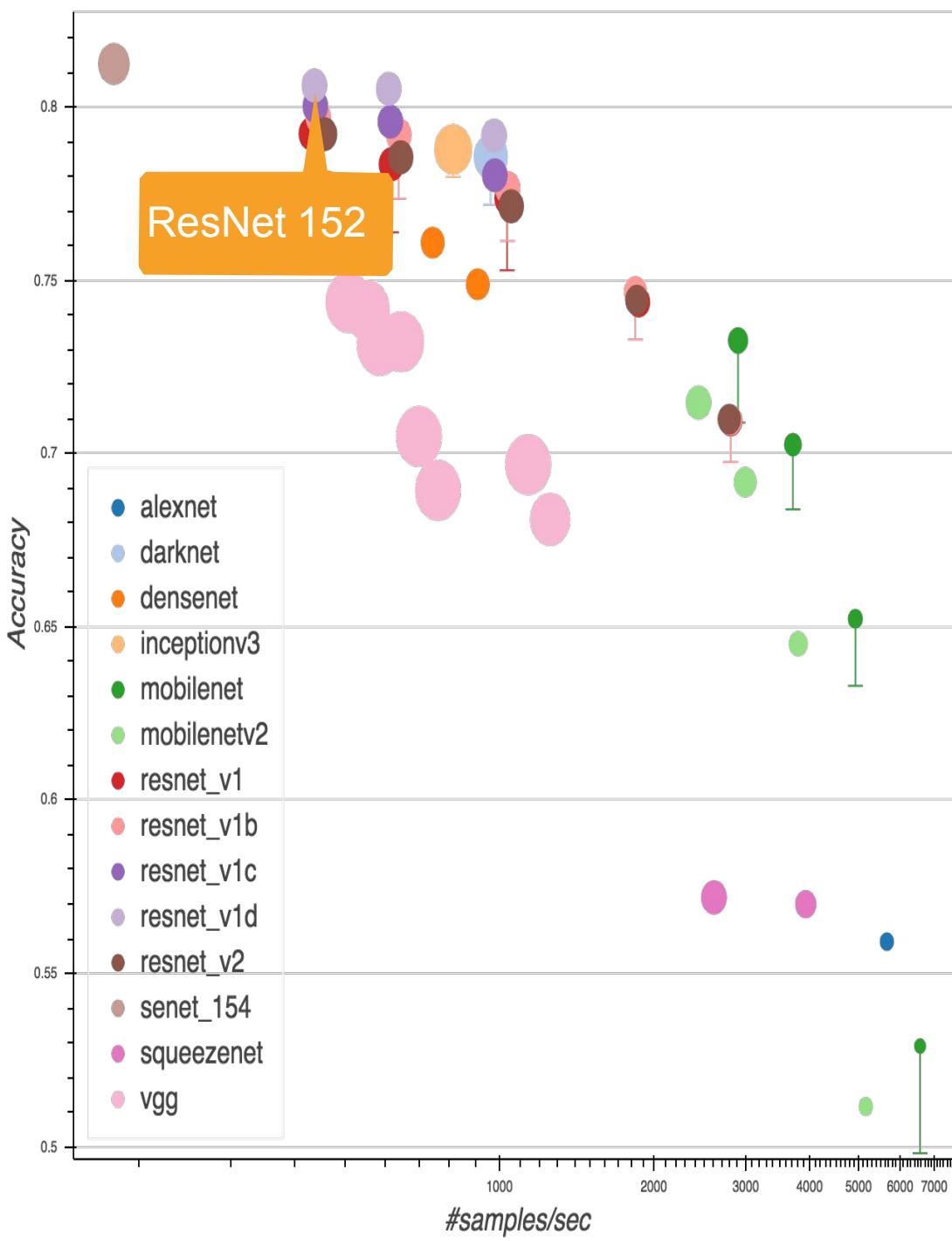
```
blk = nn.Sequential()  
for i in range(num_residuals):  
    if i == 0 and not first_block:  
  
        blk.add(Residual(num_channels,  
                         use_1x1conv=True, strides=2))  
    else:  
        blk.add(Residual(num_channels))
```

Putting it all together

- Same block structure as e.g. VGG or GoogleNet
- Residual connection to add to expressiveness
- Pooling/stride for dimensionality reduction
- Batch Normalization for capacity control

... train it at scale ...

