

Beras Tips

Use built in matrix operations in numpy (@ for matrix multiplication)

If you can't figure out a shape, play around in a REPL (i.e., run "python" from the terminal)

Many functions are "just a formula"

Forward function of linear layer: xW+b

MSE: mean((y-y_pred)**2)

get_input_gradient, returns a list of Tensors, one Tensor for each input. If one of those inputs is not a trainable parameter (i.e., the true labels in the loss function), return an array of 0's.

Others are conceptually harder...

Gradient Tape's gradient method

• Start with target (loss) gradient, backprop gradient to previous layers, push those gradients to grandparent layers, etc.

What has happened in the last 15 years?

What has changed?

- 1. Power and efficiency of compute (GPUs)
- 2. Availability of data (the internet)
- 3. New Architectures (e.g., CNNs, Transformers)



Issues with MLPs

- 1. Resource Intensive
- 2. Difficult to incorporate certain types of information
- 3. (and more)

Issues with MLPs

- 1. Resource Intensive
- 2. Difficult to incorporate certain types of information
- 3. (and more)

GPUs to the rescue!



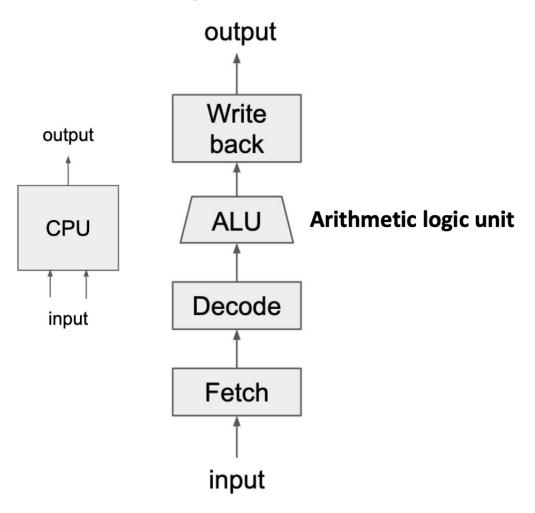


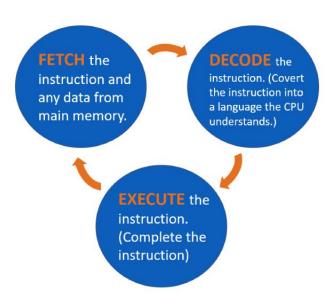
- GPUs are really good at computing mathematical operations in parallel!
- Matrix multiplication == many independent multiply and add operations

Easily parallelizable

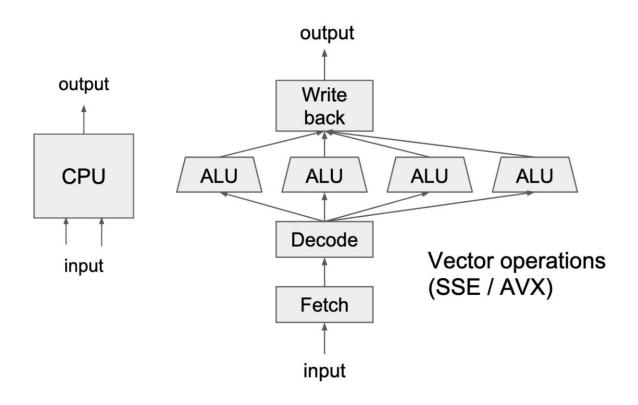
GPUs are great for this!

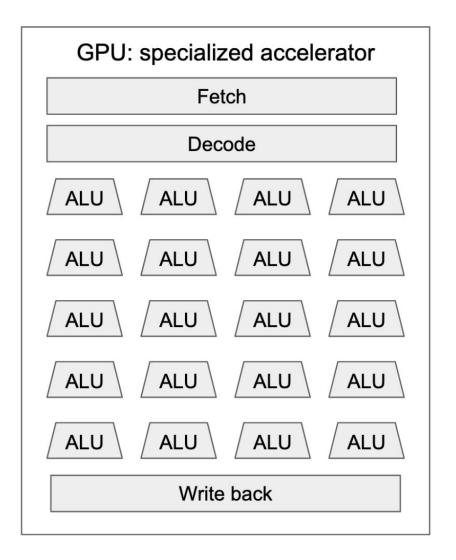
CPU v/s GPU





CPU v/s GPU





GPU-Parallel Acceleration

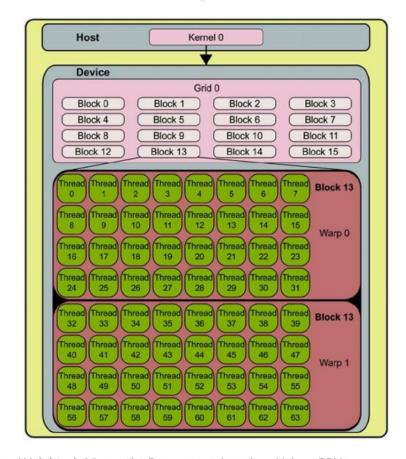
- User code (kernels) is compiled on the host (the CPU) and then transferred to the device (the GPU)
- Kernel is executed as a grid
- Each grid has multiple thread blocks
- Each thread block has multiple warps

A warp is the basic schedule unit in kernel execution

A warp consists of 32 threads

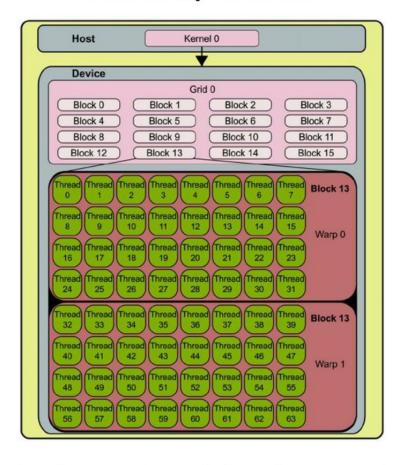
Compute Unified Device Architecture is a parallel computing platform and application programming interface (API)

CUDA compute model



GPU-Parallel Acceleration

CUDA compute model



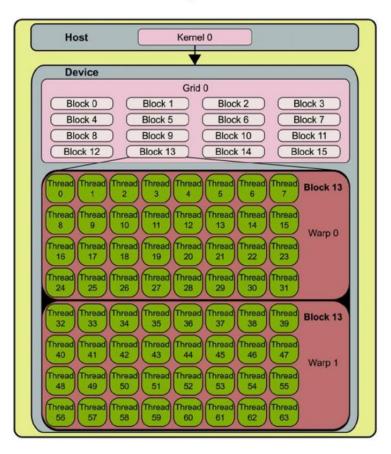
- Programmer decides how they want to parallelize the computation across grids and blocks
 - Modern deep learning frameworks take care of this for you
- CUDA compiler figures out how to schedule these units of computation on to the physical hardware

Any questions?

???

GPU-Parallel Acceleration

CUDA compute model



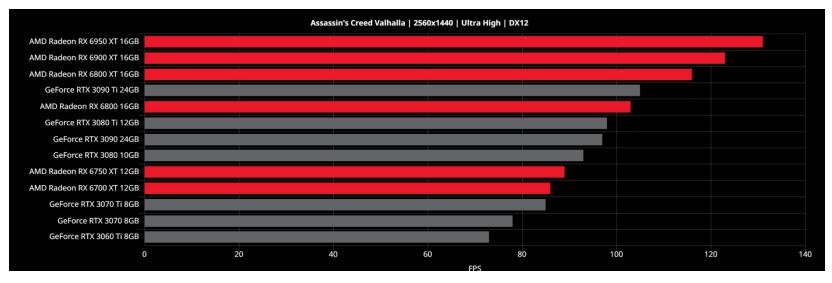
- Upshot: order of magnitude speedups!
- Example: training CNN on CIFAR-10 dataset

Device	Speed of training, examples/sec
2 x AMD Opteron 6168	440
i7-7500U	415
GeForce 940MX	1190
GeForce 1070	6500

From:

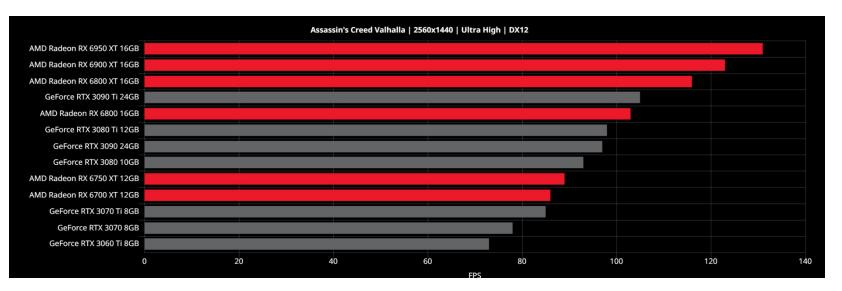
https://medium.com/@andriylazorenko/tensorflow-performance-test-cpu-vs-gpu-79fcd39170c

AMD GPUs are competitive for gaming and graphics, why not for AI?



(With a benchmarking tool made by AMD)

AMD GPUs are competitive for gaming and graphics, why not for AI?

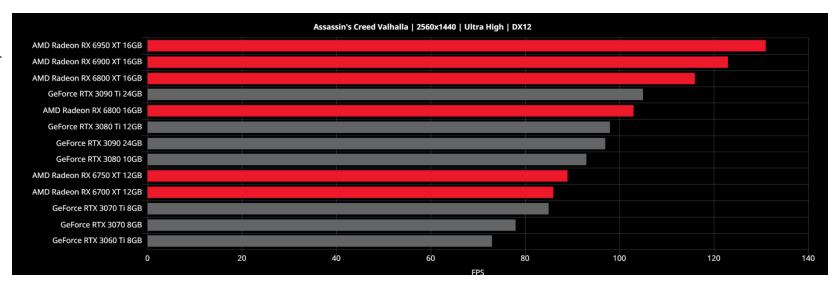


CUDA is far better than competitors (AMD)

(With a benchmarking tool made by AMD)

- Easier to use
- Better optimization
- AMD makes GPUs for graphics, NVIDIA makes GPUs for AI

AMD GPUs are competitive for gaming and graphics, why not for AI?



CUDA is far better than competitors (AMD)

(With a benchmarking tool made by AMD)

- Easier to use
- Better optimization
- AMD makes GPUs for graphics, NVIDIA makes GPUs for Al

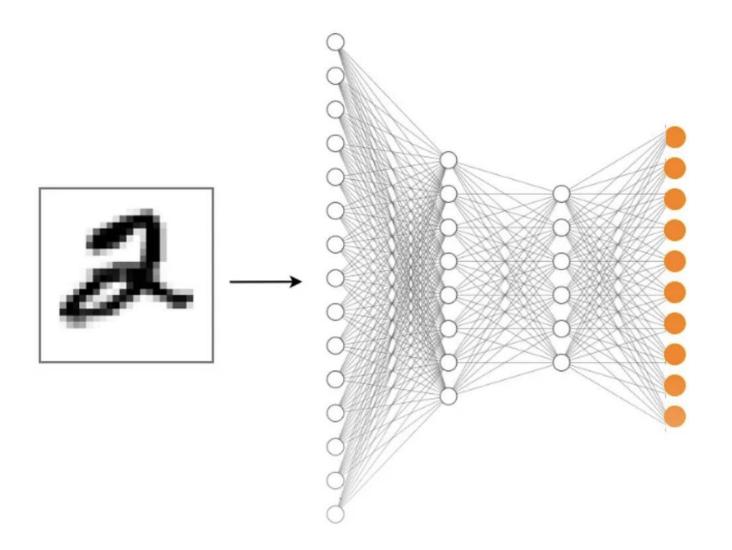
CUDA is Still a Giant Moat for NVIDIA

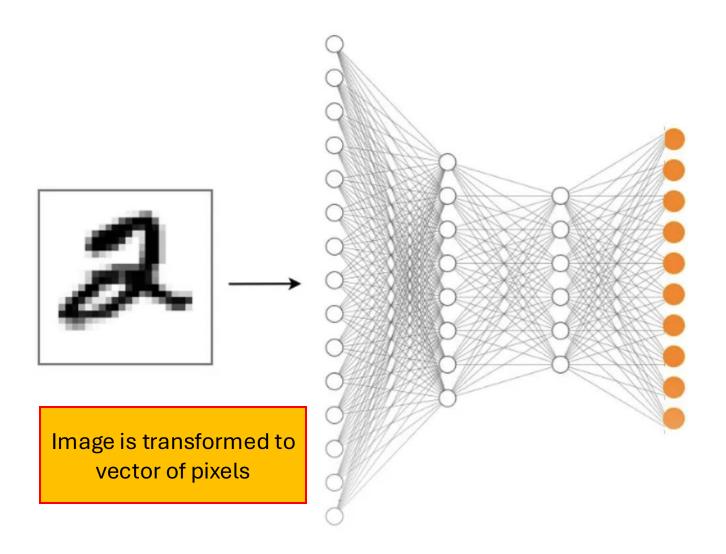
Despite everyone's focus on hardware, the software of AI is what protects NVIDIA

