

Join us for a
movie with
popcorn and
other
refreshments!

DATA SCIENCE INSTITUTE
DUG
PRESENTS

★ MOVIE NIGHT ★



FEBRUARY 24
Solomon 001



6:00 PM



CSCI 1470/2470
Spring 2023

Ritambhara Singh

February 24, 2023
Friday

Deep Learning



Recap

Convolution

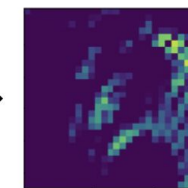
Filters/Kernels and Stride

Learning filters

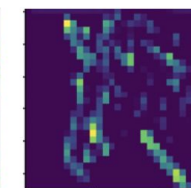
CNNs are partially connected networks



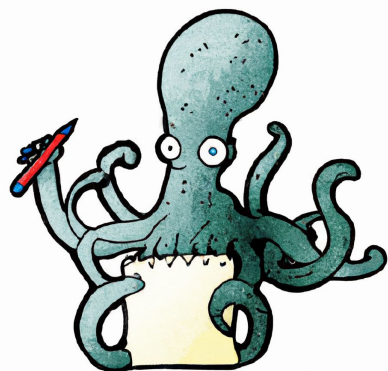
Input image



Output of filter 1



Output of filter 2



Convolution in Tensorflow

Tensorflow conv2d function

Padding

Application to MNIST/CIFAR

```
tf.nn.conv2d(input, filter, strides, padding)
```

Input Image
(4-D Tensor)

Filter/Kernel
(4-D Tensor)

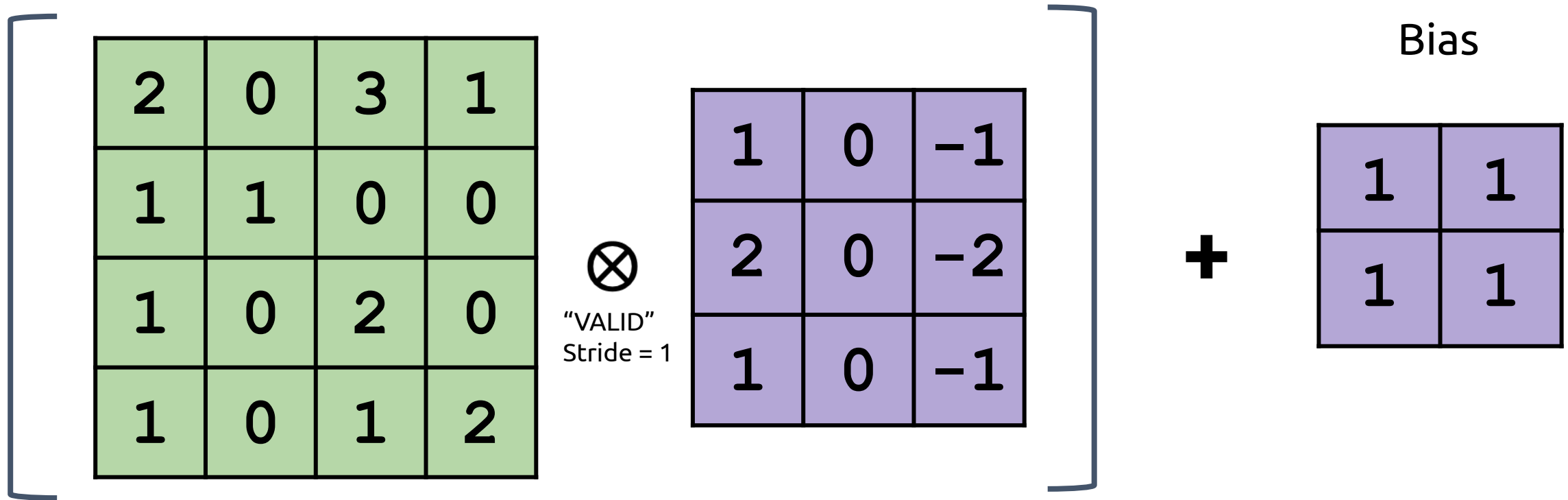
Strides along
each dimension

Type of Padding
(String "Valid" or
"Same")

Today's goal – continue to learn about CNNs

- (1) Convolutional Neural Network (CNN) architecture
- (2) First successful CNN - AlexNet
 - Pooling and translational invariance
- (3) Deeper CNNs!
 - Residual Blocks
 - Batch normalization

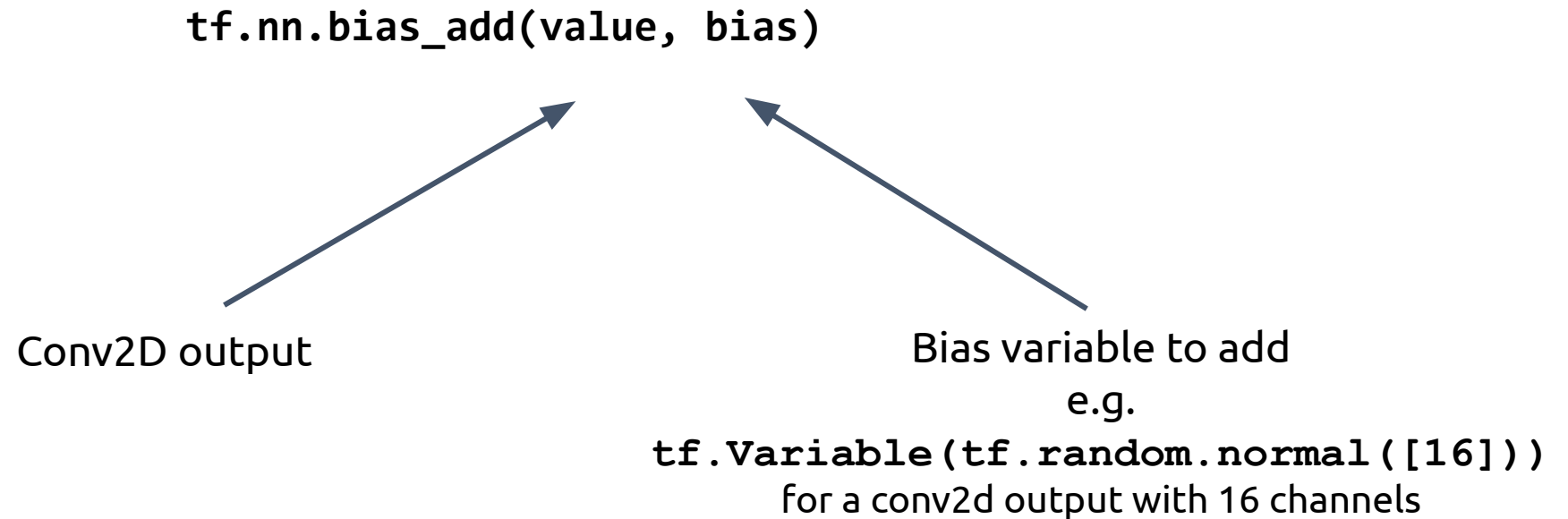
Bias Term in Convolution Layers



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:

```
tf.keras.layers.Conv2D(filters, kernel_sz, strides, padding, use_bias = True)
```

Number of filters

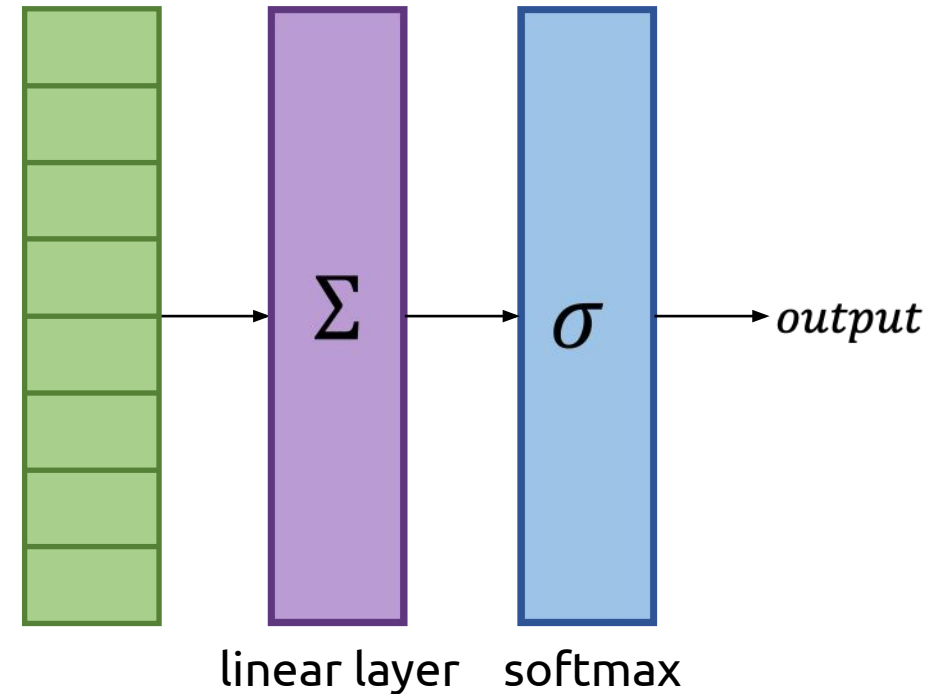
Filter Size

Strides along
each dimension

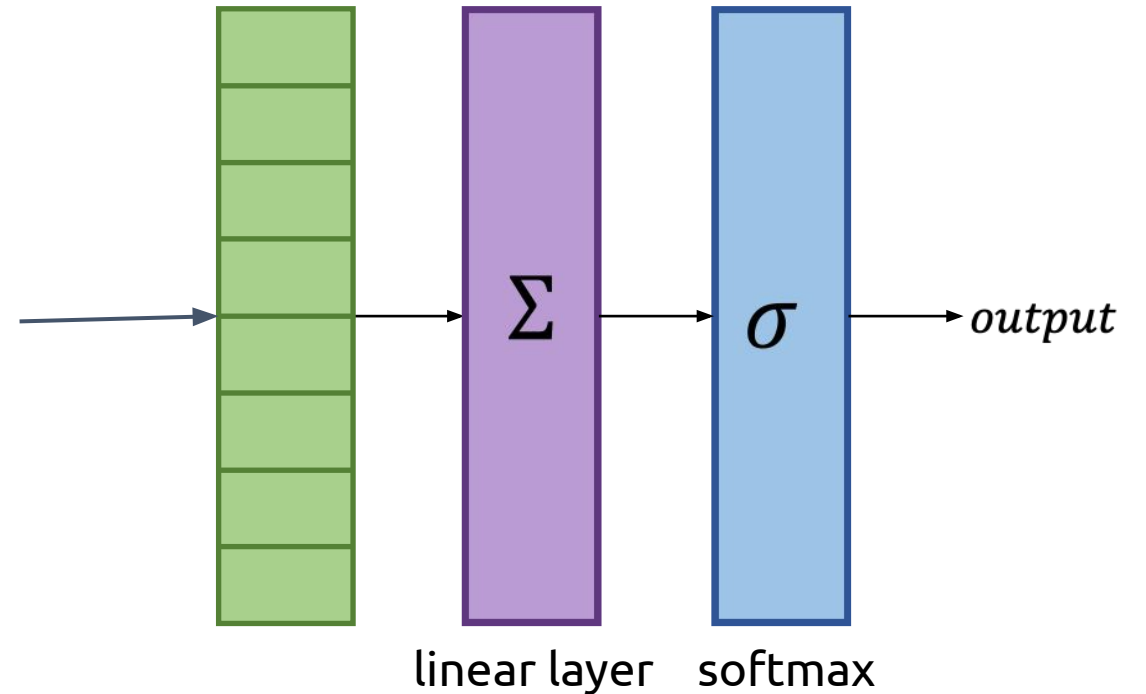
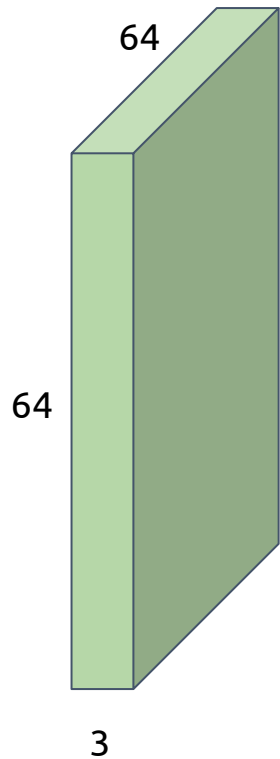
Type of Padding
(VALID or SAME)

Full documentation here: https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Conv2D

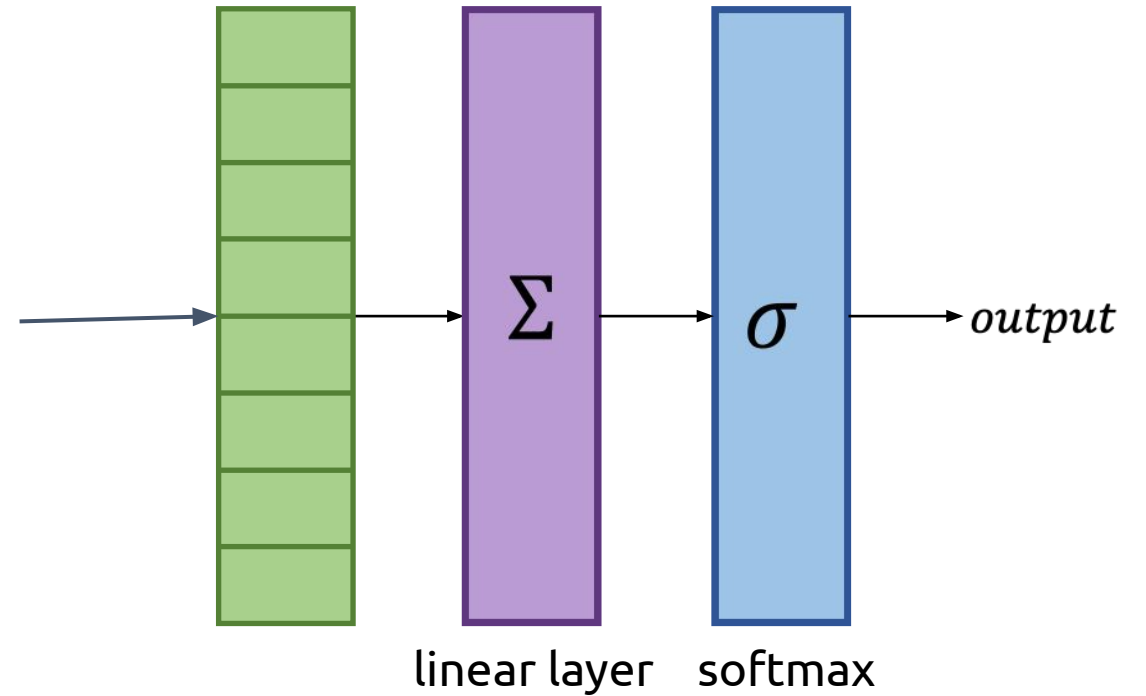
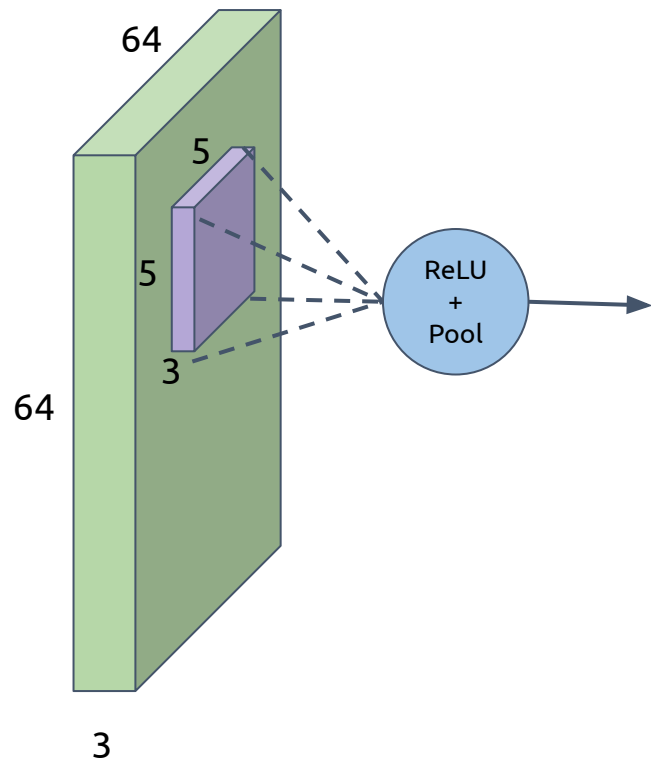
Our neural network so far



