

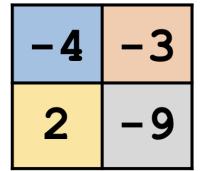
# Only certain input pixels are "connected" to certain output pixels

image

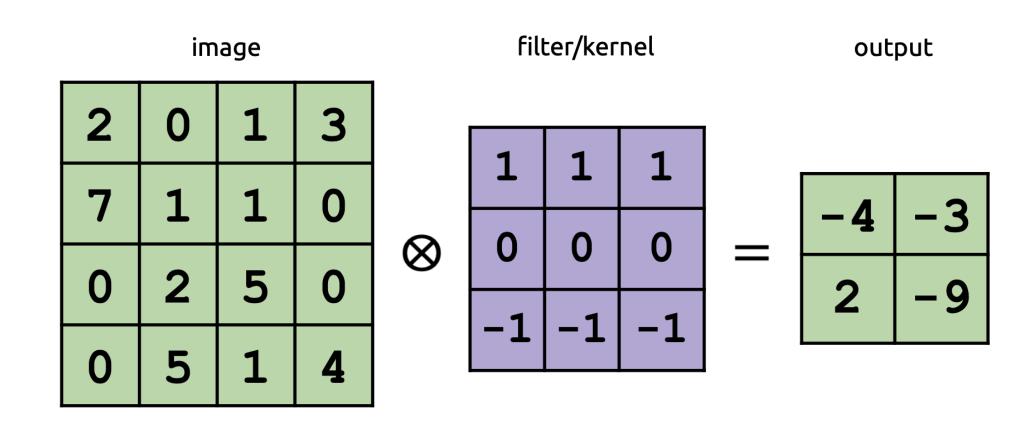
Colored dots in the input pixels represent which output pixels that input pixel contributes to

If this were fully connected, every input pixel would have all four output colors

output

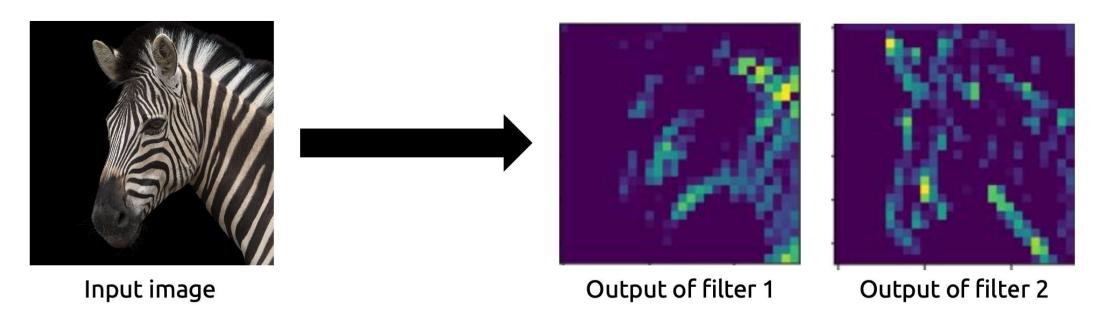


#### Key Idea 1: Filters are *Learnable*



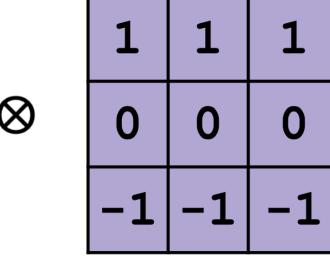
#### Key Idea 2: Learn *many* filters

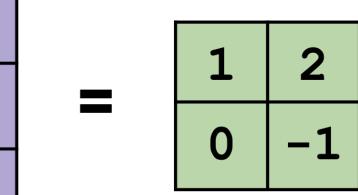
- Why are multiple filters a good idea?
  - Can learn to extract different features of the image



#### "Problem" With Convolution

2	0	1	3
0	1	1	0
0	0	2	0
0	1	1	1





- Output of convolution is always smaller than the input
- Why might we want the output size to be the same?
  - To avoid the filter "eating at the border" of the image when applying multiple conv layers

### Solution: Padding

Apply the kernel to 'imaginary' pixels surrounding the image

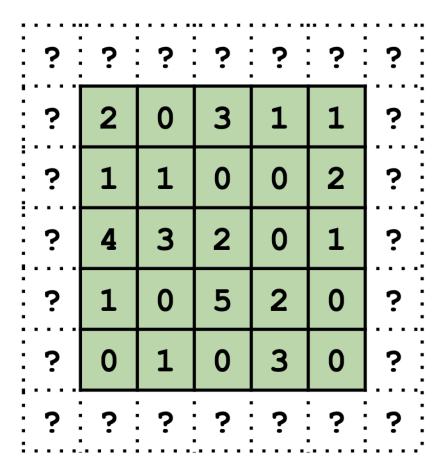
2	0	3	1	1
1	1	0	0	2
4	3	2	0	1
1	0	5	2	0
0	1	0	3	0

#### Solution: Padding

Apply the kernel to 'imaginary' pixels surrounding the image

?	?	?	?	?	?	?
?	2	0	3	1	1	?
?	1	1	0	0	2	?
?	4	3	2	0	1	?
?	1	0	5	2	0	?
?	0	1	0	3	0	?
?	?	?	?		?	?

#### What Values to Use For These Pixels?



#### What Values to Use For These Pixels?

Standard practice: fill with zeroes

0	0	0	0	0	0
2	0	3	1	1	0
1	1	0	0	2	0
4	3	2	0	1	0
1	0	5	2	0	0
0	1	0	3	0	0
0	0	0	0	0	0
	2 1 4 1	<ul><li>2 0</li><li>1 1</li><li>4 3</li><li>1 0</li><li>0 1</li></ul>	2       0       3         1       1       0         4       3       2         1       0       5         0       1       0	2       0       3       1         1       1       0       0         4       3       2       0         1       0       5       2         0       1       0       3	1       1       0       0       2         4       3       2       0       1         1       0       5       2       0         0       1       0       3       0

#### What Values to Use For These Pixels?

#### Standard practice: fill with zeroes

 Zero-valued padding pixels just result in some terms in the convolution sum being zero

$$V(x,y) = (\mathbb{I} \otimes K)(x,y) = \sum_{m} \sum_{n} I(x+m,y+n)K(m,n)$$

This is zero for a padding pixel

 End result: equivalent to a applying a 'masked' version of the filter that only covers the valid pixels

					<b>.</b>		
:	0	0	0	0	0	0	0
	0	2	0	3	1	1	0
	0	1	1	0	0	2	0
	0	4	3	2	0	1	0
	0	1	0	5	2	0	0
	0	0	1	0	3	0	0
	0	0	0	0	0	0	0

#### Padding Modes in Tensorflow

2 available options: 'VALID' and 'SAME':

#### **Valid**

Filter only slides over "Valid" regions of the data

2	0	1	3
0	1	1	0
0	0	2	0
0	1	1	1

#### Same

Filter slides over the bounds of the data, ensuring output size is the "Same" as input size (when stride = 1)

0	0	0	0	0	0
0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0

2	0	1	3
0	1	1	0
0	0	2	0
0	1	1	1

2	0	1	3
0	1	1	0
0	0	2	0
0	1	1	1

2	0	1	3
0	1	1	0
0	0	2	0
0	1	1	1

2	0	1	3
0	1	1	0
0	0	2	0
0	1	1	1

### We already tried this! (reduced output size)

2	0	3	1
1	1	0	0
1	0	2	0
1	0	1	2



1	0	-1
2	0	-2
1	0	-1

0	1
-1	-1

0	0	0	0	0	0
0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0

0	0	0	0	0	0
0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0

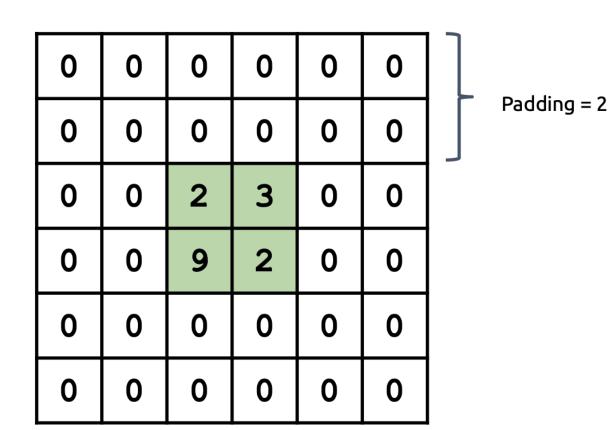
0	0	0	0	0	0
0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0

0	0	0	0	0	0
0	2	0	1	თ	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0

#### Output Size of a Convolution Layer

The output size of a convolution layer depends on 4 Hyperparameters:

- Number of filters, N
- The size of these filters, F
- The stride, S
- The amount of padding, P



#### Output Size of a Convolution Layer

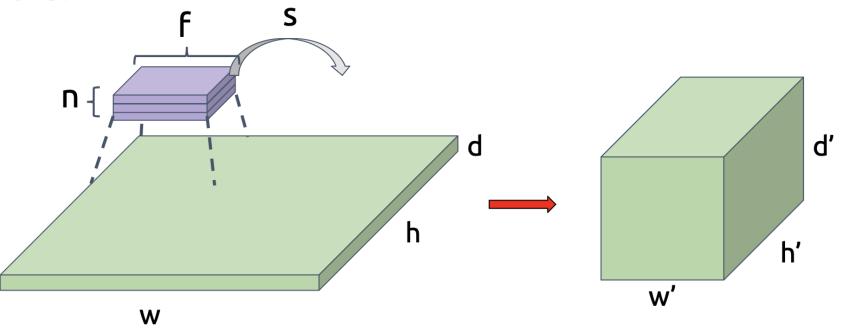
Suppose we know the number of filters, their size, the stride, and padding (n,f,s,p).

Then for a convolution layer with input dimension w x h x d, the output dimensions w' x h' x d' are:

$$w' = \frac{w - f + 2p}{s} + 1$$

$$h' = \frac{h - f + 2p}{s} + 1$$

$$d' = n$$



$$w' = \frac{w - f + 2p}{s} + 1$$

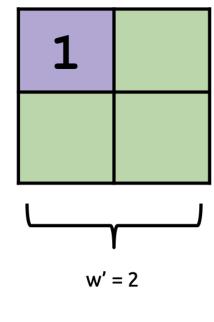
Let 
$$w = 4$$

num filters 
$$n = 1$$
  
filter size  $f = 3$   
stride  $s = 1$   
padding  $p = 0$ 

$$w' = \frac{4 - 3 + 2 \cdot 0}{1} + 1$$
$$= 1 + 1 = 2$$

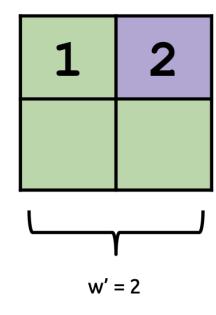
$$w' = \frac{w - f + 2p}{s} + 1$$

2	0	1	3		
0	1	1	0		
0	0	2	0		
0	1	1	1		
w = 4					

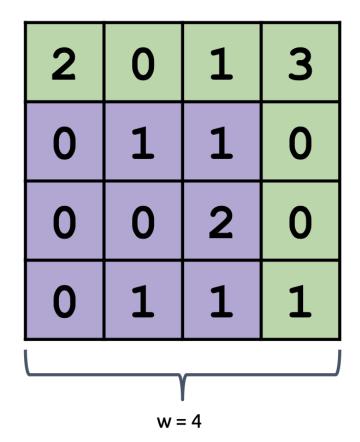


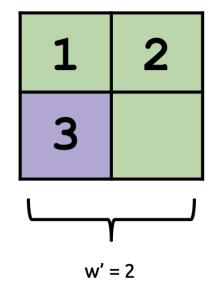
$$w' = \frac{w - f + 2p}{s} + 1$$

2	0	1	თ			
0	1	1	0			
0	0	2	0			
0	1	1	1			
<b>y</b> w = 4						

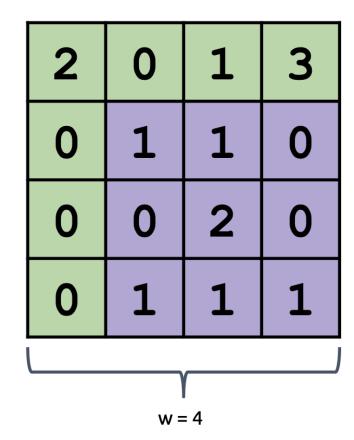


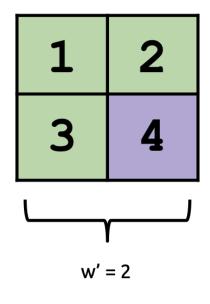
$$w' = \frac{w - f + 2p}{s} + 1$$





$$w' = \frac{w - f + 2p}{s} + 1$$





$$w' = \frac{w - f + 2p}{s} + 1$$

num filters 
$$n = 1$$
  
filter size  $f = 3$   
stride  $s = 1$   
padding  $p = 1*$ 

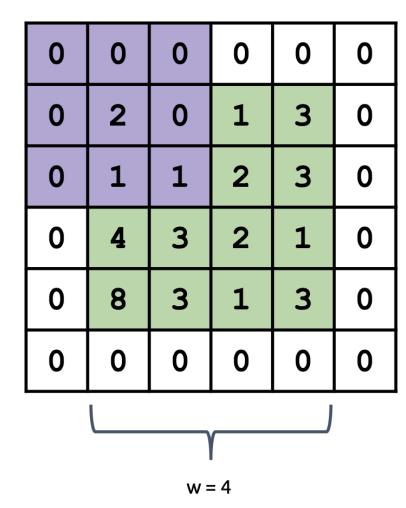
Let 
$$w = 4$$

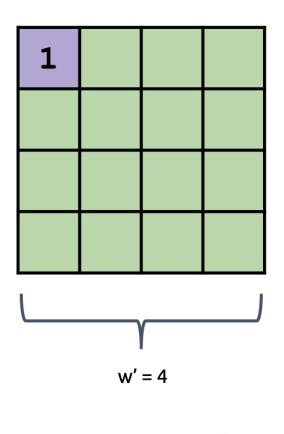
$$w' = \frac{4 - 3 + 2 \cdot 1}{1} + 1$$