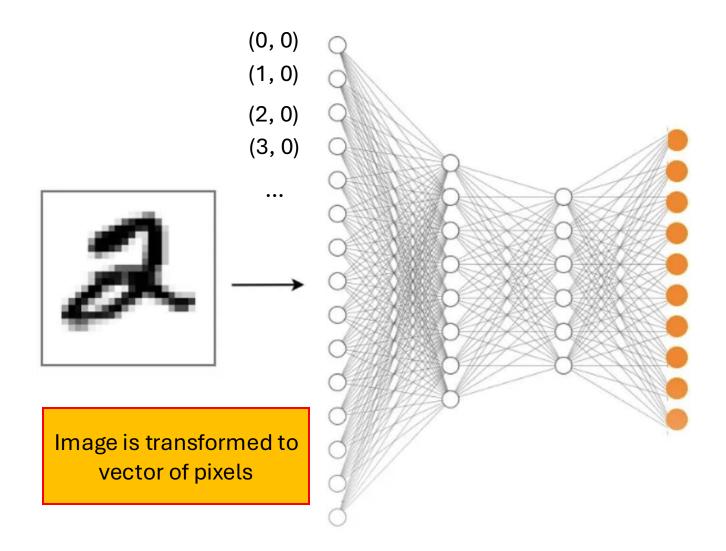


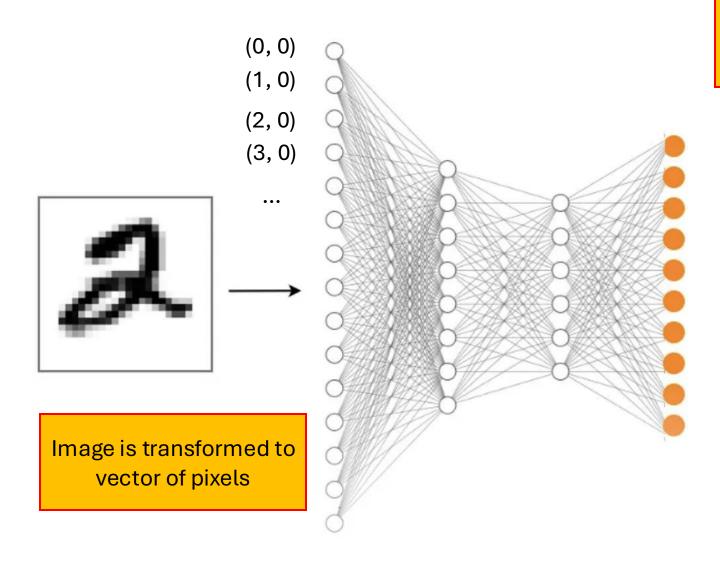
Issues with MLPs

- 1. Resource Intensive
- 2. Difficult to incorporate certain types of information

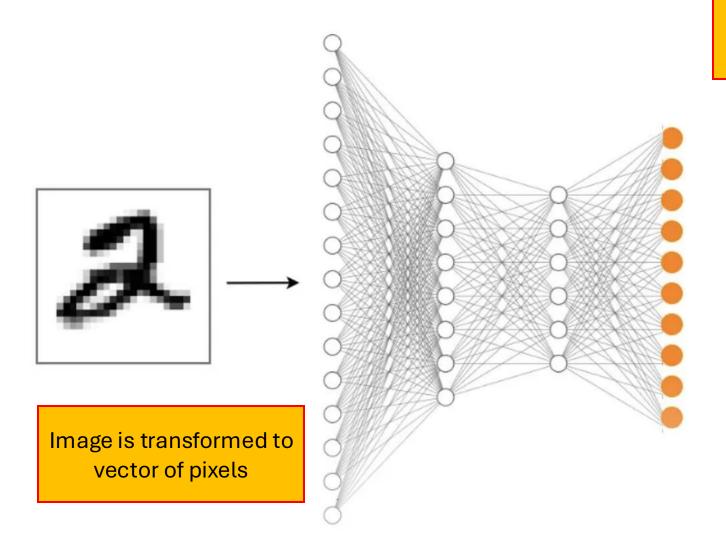




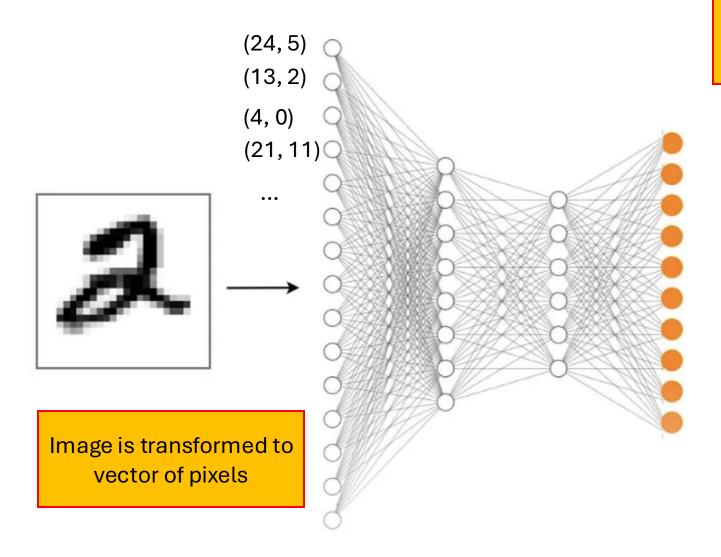




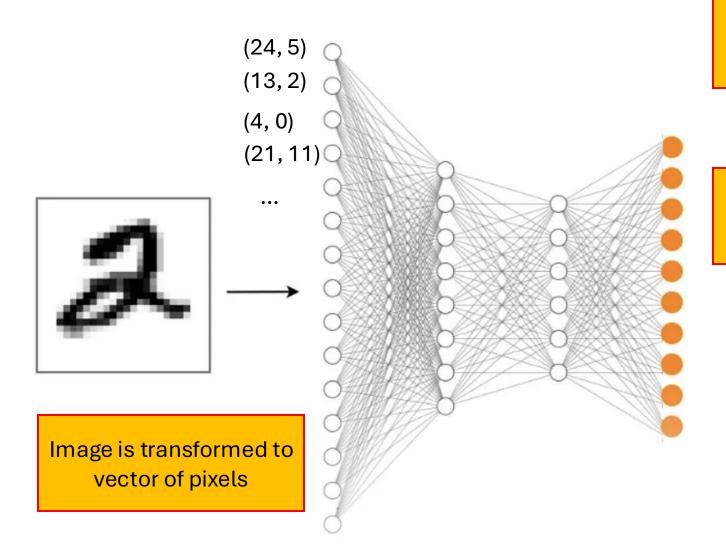
What would happen if we permuted the ordering of the pixels?



What would happen if we permuted the ordering of the pixels?

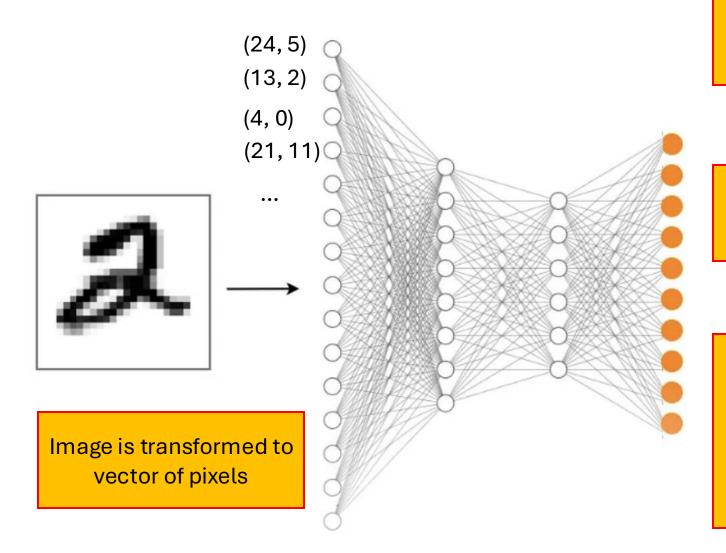


What would happen if we permuted the ordering of the pixels?



What would happen if we permuted the ordering of the pixels?

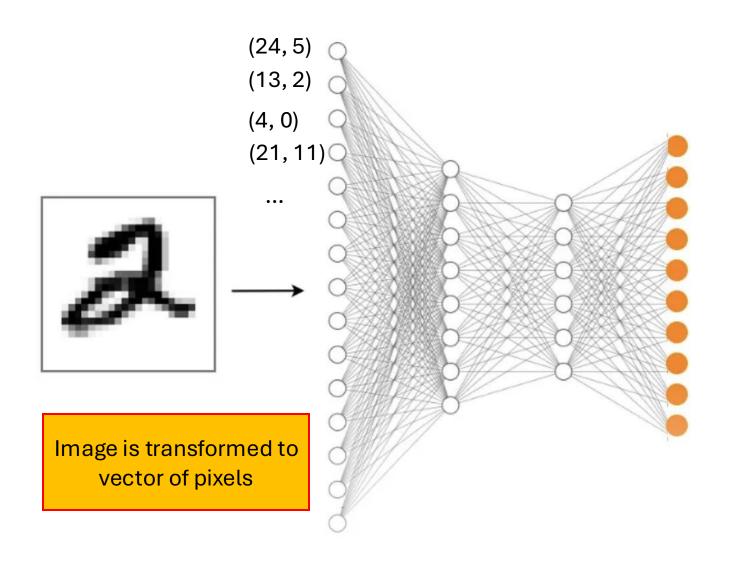
Will the training of the neural network differ?



What would happen if we permuted the ordering of the pixels?

Will the training of the neural network differ?

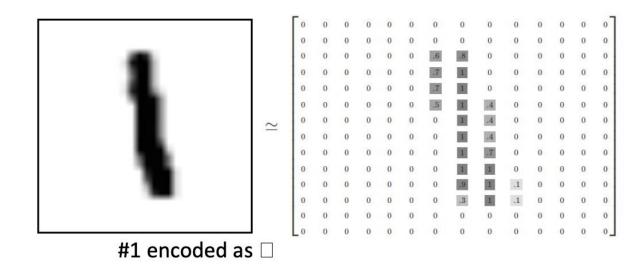
No! MLPs do not use spatial information, it does not matter which order the pixels are fed in so long as it is the same ordering for every input



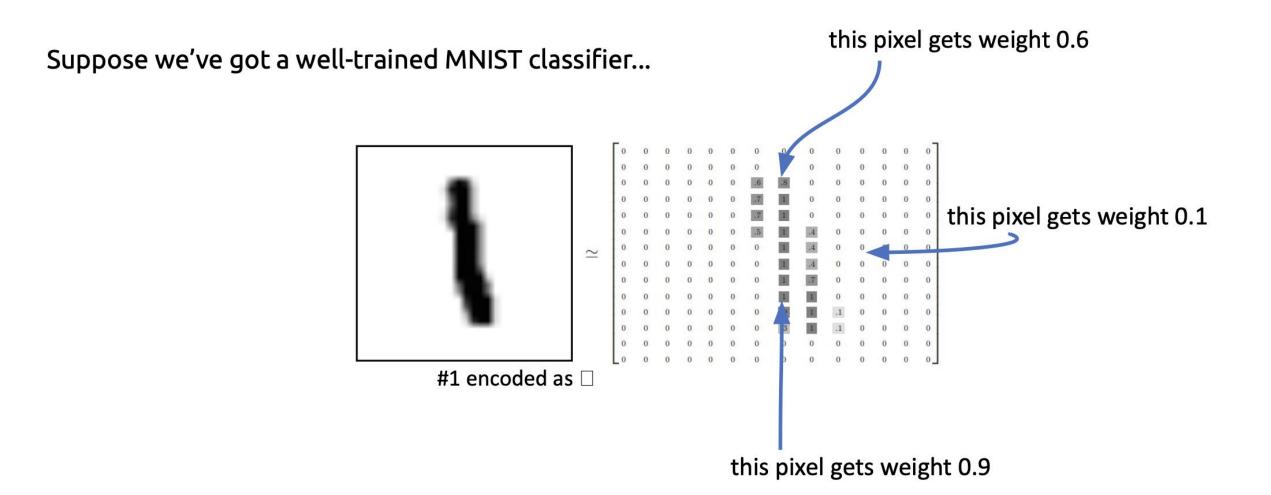
Isn't this actually a hard problem that we are trying to learn?

Limitations of Full Connections for MNIST

Suppose we've got a well-trained MNIST classifier...