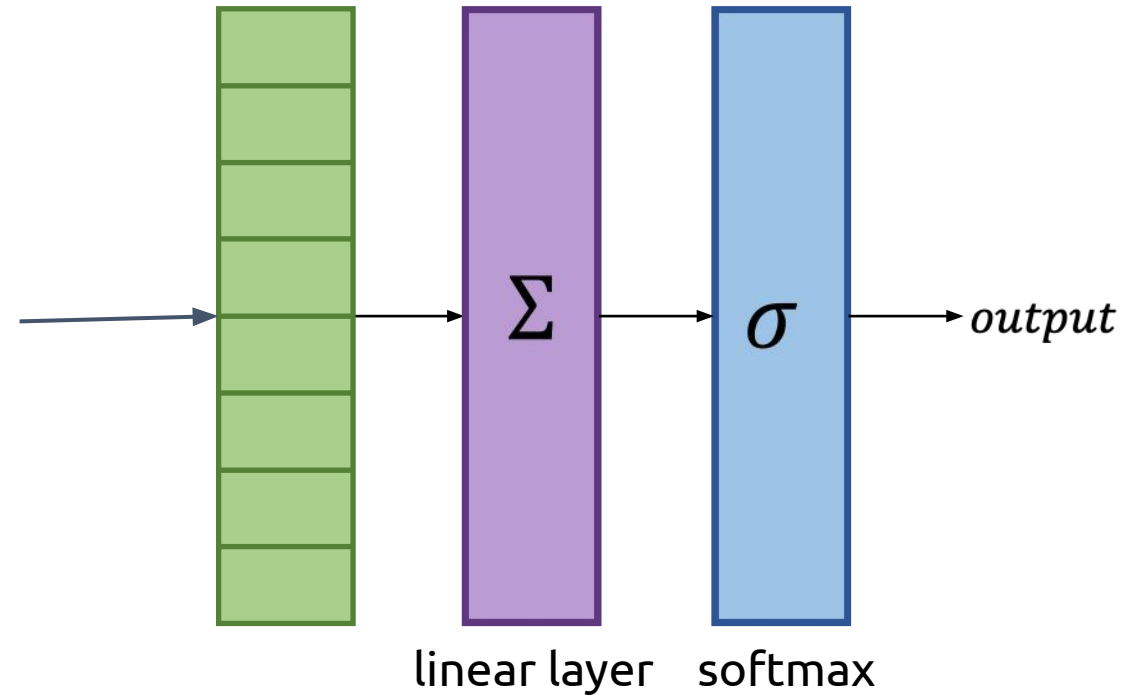
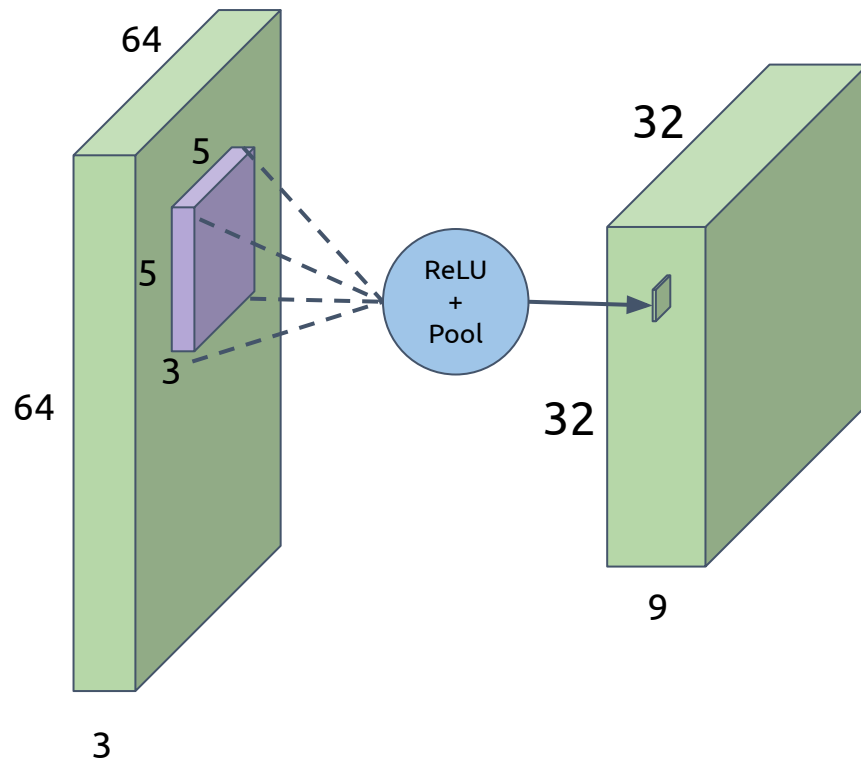
Convolutional Neural Network Architecture