$$= 3 + 1 = 4$$

$$w' = \frac{w - f + 2p}{s} + 1$$

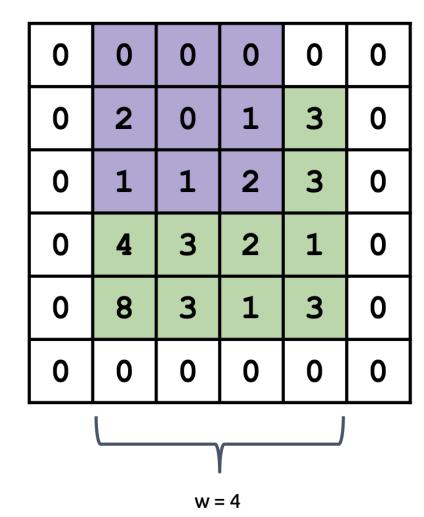
num filters 
$$n = 1$$
  
filter size  $f = 3$   
stride  $s = 1$   
padding  $p = 1*$ 

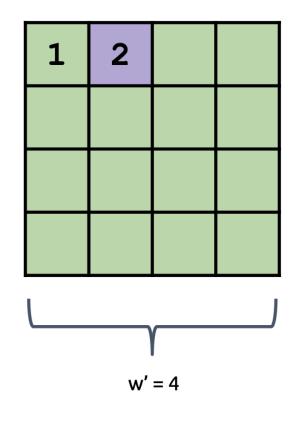




$$w' = \frac{w - f + 2p}{s} + 1$$

num filters 
$$n = 1$$
  
filter size  $f = 3$   
stride  $s = 1$   
padding  $p = 1*$ 

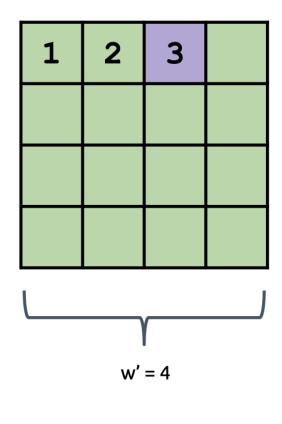




$$w' = \frac{w - f + 2p}{s} + 1$$

num filters 
$$n = 1$$
  
filter size  $f = 3$   
stride  $s = 1$   
padding  $p = 1*$ 

0	0	0	0	0	0
0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0
Υ					
w = 4					





$$w' = \frac{w - f + 2p}{s} + 1$$

num filters n = 1filter size f = 3stride s = 1padding p = 1\*

0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0
	Ī			J	

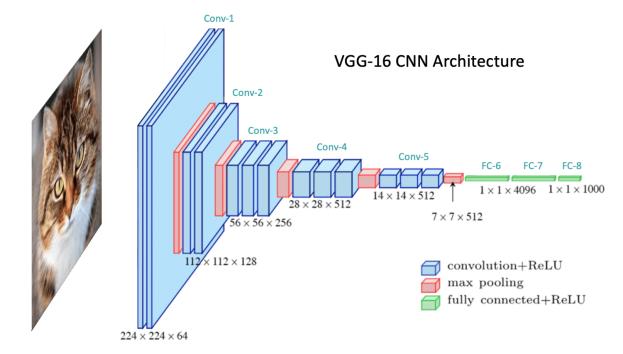
1	2	3	4

w' = 4

#### Getting network output

Remaining Question: If the convolution creates another [h x w x d] tensor, how do we actually get an output?

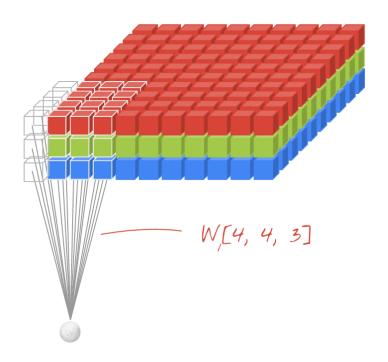
How can we turn use convolutions for classification?



https://learnopencv.com/understa nding-convolutional-neuralnetworks-cnn/

#### Color images...

What if our input has multiple channels (colors)? Do we apply filters to each individual color matrix? Or in some other way?



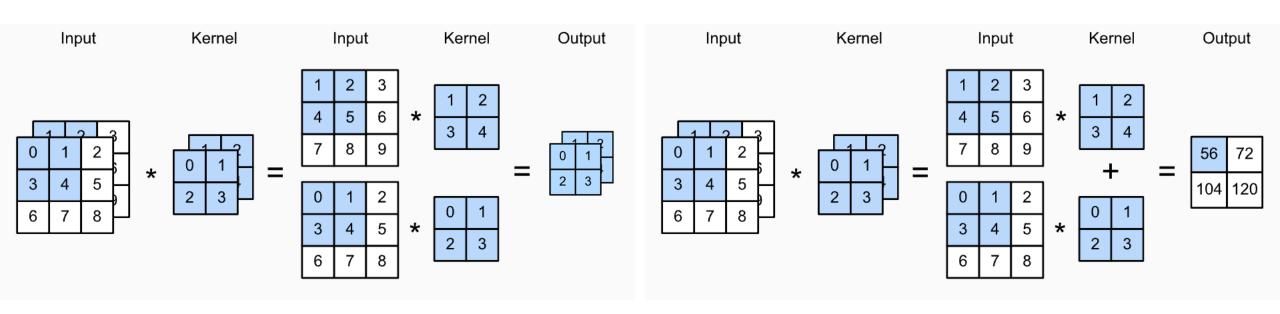
Source: Martin Görner

#### Multi-Channel Input

Which makes more sense?

Option #1 n channels to n outputs

Option #2 N channels to 1 output

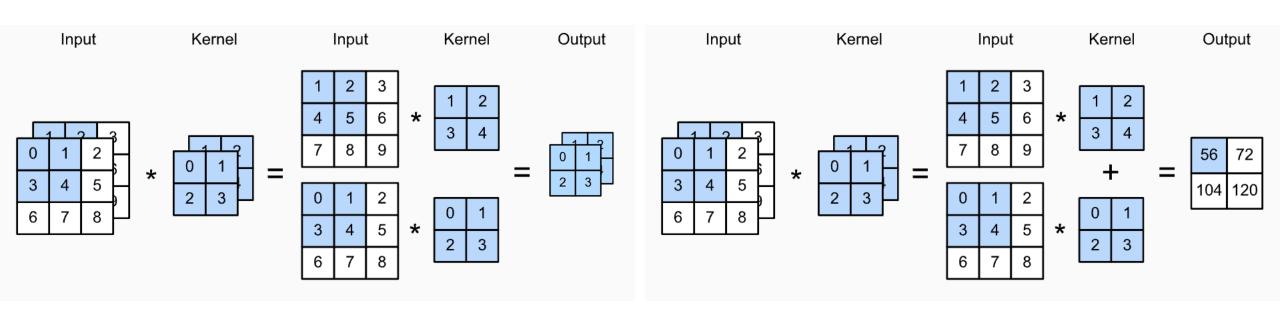


#### Multi-Channel Input

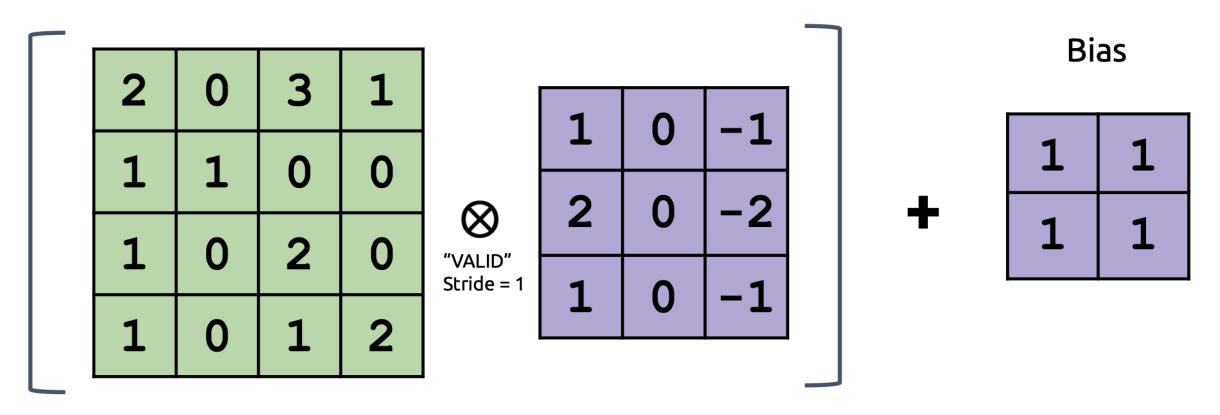
N-channels to 1 output allows information from separate channels to be used together

Option #1 n channels to n outputs

Option #2 N channels to 1 output



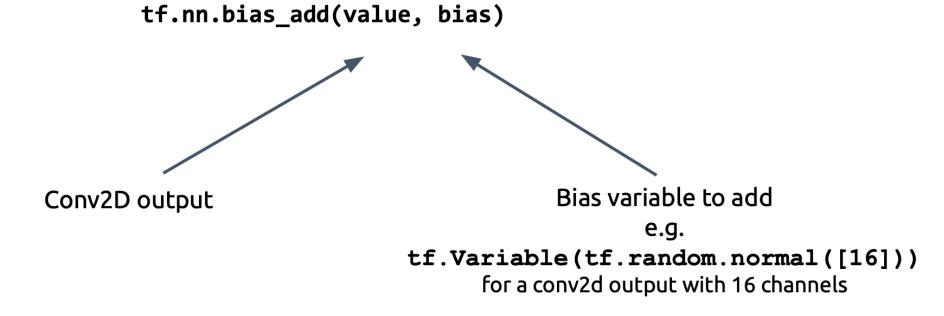
# Bias Term in Convolution Layers



Just like a fully connected layer, we can have a learnable additive bias for convolution.

# Adding a Bias in Tensorflow

If you use tf.nn.conv2d, bias can be added with:

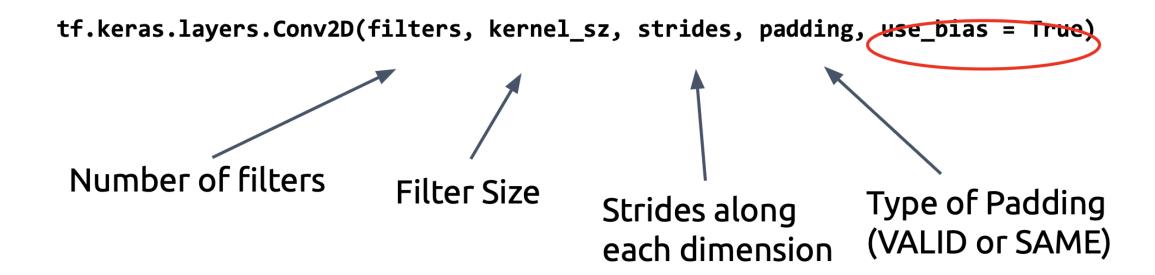


Full documentation here:

https://www.tensorflow.org/api\_docs/python/tf/nn/bias\_add

# Adding a Bias in Tensorflow

If you are using keras layers, bias is included by default:



Full documentation here:

https://www.tensorflow.org/versions/r2.0/api\_docs/python/tf/keras/layers/Conv2D

# Today's Goals

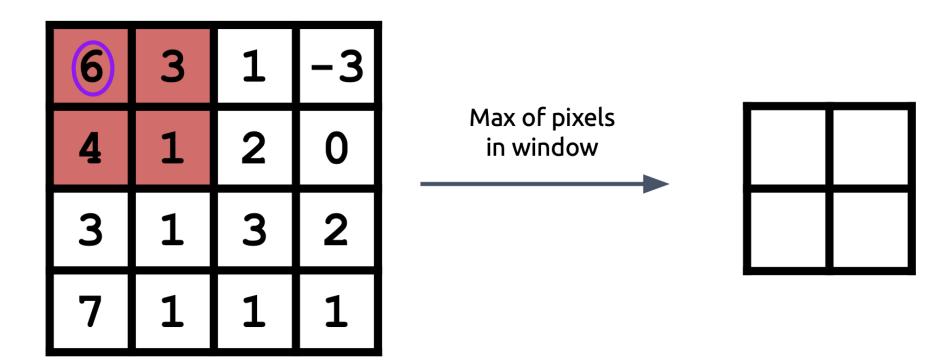
- (1) What non-linear activation functions are available to us?
- (2) Learn about Convolutional Architectures
  - (1) Many more decisions to make about structure of network than MLPs