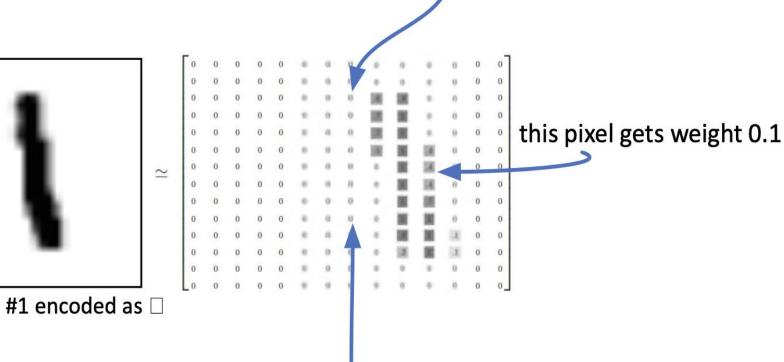


Limitations of Full Connections for MNIST



Limitations of Full Connections for MNIST

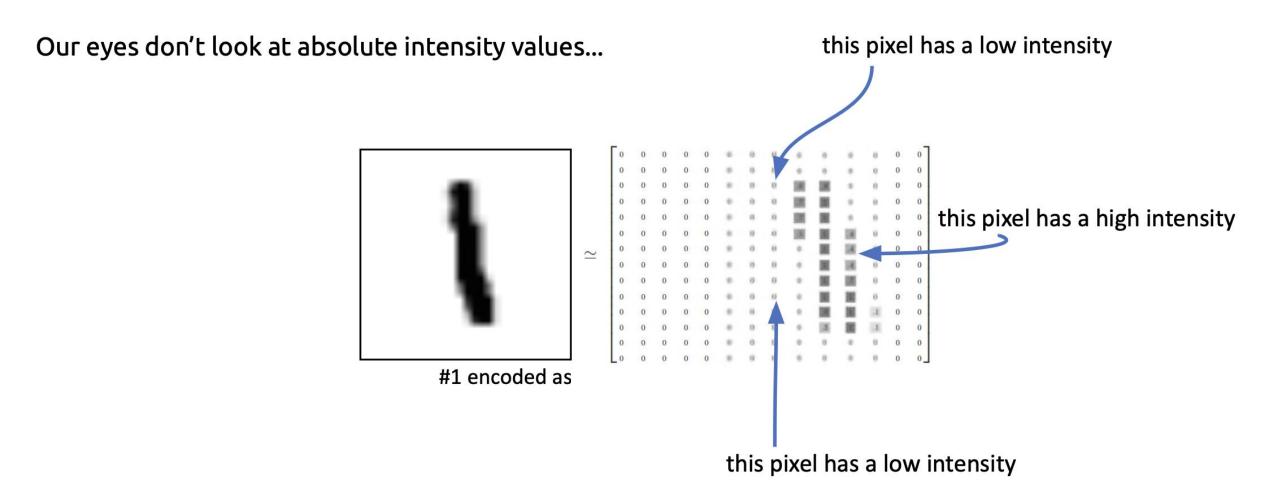
If we shift the digit to the right, then a different set of weights becomes relevant → network might have trouble classifying this as a 1...



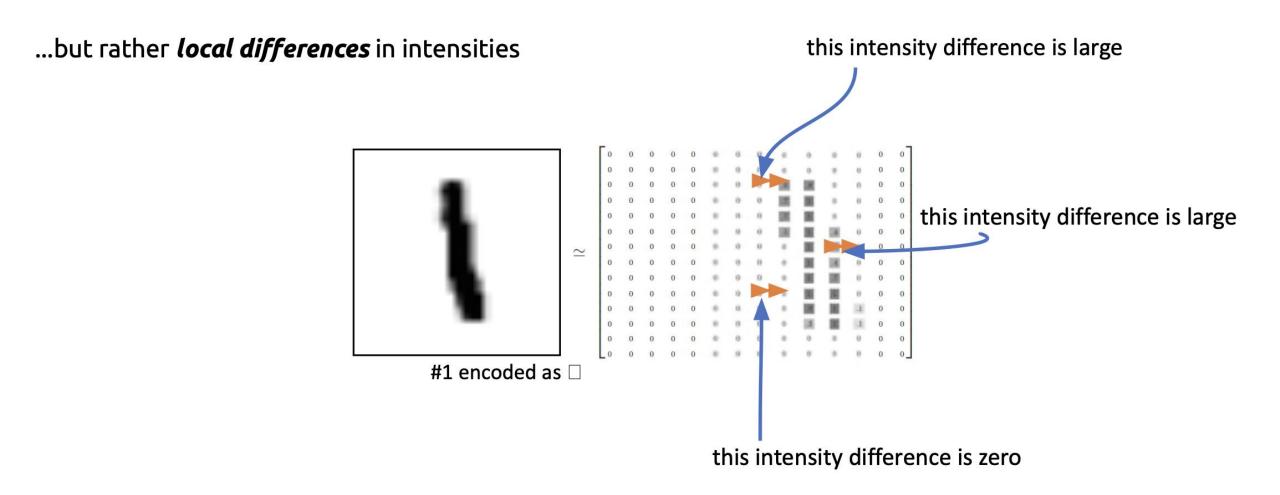
this pixel gets weight 0.9

this pixel gets weight 0.6

This would **not** be a problem for the human visual system

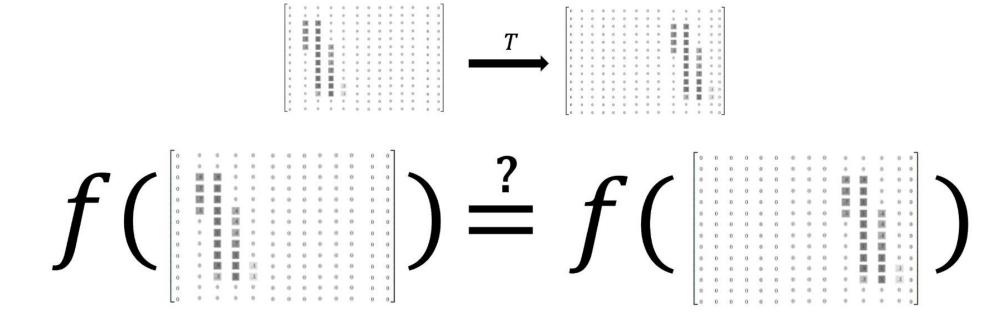


This would **not** be a problem for the human visual system

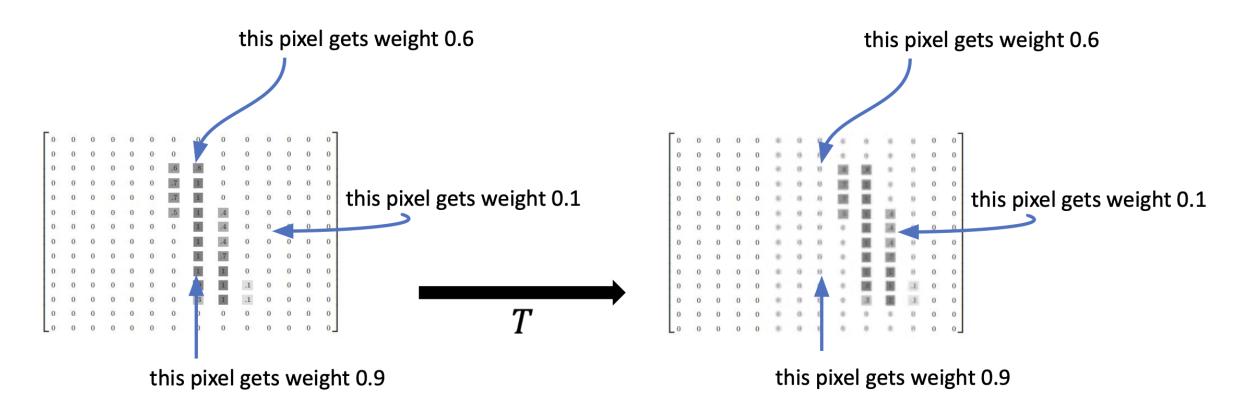


Translational Invariance

- 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)



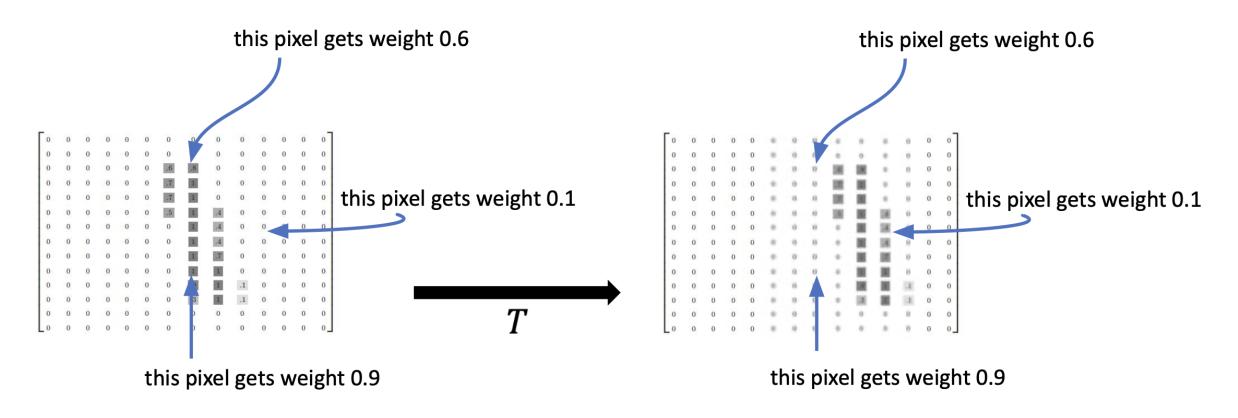
Fully Connected Nets are *not*Translationally Invariant



Sum of these three: $0.6 \cdot 0.8 + 0.1 \cdot 0 + 0.9 \cdot 1 = 1.38$

Sum of these three: $0.6 \cdot 0 + 0.1 \cdot 0.4 + 0.9 \cdot 0 = 0.4$

Fully Connected Nets are *not*Translationally Invariant



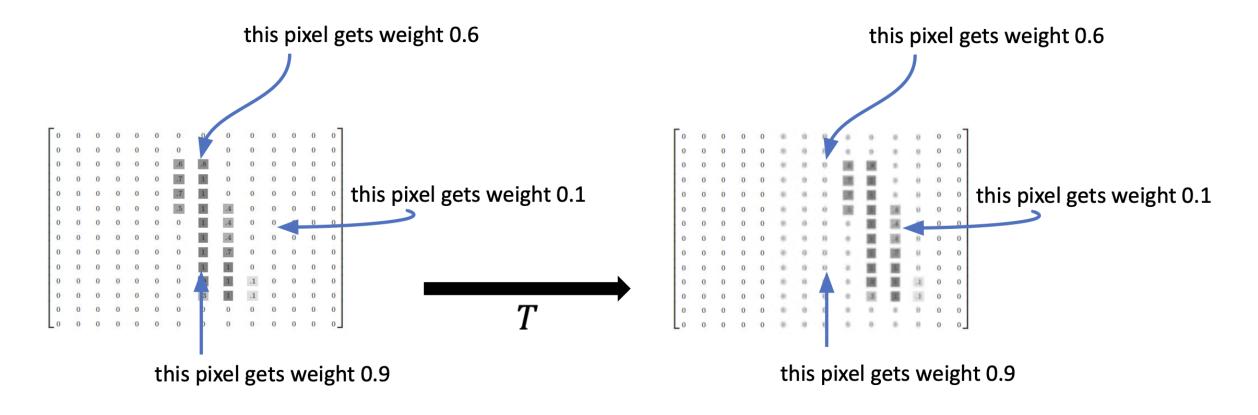
Sum of these three: $0.6 \cdot 0.8 + 0.1 \cdot 0 + 0.9 \cdot 1 = 1.38$

Sum of these three: $0.6 \cdot 0 + 0.1 \cdot 0.4 + 0.9 \cdot 0 = 0.4$

Fully Connected Nets are *not*Translationally Invariant

How to make the network translationally invariant?