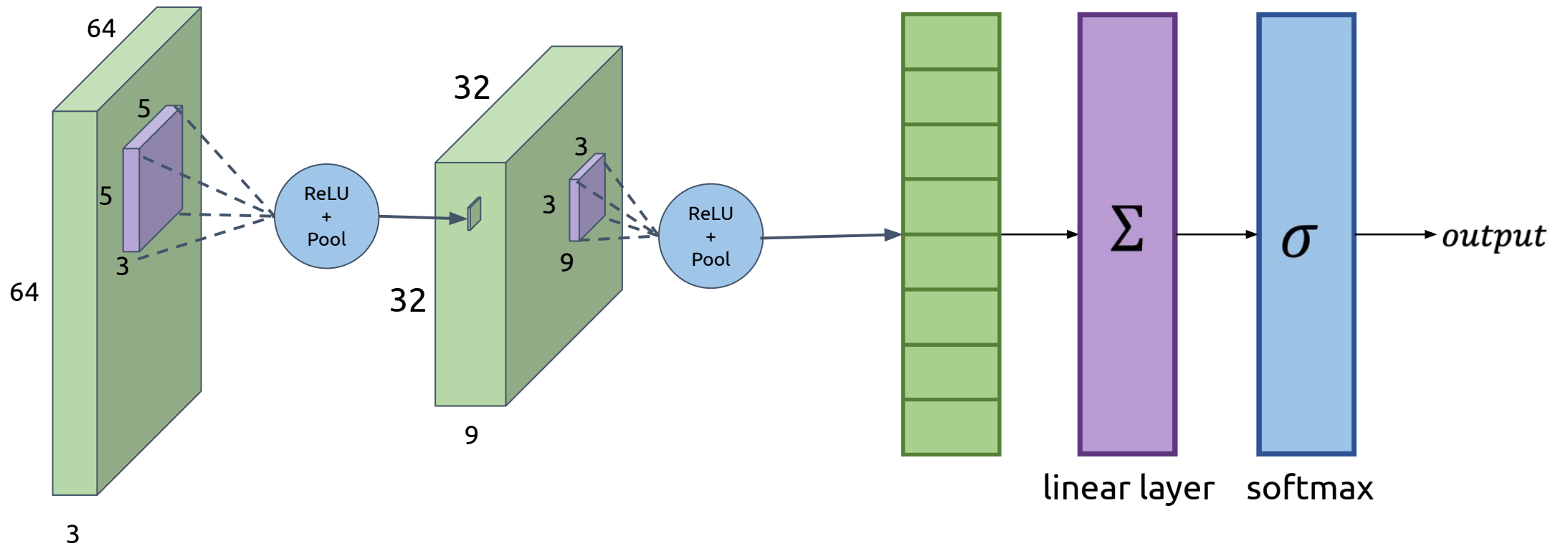
CNN Architecture



CNN Architecture

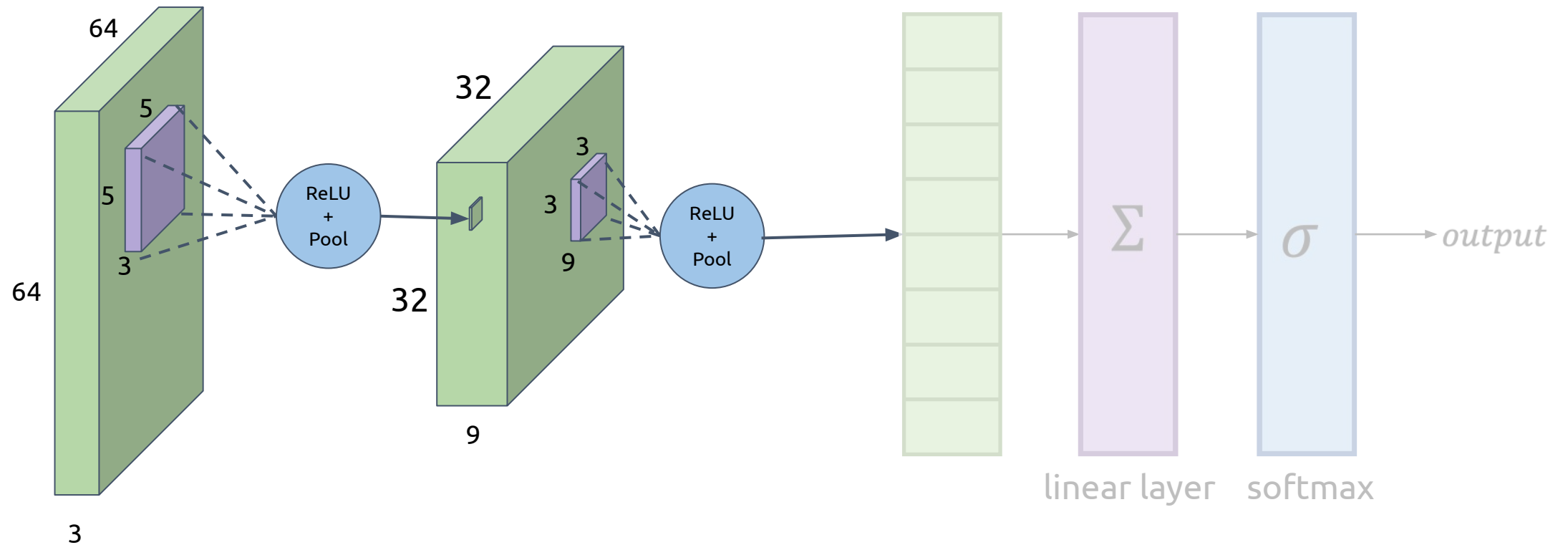


CNN Architecture

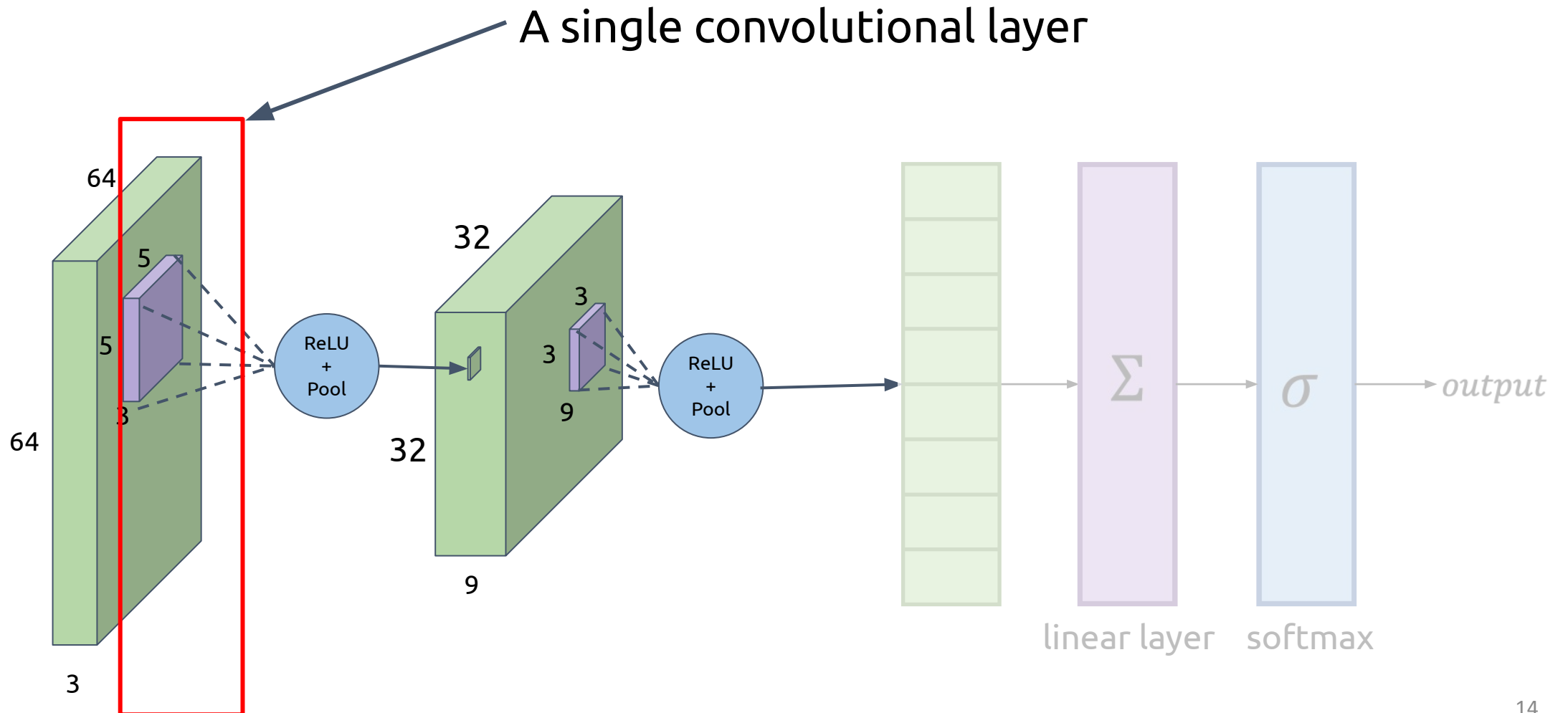


CNN Architecture

This part learns to extract **features** from the image

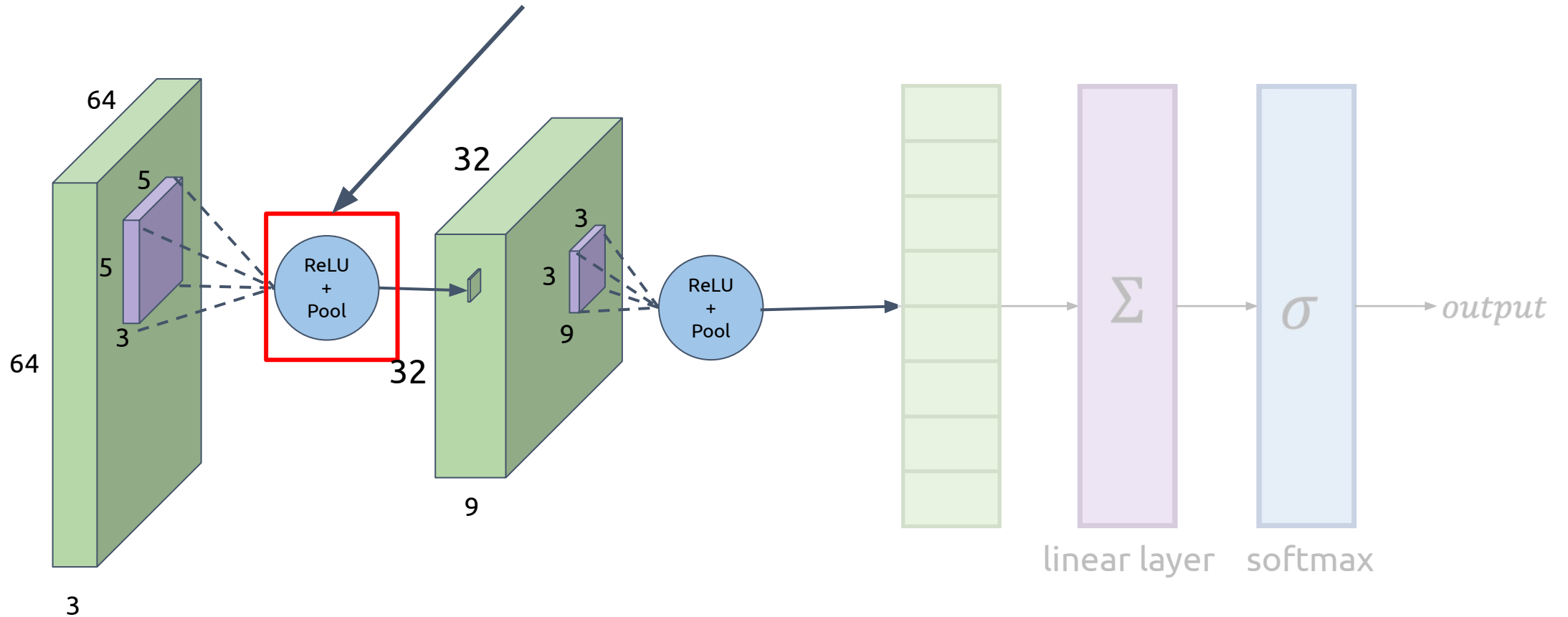


CNN Architecture



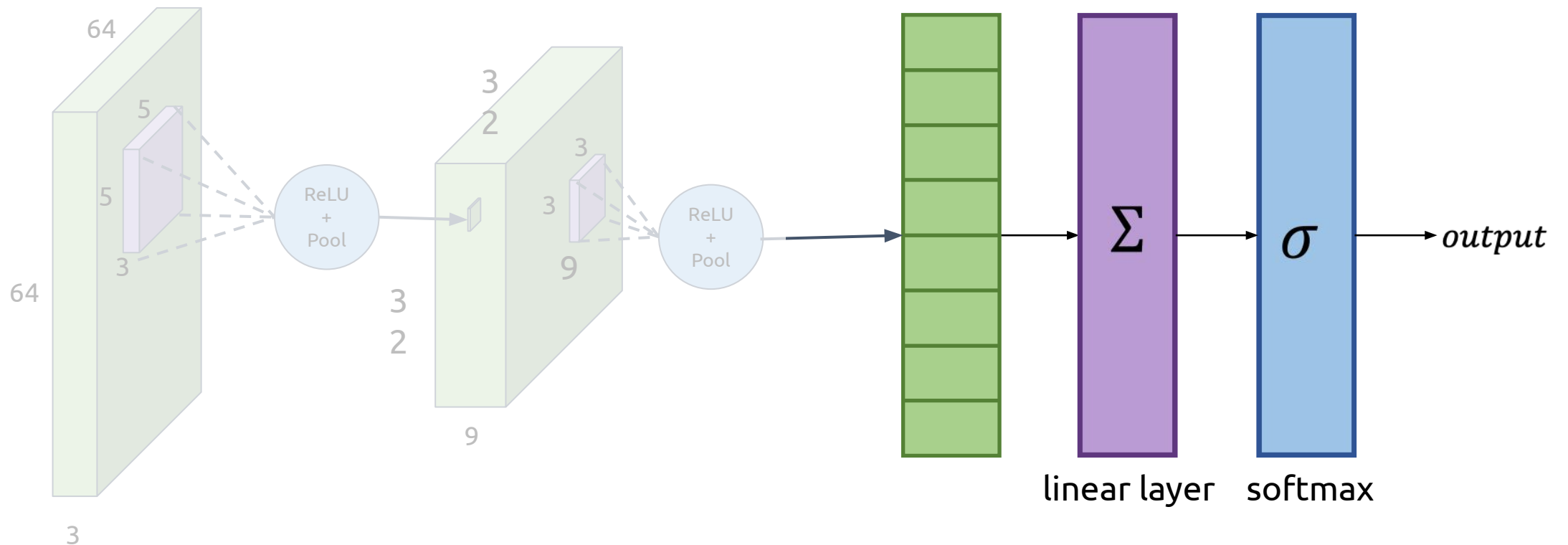
CNN Architecture

Activation after filter passes over image

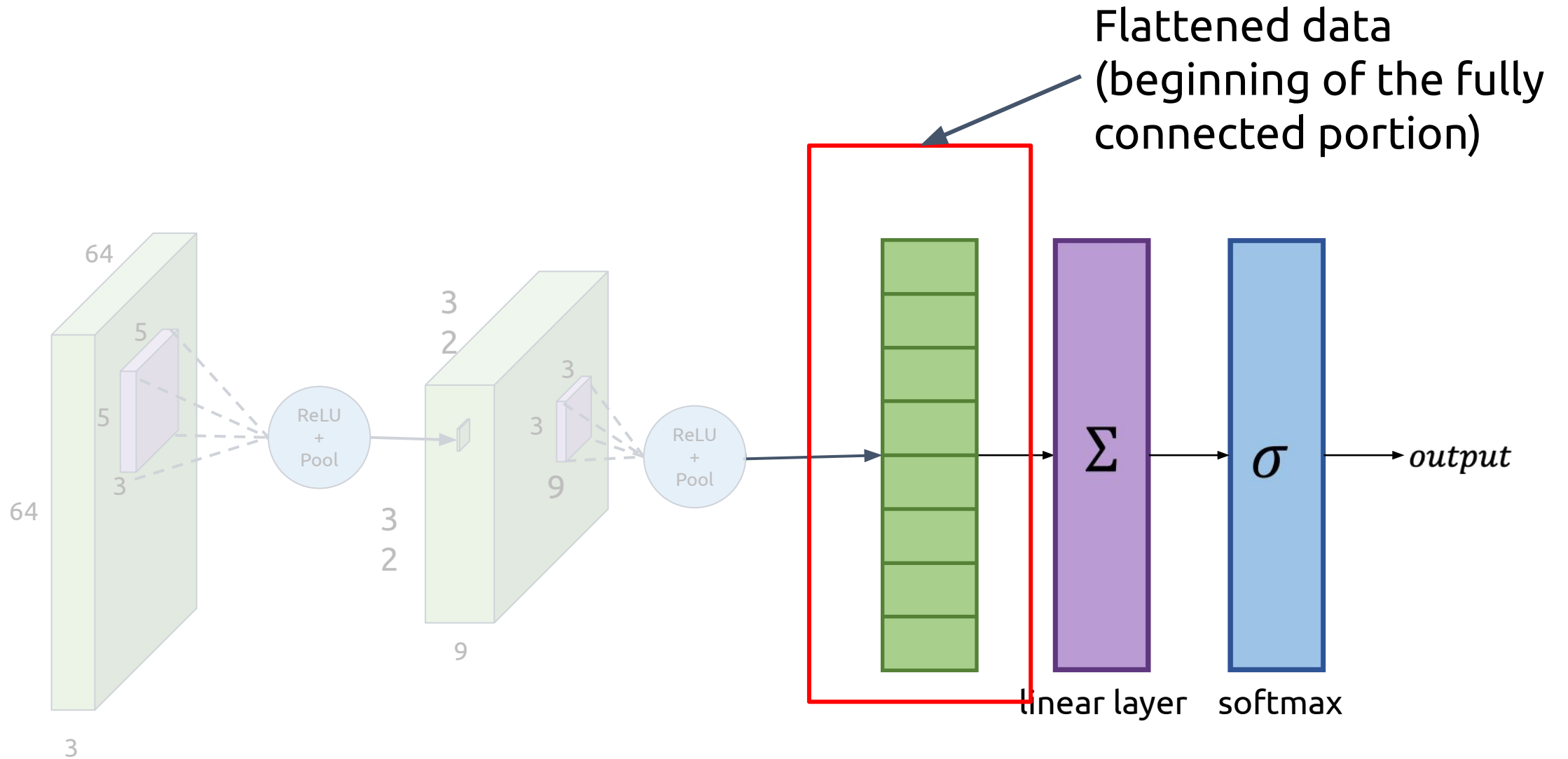


CNN Architecture

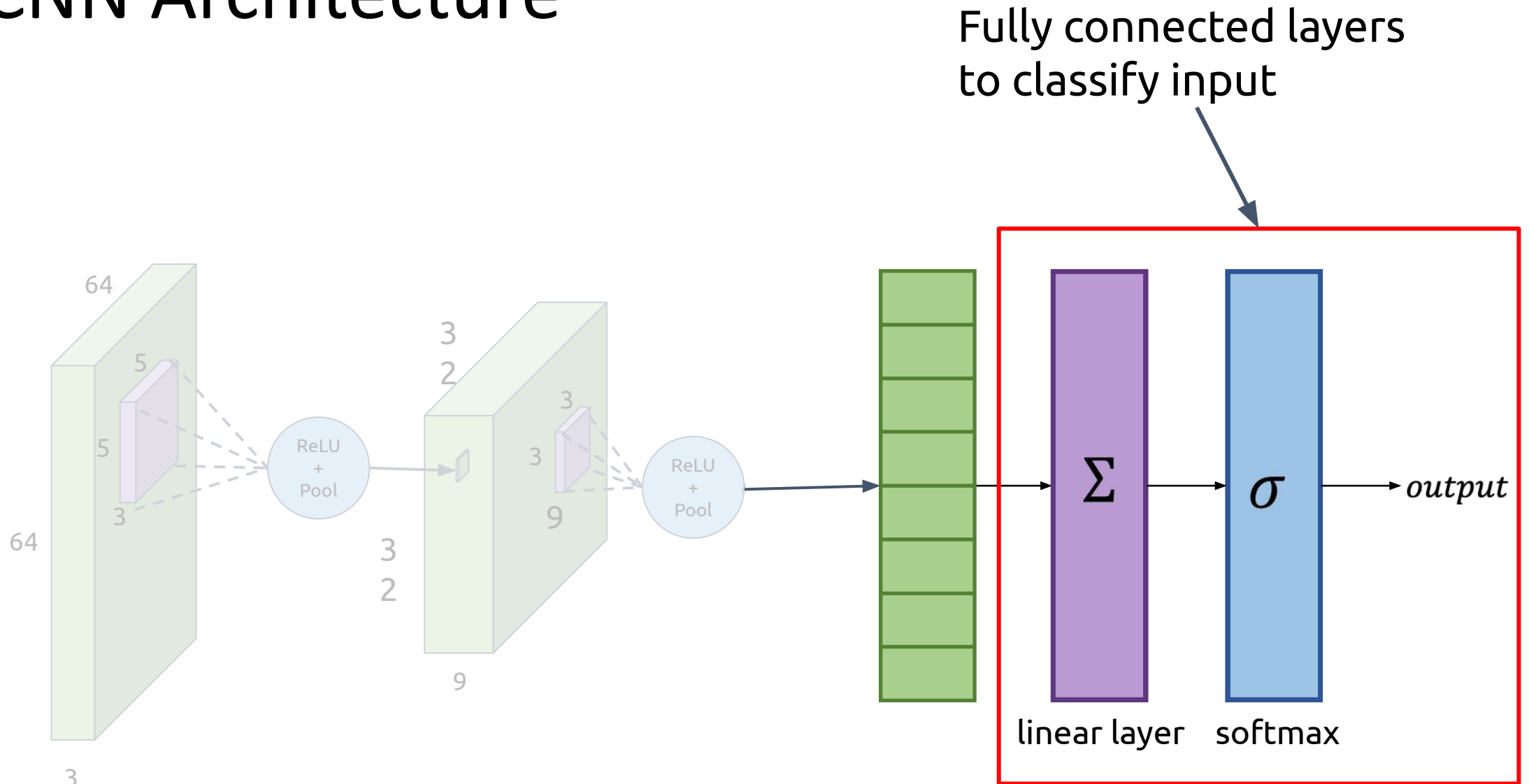
This part learns to perform a specific task
(e.g. classification) using those features