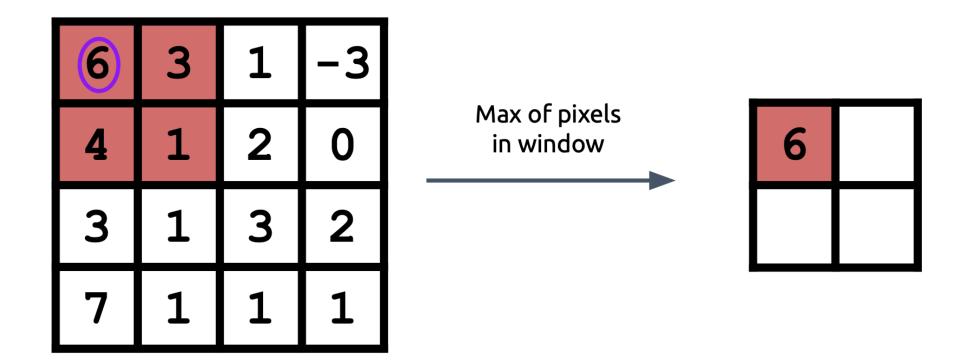
## **Activation Functions**

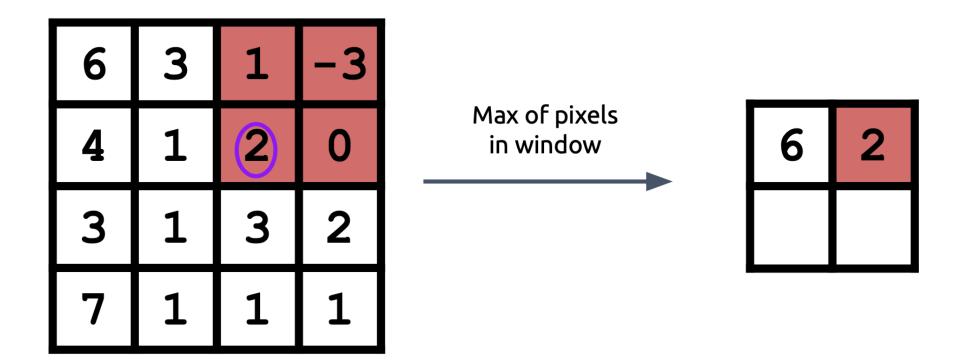
Remember, a linear combination of features, even if repeated many times, will always be linear.

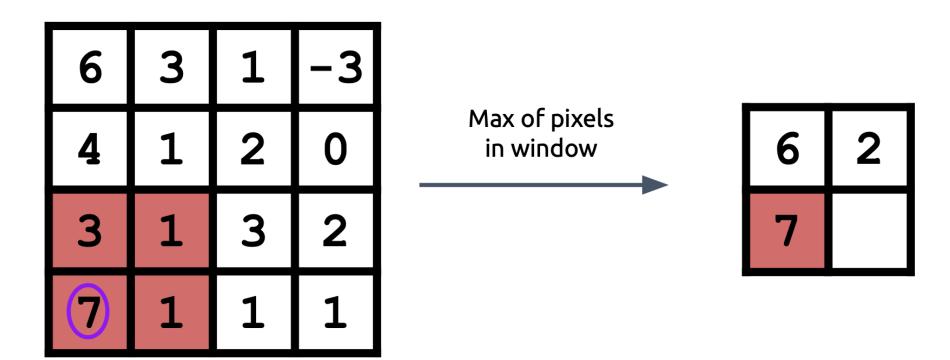
Still need some type of non-linear activation (e.g., ReLUs)

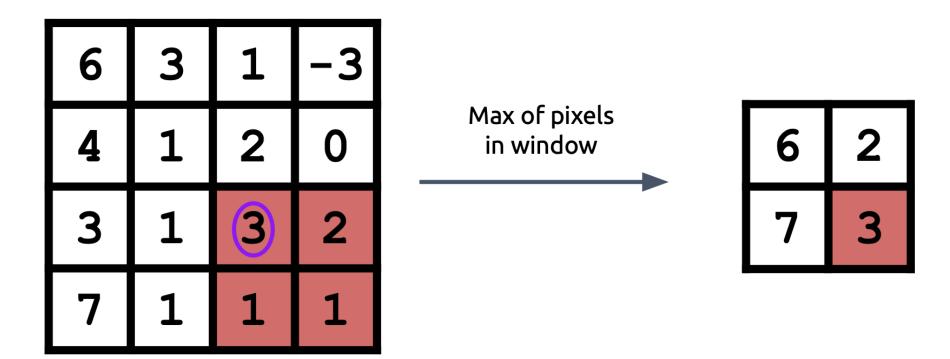
We also have other convolution-specific activation functions called "pooling" operations



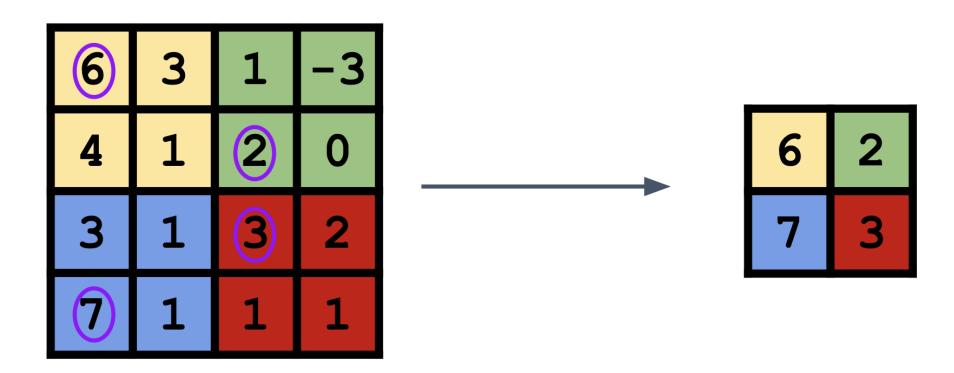








Max pooling with stride 2 and 2x2 filters



Why use Max Pooling?

# **Pooling: Motivation**

#### **Max Pooling**

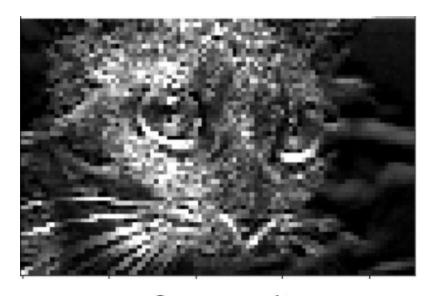
- Keeps track of regions with highest activations, indicating object presence
- Controllable way to lower (coarser) resolution (down sample the convolution output)



Original Image

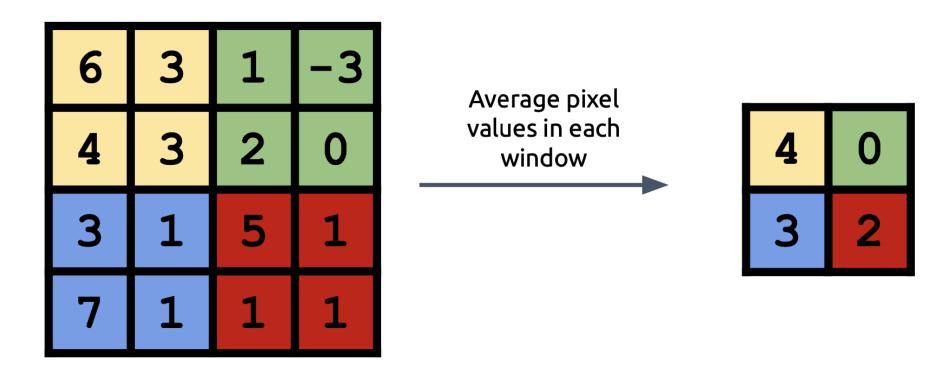


**Convolution Output** 



After Pooling

# Other Pooling Techniques

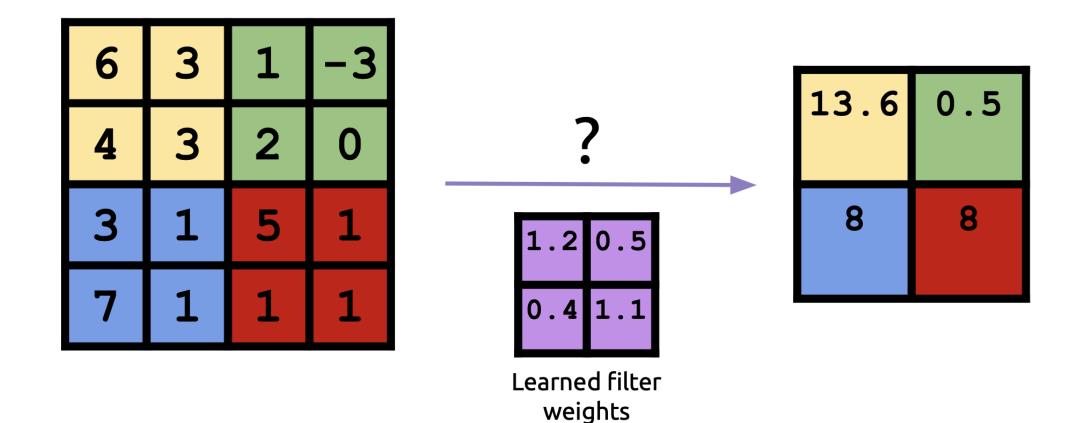


#### Any questions?

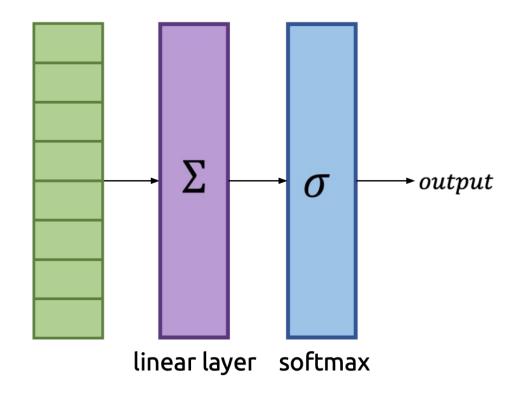


# Learning a Pooling Function

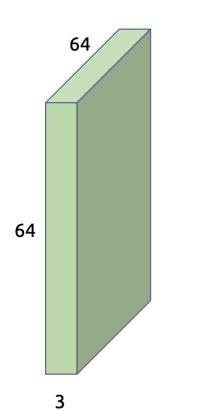
- The network can learn its own pooling function
- · Implement via a strided convolution layer

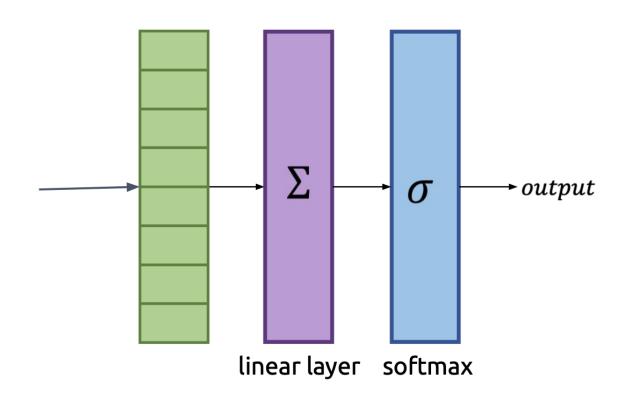


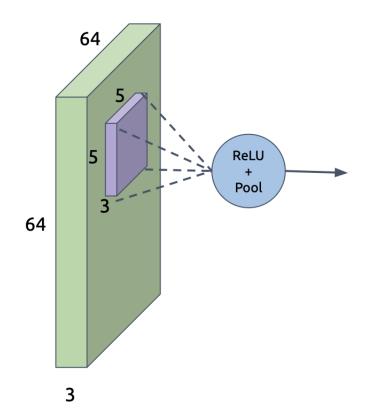
## Our neural network so far

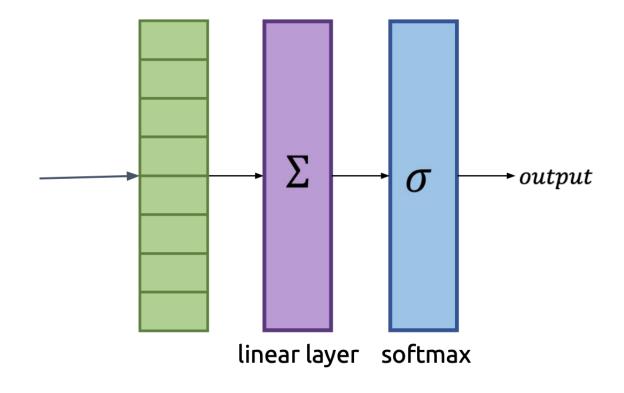


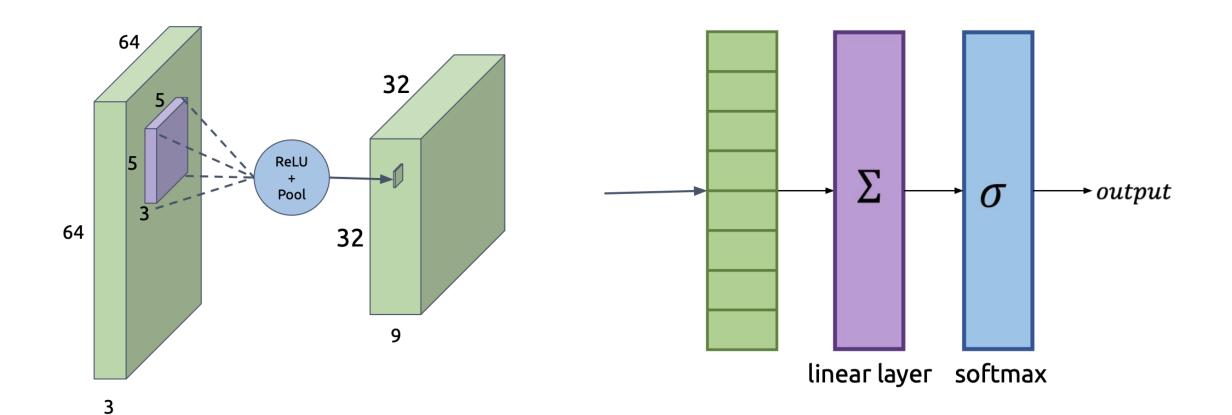
## Convolutional Neural Network Architecture