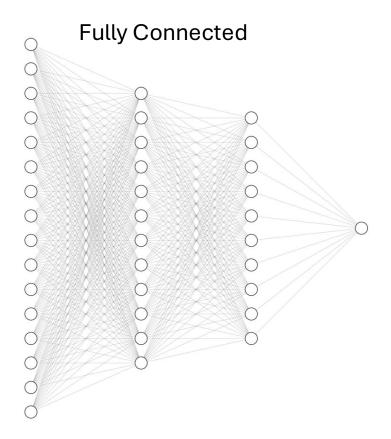
Focus on local differences/patterns



Sum of these three: $0.6 \cdot 0.8 + 0.1 \cdot 0 + 0.9 \cdot 1 = 1.38$

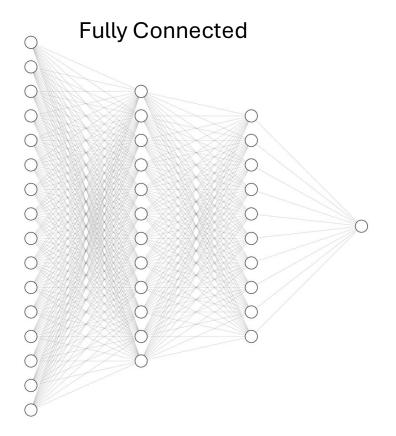
Sum of these three: $0.6 \cdot 0 + 0.1 \cdot 0.4 + 0.9 \cdot 0 = 0.4$

MLPs (also called fully-connected networks) have weights from every pixel to every neuron



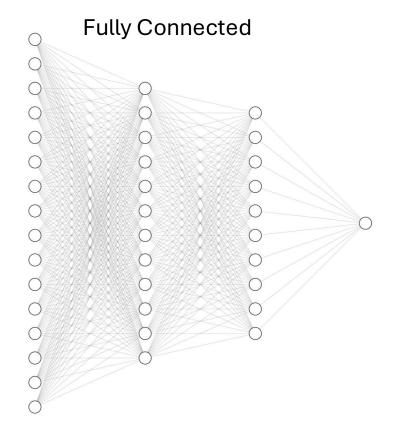
How can we change a fullyconnected network to account for spatial information?

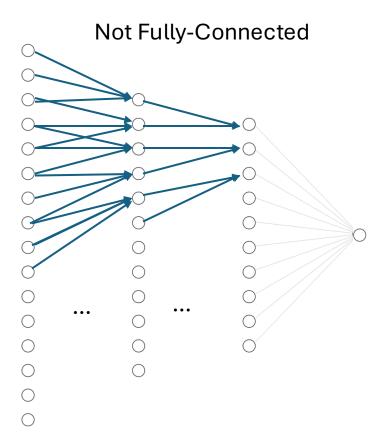
MLPs (also called fully-connected networks) have weights from every pixel to every neuron



How can we change a fullyconnected network to account for spatial information?

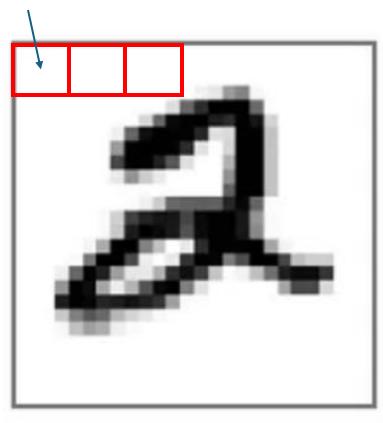
MLPs (also called fully-connected networks) have weights from every pixel to every neuron





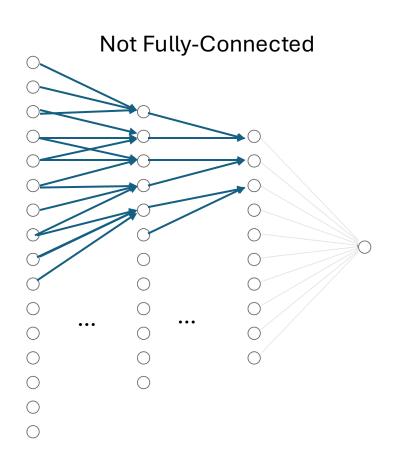
MLPs and Spatial Reasoning

Patches: Pixels close to each other



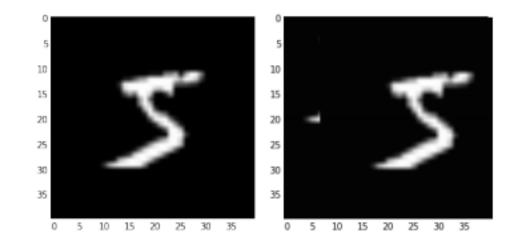
Advantages of Not Fully Connected Layers

- Fewer weights → Faster?
- The outputs of neurons are "features" for local "patches"
- Incorporates spatial information (pixels that are close together matter)



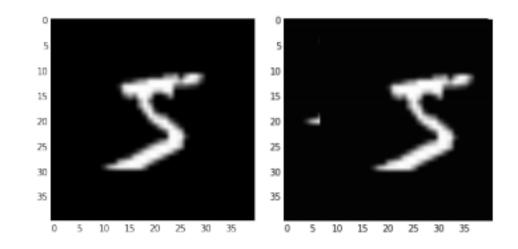
Disadvantages of Not Fully Connected Layers

- What happens if the image is Translated?
- The patches on the right side were never trained with 5's in that side.



Disadvantages of Not Fully Connected Layers

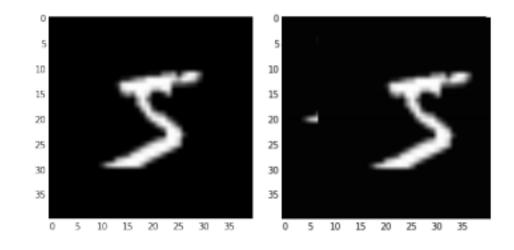
- What happens if the image is Translated?
- The patches on the right side were never trained with 5's in that side.



Even though we include spatial *information*, we still don't have spatial *reasoning*. (Can't recognize a shifted 5 is still a 5)

Disadvantages of Not Fully Connected Layers

- What happens if the image is Translated?
- The patches on the right side were never trained with 5's in that side.



Even though we include spatial *information*, we still don't have spatial *reasoning*. (Can't recognize a shifted 5 is still a 5)

What if we used the same weights for each patch? (Weight Sharing)

The Main Building Block: Convolution

Convolution is an operation that takes two inputs:

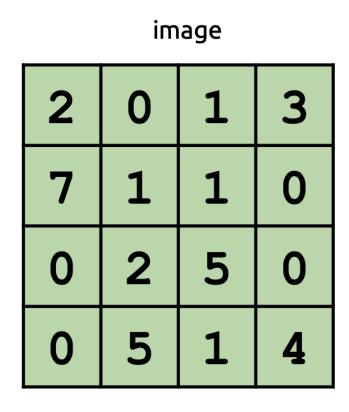
(1) An image (2D - B/W)

(2) A filter (also called a kernel)

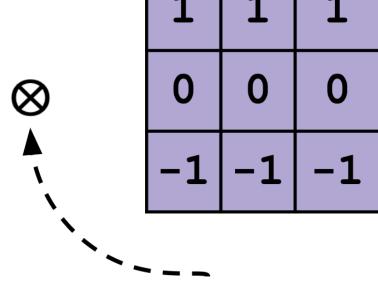


1	1	1		
0	0	0		
-1	-1	-1		

2D array of numbers; could be any values

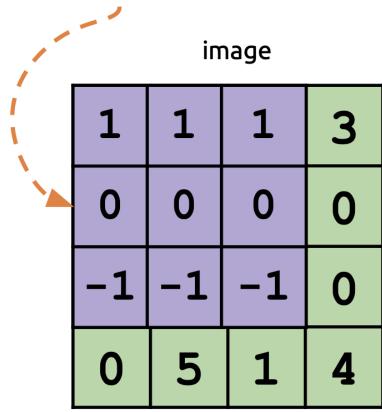


filter/kernel

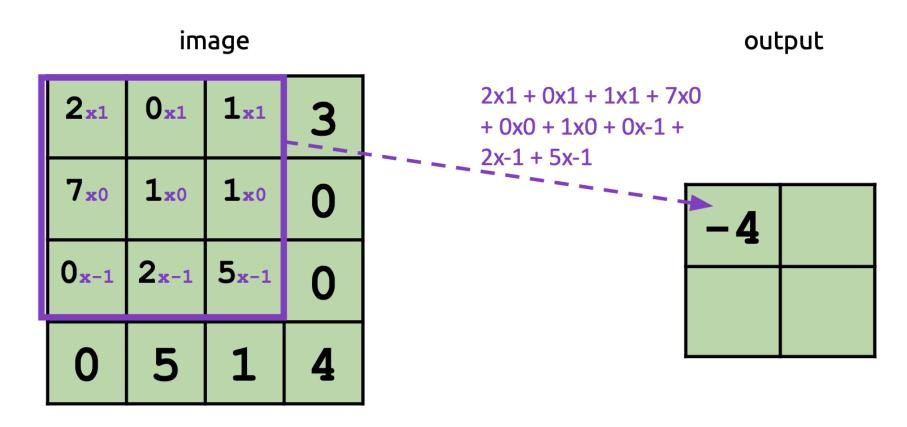


(We use this symbol for convolution) (The verb form is "convolve")

Overlay the filter on the image



Sum up multiplied values to produce output value

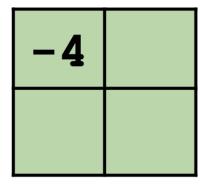


Move the filter over by one pixel

image

1	1	1	3	
0	0	0	0	
-1	-1	-1	0	
0	5	1	4	

output

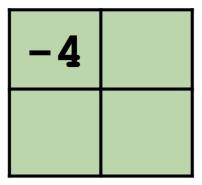


Move the filter over by one pixel

image

2	1	1	1
7	0	0	0
0	-1	-1	-1
0	5	1	4

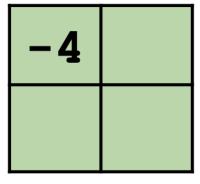
output



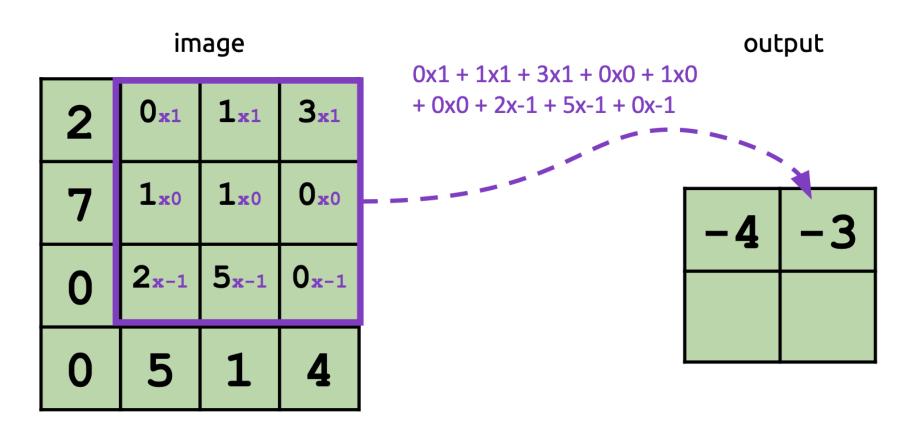
Repeat (multiply, sum up)

image

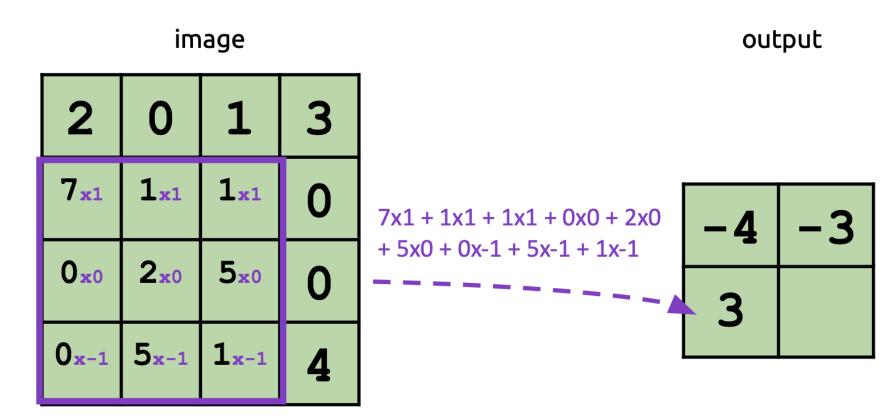
 output



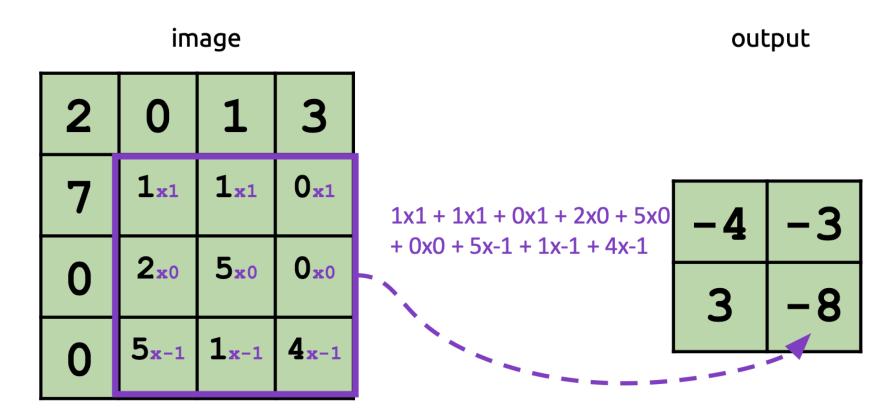
Repeat (multiply, sum up)



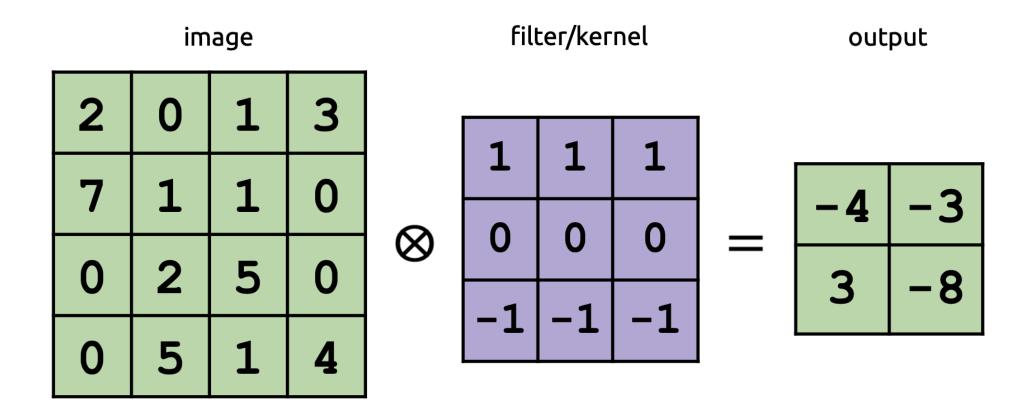
Repeat...



