# Deep Learning

CSCI 1470
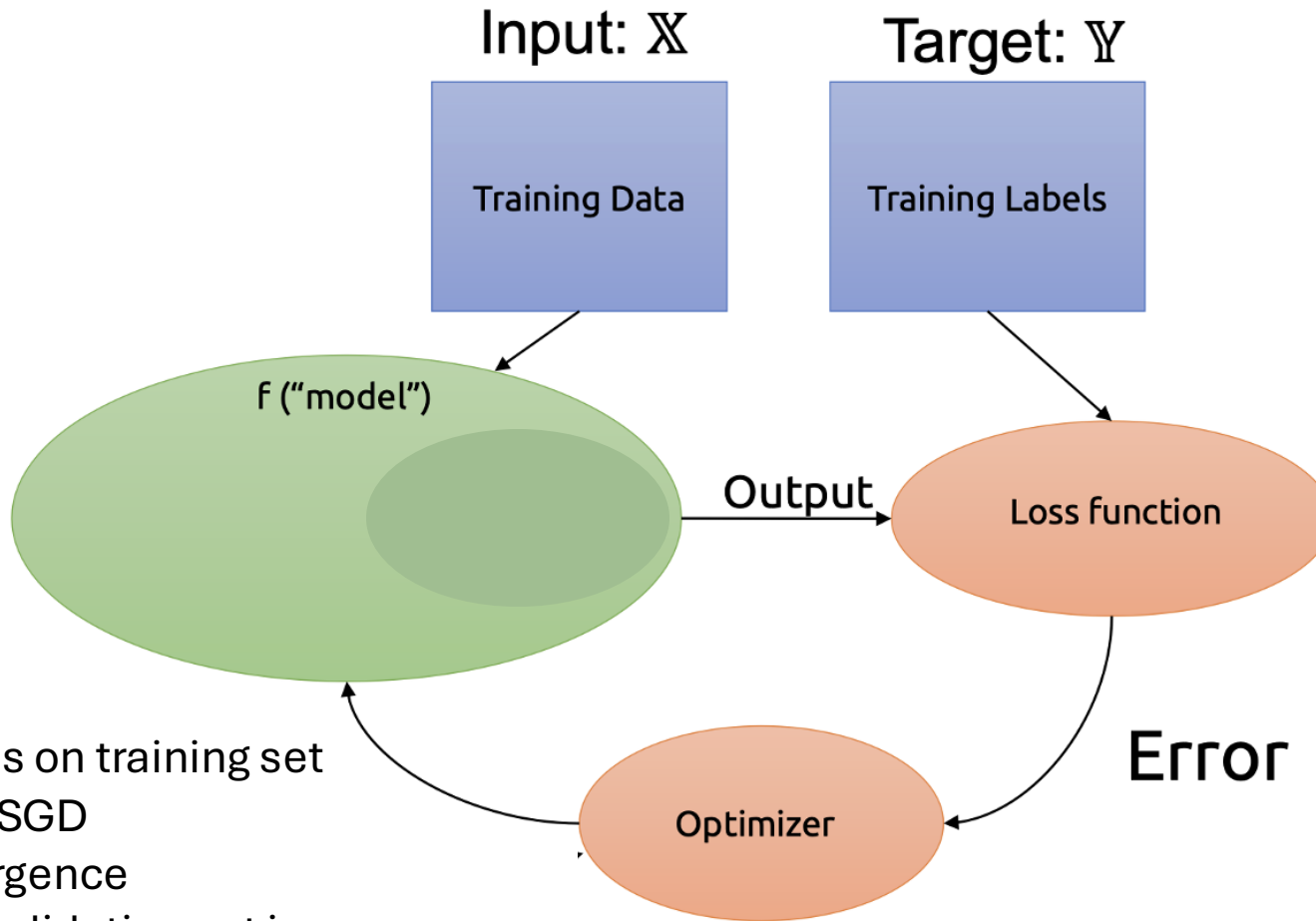
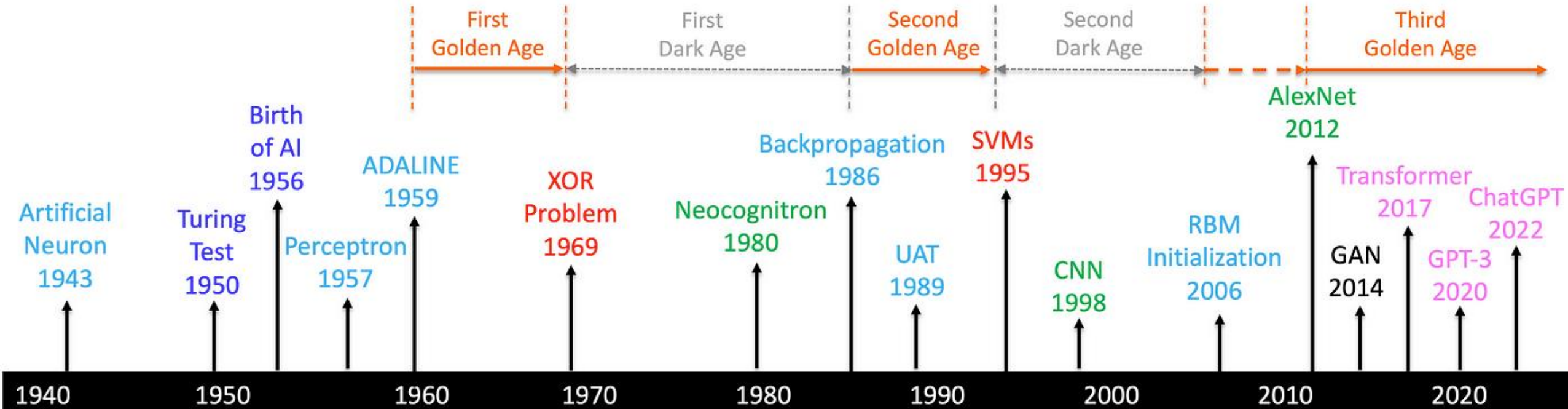Eric Ewing