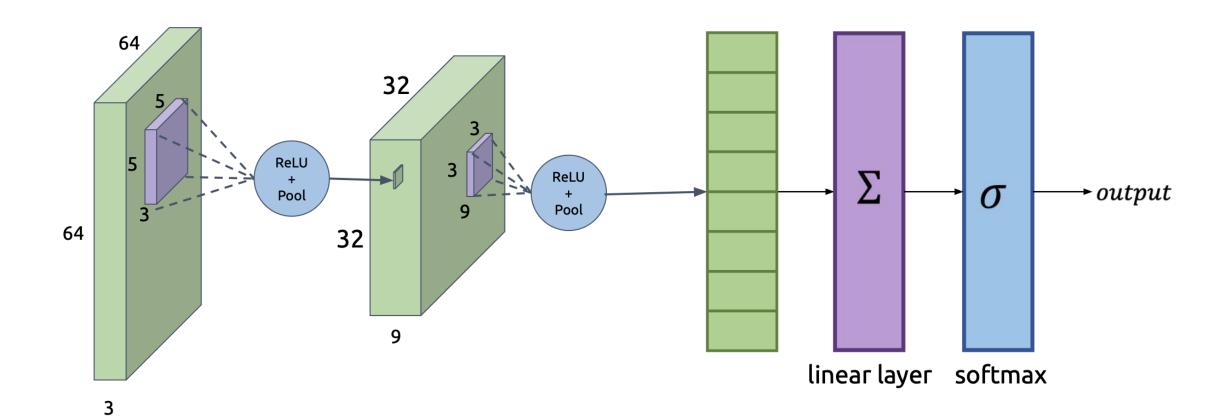




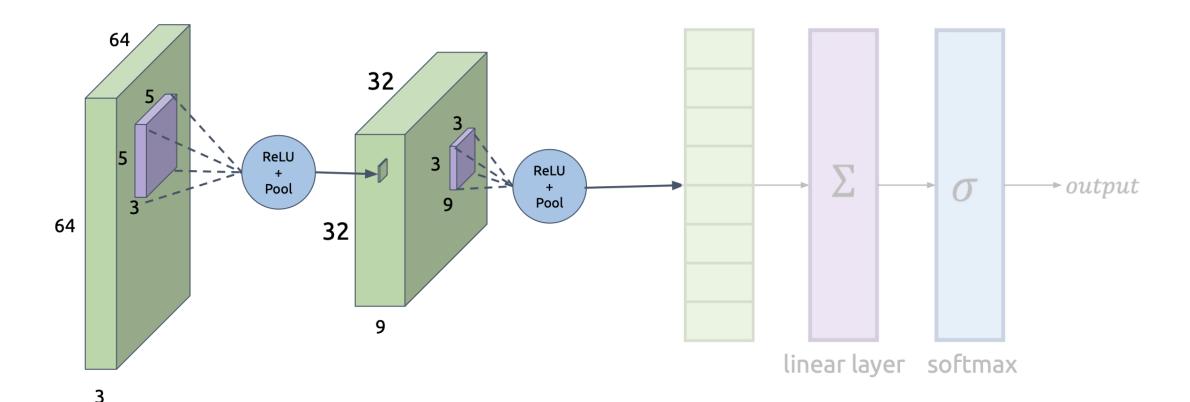


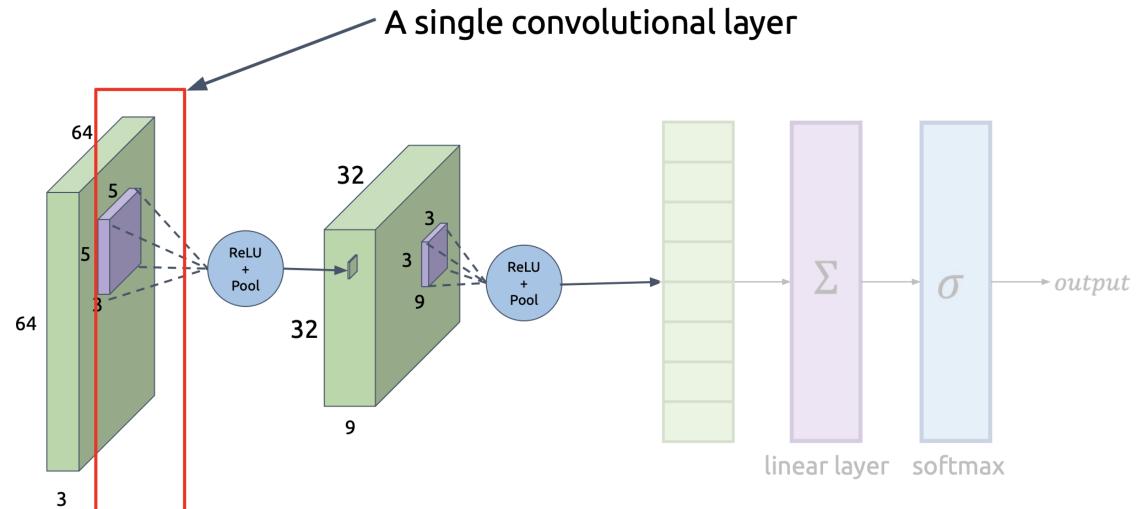




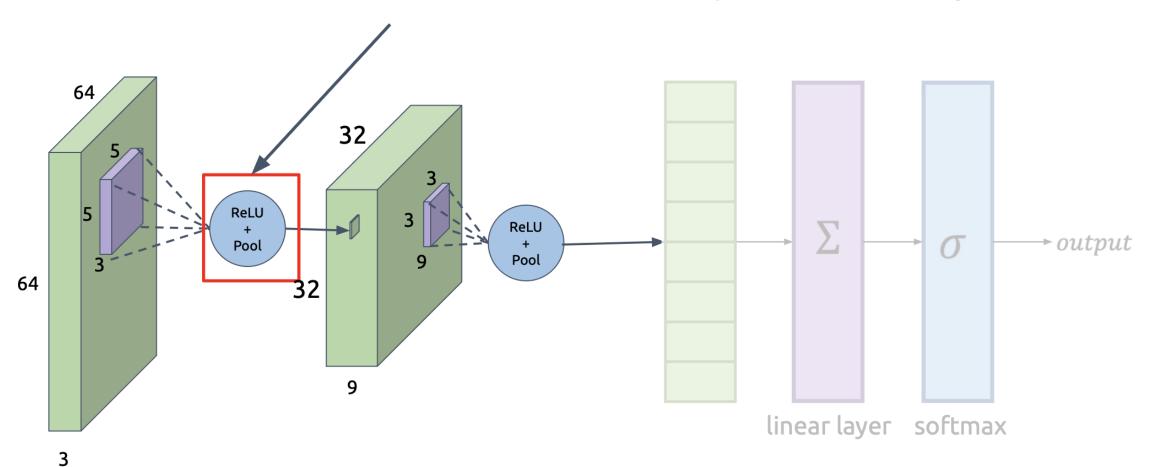


This part learns to extract *features* from the image

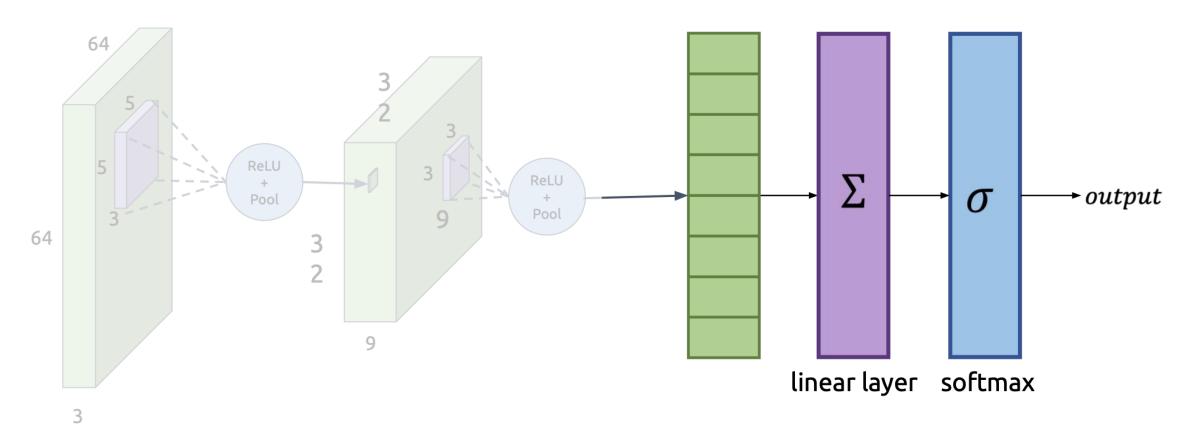




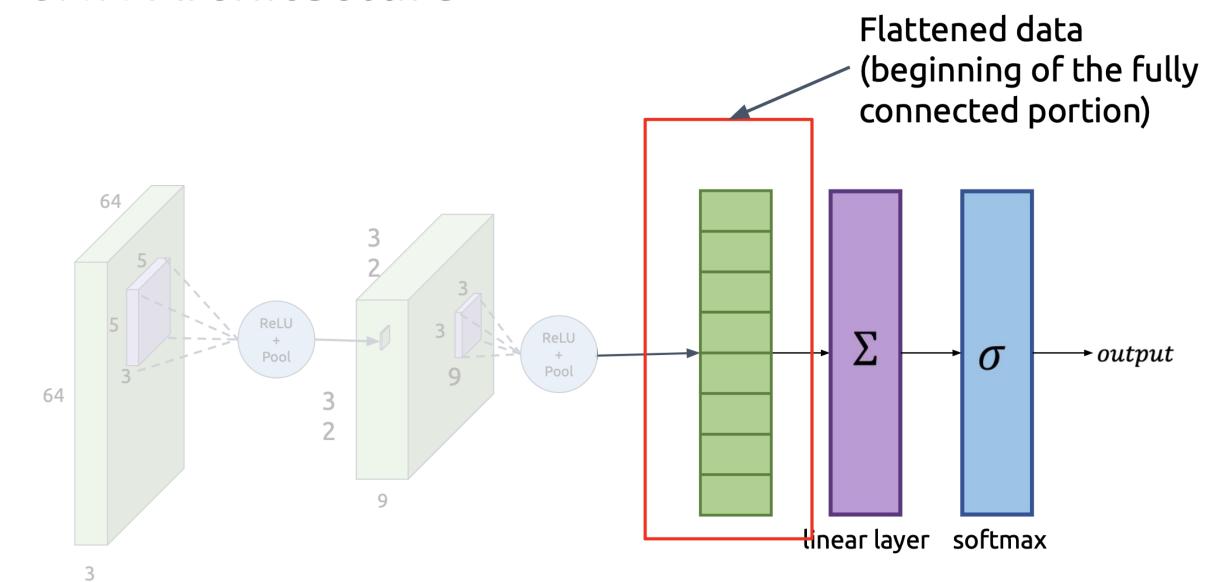
#### Activation after filter passes over image

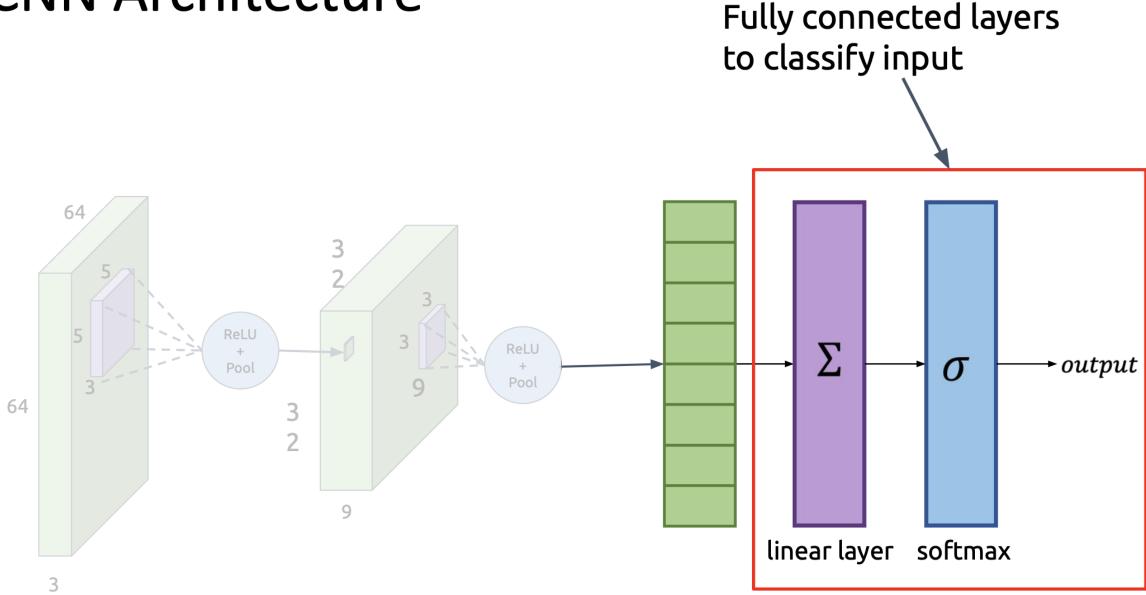


This part learns to perform a specific task (e.g. classification) using those features