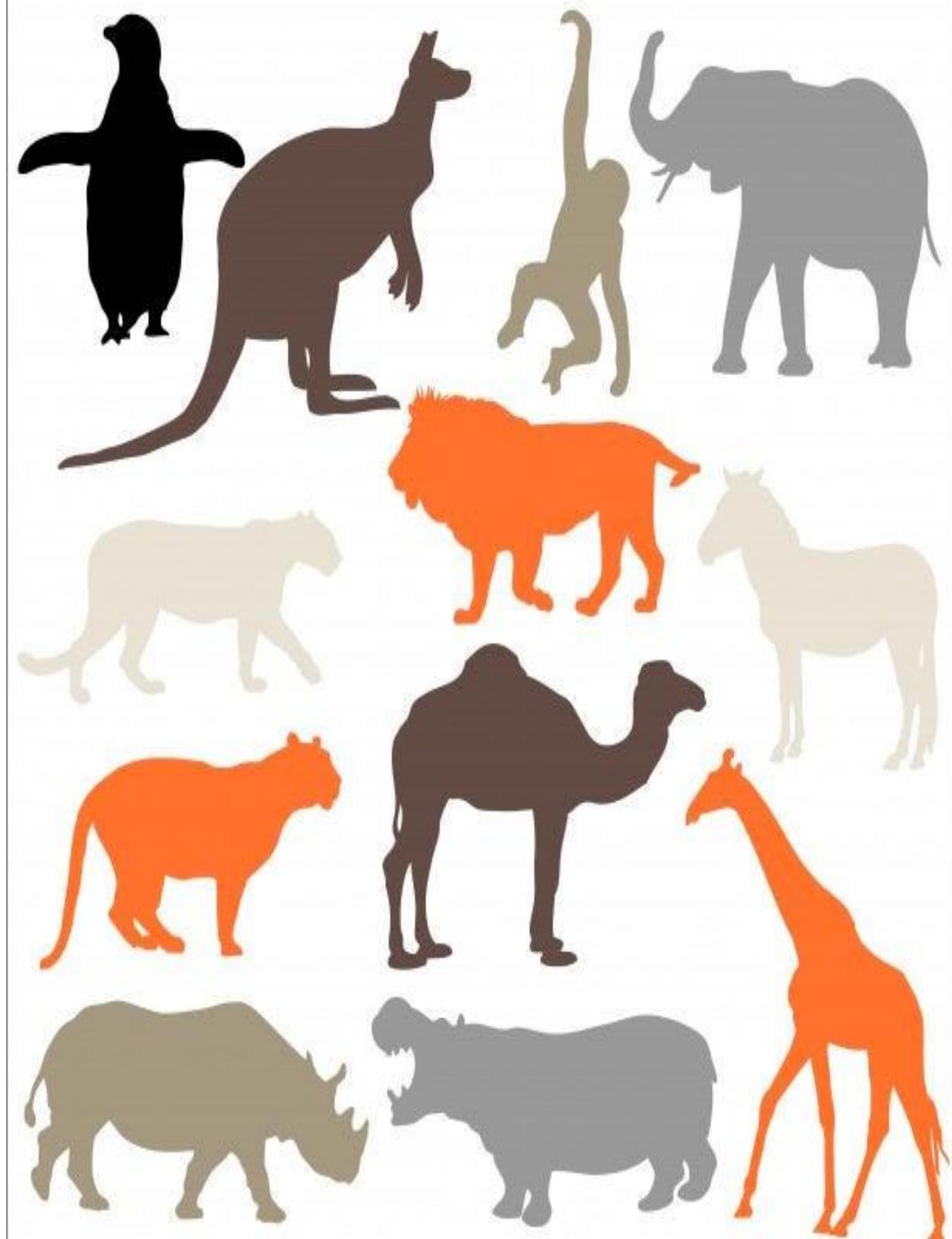




GluonCV Model Zoo
https://gluon-cv.mxnet.io/model_zoo/classification.html

Jupyter Notebook

More Ideas



DenseNet (Huang et al., 2016)