CNN Architecture

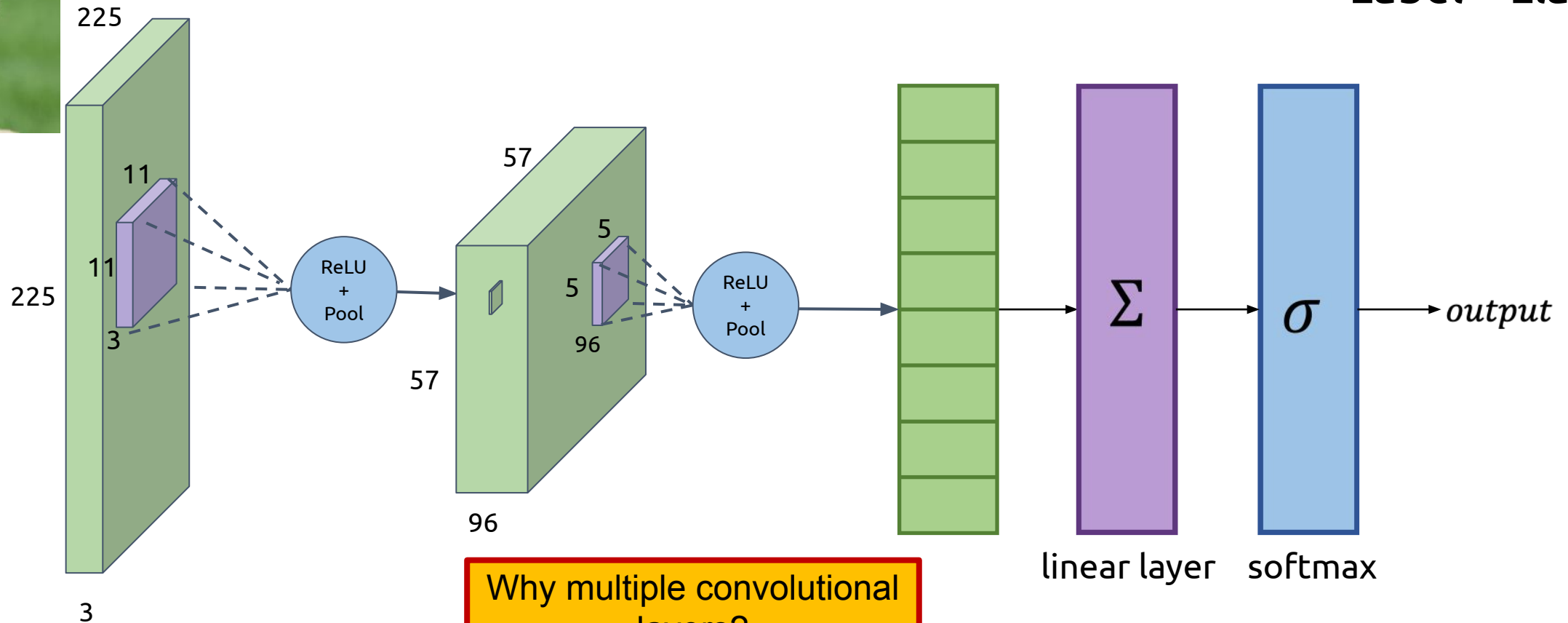


CNN Architecture



CNN Architecture

Input

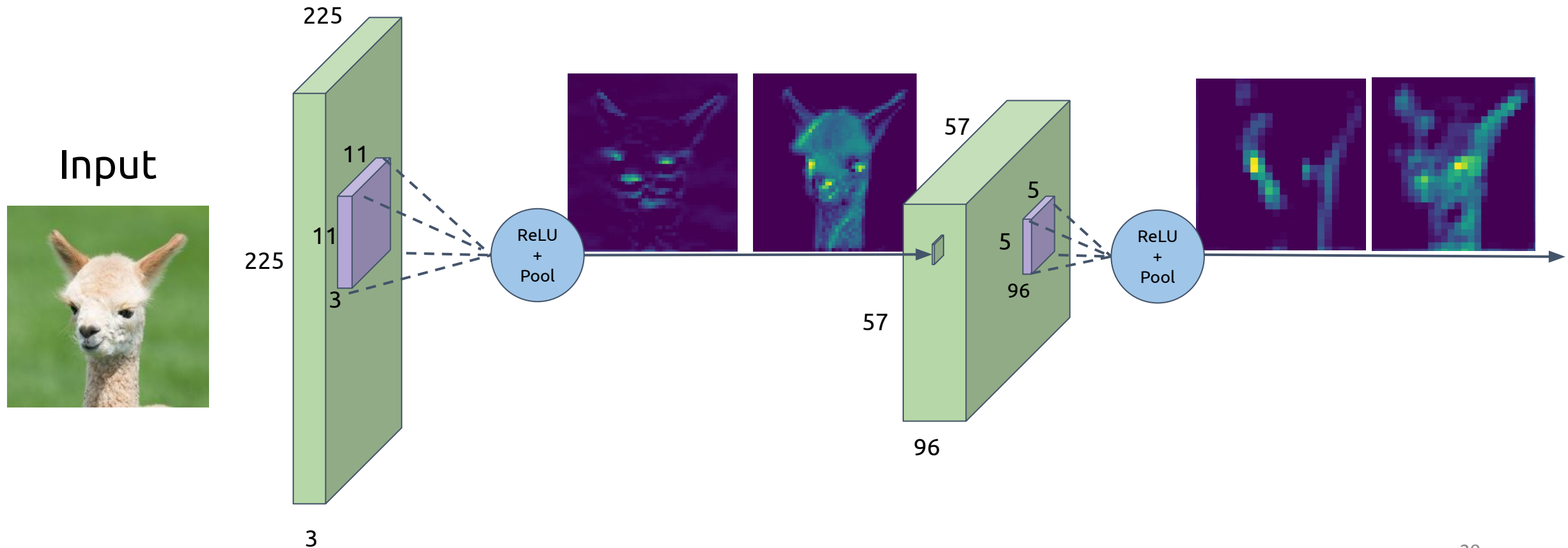


Why multiple convolutional layers?

Feature Extraction using multiple convolution layers

Hierarchy of features

Sequence of layers detect broader and broader features

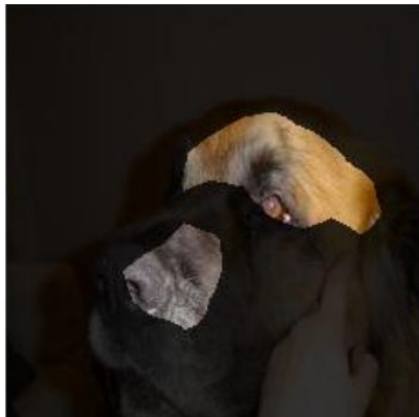




Example: Network Dissection

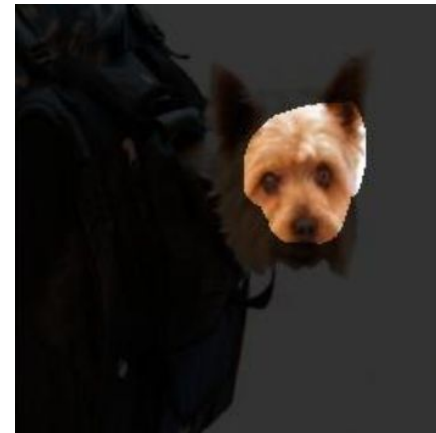
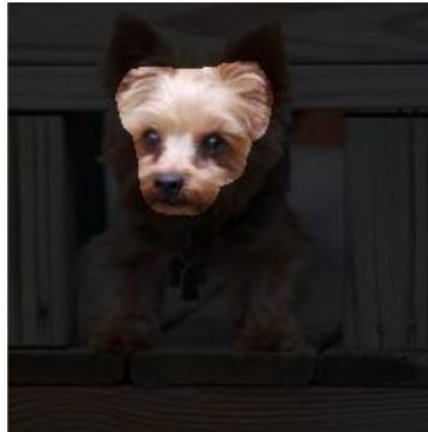
<http://netdissect.csail.mit.edu/>

Layer 3 active regions



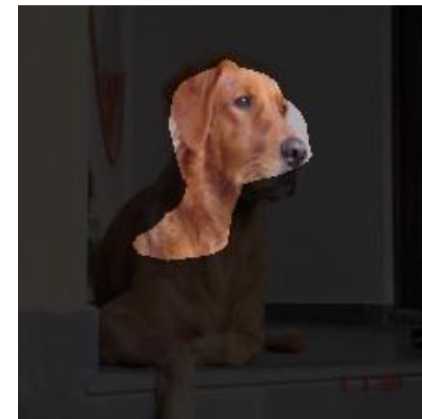
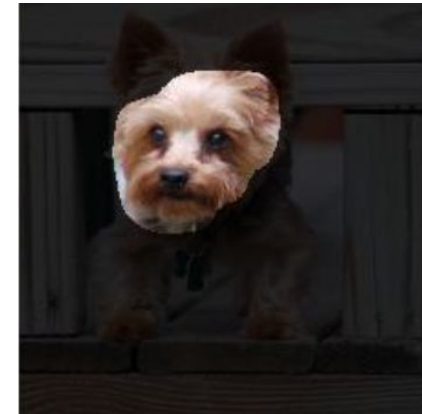
"Eye Detector"

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



https://commons.wikimedia.org/wiki/File:Common_zebra_1.jpg

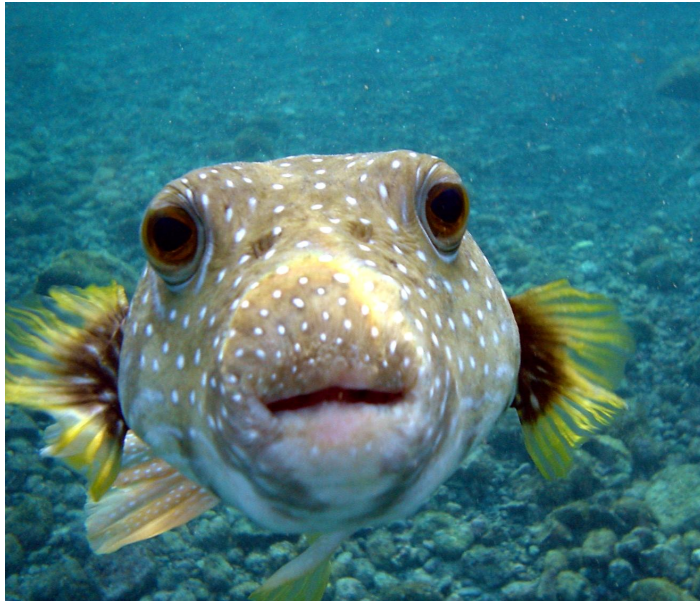
Predictions:

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



ILSVRC 2012

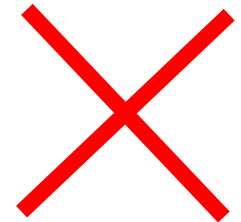
Percentage that model fails to classify is known as ***Top 5 Error Rate***



https://commons.wikimedia.org/wiki/File:Puffer_Fish_DSC01257.JPG

Predictions:

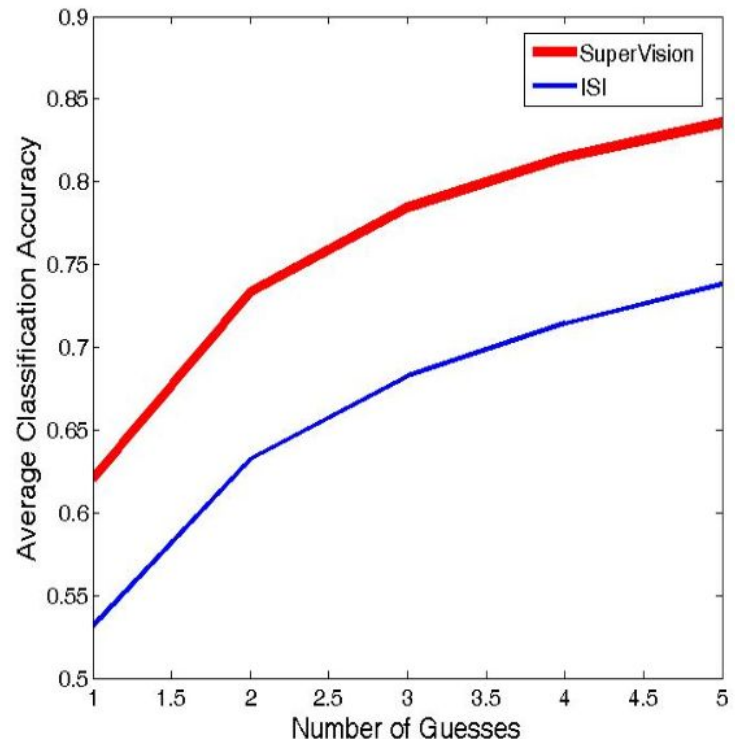
1. Sponge
2. Person
3. Llama
4. Flower
5. Boat