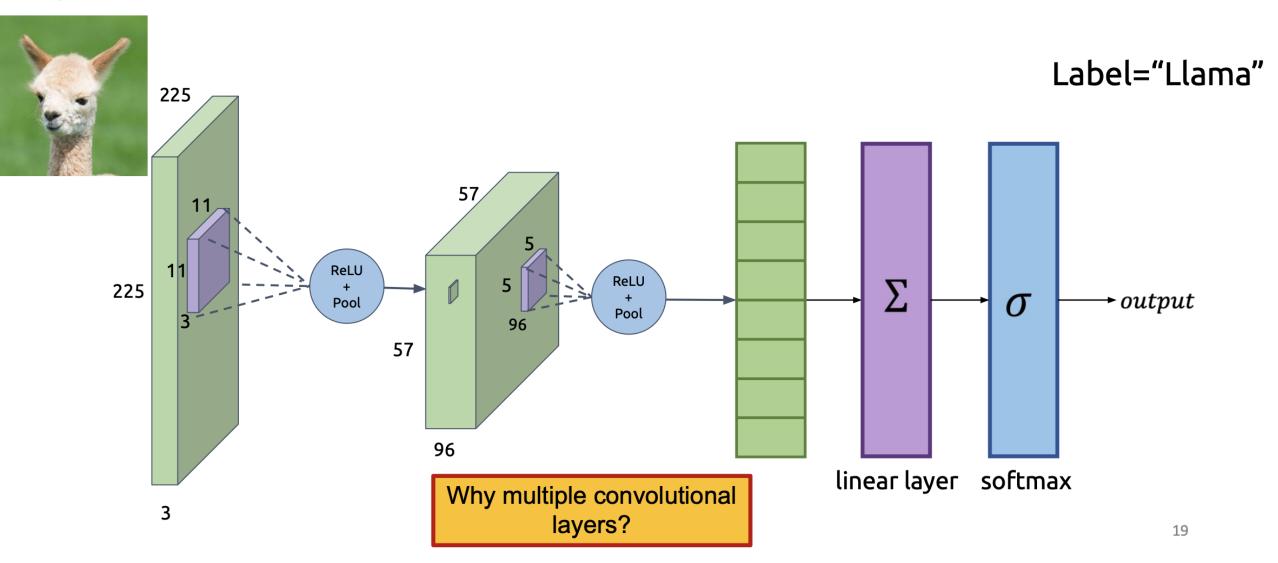


16



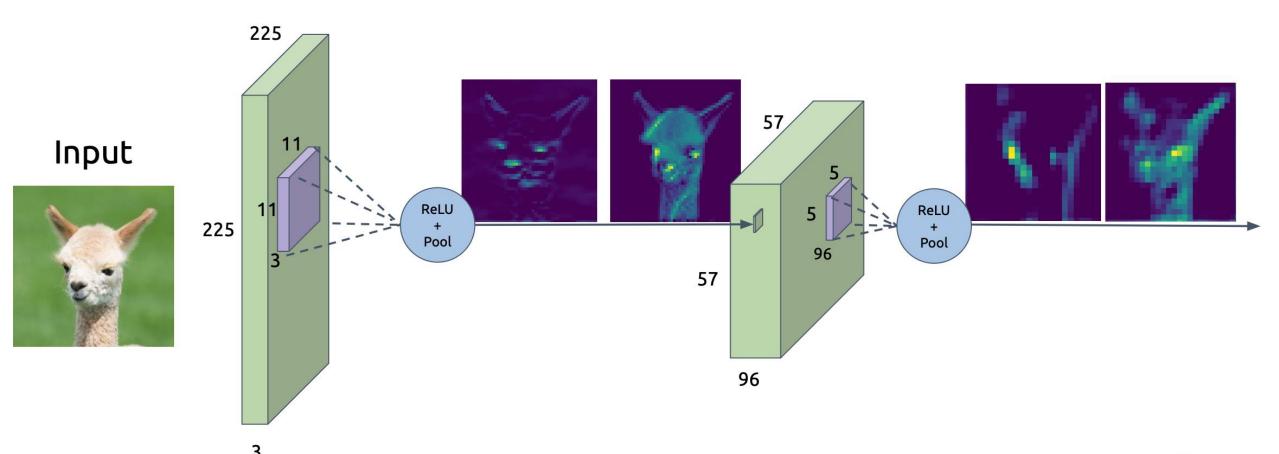


#### Input



# Feature Extraction using multiple convolution layers Hierarchy of features

Sequence of layers detect broader and broader features





# **Example: Network Dissection**

http://netdissect.csail.mit.edu/







"Eye Detector"

Layer 4 active regions





"Eyes and Nose Detector"

#### Layer 5 active regions





"Dog Face Detector"

## **ILSVRC 2012**

(ImageNet Large Scale Visual Recognition Challenge)

The classification task on ImageNet:

For each image, assign 5 labels in order of decreasing confidence. one of these labels matches the ground truth

Success if



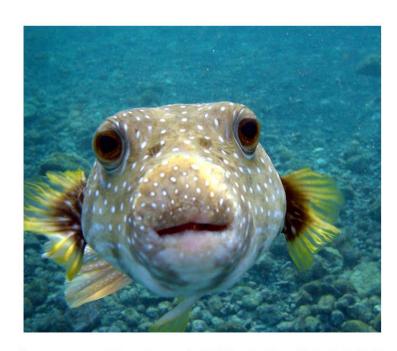
#### **Predictions:**

- 1. Carpet
- 2. Zebra
- 3. Llama
- 4. Flower
- 5. Horse



## **ILSVRC 2012**

### Percentage that model fails to classify is known as *Top 5 Error Rate*



https://commons.wikimedia.org/wiki/File:Puffer\_Fish\_DSC01257.JPG

#### **Predictions:**

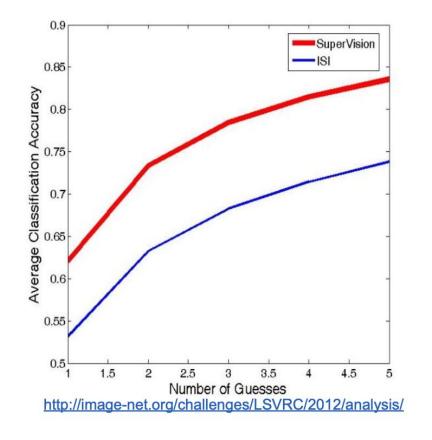
- 1. Sponge
- 2. Person
- 3. Llama
- 4. Flower
- 5. Boat



# AlexNet: Why CNNs Are a Big Deal

Major performance boost on ImageNet at ILSCRV 2012

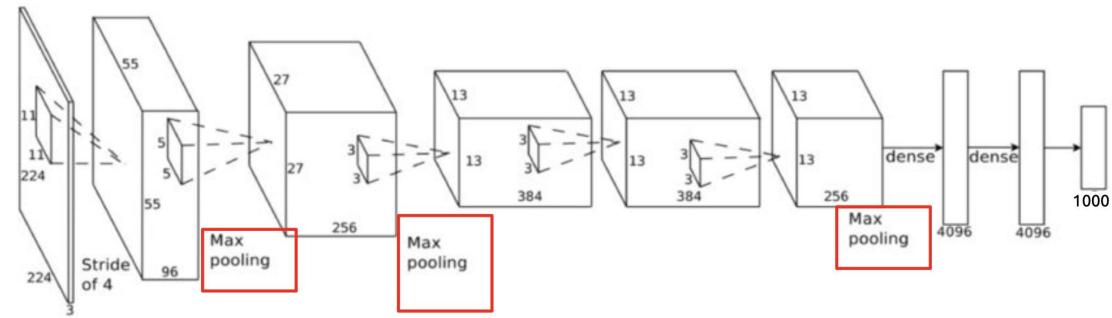
Top 5 error rate of 15.3% compared to 26.2% achieved by 2nd place



Note: SuperVision is the name of Alex's team

## AlexNet

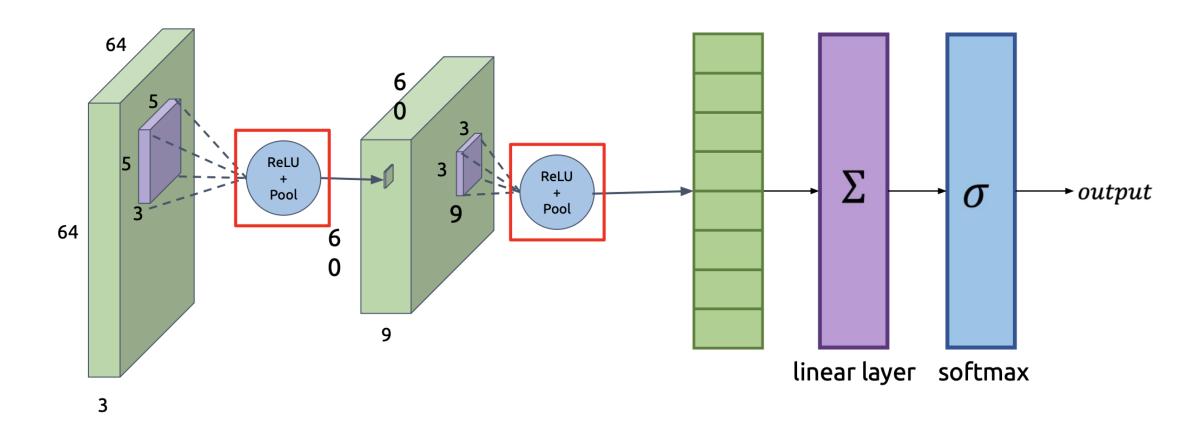
- 60 million parameters
- 5 Convolutional Layers
- 3 Fully Connected Layers



[Alex Krizhevsky et al. 2012]

https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

# Pooling



# So...did we achieve our goal of translational invariance?

# What was Translational Invariance again?

- To make a neural net f robust in this same way, it should ideally satisfy **translational invariance**: f(T(x)) = f(x), where
  - x is the input image
  - T is a translation (i.e. a horizonal and/or vertical shift)

## **Are CNNs Translation Invariant?**

- Convolution is translation equivariant
  - A translated input results in an output translated by the same amount

• 
$$f(T(I)) = T(f(I))$$

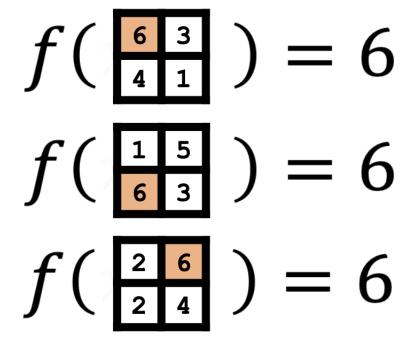
• 
$$(T(I) \otimes K)(x, y) = T(I \otimes K)(x, y)$$

$$f(|||) = |||$$

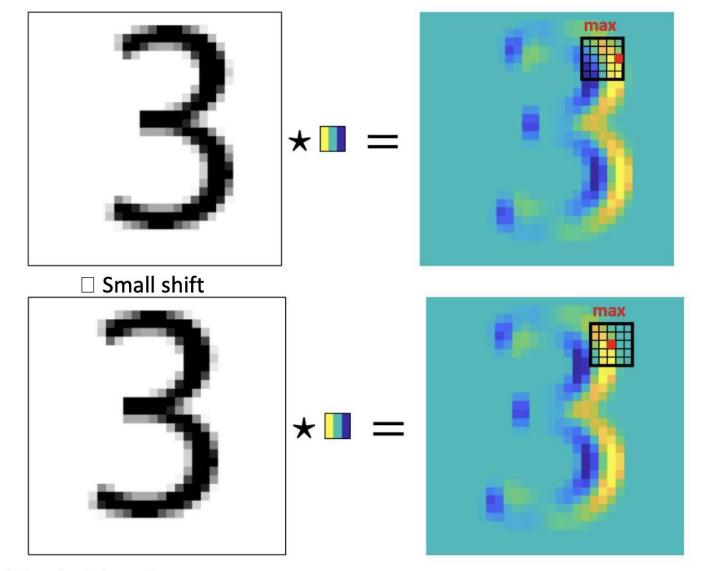
$$||T||$$

$$f(|||) = ||$$

- Max pooling is intended to give invariance to small translations
  - The highest activation pixel can shift around within the pooling window, and the output does not change



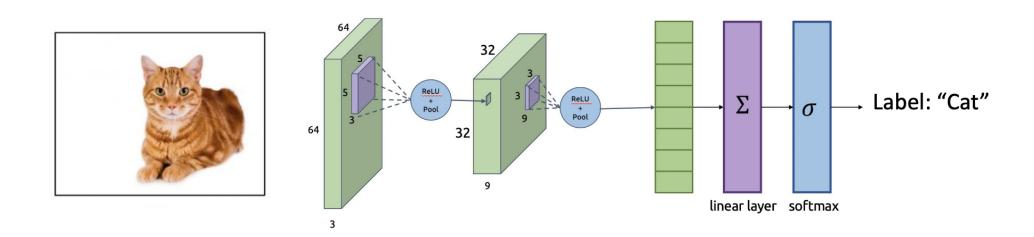
## So how does it all come together?