Wednesday, 2/12/25

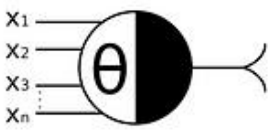Day 10: Introduction to Convolutions

# Recap: MLPs



1. Compute Error/Loss on training set
2. Run Backprop and SGD
3. Repeat until convergence
4. If performance on validation set is acceptable, terminate, else try new hyperparameters
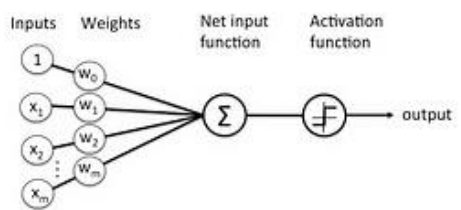
# A Brief History of AI with Deep Learning

# 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!



- *Graphics* Processing Units

- 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



output

Write back

ALU    **Arithmetic logic unit**

Decode

Fetch

input

output

CPU

input

**FETCH** the instruction and any data from main memory.

**DECODE** the instruction. (Covert the instruction into a language the CPU understands.)

**EXECUTE** the instruction. (Complete the instruction)

# CPU v/s GPU



Vector operations (SSE / AVX)

# 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

**CUDA compute model**



https://www.researchgate.net/publication/236666656_Accelerating_Fibre_Orientation_Estimation_from_Diffusion_Weighted_Magnetic_Resonance_Imaging_Using_GPUs

14

# 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

# GPU-Parallel Acceleration

Any questions?

**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

https://www.researchgate.net/publication/236666656_Accelerating_Fibre_Orientation_Estimation_from_Diffusion_Weighted_Magnetic_Resonance_Imaging_Using_GPUs

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



Assassin's Creed Valhalla | 2560x1440 | Ultra High | DX12

- 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

# CUDA is Still a Giant Moat for NVIDIA

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

JAMES WANG
MAR 23, 2024

# Issues with MLPs

1.  Resource Intensive

2.  Difficult to incorporate certain types of information

# MLPs and Spatial Reasoning



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

(0, 0)
(1, 0)
(2, 0)
(3, 0)
...

Image is transformed to vector of pixels

# MLPs and Spatial Reasoning

(24, 5)
(13, 2)
(4, 0)
(21, 11)
...

Image is transformed to vector of pixels

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

# MLPs and Spatial Reasoning



(24, 5)
(13, 2)
(4, 0)
(21, 11)
...

Image is transformed to vector of pixels

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...



#1 encoded as □

# Limitations of Full Connections for MNIST

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

this pixel gets weight 0.6

this pixel gets weight 0.1



#1 encoded as ▢

this pixel gets weight 0.9

# Limitations of Full Connections for MNIST

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

this pixel gets weight 0.6

this pixel gets weight 0.1



#1 encoded as ☐

this pixel gets weight 0.9

Can you tell this is a 1?