AlexNet: Why CNNs Are a Big Deal

Major performance boost on ImageNet at ILSVRC 2012

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

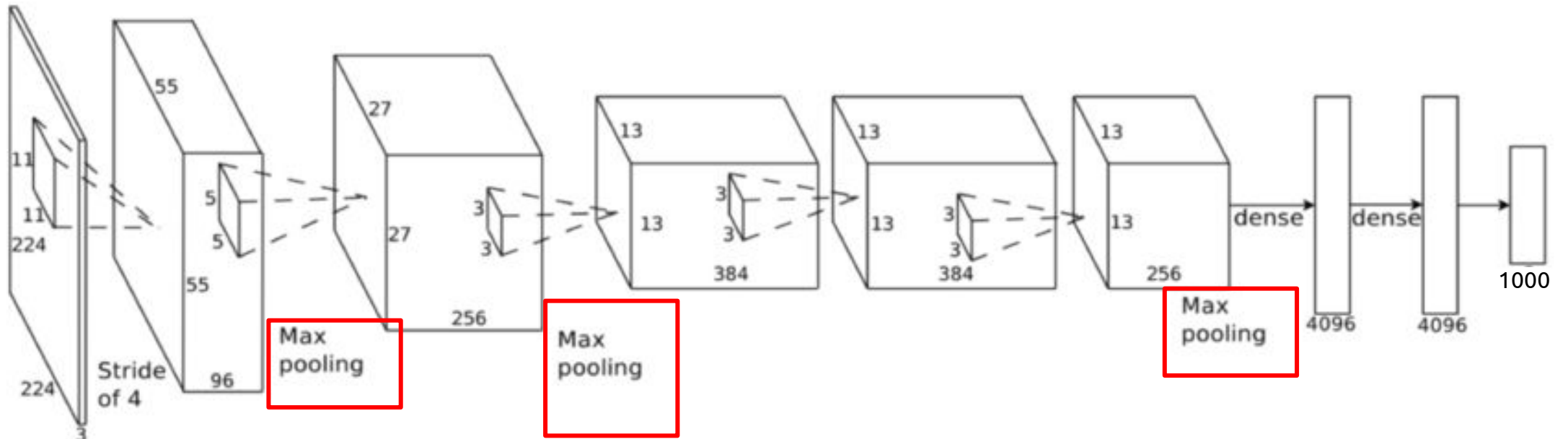


Note: SuperVision is the name of Alex's team

<http://image-net.org/challenges/LSVRC/2012/analysis/>

AlexNet

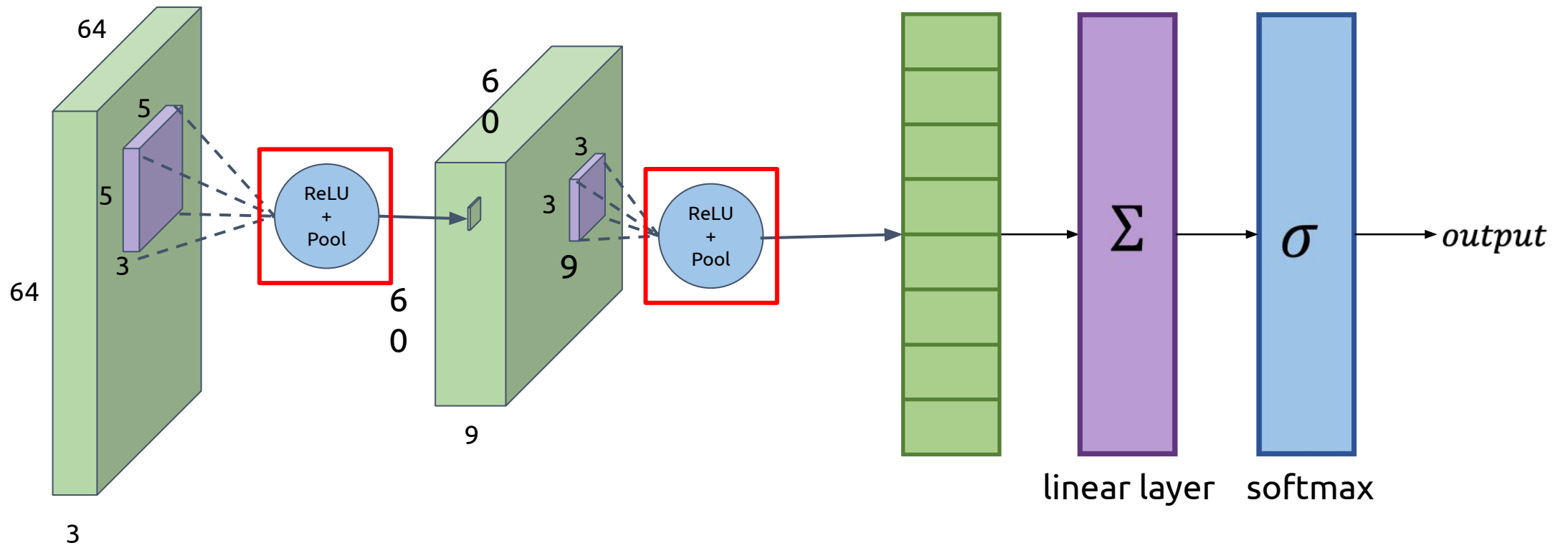
- 60 million parameters
- 5 Convolutional Layers
- 3 Fully Connected Layers



[Alex Krizhevsky et al. 2012]

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

Pooling



Max Pooling

Max pooling with stride 2 and 2x2 filters

6	3	1	-3
4	1	2	0
3	1	3	2
7	1	1	1

Max of pixels
in window



Max Pooling

Max pooling with stride 2 and 2x2 filters

6	3	1	-3
4	1	2	0
3	1	3	2
7	1	1	1

Max of pixels
in window



6	

Max Pooling

Max pooling with stride 2 and 2x2 filters

6	3	1	-3
4	1	2	0
3	1	3	2
7	1	1	1

Max of pixels
in window



6	2

Max Pooling

Max pooling with stride 2 and 2x2 filters

6	3	1	-3
4	1	2	0
3	1	3	2
7	1	1	1

Max of pixels
in window



6	2
7	

Max Pooling

Max pooling with stride 2 and 2x2 filters

6	3	1	-3
4	1	2	0
3	1	3	2
7	1	1	1

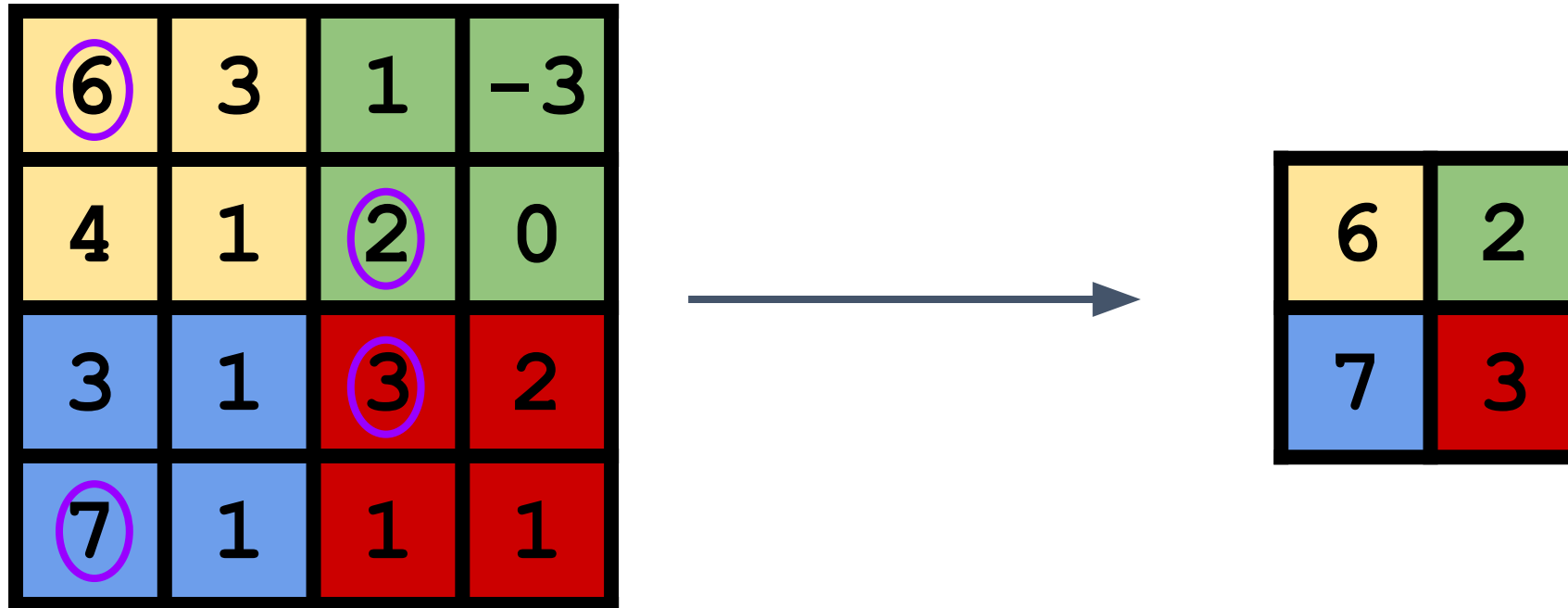
Max of pixels
in window



6	2
7	3

Max Pooling

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



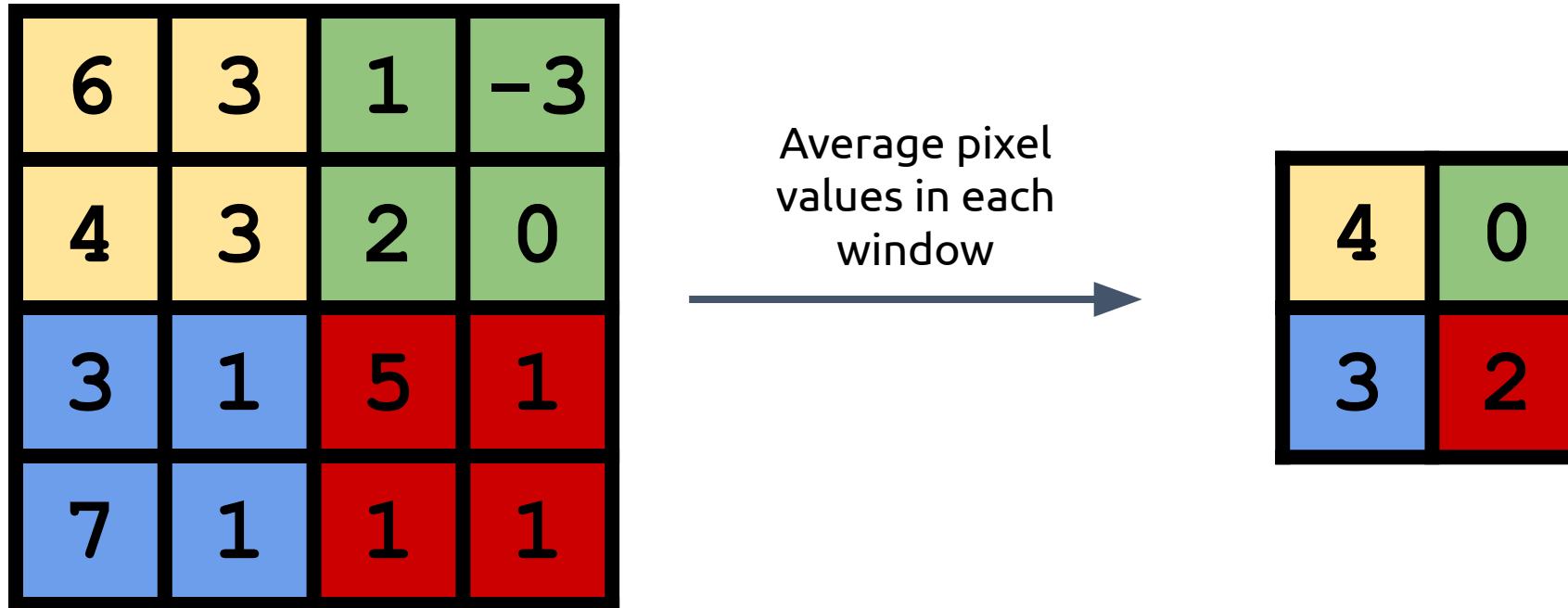
Convolution Output



After Pooling

Other Pooling Techniques

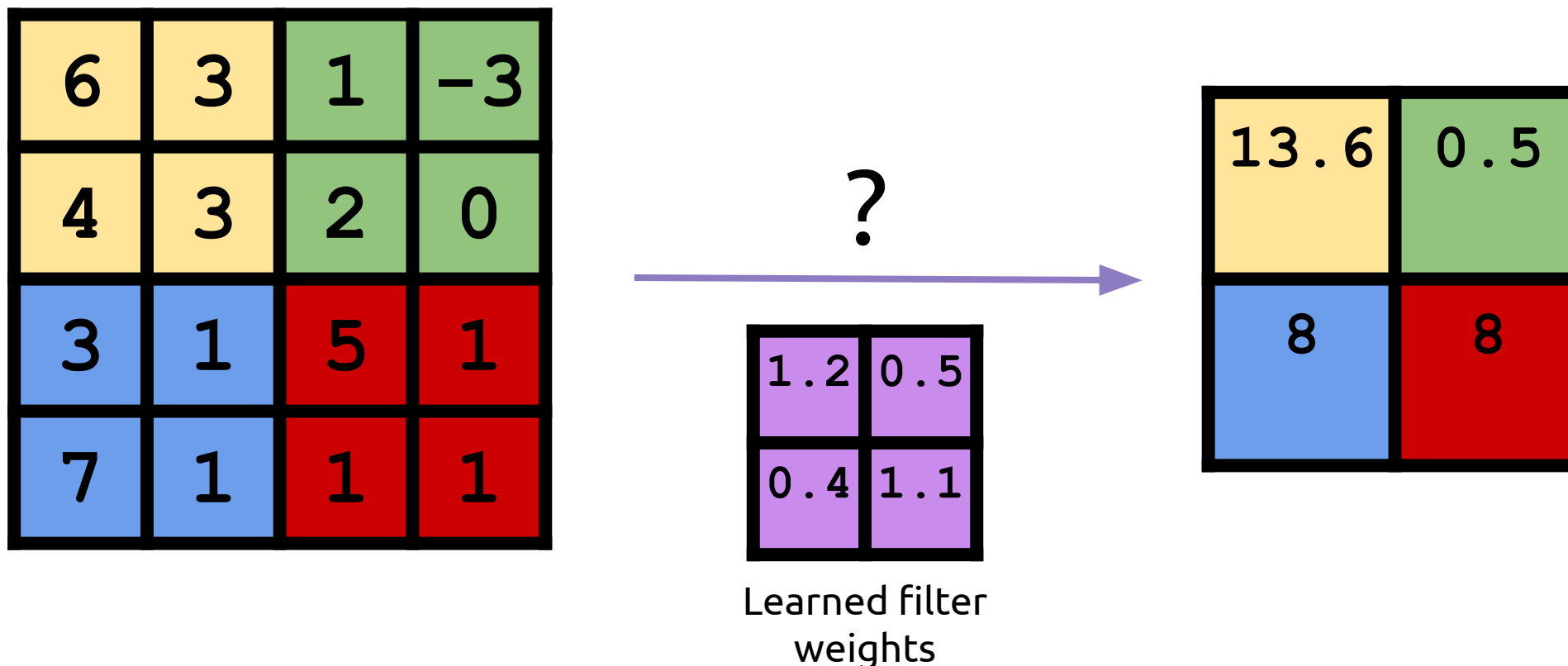
Average pooling with stride 2 and 2x2 filters





Learning a Pooling Function

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

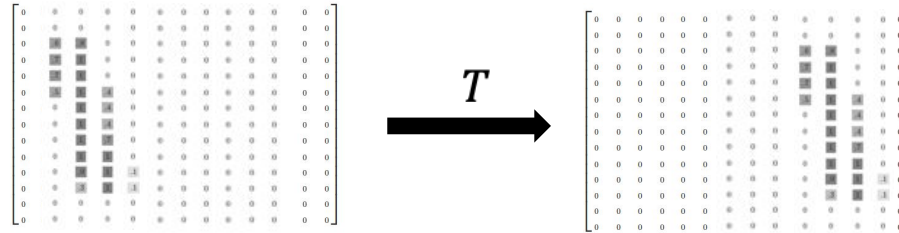


So...did we achieve our goal of translational invariance?