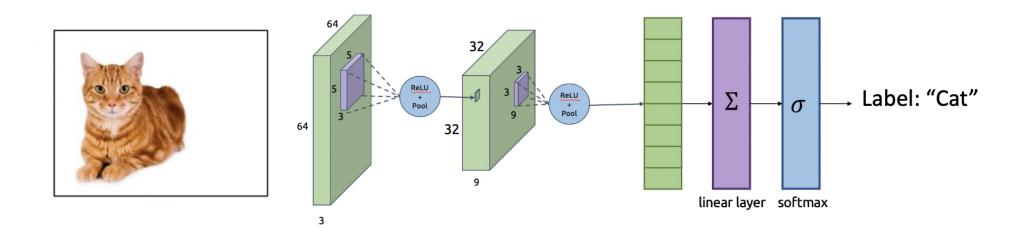
Convolution is translation equivariant

Max pooling gives invariance to small translations

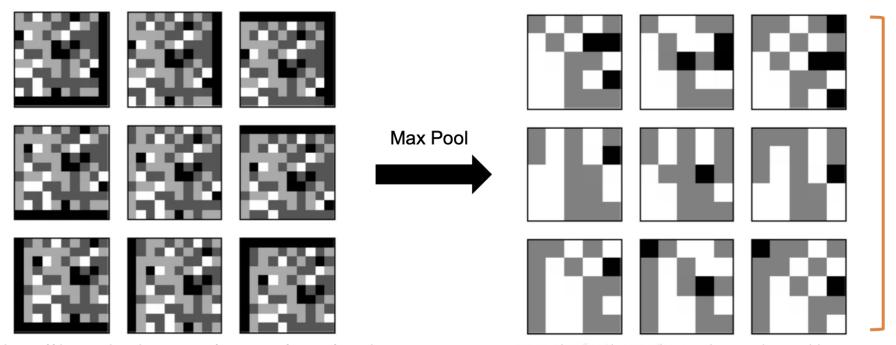
- Answer: CNNs are "sort of" translation invariant
  - Shifting the content of the image around tends not to drastically effect the output classification probabilities...



- Answer: CNNs are "sort of" translation invariant
  - Shifting the content of the image around tends not to drastically effect the output classification probabilities...



- Answer: CNNs are "sort of" translation invariant
  - Shifting the content of the image around tends not to drastically effect the output classification probabilities...
  - ...but they are *not*, strictly speaking, translation invariant



These are **not** all the same!

























Size Invariance



















Size Invariance







Illumination Invariance



















Size Invariance







Illumination Invariance







- All are desirable properties!
- How do CNNs fare?
  - Max pooling gives some amount of size and translational invariance
  - But in general, CNNs do not fare well with large changes in lighting or scale.
- Consequences of not having these invariances?
  - Require lots of training data
  - Have to show network many examples of lighting changes, scale changes, etc.

#### Rotation/Viewpoint Invariance













Size Invariance







Illumination Invariance







- All are desirable properties!
- How do CNNs fare?
  - Max pooling gives some amount of size and translational invariance
  - But in general, CNNs do not fare well with large changes in lighting or scale.
- Consequences of not having these invariances?
  - Require lots of training data
  - Have to show network many examples of lighting changes, scale changes, etc.

Can we address these concerns without collecting additional data?

#### Rotation/Viewpoint Invariance













Size Invariance







Illumination Invariance







Data Augmentation! Use rotated/scaled/shifted images from your dataset to train

- All are desirable properties!
- How do CNNs fare?
  - Max pooling gives some amount of size and translational invariance
  - But in general, CNNs do not fare well with large changes in lighting or scale.
- Consequences of not having these invariances?
  - Require lots of training data
  - Have to show network many examples of lighting changes, scale changes, etc.

Can we address these concerns without collecting additional data?



If this is a cat in our dataset, it is an image with a label (cat)

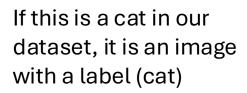




If this is a cat in our dataset, it is an image with a label (cat)

This is also a cat







This is also a cat

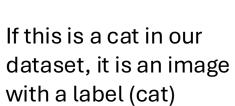


This is also a cat



This is also a cat







This is also a cat

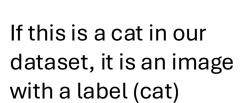


This is also a cat



This is also a cat





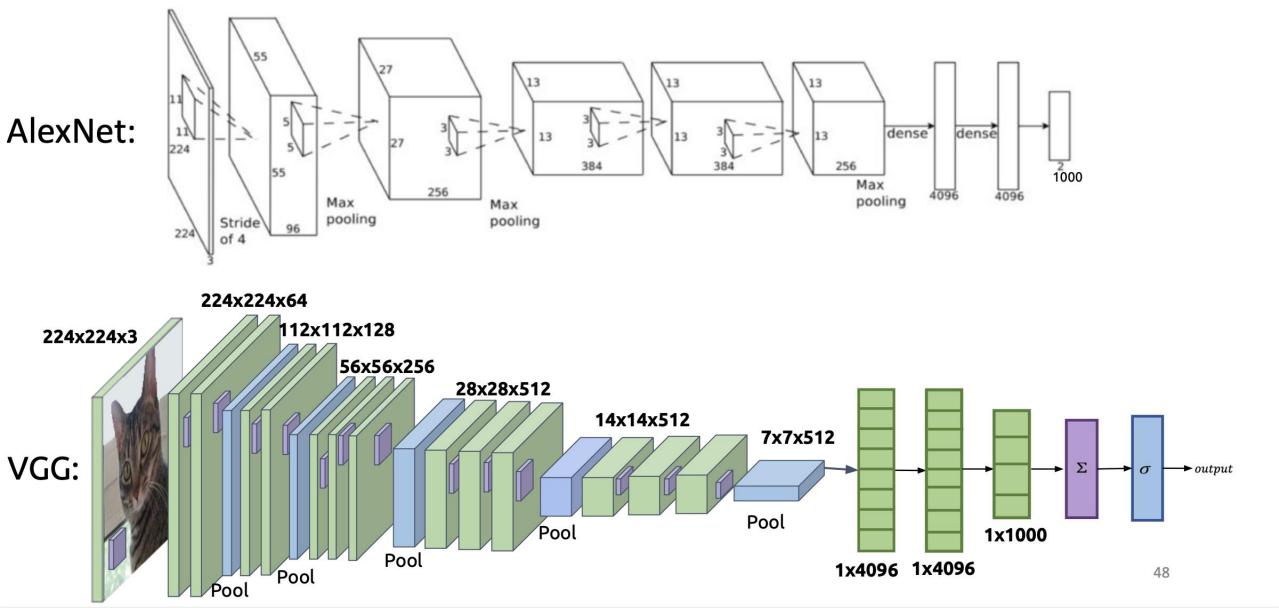


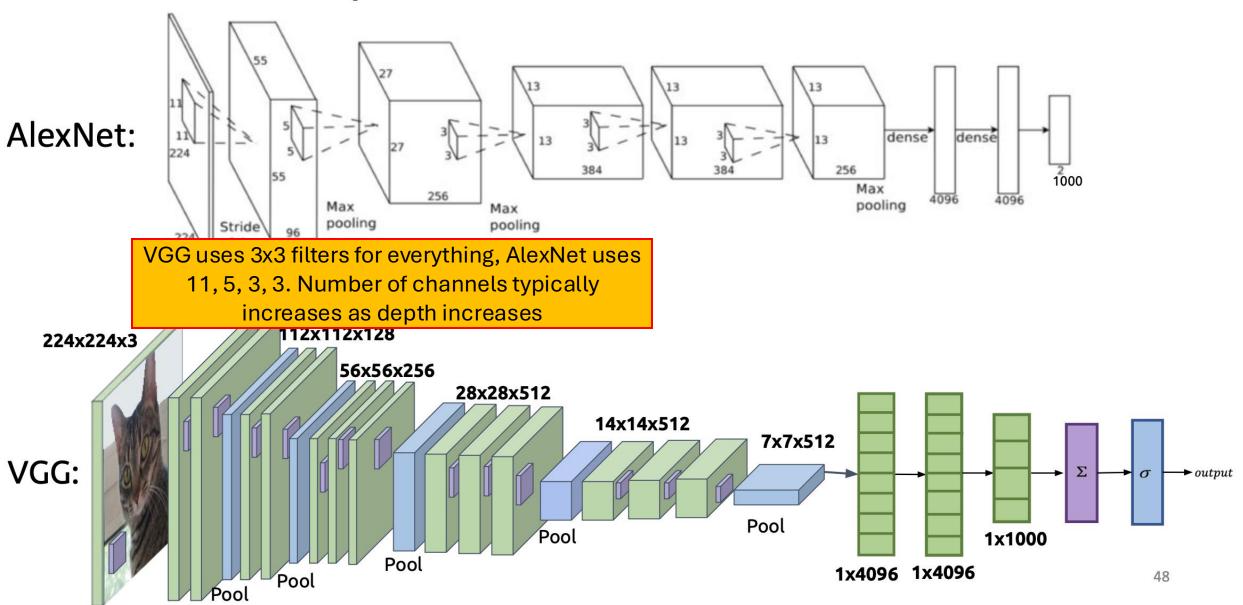
This is also a cat

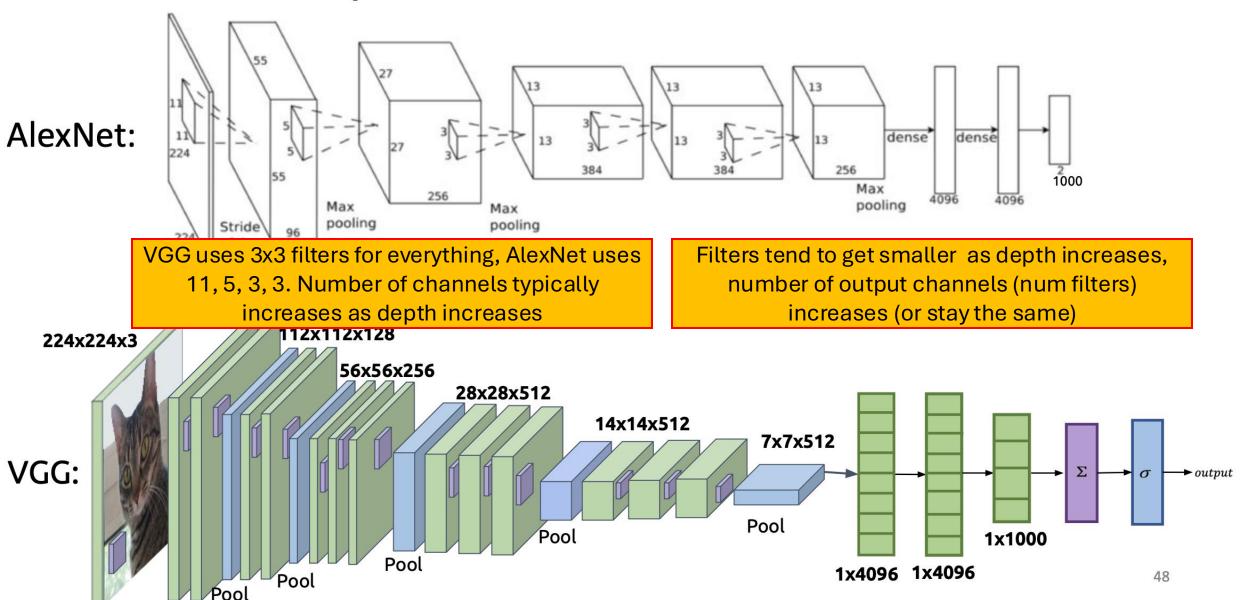


This is also a cat





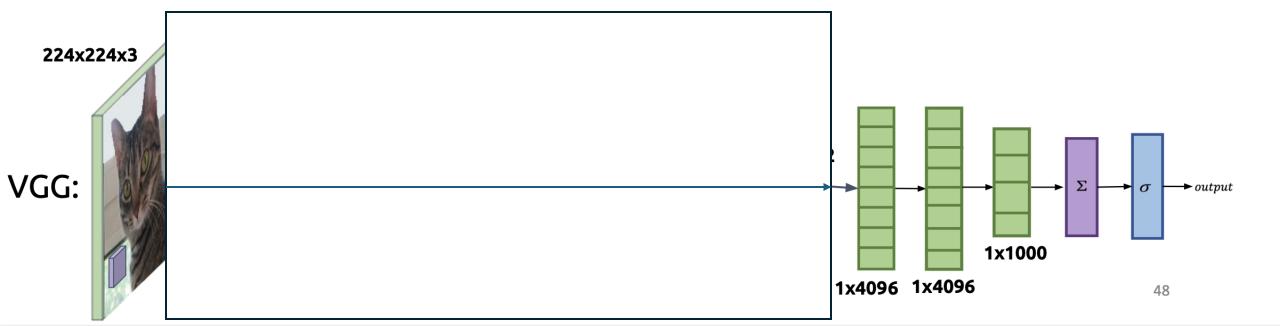




### What if we didn't use a convolution?

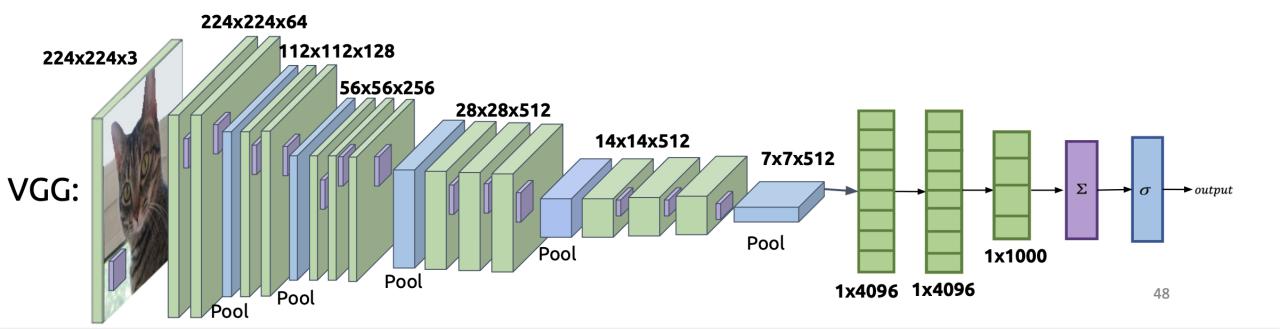
How many weights would there be if we have an input image of 224x224x3 and want to go to a hidden layer size of 4096?