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

Our eyes don't look at absolute intensity values...

this pixel has a low intensity

this pixel has a high intensity

#1 encoded as

this pixel has a low intensity

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

...but rather *local differences* in intensities



#1 encoded as □

this intensity difference is large

this intensity difference is large

this intensity difference is zero

# 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)



$$ f\left( \quad \right) \overset{?}{=} f\left( \quad \right) $$

# Fully Connected Nets are *not* Translationally Invariant



this pixel gets weight 0.6

this pixel gets weight 0.1

this pixel gets weight 0.9

$T$

this pixel gets weight 0.6

this pixel gets weight 0.1

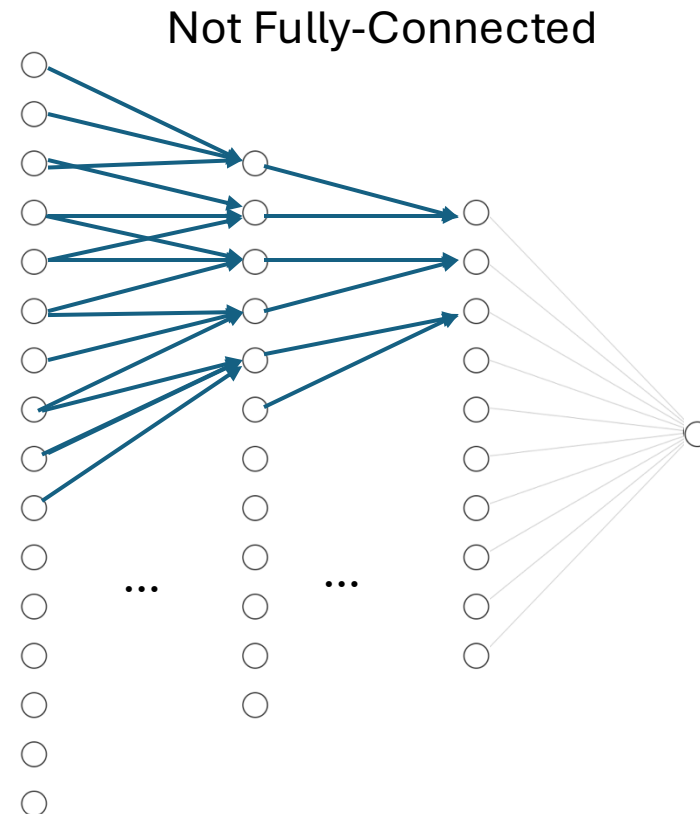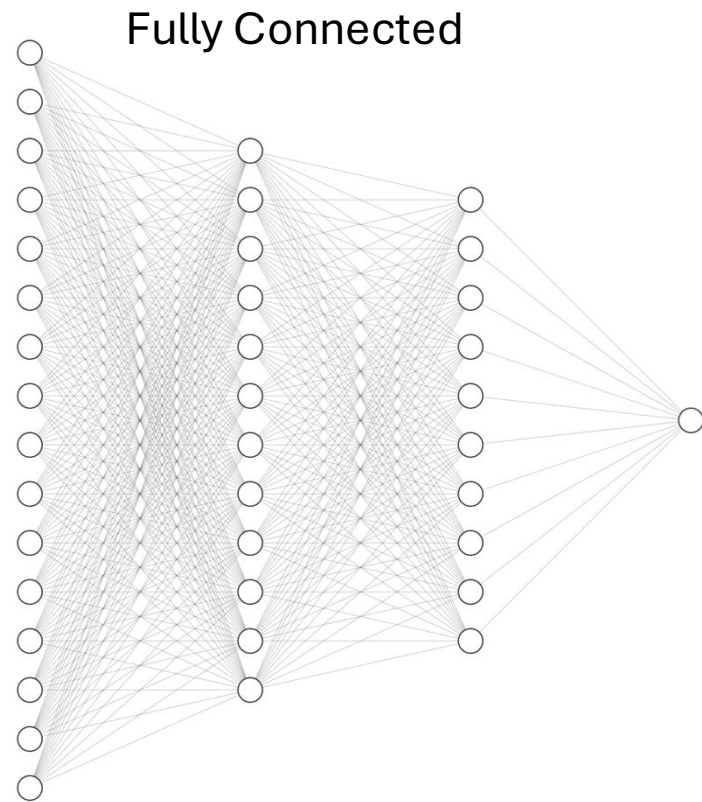this pixel gets weight 0.9