Repeat...



In summary:



Handmade Kernels and Filters

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

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

$$\left[egin{array}{ccc} 0 & -1 & 0 \ -1 & 5 & -1 \ 0 & -1 & 0 \ \end{array}
ight]$$

Identity kernel

Edge detection

Sharpen kernel

$$\frac{1}{9} \left[\begin{array}{rrr} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{array} \right]$$

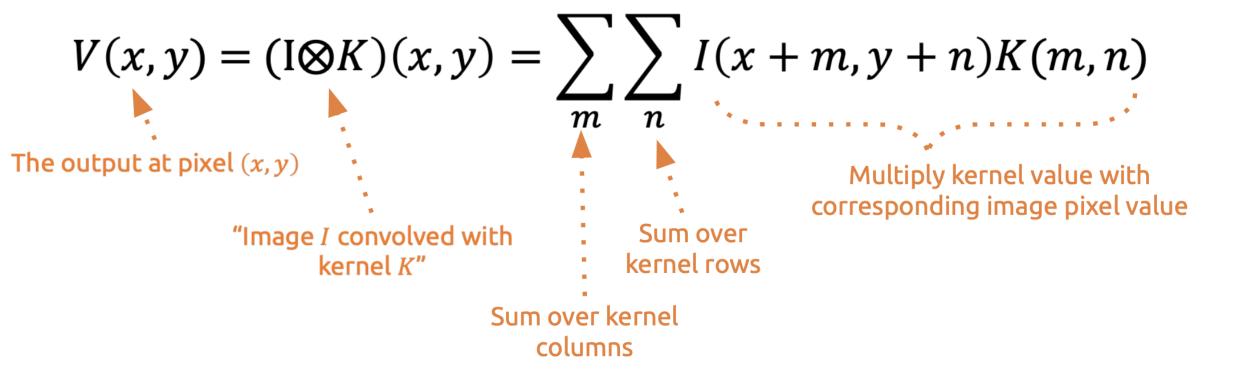
Box blur

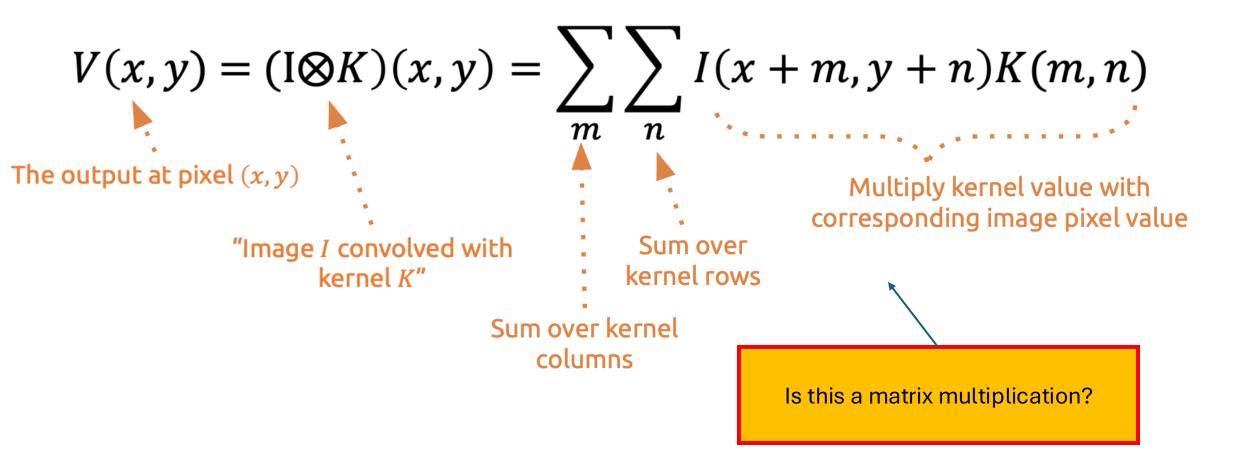
$$\frac{1}{256} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$$

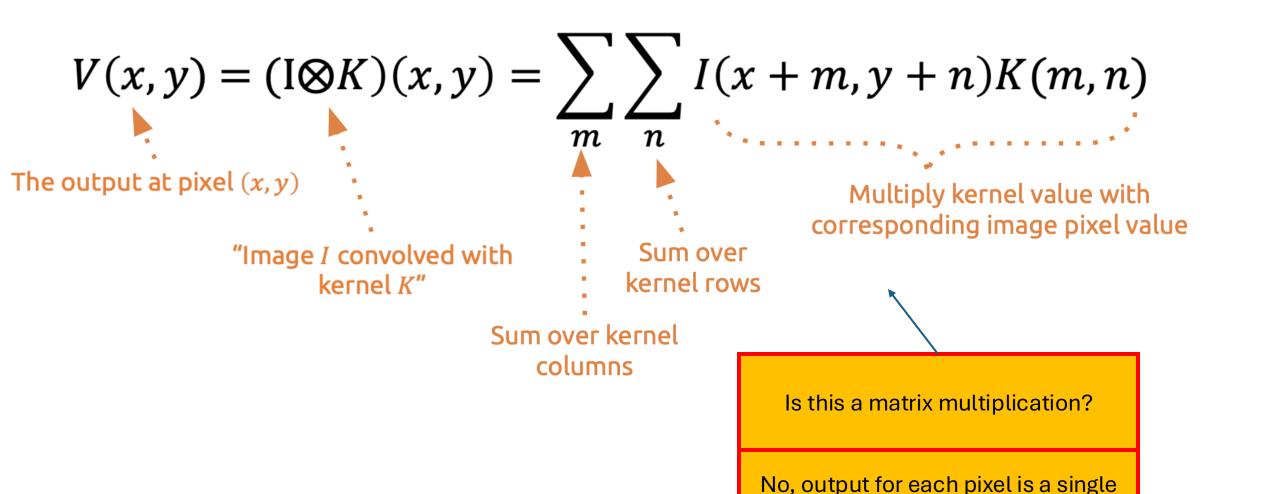
Gaussian blurr kernel

Operation Image result g(x,y) Kernel ω Identity Ridge or edge detection Sharpen **Box blur** (normalized) Gaussian blur 3 x 3 (approximation)

Source: Wikipedia



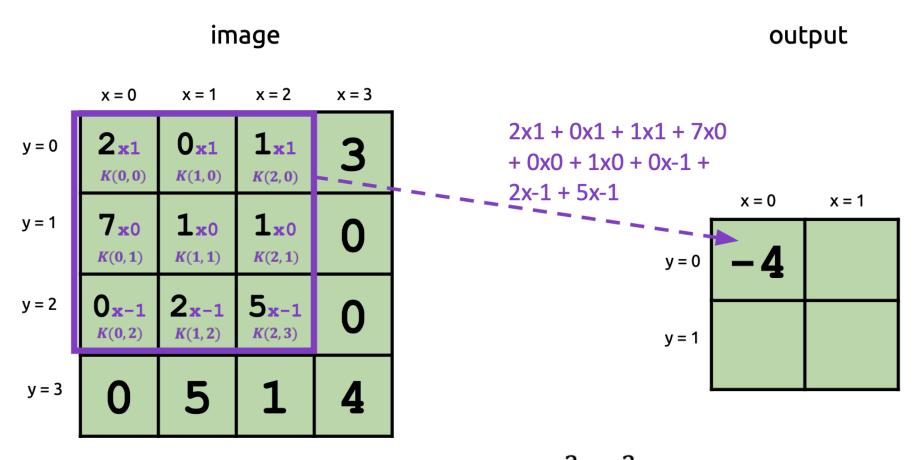




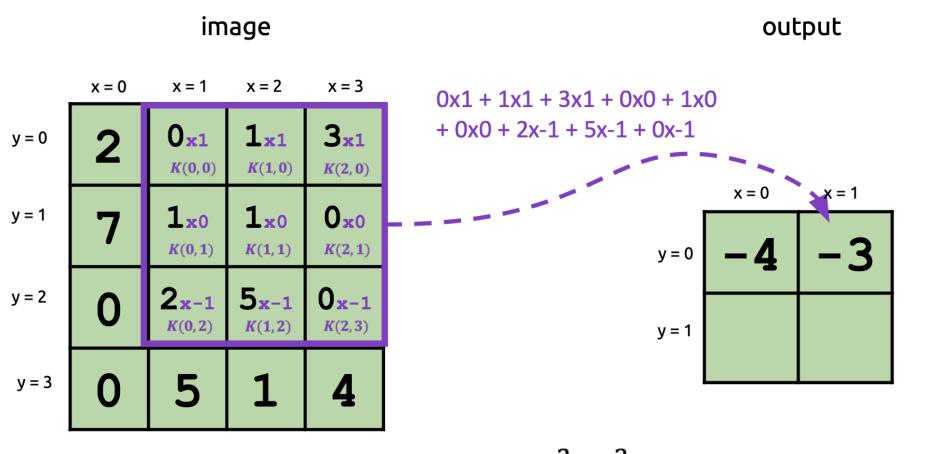
value, not a matrix.

image

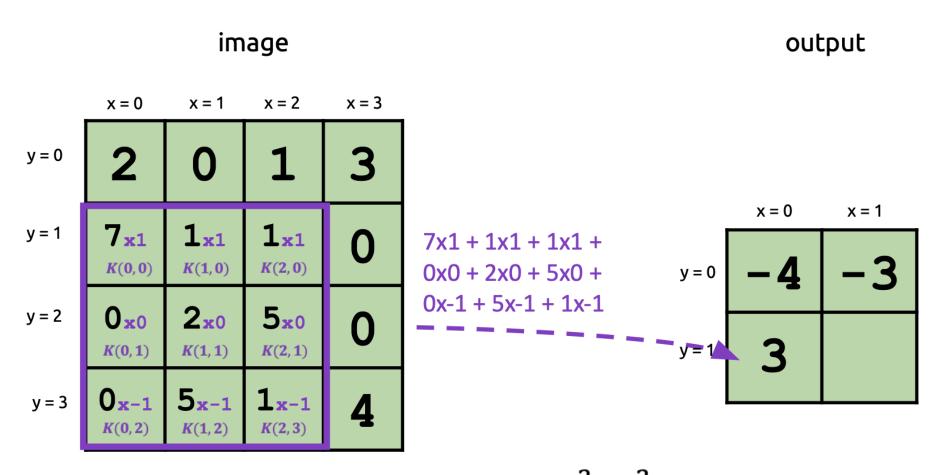
	x = 0	x = 1	x = 2	x = 3	_	fill	er/ker		out	put	
y = 0	2	0	1	3		m = 0	m = 1	m = 2	•		
					n = 0	1	1	1		x = 0	x = 1
y = 1	7	1	1	0		1	•	1		A	
				0	n=1)	0	0	y = 0	-4	-3
y = 2	0	2	5	0	\otimes	0	U	U	_		
)	4	ว	כ	n = 2	7	7	1	y = 1	3	-8
y = 3	>	П	1	A				_T			
	0	כ		4	'						