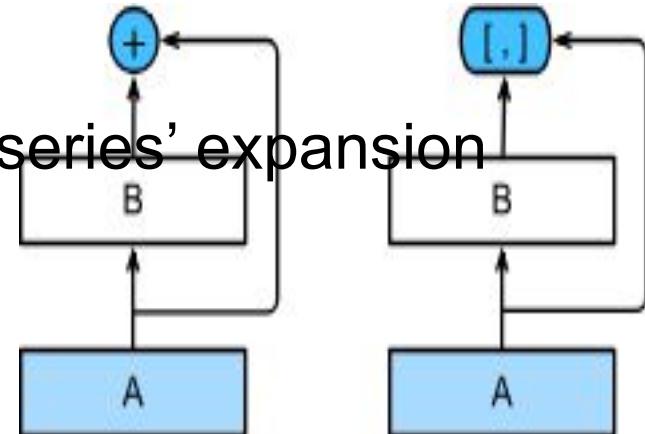
- ResNet combines x and $f(x)$
- DenseNet uses higher order ‘Taylor series’ expansion

$$x_{i+1} = [x_i, f_i(x_i)]$$

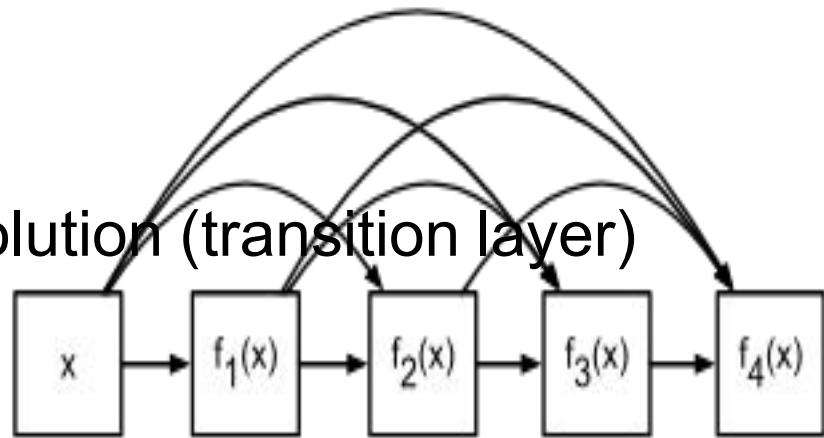
$$x_1 = x$$

$$x_2 = [x, f_1(x)]$$

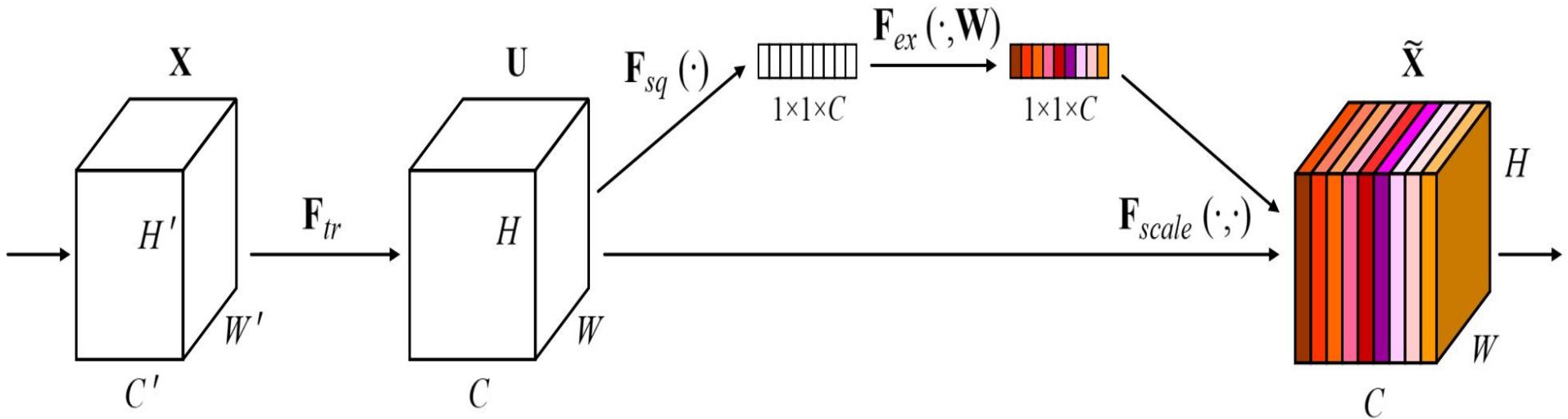
$$x_2 = [x, f_1(x), f_2([x, f_1(x)])]$$