What was Translational Invariance again?

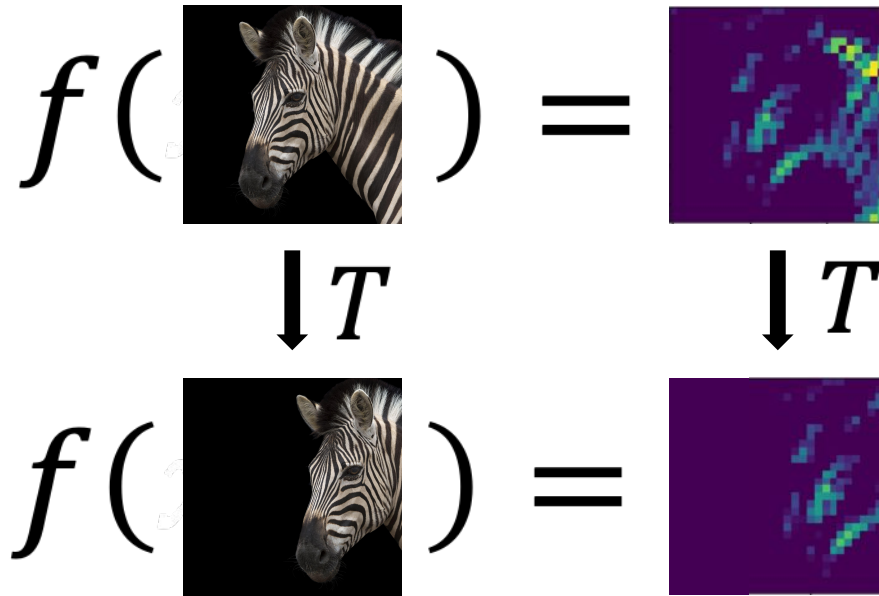
- To make a neural net f robust in this same way, it should ideally satisfy **translational invariance**: $f(T(x)) = f(x)$, where
 - x is the input image
 - T is a translation (i.e. a horizontal and/or vertical shift)



$$f\left(\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}\right) \stackrel{?}{=} f\left(\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}\right)$$

Are CNNs Translation Invariant?

- Convolution is ***translation equivariant***
 - A translated input results in an output translated by the same amount
 - $f(T(I)) = T(f(I))$
 - $(T(I) \otimes K)(x, y) = T(I \otimes K)(x, y)$



* Here, $(I \otimes K)(x, y) = \sum_m \sum_n I(x + m, y + n)K(m, n)$

Are CNNs Translation Invariant?

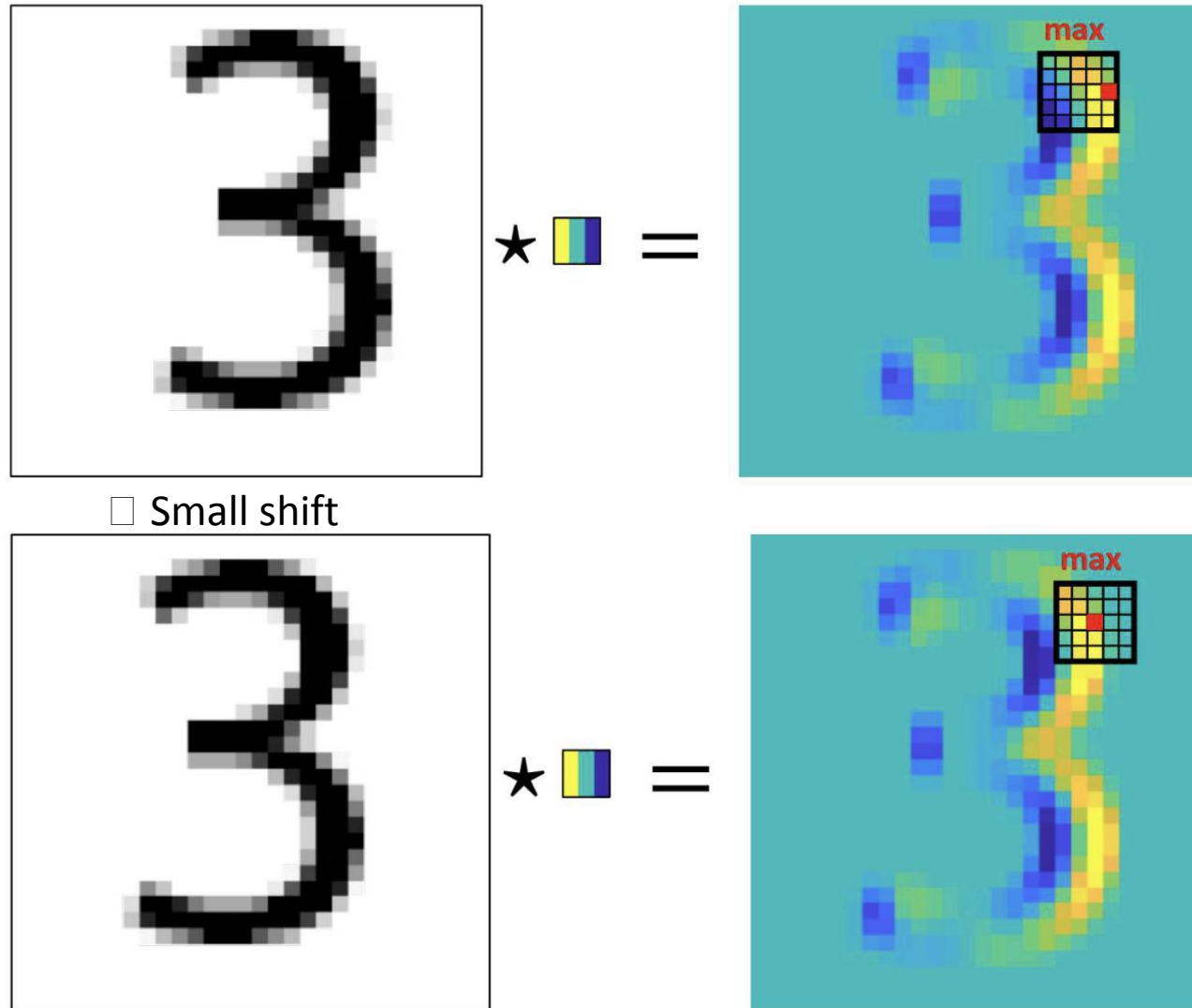
- Max pooling is intended to give invariance to small translations
 - The highest activation pixel can shift around within the pooling window, and the output does not change

$$f\left(\begin{array}{|c|c|}\hline 6 & 3 \\ \hline 4 & 1 \\ \hline\end{array}\right) = 6$$

$$f\left(\begin{array}{|c|c|}\hline 1 & 5 \\ \hline 6 & 3 \\ \hline\end{array}\right) = 6$$

$$f\left(\begin{array}{|c|c|}\hline 2 & 6 \\ \hline 2 & 4 \\ \hline\end{array}\right) = 6$$

So how does it all come together?

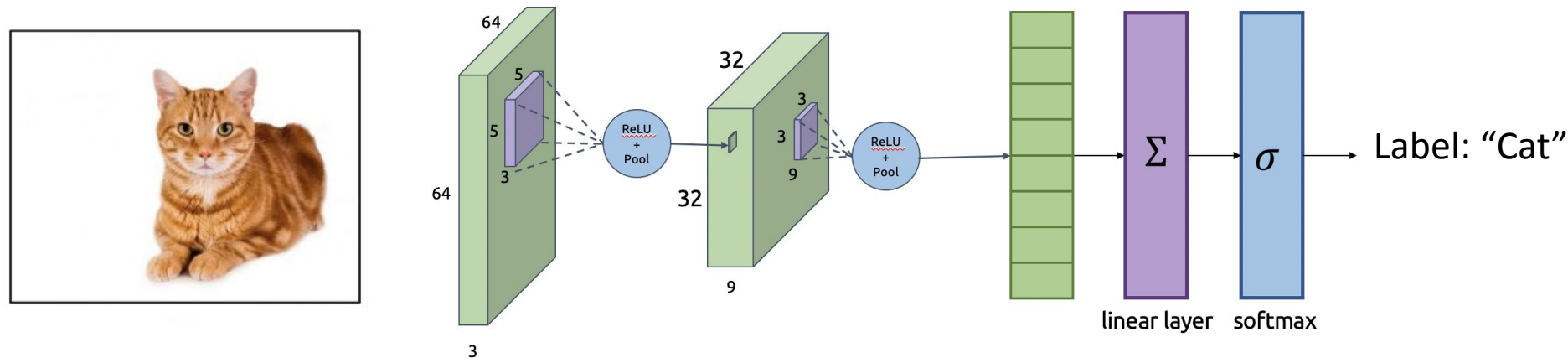


Convolution is
***translation
equivariant***

Max pooling gives
invariance to
small translations

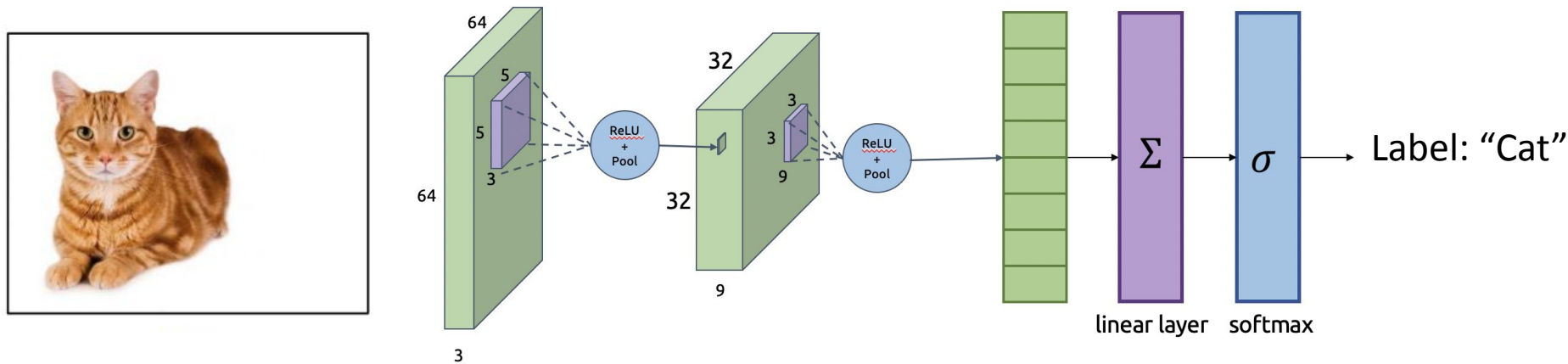
Are CNNs Translation Invariant?

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



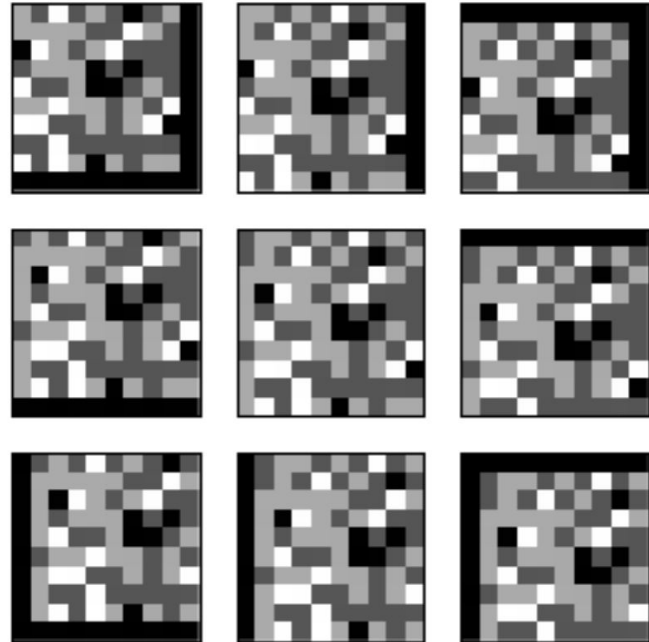
Are CNNs Translation Invariant?

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

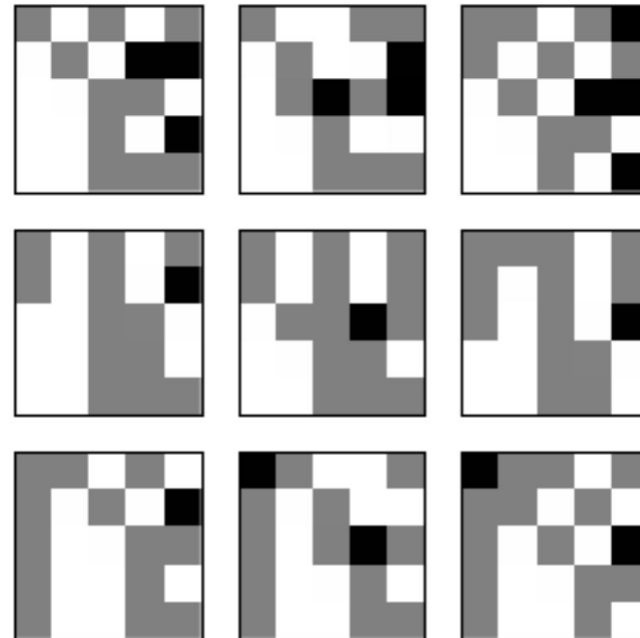


Are CNNs Translation Invariant?

- 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



Max Pool
→



These are
not all the
same!

Are CNNs Translation Invariant?