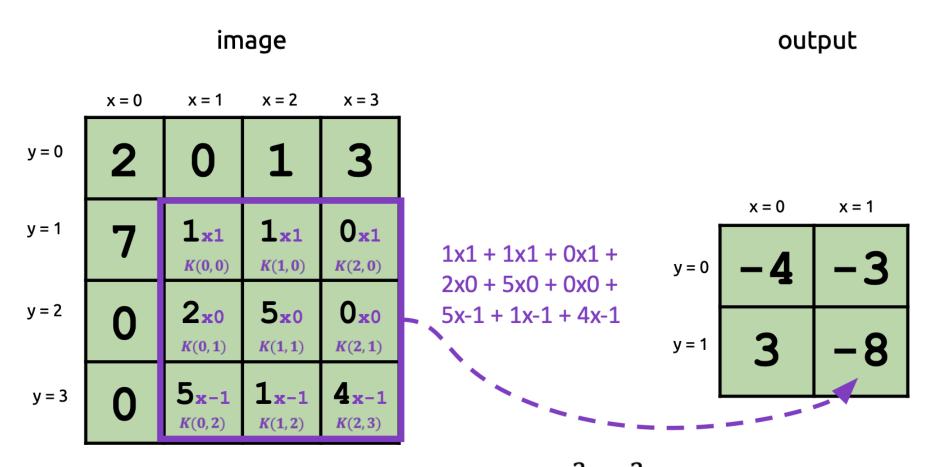
$$V(0,0) = (I \otimes K)(0,0) = \sum_{m=0}^{2} \sum_{n=0}^{2} I(0+m,0+n)K(m,n)$$



$$V(1,0) = (I \otimes K)(1,0) = \sum_{m=0}^{2} \sum_{n=0}^{2} I(1+m,0+n)K(m,n)$$



$$V(0,1) = (I \otimes K)(0,1) = \sum_{m=0}^{2} \sum_{n=0}^{2} I(0+m,1+n)K(m,n)$$



$$V(1,1) = (I \otimes K)(1,1) = \sum_{m=0}^{2} \sum_{n=0}^{2} I(1+m,1+n)K(m,n)$$

What Convolution Does (In Code)

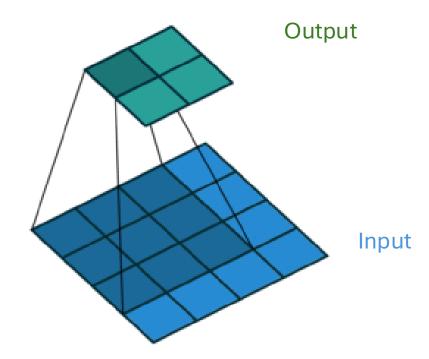
```
// Input: Image I, Kernel K, Output V, pixel index x,y
// Assumes K is 3x3
function apply_kernel(I, K, V, x, y)
  for m = 0 to 2:
    for n = 0 to 2:
        V(x,y) += K(m,n) * I(m+x, n+y)
```

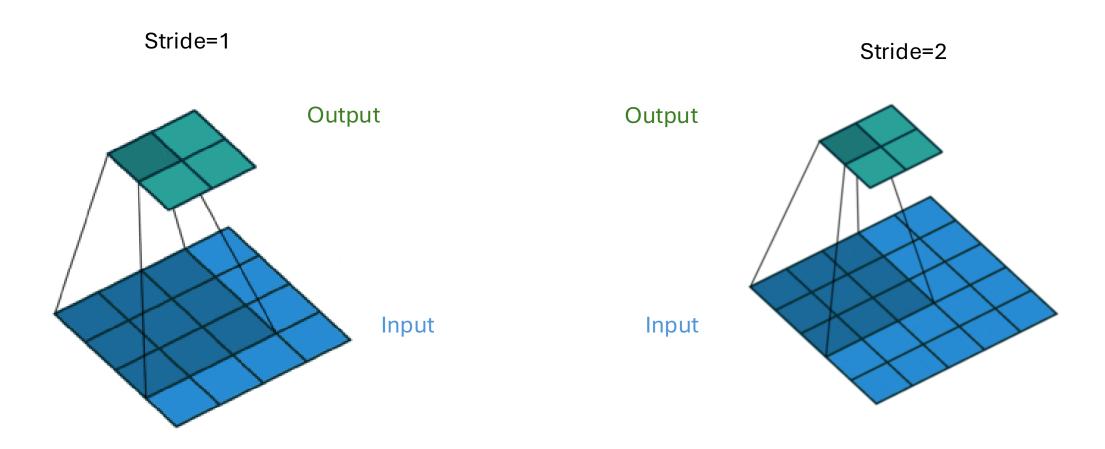
Equation:
$$V(x,y) = (I \otimes K)(x,y) = \sum_{m} \sum_{n} I(x+m,y+n)K(m,n)$$

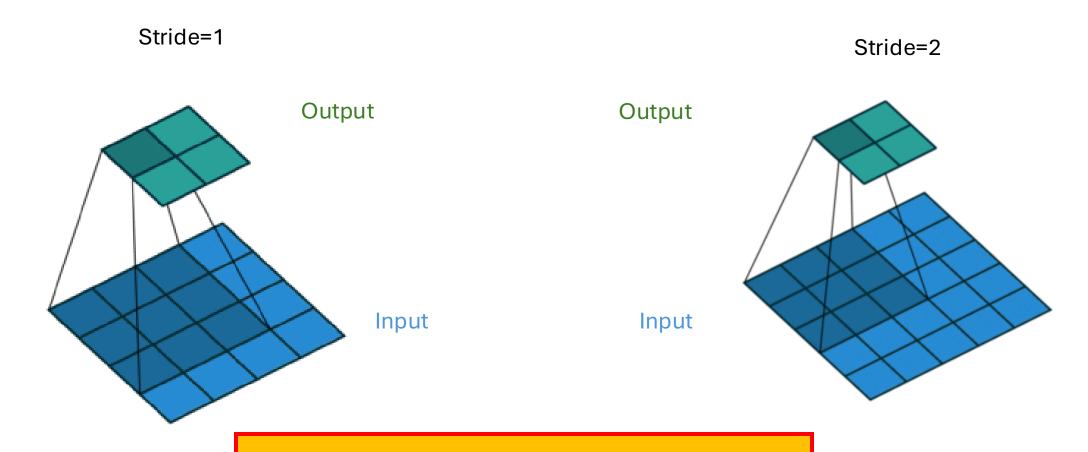
- We don't just have to slide the filter by one pixel every time
- The distance we slide a filter by is called *stride*
 - All the examples we've seen thus far have been stride = 1

1	1	1	3		2	1	1	1
0	0	0	0	stride = 1	7	0	0	0
-1	-1	-1	0		0	-1	-1	-1
0	5	1	4		0	5	1	4

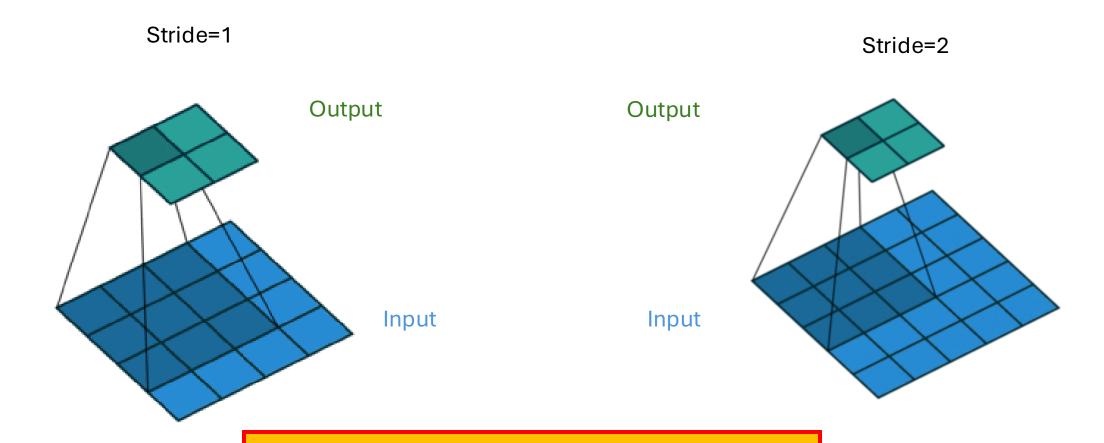
Stride=1



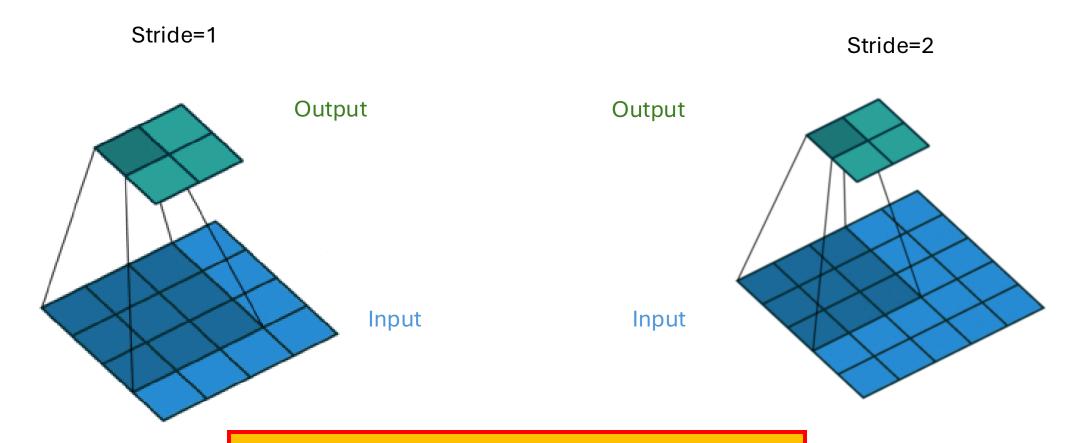




Any connection between input and output size?

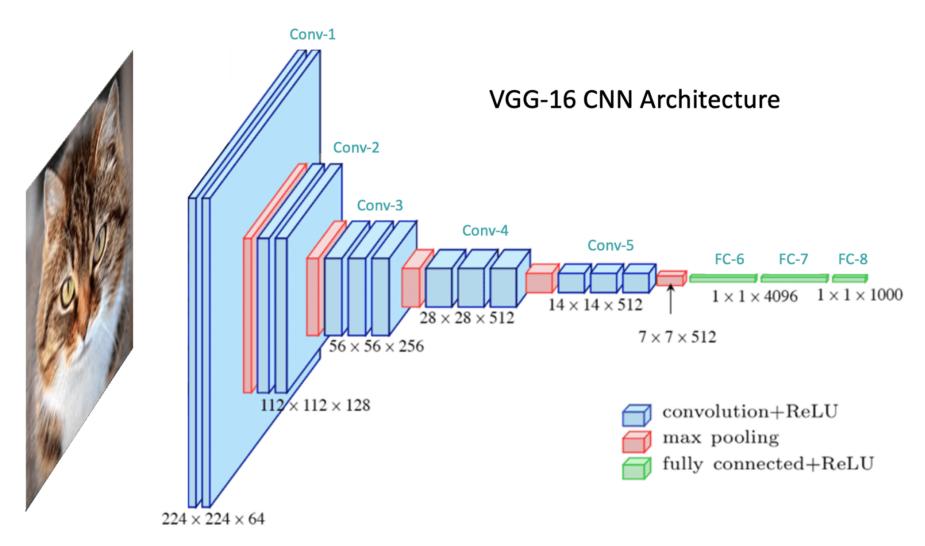


Larger stride turns larger input into same sized output



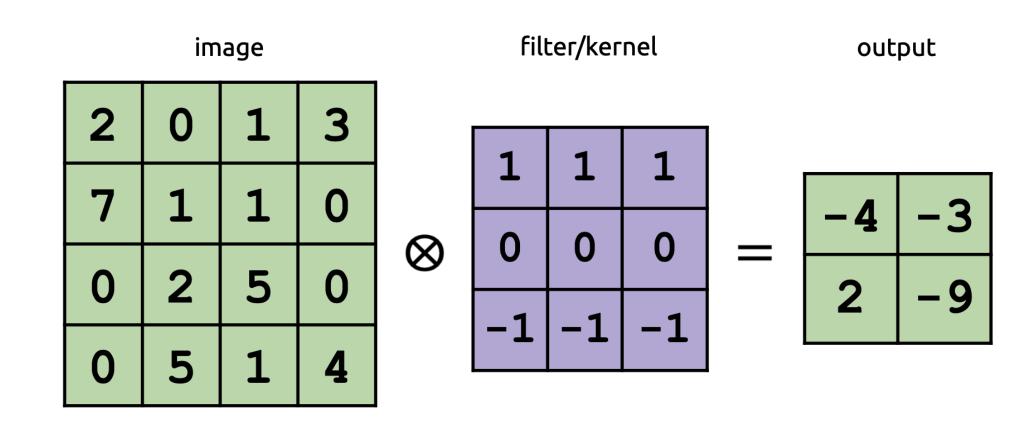
Larger stride turns **same** sized input into **smaller** sized output