- Occasionally need to reduce resolution (transition layer)



Squeeze-Excite Net (Hu et al., 2017)

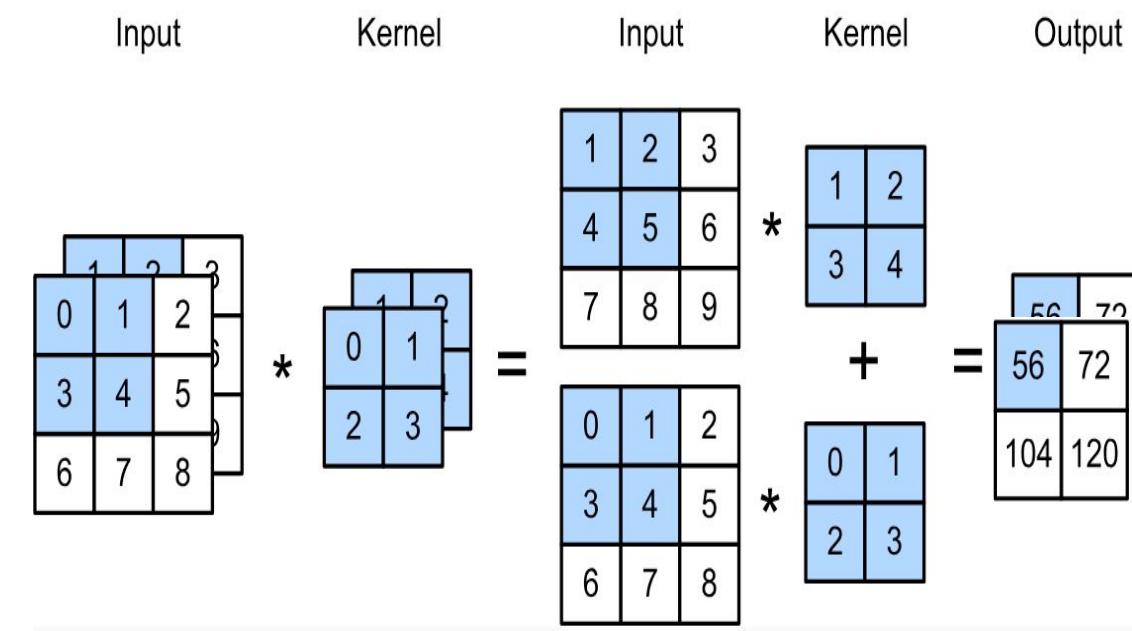


- Learn global weighting function per channel
- Allows for fast information transfer between pixels in different locations of the image