What is the size of the Jacobian  $\frac{\partial z}{\partial W}$ ?



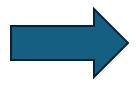
#### With Convolutions

VGG uses 3x3 convolutions, how many weights are in the first filter bank to go from 224x224x3 to 224x224x64?



## Convolutions and Depth

Convolutions are much faster to run than a linear layer on the same size input

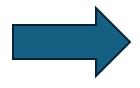


We can add more layers to CNNs than MLPs with the same inference time

Theory: Having more layers gives better performance with the same number of total weights (with lots of caveats)

## Convolutions and Depth

Convolutions are much faster to run than a linear layer on the same size input

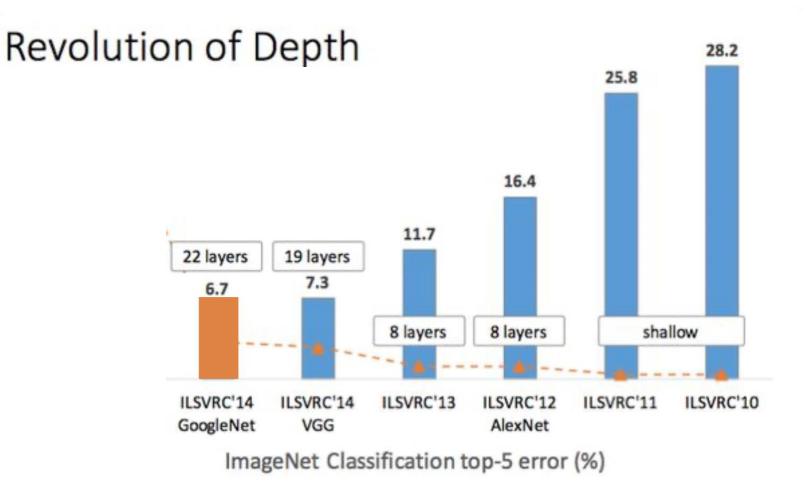


We can add more layers to CNNs than MLPs with the same inference time

Theory: Having more layers gives better performance with the same number of total weights (with lots of caveats)

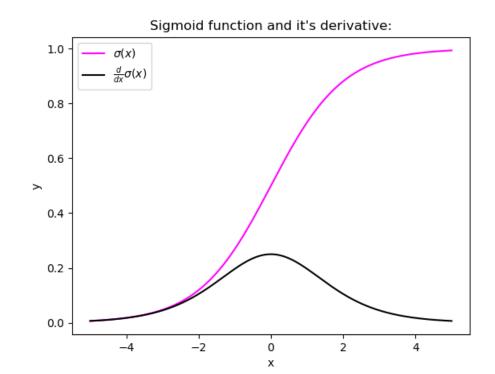
But we start to run into other issues as the depth of our neural networks increase...

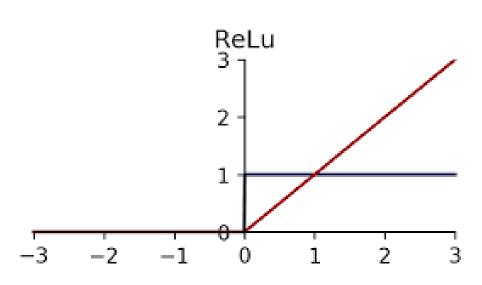
# What's the biggest limitation in increasing depth?

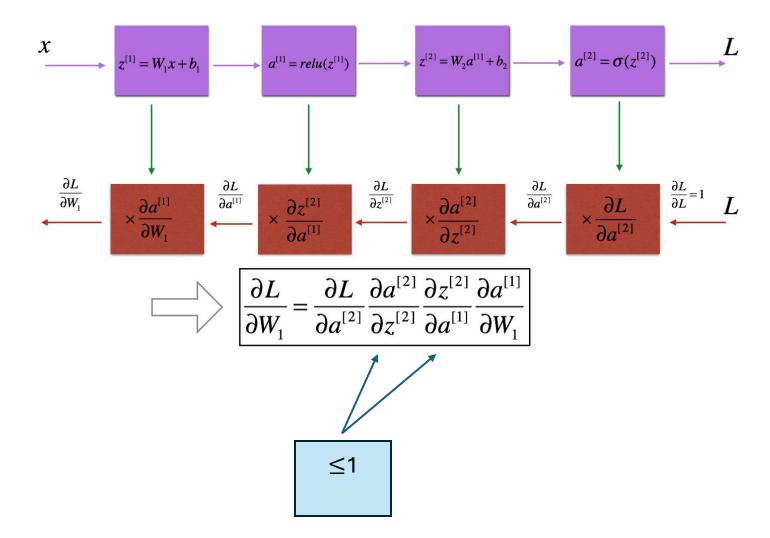


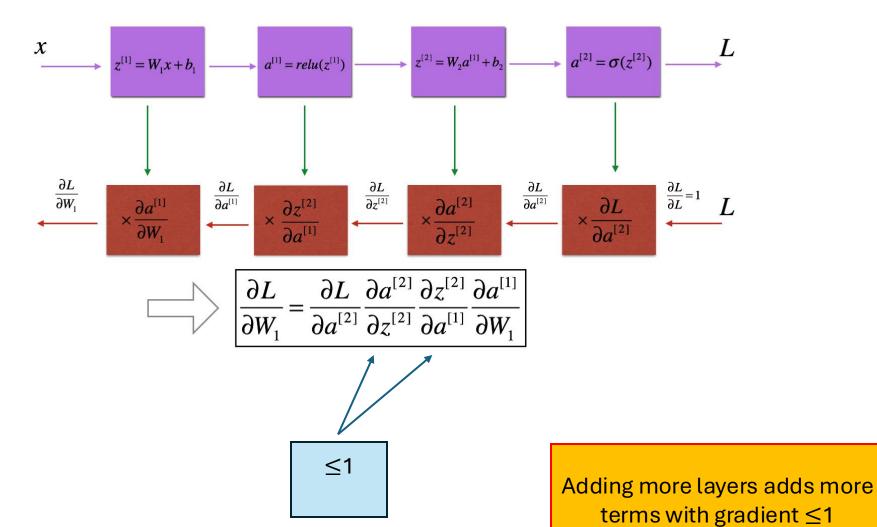
#### The Return of Gradients

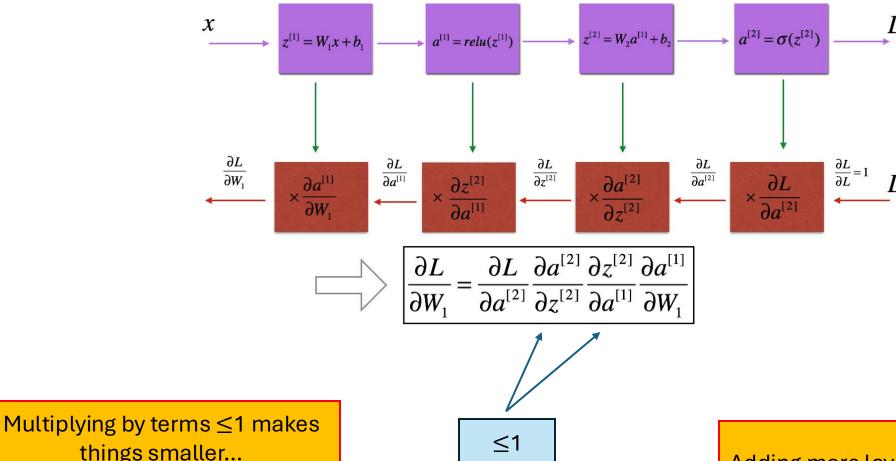
Common activation functions typically have a derivative smaller than 1 (or at least not more than 1)











things smaller...
Gradients earlier in the network tend to "Vanish"

Adding more layers adds more terms with gradient ≤1

# Could we fix it by making everything "steeper"

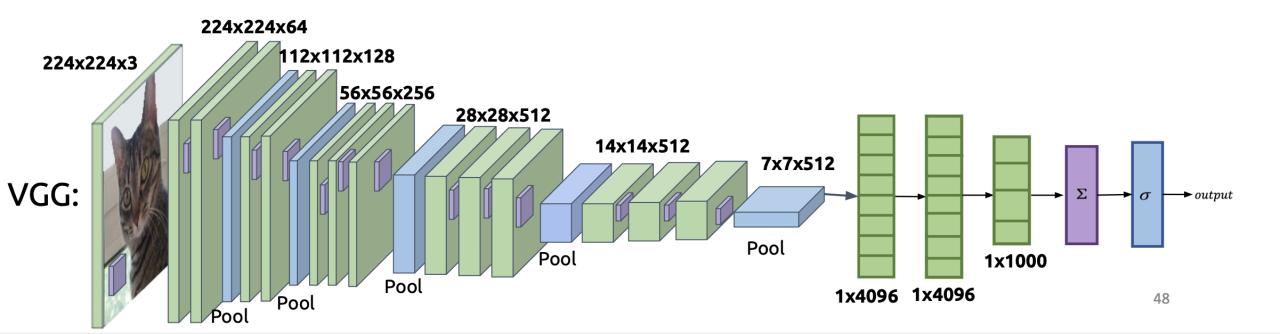
- Vanishing gradients are caused by the repeated multiplication of numbers smaller than 1
- If we make those numbers larger than 1, we have a separate problem...

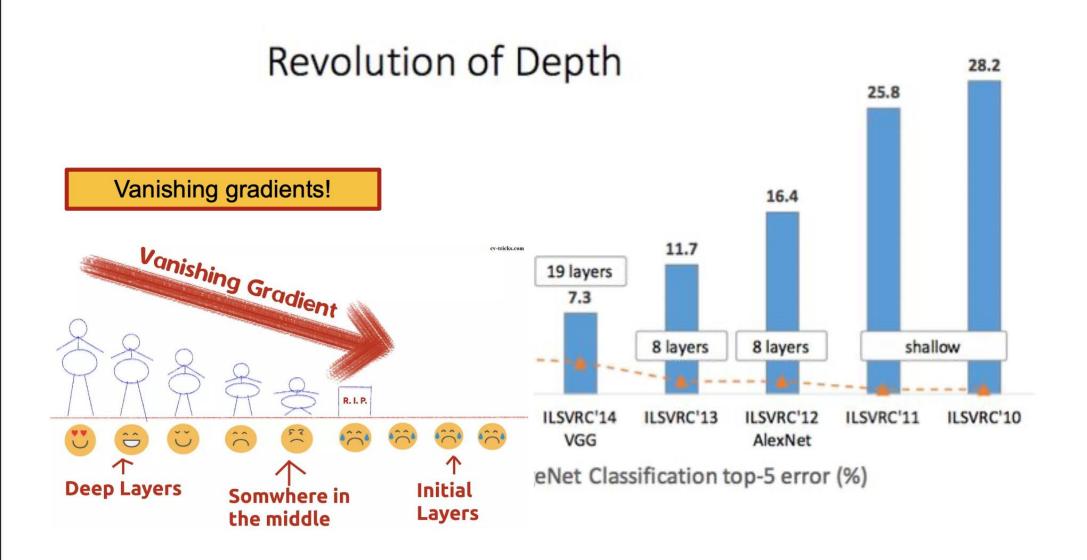
# Could we fix it by making everything "steeper"

- Vanishing gradients are caused by the repeated multiplication of numbers smaller than 1
- If we make those numbers larger than 1, we have a separate problem...

**Exploding Gradients** 

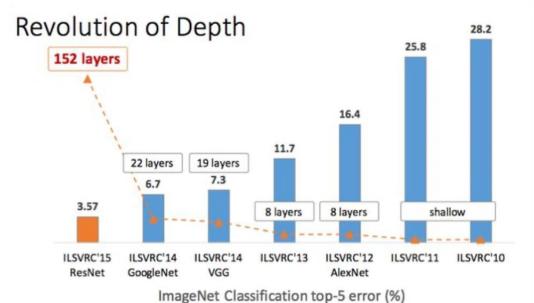
If you could make one change to a weight to have the biggest change on output, which weight would you pick?