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 and Spatial Reasoning

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

Fully Connected

Not Fully-Connected

# 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)

Not Fully-Connected

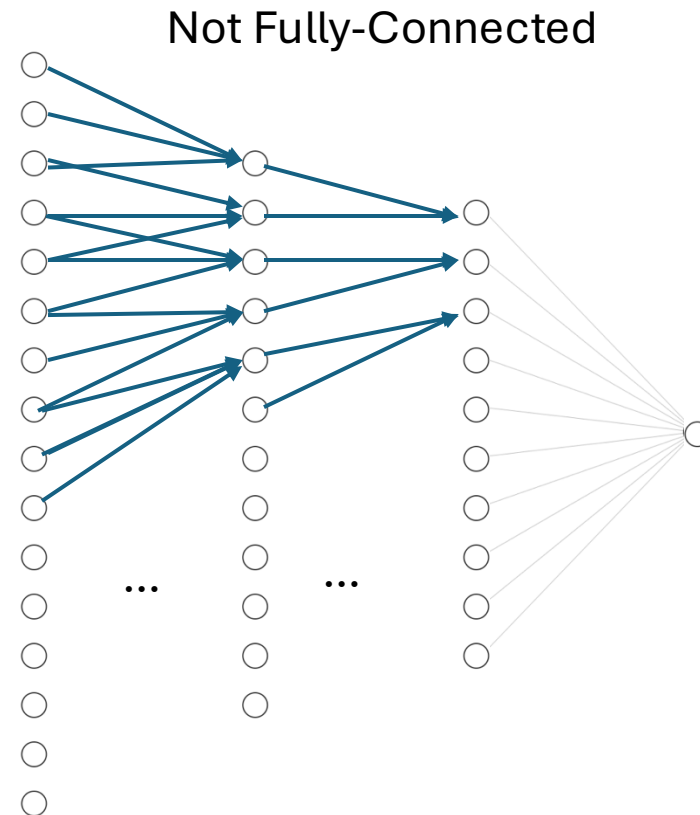# Disadvantages of Not Fully Connected Layers

- 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

# What Convolution Does (Visually)

image

| 2 | 0 | 1 | 3 |
|---|---|---|---|
| 7 | 1 | 1 | 0 |
| 0 | 2 | 5 | 0 |
| 0 | 5 | 1 | 4 |