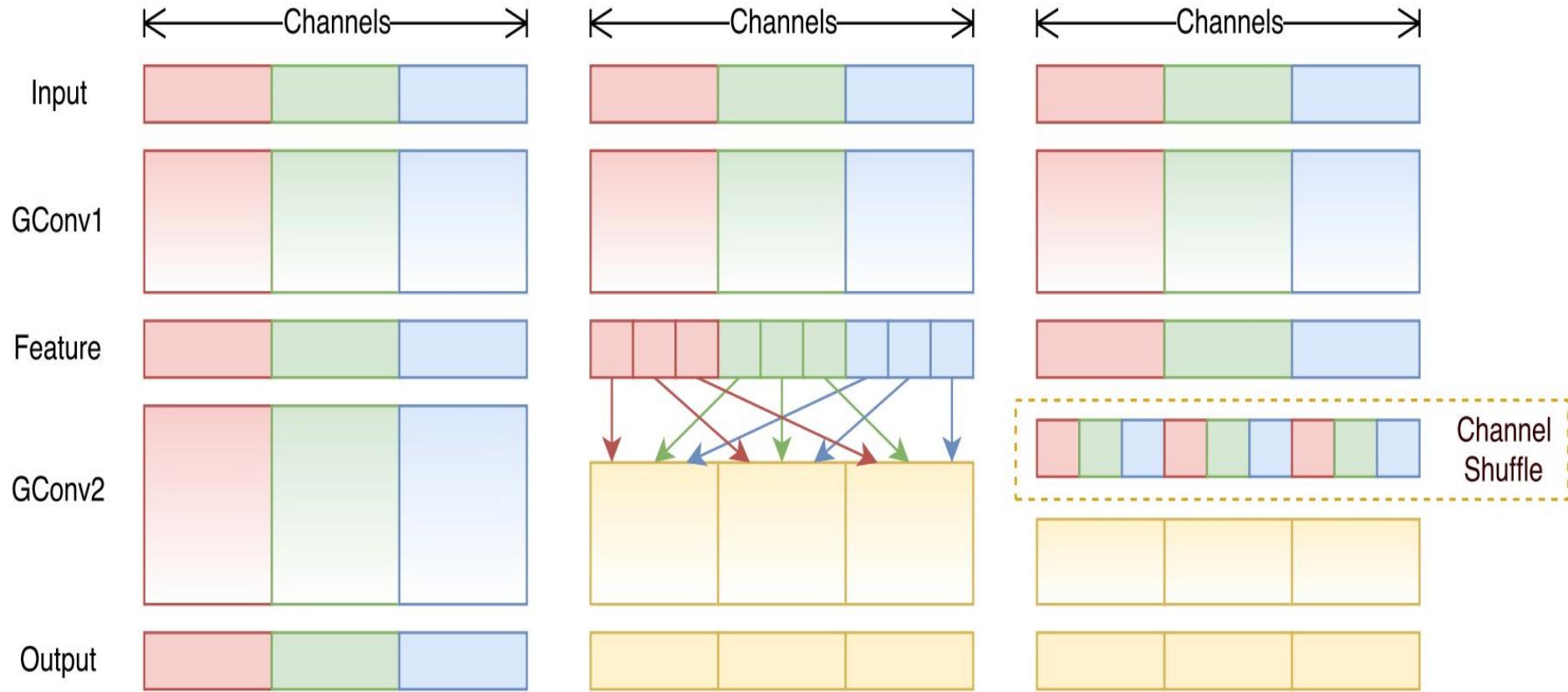
Separable Convolutions - all channels separate

- Parameters $k_h \cdot k_w \cdot c_i$
- Computation $c_o m_h \cdot m_w \cdot k_h \cdot k_w$
- Break up channels to the extreme
No mixing between channels

$$m_h \cdot m_w \cdot k_h \cdot k_w \cdot c$$



ShuffleNet (Zhang et al., 2018)



- ResNext breaks convolution into channels
- ShuffleNet mixes by grouping (very efficient for mobile)

Outline

- **GPUs**
- **Convolutions**
- **Pooling, Padding and Stride**
- **Convolutional Neural Networks (LeNet)**
- **Deep ConvNets (AlexNet)**
- **Networks using Blocks (VGG)**
- **Residual Neural Networks (ResNet)**