Convolutional Neural Networks

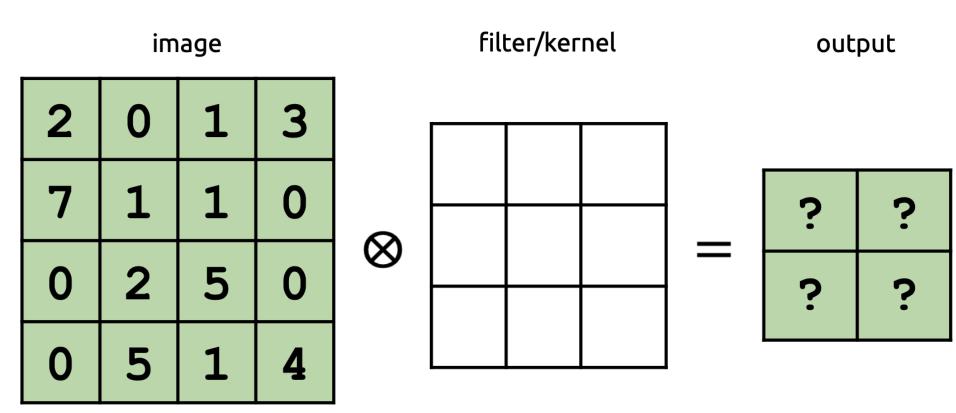


https://learnopencv.com/understanding-convolutional-neural-networks-cnn/

Key Idea 1: Filters are *Learnable*

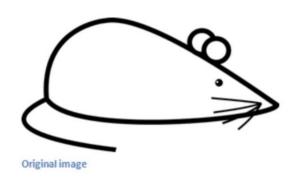


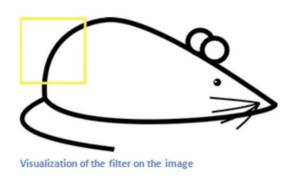
Key Idea 1: Filters are *Learnable*



 $k_{i,j}$ are learnable parameters

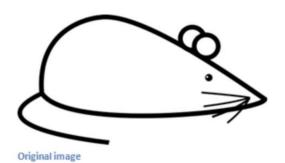
Key Idea 1: Filters are *Learnable*

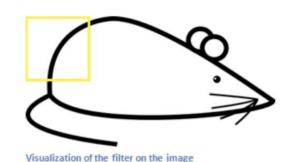


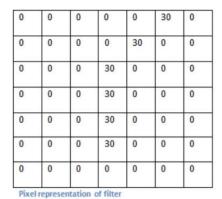


Label="Mouse"

Detecting patterns using learned filters









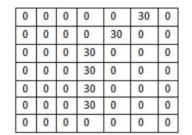
Visualization of a curve detector filter



Visualization of the receptive field

0	0	0	0	0	0	30
0	0	0	0	50	50	50
0	0	0	20	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0

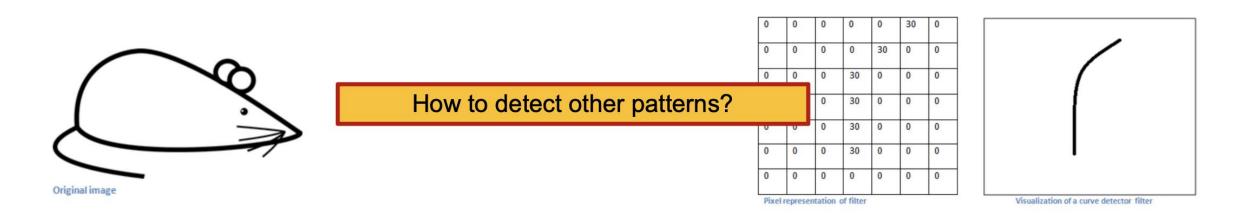
Pixel representation of the receptive field



Pixel representation of filter

Multiplication and Summation = (50*30)+(50*30)+(50*30)+(50*30)+(50*30)=6600 (A large number!)

Detecting patterns using learned filters





0	0	0	0	0	0	0
0	40	0	0	0	0	0
40	0	40	0	0	0	0
40	20	0	0	0	0	0
0	50	0	0	0	0	0
0	0	50	0	0	0	0
25	25	0	50	0	0	0



0	0	0	0	0	30	0
0	0	0	0	30	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	0	0	0	0

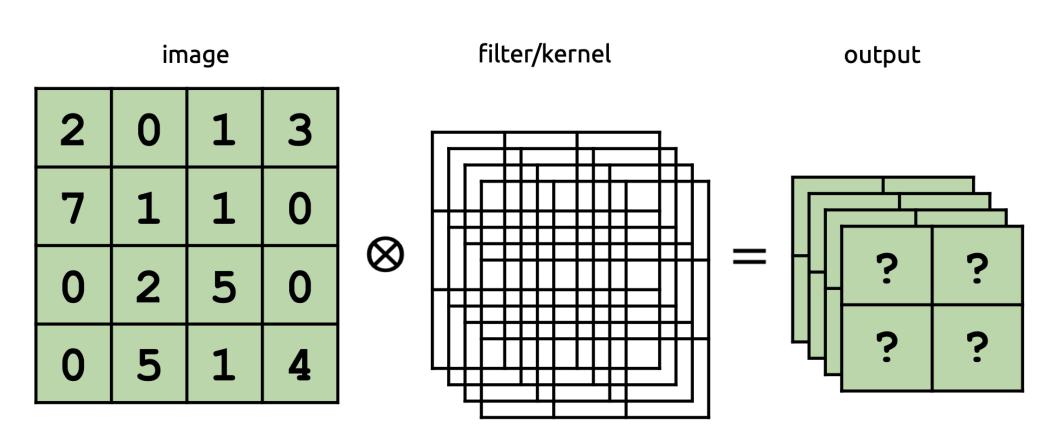
Visualization of the filter on the image

Pixel representation of receptive field

Pixel representation of filter

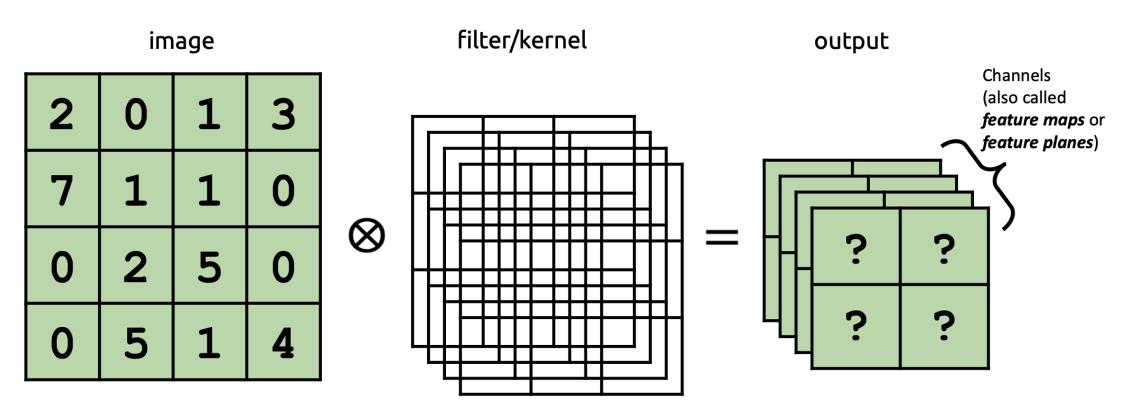
Multiplication and Summation = 0

Key Idea 2: Learn *many* filters



This block of filters is called a *filter bank*

Key Idea 2: Learn *many* filters

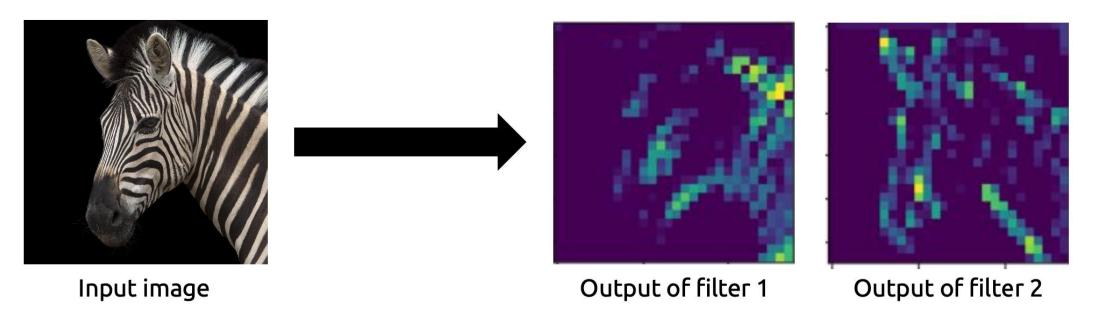


The output is now a multi-channel image



Key Idea 2: Learn *many* filters

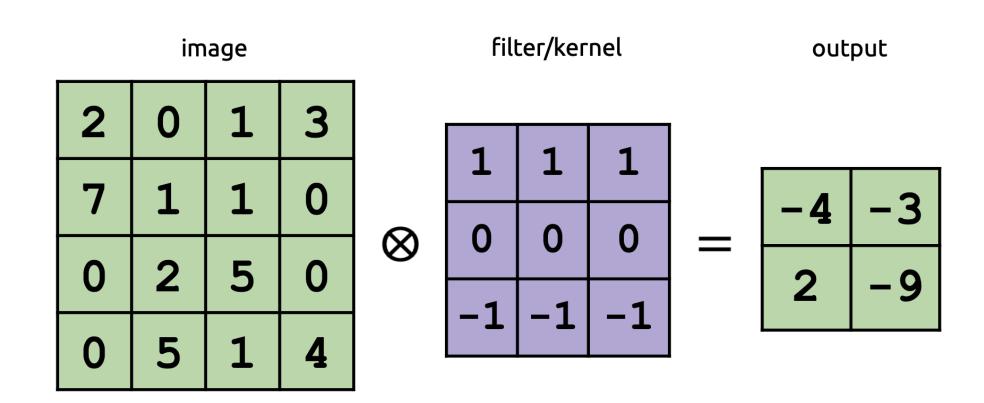
- Why are multiple filters a good idea?
 - Can learn to extract different features of the image



How is convolution "partially connected?"

Fully Connected Partially Connected

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



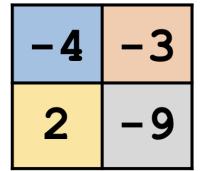
Only certain input pixels are "connected" to certain output pixels

image

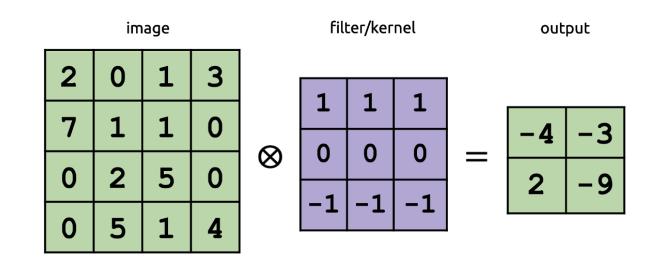
Colored dots in the input pixels represent which output pixels that input pixel contributes to

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

output



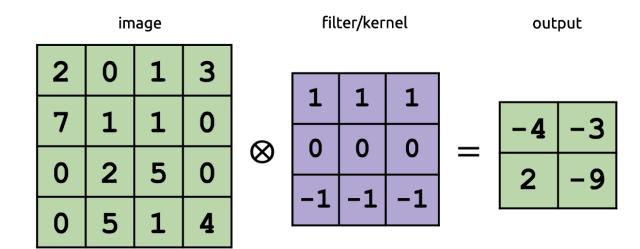
tf.nn.conv2d(input, filter, stride, padding)



Documentation: https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/nn/conv2d

tf.nn.conv2d(input, filter, stride, padding)

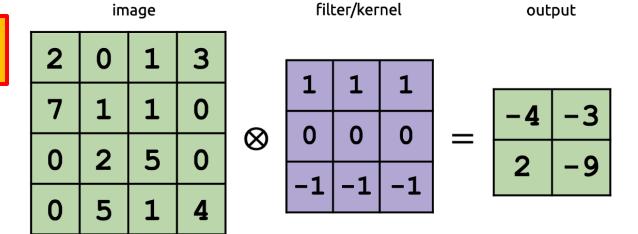




tf.nn.conv2d(input, filter, stride, padding)

Input "Image" Kernel/filter

What is the shape of the input?
Tensor of:[# items in batch, width, height, # channels]



tf.nn.conv2d(input, filter, stride, padding)



What is the shape of the input?
Tensor of:[# items in batch, width, height, # channels]

Channels: Number of input "colors"

image			filter/kernel			out	put		
2	0	1	3						
7	1	1	0		1	1	1	-4	-3
0	2	5	0	\otimes	0	0	0	2	-9
0	5	1	4		_T	_T	_T		

Documentation: https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/nn/conv2d

tf.nn.conv2d(input, filter, stride, padding)

Input "Image" Kernel/filter

What is the shape of the input?
Tensor of:[# items in batch, width, height, # channels]

Channels: Number of input "colors"

 2
 0
 1
 3

 7
 1
 1
 0

 0
 2
 5
 0

 0
 5
 1
 4

image

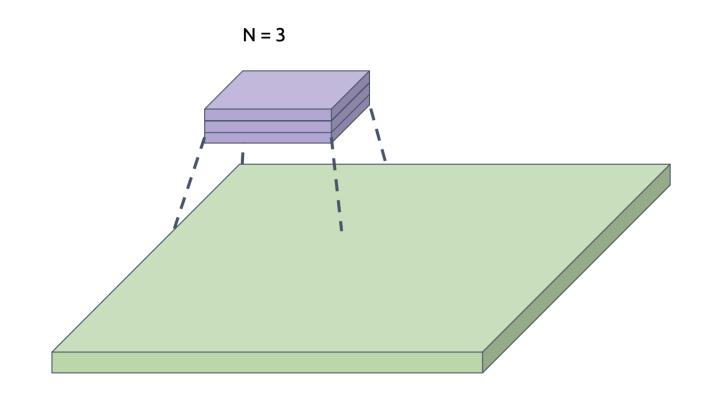
output

filter/kernel

How can we determine the output size of a convolution?