- Is it possible to build a truly translation invariant CNN?
 - Yes!
 - Have to properly “pre-filter” images before pooling them
 - Comes from signal processing theory (The Sampling Theorem)
 - Take CS 1230 (Computer Graphics) if you want to learn about this!
- One effort to make a translation-invariant CNN:
<https://arxiv.org/pdf/1904.11486.pdf>

Other Invariances

Rotation/Viewpoint Invariance



Other Invariances

Rotation/Viewpoint Invariance



Size Invariance





Other Invariances

Rotation/Viewpoint Invariance



Size Invariance



Illumination Invariance



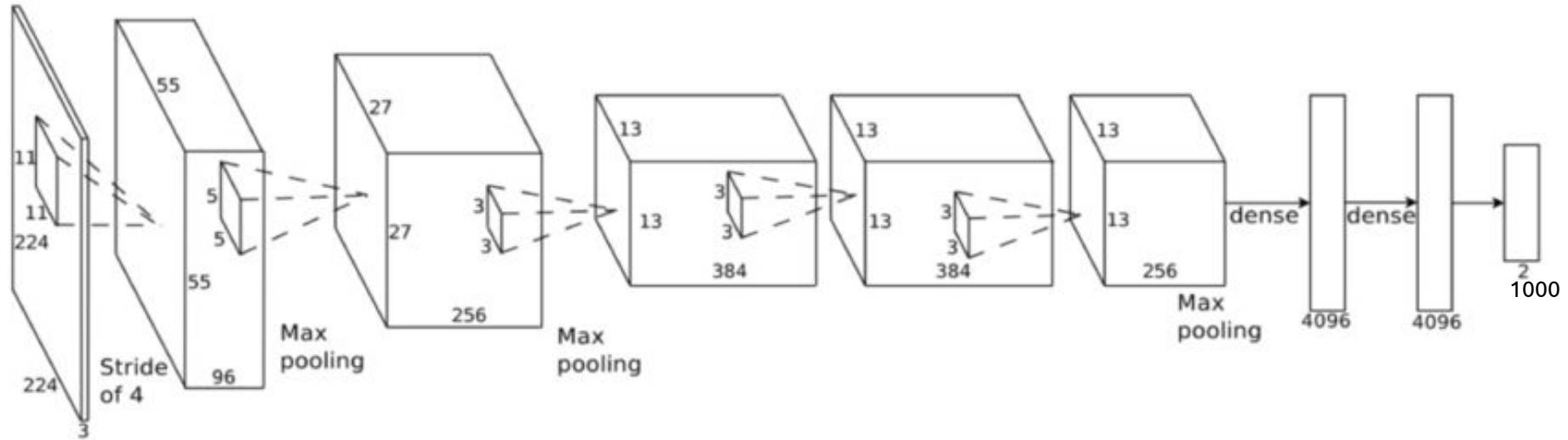
- All of these are desirable
- How do CNNs fare?
 - Max pooling gives some small amount of size invariance...
 - ...but in general, CNNs don't do well with big changes in size, pose, or lighting

What should we we do?

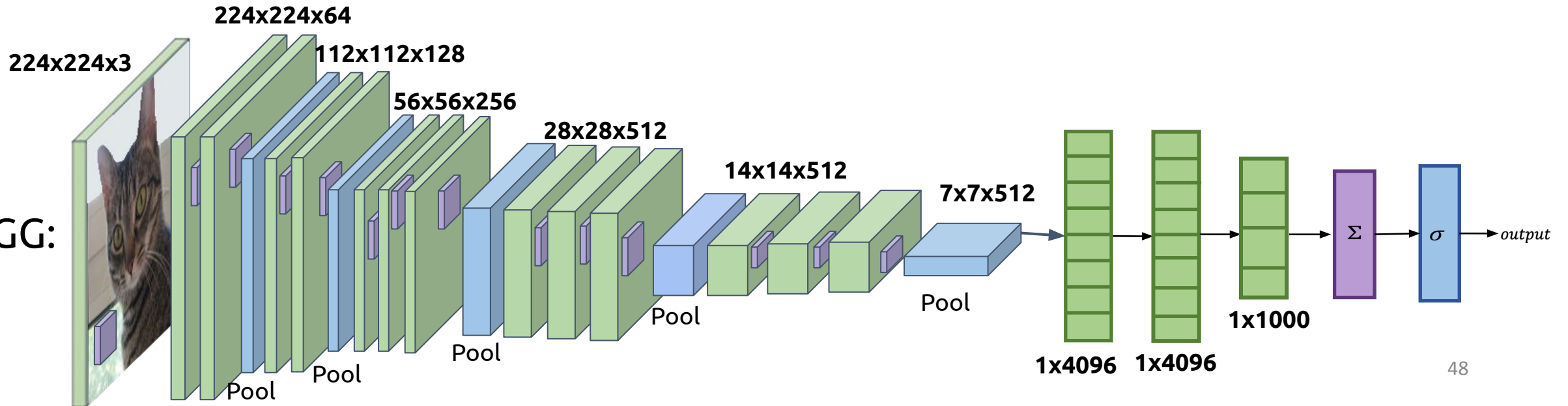
- Consequence of not having these invariances?
 - Need **lots** of training data
 - Have to show the network examples of everything under different poses, lighting, etc.
 - Data Augmentation

More Complicated Networks

AlexNet:



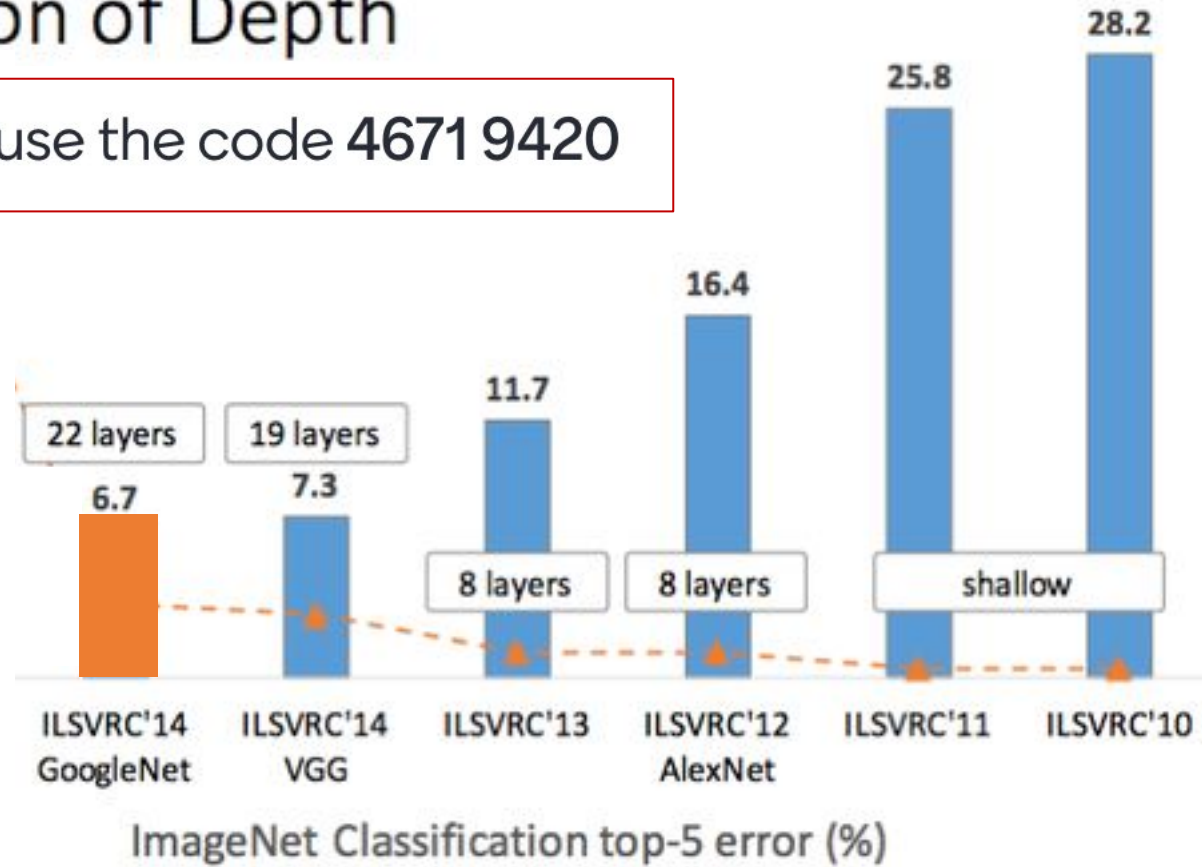
VGG:



Can you guess what was
the biggest bottleneck to
adding more layers?

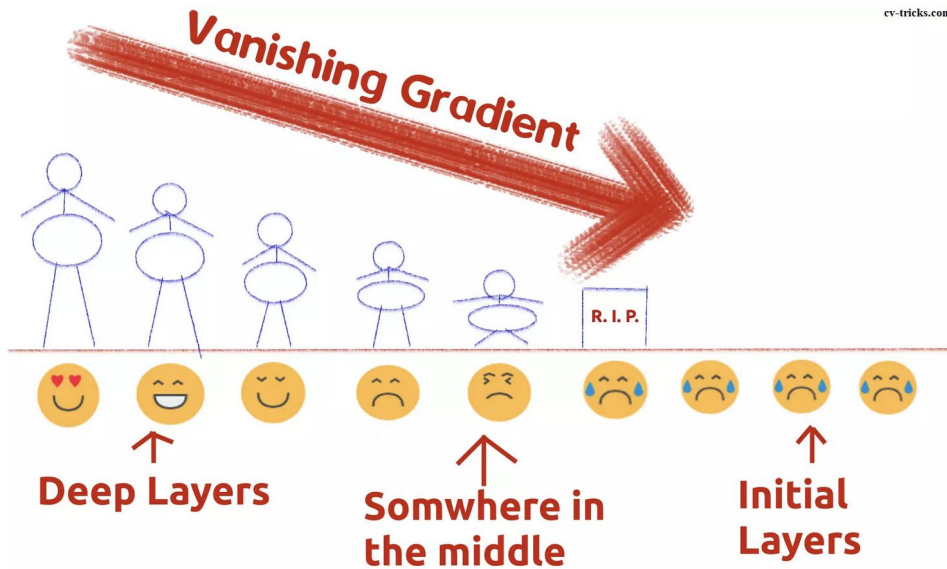
Revolution of Depth

Go to www.menti.com and use the code 4671 9420

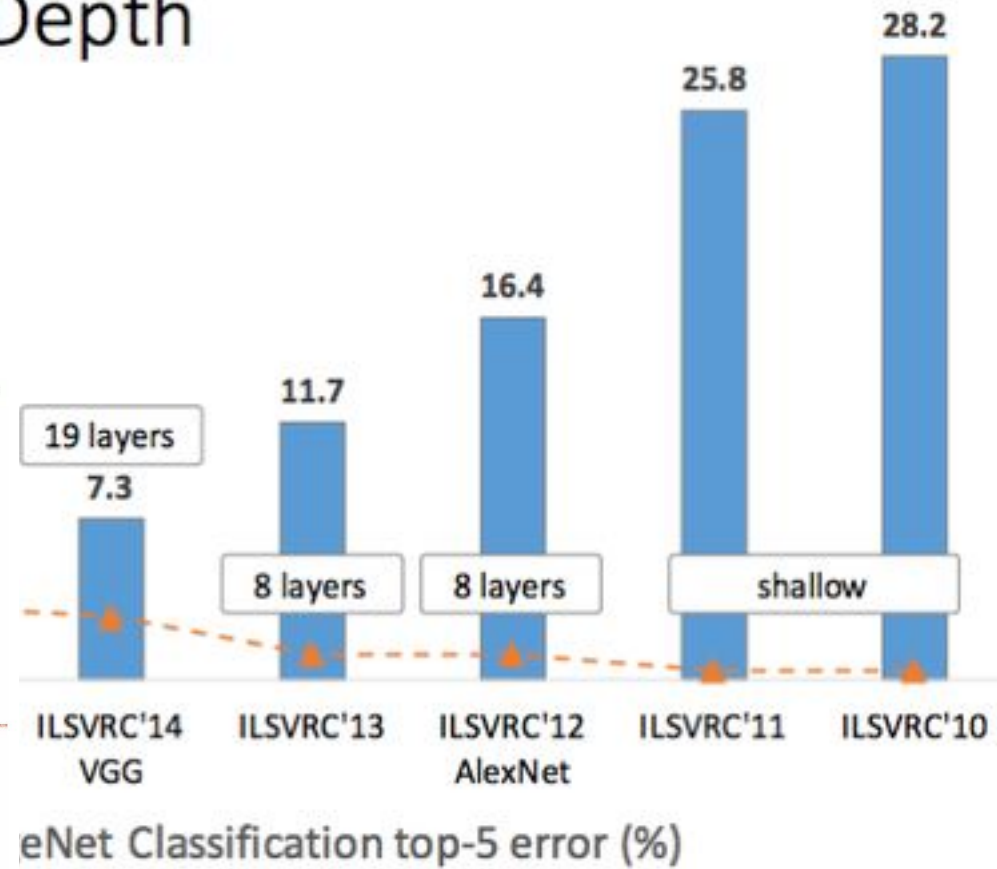


Revolution of Depth

Vanishing gradients!



cv-tricks.com

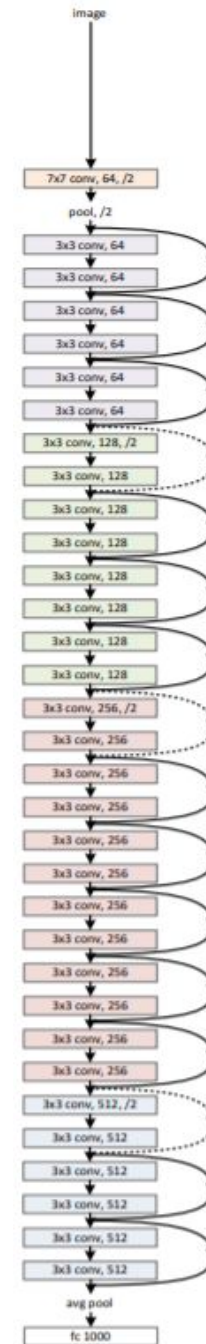
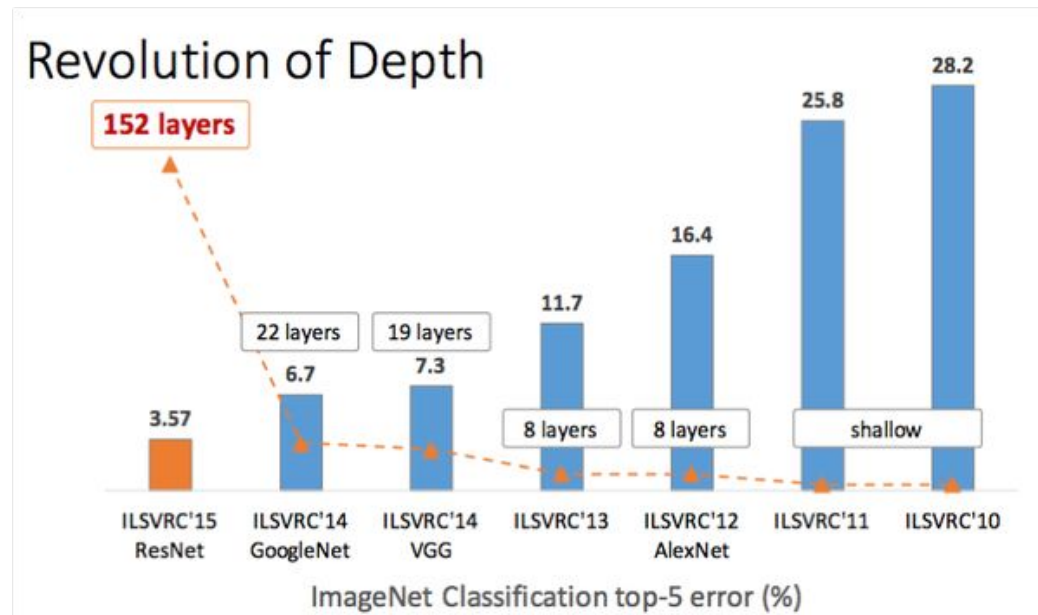


More Complicated Networks

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.