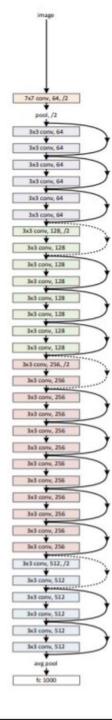


#### **ResNet:**

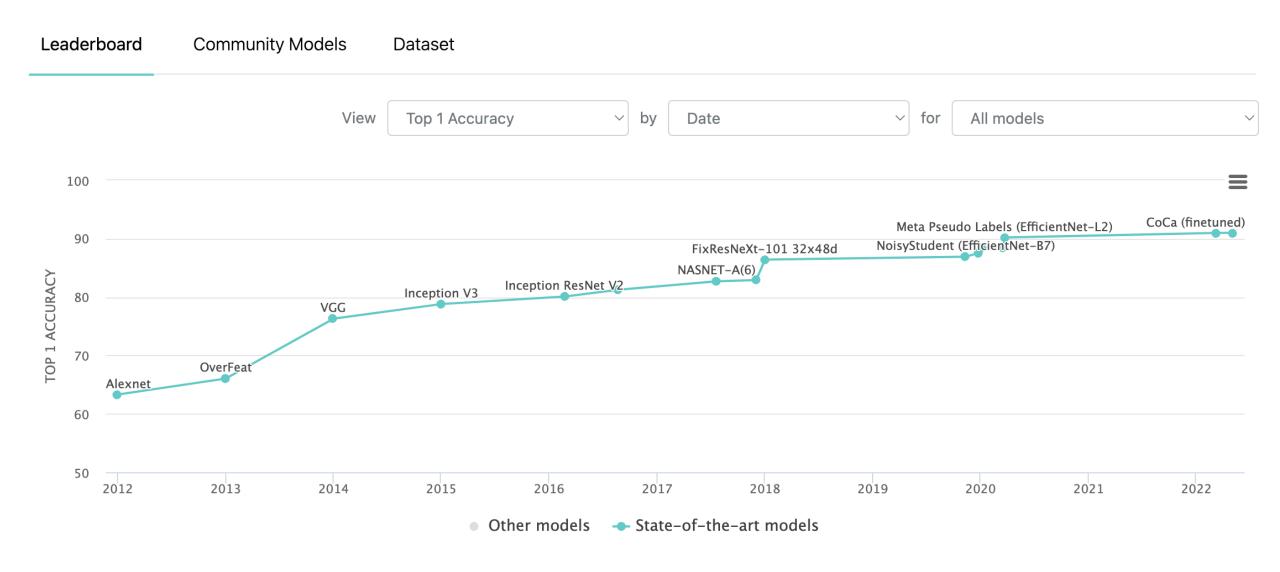
Lots of layers, tons of learnable parameters Avoids Vanishing Gradient problem but how?



K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.

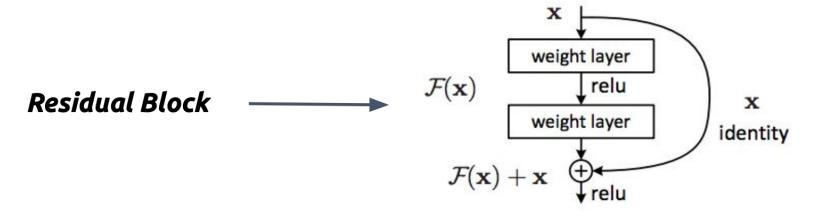


## Image Classification on ImageNet

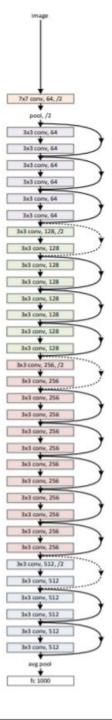


#### ResNet:

Lots of layers, tons of learnable parameters Avoids Vanishing Gradient problem

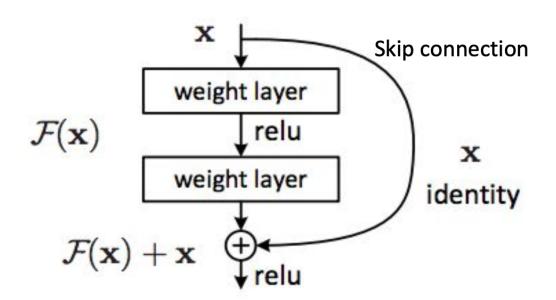


K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.



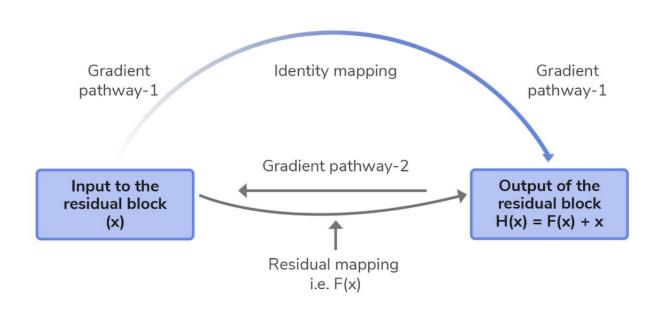
### Residual Blocks

- In very deep nets, each layer often needs to learn just a small transformation of the preceding layer (identity + change)
- Idea: explicitly design the network such that the output of each layer is the identity
   + some deviation from it
  - Deviation is known as a residual



## Residual Blocks

- In very deep nets, each layer often needs to learn just a small transformation of the preceding layer (identity + change)
- Idea: explicitly design the network such that the output of each layer is the identi + some deviation from it
  - Deviation is known as a residual
- Allows gradient to flow through two pathways
- Significantly stabilizes training of very deep networks



## Tensorflow

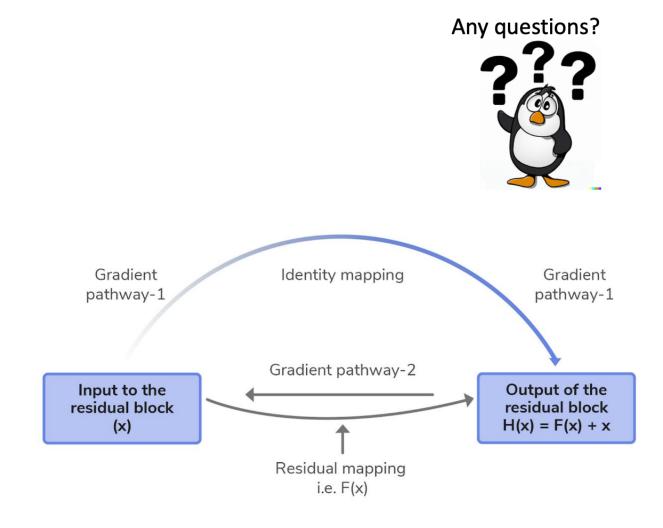
Option #1: Residual Block tfm.vision.layers.ResidualBlock(filters, strides) Option #2:

```
# Residual Block
def ResBlock(inputs):
    x = layers.Conv2D(64, 3, padding="same", activation="relu")(inputs)
    x = layers.Conv2D(64, 3, padding="same")(x)
    x = layers.Add()([inputs, x])
    return x
Original Input Intermediate Output
```

https://keras.io/examples/vision/edsr/

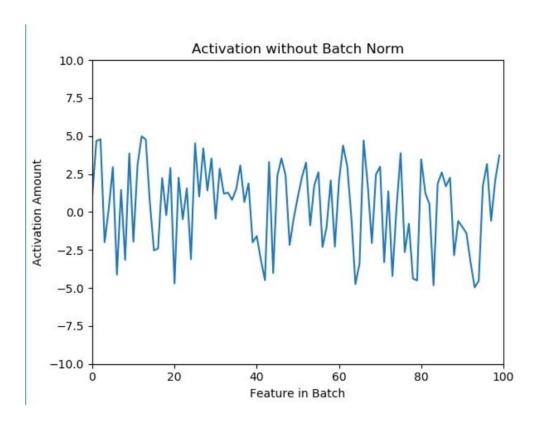
### Residual Blocks

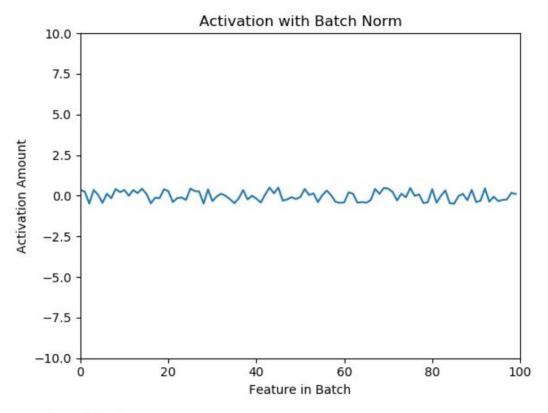
- In very deep nets, each layer often needs to learn just a small transformation of the preceding layer (identity + change)
- Idea: explicitly design the network such that the output of each layer is the identi + some deviation from it
  - Deviation is known as a residual
- Allows gradient to flow through two pathways
- Significantly stabilizes training of very deep networks



## Batch Normalization (stabilizing training)

Idea: normalize the activations for each feature at each layer



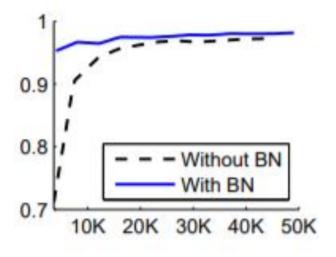


Why might we want to do this?

## **Batch Normalization: Motivation**

More stable inputs = faster training

MNIST test accuracy vs number of training steps



https://arxiv.org/pdf/1502.03167.pdf

## Batch Normalization: Implementation

For each feature x, Start by calculating the batch mean and standard deviation for each feature:

$$\mu_{batch} = \frac{\sum_{i=0}^{batch\_size} x_i}{batch\_size}$$

$$\sigma_{batch} = \sqrt{\frac{\sum_{i=0}^{batch\_size} (x_i - \mu_{batch})^2}{batch_{size}}}$$

## Batch Normalization: Implementation

Normalize by subtracting feature x's batch mean, then divide by batch standard deviation.

$$x' = \frac{x - \mu_{batch}}{\sigma_{batch}}$$

Feature x now has mean 0 and variance 1 along the batch

## **Batch Normalization in Tensorflow**

tf.keras.layers.BatchNormalization(input)

Documentation: <a href="https://www.tensorflow.org/versions/r2.0/api\_docs/python/tf/keras/layers/BatchNormalization">https://www.tensorflow.org/versions/r2.0/api\_docs/python/tf/keras/layers/BatchNormalization</a>

## Motivation of BatchNorm

- Reduce "internal co-variate shift"
- Neural networks are trained on a certain distribution of data and are expected to be tested on the same distribution
- If we were to scale the colors of an image significantly at test time, we wouldn't expect a neural network to do well
- The same can be said for our intermediate layers
  - They expect a certain distribution of inputs, if that changes significantly from example to example, it will be hard to learn
- (Most commonly cited reason for using BatchNorm)

## The only issue is that controlling internal covariate shift does not matter that much...

### **How Does Batch Normalization Help Optimization?**

Shibani Santurkar\* MIT shibani@mit.edu

.edu tsip

Dimitris Tsipras\* An MIT tsipras@mit.edu aily

Andrew Ilyas\*
MIT
ailyas@mit.edu

Aleksander Mądry MIT madry@mit.edu

#### **Abstract**

Batch Normalization (BatchNorm) is a widely adopted technique that enables faster and more stable training of deep neural networks (DNNs). Despite its pervasiveness, the exact reasons for BatchNorm's effectiveness are still poorly understood. The popular belief is that this effectiveness stems from controlling the change of the layers' input distributions during training to reduce the so-called "internal covariate shift". In this work, we demonstrate that such distributional stability of layer inputs has little to do with the success of BatchNorm. Instead, we uncover a more fundamental impact of BatchNorm on the training process: it makes the optimization landscape significantly smoother. This smoothness induces a more predictive and stable behavior of the gradients, allowing for faster training.

# BatchNorm makes the loss landscape smoother with fewer local minima, saddle points, and other problematic areas for gradient descent

### **How Does Batch Normalization Help Optimization?**

Shibani Santurkar\* MIT shibani@mit.edu Dimitris Tsipras\* MIT tsipras@mit.edu Andrew Ilyas\*
MIT
ailyas@mit.edu

Aleksander Mądry MIT madry@mit.edu

#### **Abstract**

Batch Normalization (BatchNorm) is a widely adopted technique that enables faster and more stable training of deep neural networks (DNNs). Despite its pervasiveness, the exact reasons for BatchNorm's effectiveness are still poorly understood. The popular belief is that this effectiveness stems from controlling the change of the layers' input distributions during training to reduce the so-called "internal covariate shift". In this work, we demonstrate that such distributional stability of layer inputs has little to do with the success of BatchNorm. Instead, we uncover a more fundamental impact of BatchNorm on the training process: it makes the optimization landscape significantly smoother. This smoothness induces a more predictive and stable behavior of the gradients, allowing for faster training.

## Theory, intuition, and experimental results can all tell you different things

Why does BatchNorm work so well? Intuition: If normalizing input data works so well for training, why not normalize input features to intermediate layers?

Theory/experiments: Makes gradients of loss function "better"

Why do CNNs work so well?
Intuition: Looking for a way to get
"spatial reasoning" or translational
invariance

Theory/experiments: Maybe it's just that using fewer weights lets us go deeper and deep networks learn better (and also they have spatial reasoning)

## Recap

**Translations** 

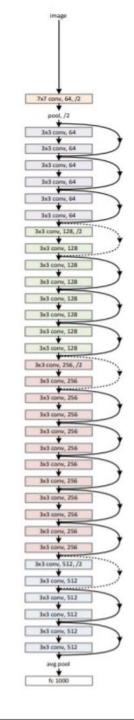
Convolutions and Pooling give us translationally equivariant layers in our network

Small translations in input cause translations in output

Convolutions let us train train deeper networks than MLPs

Adding significantly more depth presents new challenges (vanishing/exploding gradients)

Residual layers and batch norm can help reduce those effects



**CNN Architectures**