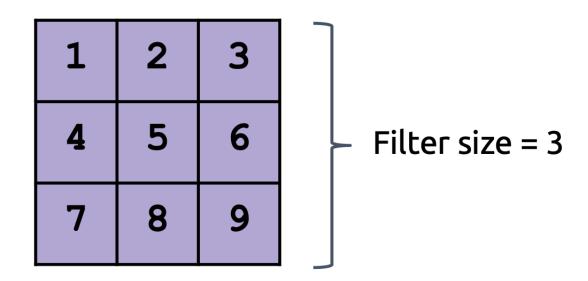
The output size of a convolution layer depends on 4 Hyperparameters:

Number of filters, N



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

- Number of filters, N
- The size of these filters, F



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

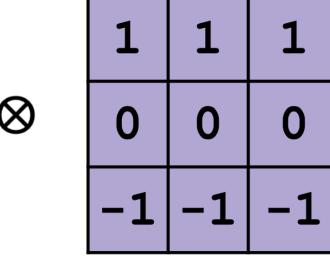
- Number of filters, N
- The size of these filters, F
- The stride, **S**

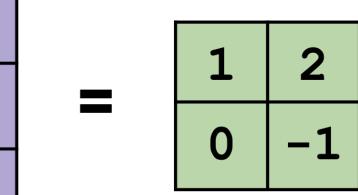
2	0	3	1	0
2	4	5	2	3
0	0	3	3	1
2	9	9	7	8
3	4	7	2	1

2	0	3	1	0
2	4	5	2	3
0	0	3	3	1
2	9	9	7	8
3	4	7	2	1

"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

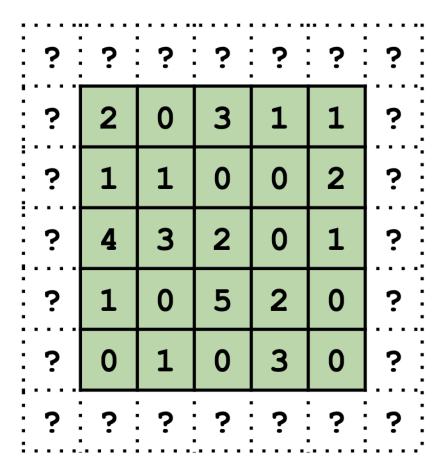
2	0	3	1	1
1	1	0	0	2
4	3	2	0	1
1	0	5	2	0
0	1	0	3	0

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	2 01 14 31 00 1	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

					.		
:	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

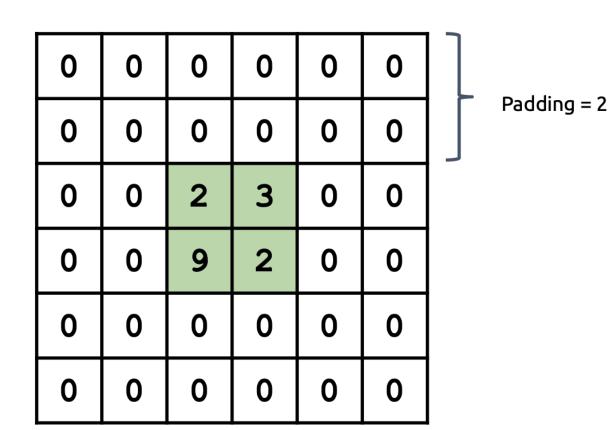
0	0	0	0	0	0
0	2	0	1	3	0
0	1	1	2	3	0
0	4	3	2	1	0
0	8	3	1	3	0
0	0	0	0	0	0

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

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



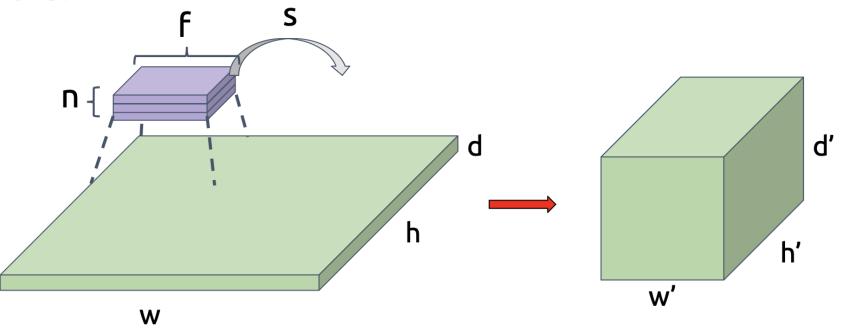
Suppose we know the number of filters, their size, the stride, and padding (n,f,s,p).

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
$$w = 4$$

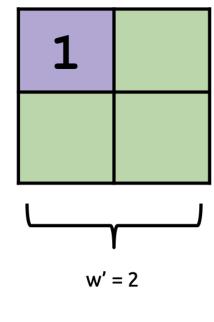
num filters
$$n = 1$$

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

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

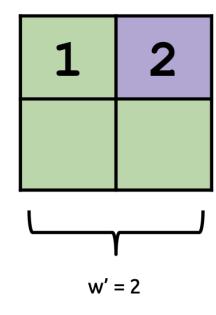
$$w' = \frac{w - f + 2p}{s} + 1$$

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

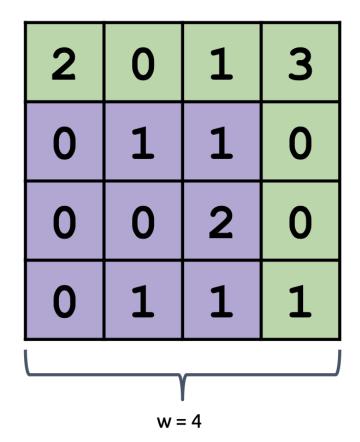


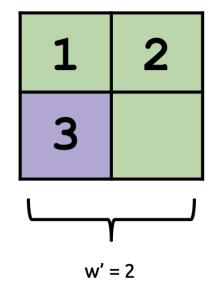
$$w' = \frac{w - f + 2p}{s} + 1$$

2	0	1	თ	
0	1	1	0	
0	0	2	0	
0	1	1	1	
v = 4				

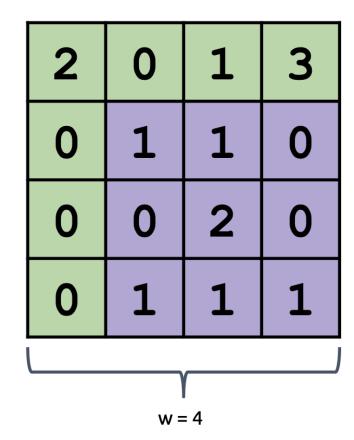


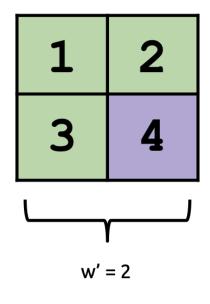
$$w' = \frac{w - f + 2p}{s} + 1$$





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





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

num filters
$$n = 1$$

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

Let
$$w = 4$$

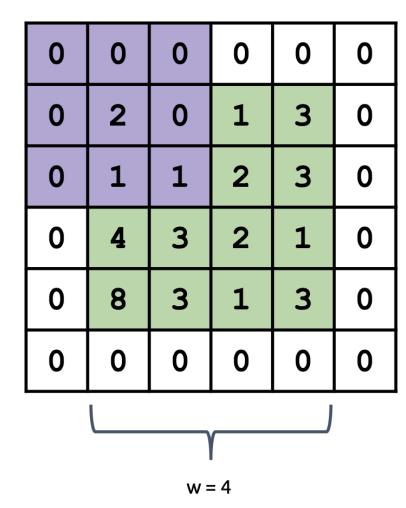
$$w' = \frac{4 - 3 + 2 \cdot 1}{1} + 1$$

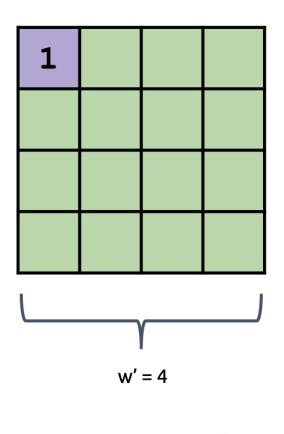
$$= 3 + 1 = 4$$

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

num filters
$$n = 1$$

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

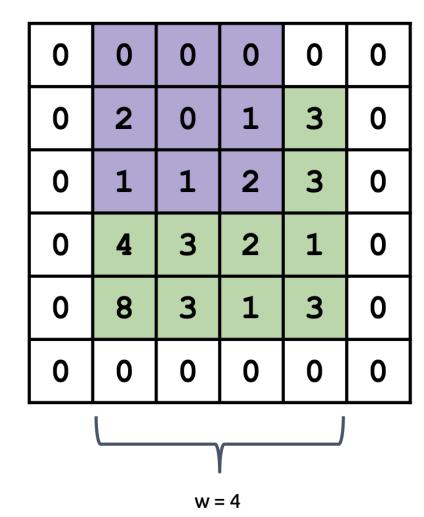


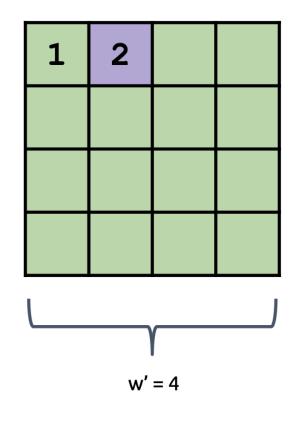


$$w' = \frac{w - f + 2p}{s} + 1$$

num filters
$$n = 1$$

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



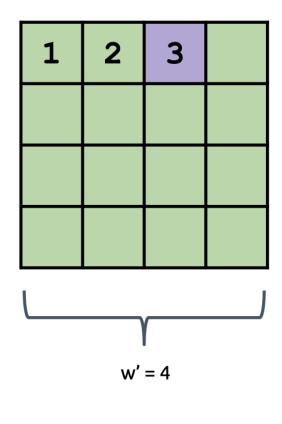


$$w' = \frac{w - f + 2p}{s} + 1$$

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
<u> </u>					
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

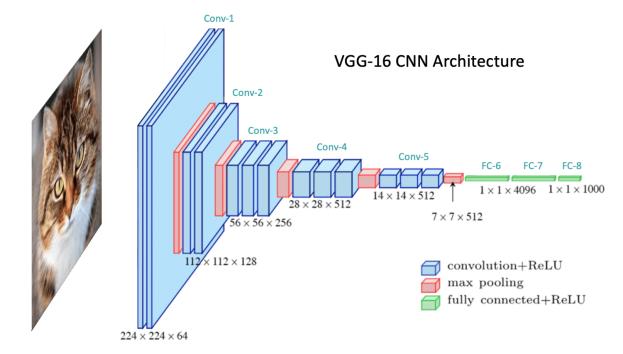
1	2	3	4	

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?

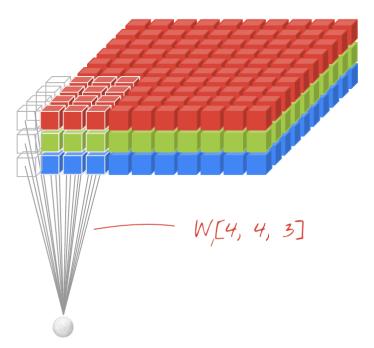


https://learnopencv.com/understa nding-convolutional-neuralnetworks-cnn/

Color images...

Remaining Question: 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