filter/kernel

| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -1 | -1 |

$\otimes$

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

# What Convolution Does (Visually)
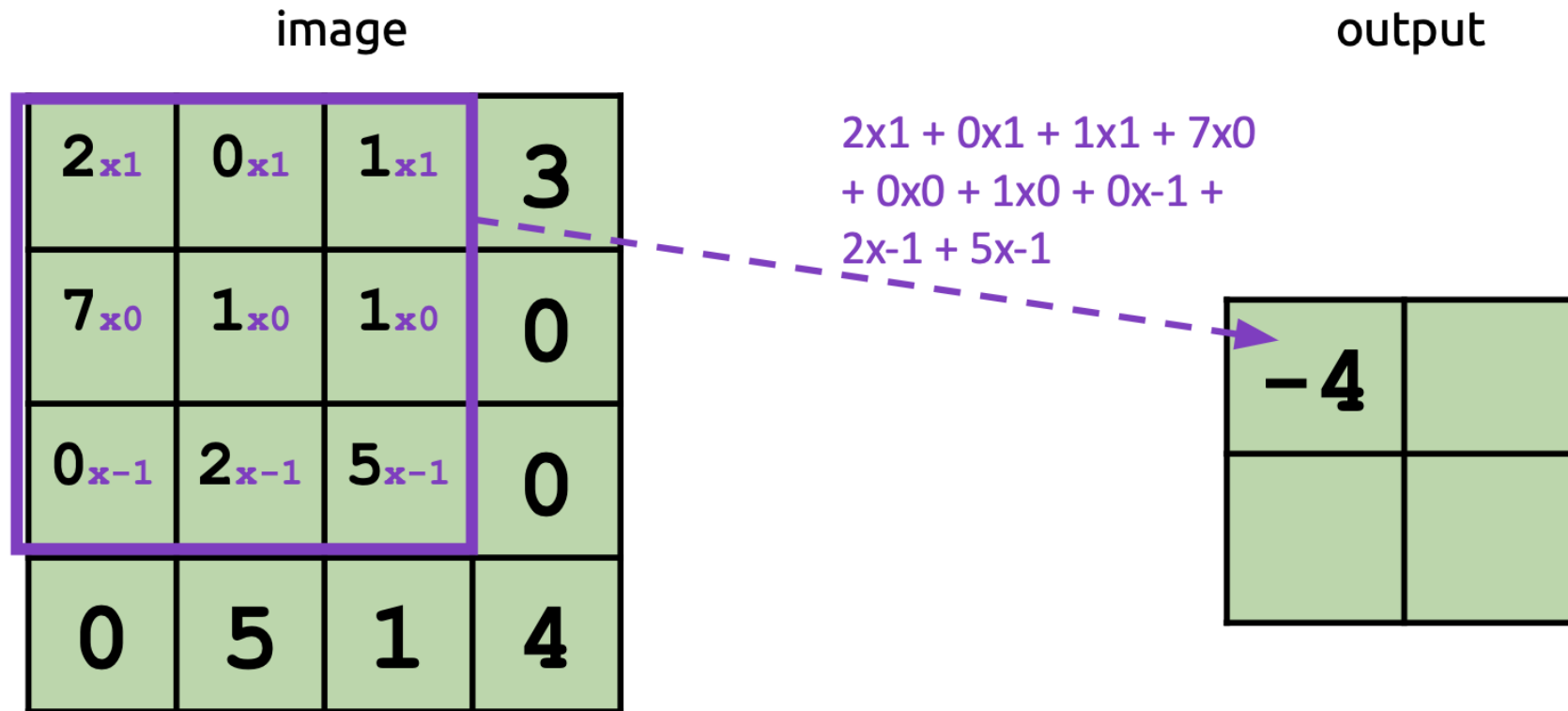
Overlay the filter on the image

image

# What Convolution Does (Visually)

Sum up multiplied values to produce output value

image



$2x1 + 0x1 + 1x1 + 7x0$
$+ 0x0 + 1x0 + 0x{-1} +$
$2x{-1} + 5x{-1}$

output

# What Convolution Does (Visually)

Move the filter over by one pixel

image



output

# What Convolution Does (Visually)

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

| | |
|---|---|
| -4 | |
| | |

# What Convolution Does (Visually)

Repeat (multiply, sum up)

image

| | | | |
|---|---|---|---|
| 2 | $0_{x1}$ | $1_{x1}$ | $3_{x1}$ |
| 7 | $1_{x0}$ | $1_{x0}$ | $0_{x0}$ |
| 0 | $2_{x-1}$ | $5_{x-1}$ | $0_{x-1}$ |
| 0 | 5 | 1 | 4 |

output

| | |
|---|---|
| -4 | |
| | |

# What Convolution Does (Visually)

Repeat (multiply, sum up)

image

$$0x1 + 1x1 + 3x1 + 0x0 + 1x0 + 0x0 + 2x{-}1 + 5x{-}1 + 0x{-}1$$

output

# What Convolution Does (Visually)

Repeat...

image

output

| 2 | 0 | 1 | 3 |
|---|---|---|---|
| $7_{x1}$ | $1_{x1}$ | $1_{x1}$ | 0 |
| $0_{x0}$ | $2_{x0}$ | $5_{x0}$ | 0 |
| $0_{x-1}$ | $5_{x-1}$ | $1_{x-1}$ | 4 |

$7x1 + 1x1 + 1x1 + 0x0 + 2x0$
$+ 5x0 + 0x-1 + 5x-1 + 1x-1$

| -4 | -3 |
|----|----|
| 3 |  |

# What Convolution Does (Visually)

Repeat...

image

output

| 2 | 0 | 1 | 3 |
|---|---|---|---|
| 7 | $1_{x1}$ | $1_{x1}$ | $0_{x1}$ |
| 0 | $2_{x0}$ | $5_{x0}$ | $0_{x0}$ |
| 0 | $5_{x-1}$ | $1_{x-1}$ | $4_{x-1}$ |

$1x1 + 1x1 + 0x1 + 2x0 + 5x0 + 0x0 + 5x{-1} + 1x{-1} + 4x{-1}$

| -4 | -3 |
|----|----|
| 3 | -8 |

# What Convolution Does (Visually)

In summary:

image

| 2 | 0 | 1 | 3 |
|---|---|---|---|
| 7 | 1 | 1 | 0 |
| 0 | 2 | 5 | 0 |
| 0 | 5 | 1 | 4 |

$\otimes$

filter/kernel

| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -1 | -1 |

$=$

output

| -4 | -3 |
|----|----|
| 3  | -8 |

# Handmade Kernels and Filters

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

**Identity kernel**

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

**Edge detection**

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

**Sharpen kernel**

$$\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

**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 | Kernel ω | Image result g(x,y) |
|---|---|---|
| **Identity** | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ | |
| **Ridge** or **edge detection** | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ | |
| **Sharpen** | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | |
| **Box blur** (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |
| **Gaussian blur** 3 × 3 (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ | |

# What Comes Next?

Can we learn a filter for our images rather than "hand crafting" one?

# Recap

Participation Quiz #3 is up!

Fully Connected Neural Networks for images lack spatial information

We can add spatial reasoning by connecting pixels that are "spatially" close

Convolutions/Filters/Kernels are a technique from image processing that combine close pixels with a linear transformation