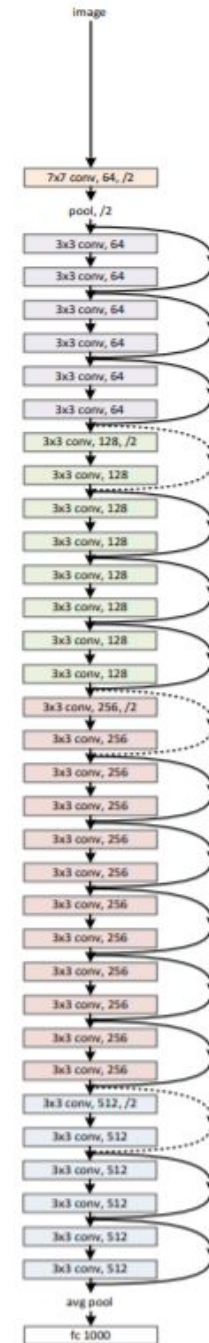
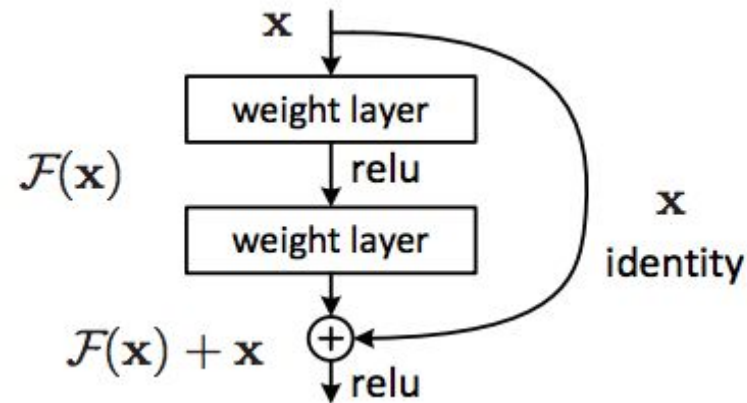
More Complicated Networks

ResNet:

Lots of layers, tons of learnable parameters

Avoids Vanishing Gradient problem

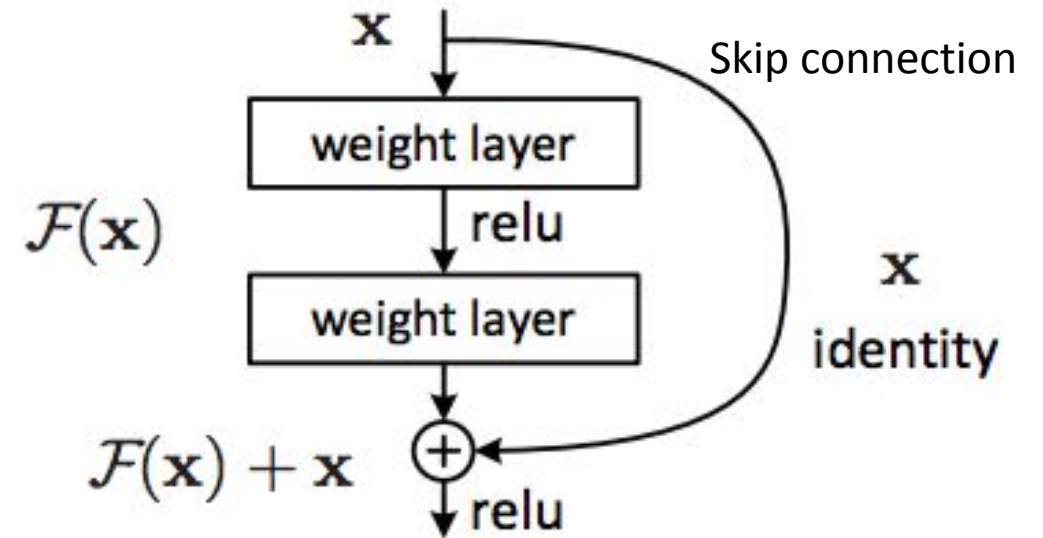
Residual Block



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

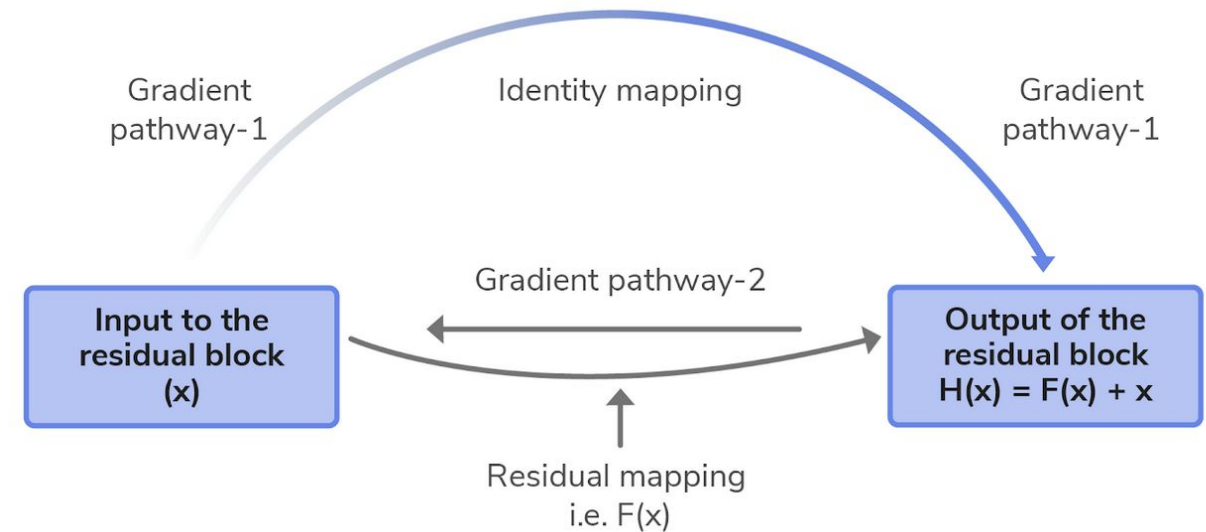


Any questions?



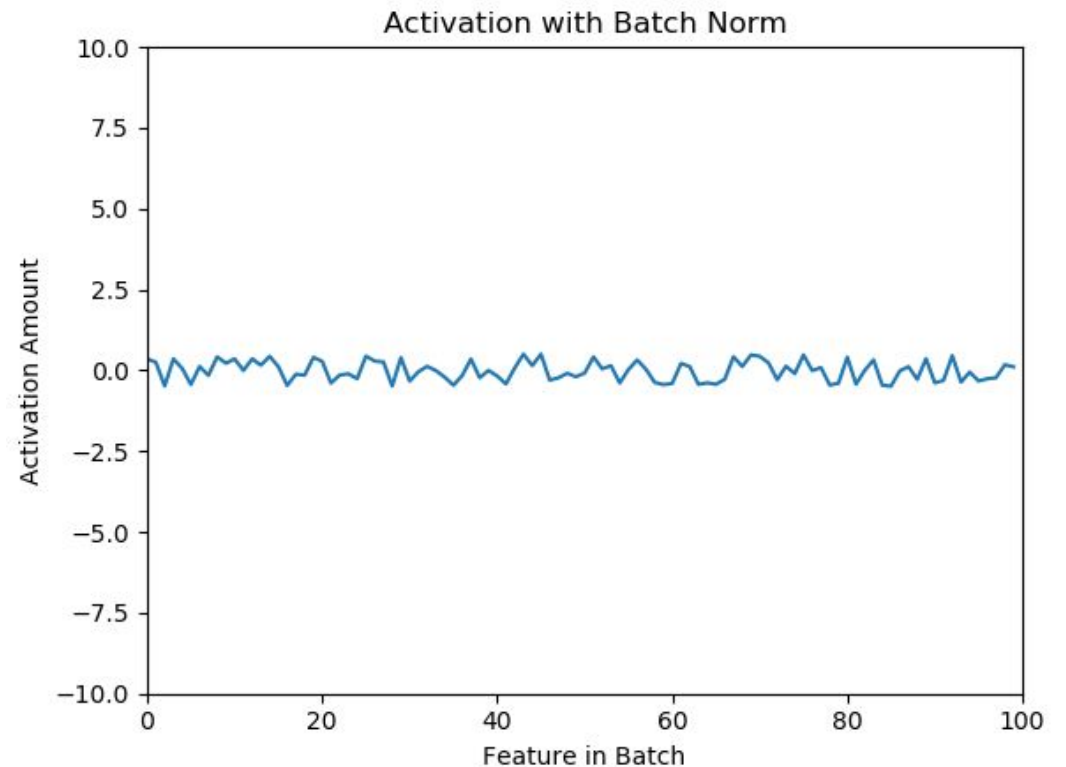
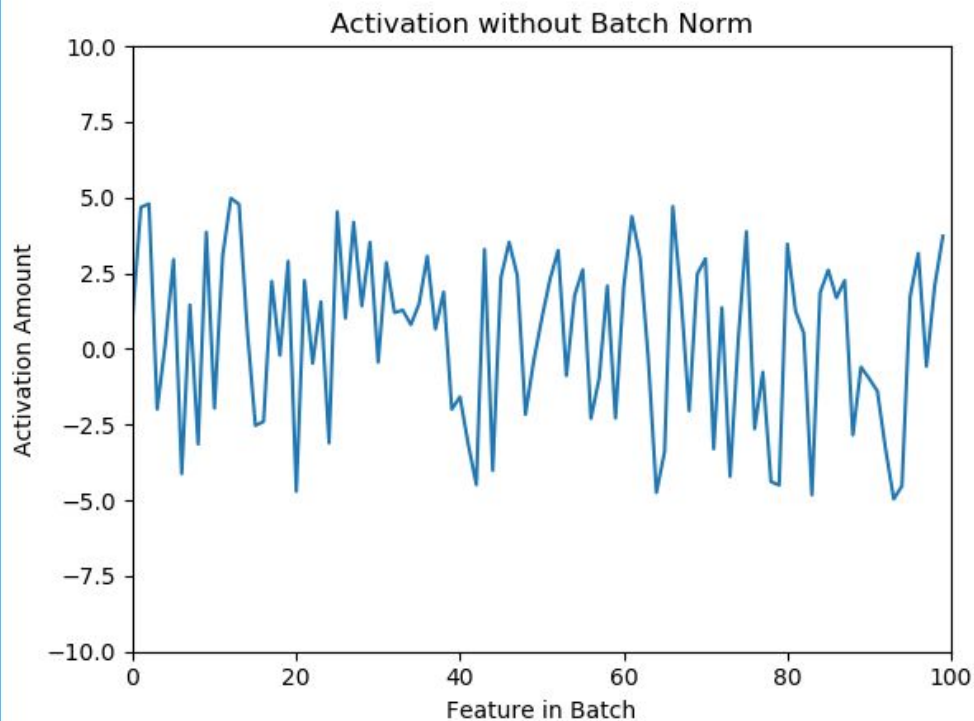
Residual Blocks

- 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
- 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

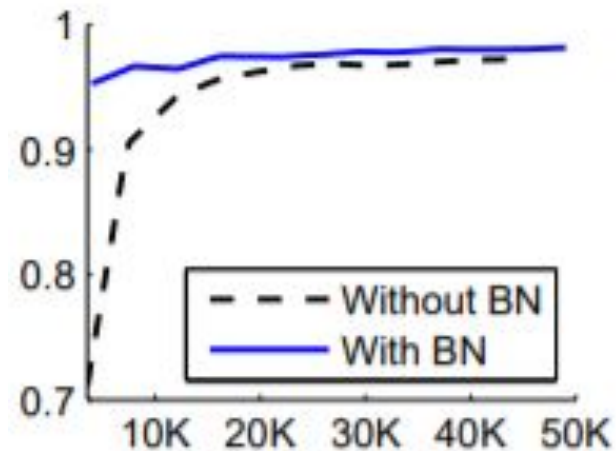


Why might we want to do this?

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: https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/BatchNormalization

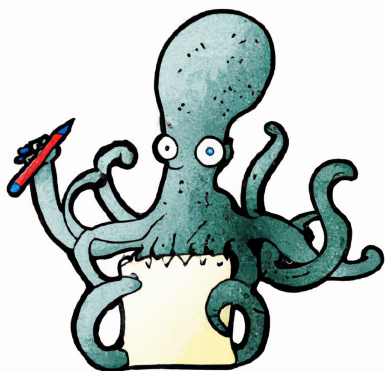
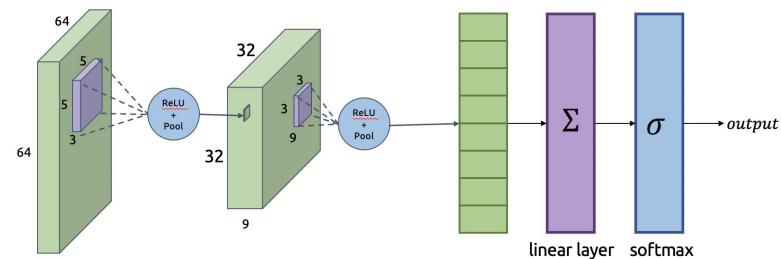
Recap

CNNs

Architecture

AlexNet + Pooling

CNNs are “sort of” translationally invariant



Deeper CNNs

Many layers = vanishing gradient

ResNet + Residual blocks

Batch normalization

Revolution of Depth

