

# Uncertainty in Deep Learning for Image Classification

**Method Seminar HS2018**

**Elvis Murina**

**Institute of Data Analysis and Process Design  
Zurich University of Applied Sciences**

**19.11.2018**

# Outline

- Definitions and History
  - AI, Machine Learning, Deep Learning
- Image Classification
  - Motivation, challenges
  - Traditional vs Deep Learning
  - Imagenet challenge
- Artificial neural networks
  - Fully connected vs Convolution neural networks
  - Stochastic gradient descent
  - Backpropagation and the chain rule
- Tricks to fight overfitting
  - Data Augmentation, Dropout
- Uncertainty
  - MC-Dropout and uncertainty in Deep Learning
- HCS Application
- Questions

# What is Deep Learning?



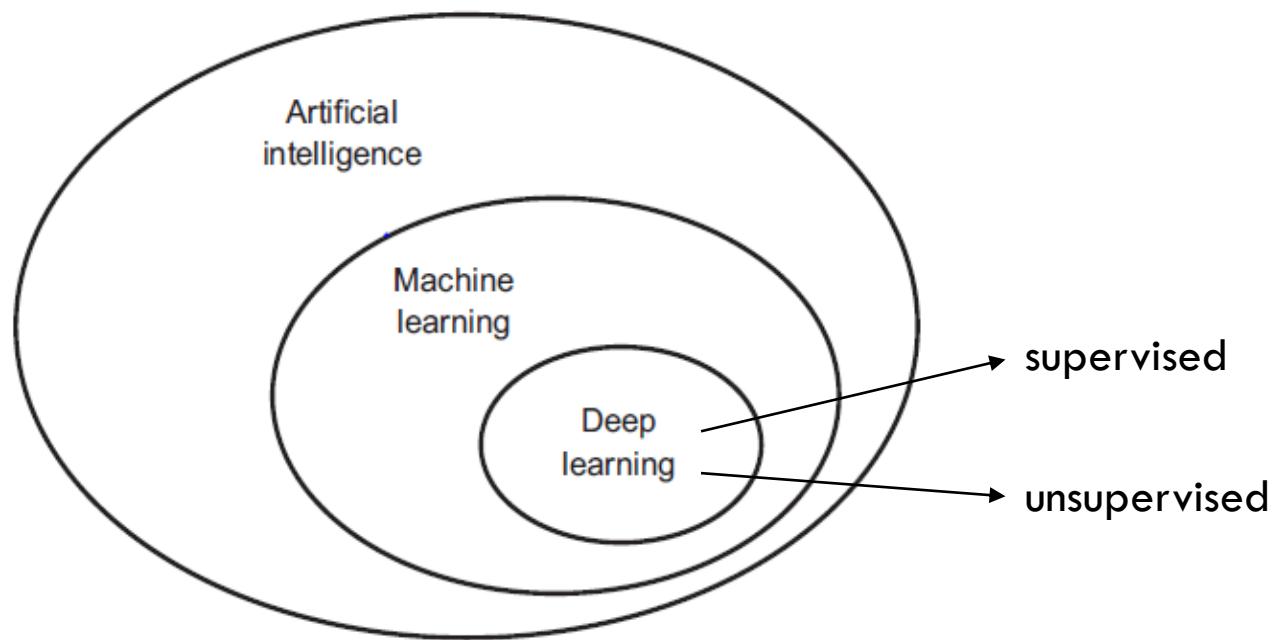


# Definitions

# Definitions

- AI
  - The effort to automate intellectual tasks normally performed by humans (François Chollet)
- Machine Learning
  - Field of study that gives computers the ability to learn without being explicitly programmed. (Arthur Samuel)
  - A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E. (Tom M. Mitchell)
- Deep Learning
  - A subset of ML which uses deep artificial neural networks as models and does not require feature engineering.

# Definitions





# History of Neural Networks

# History

## NEURAL NETWORKS ARE AN OLD IDEA

One of the very first ideas in machine learning and artificial intelligence

- Date back to 1940s
- Many cycles of boom and bust
- Repeated promises of "true AI" that were unfulfilled followed by "AI winters"



*"[The perceptron is] the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."*

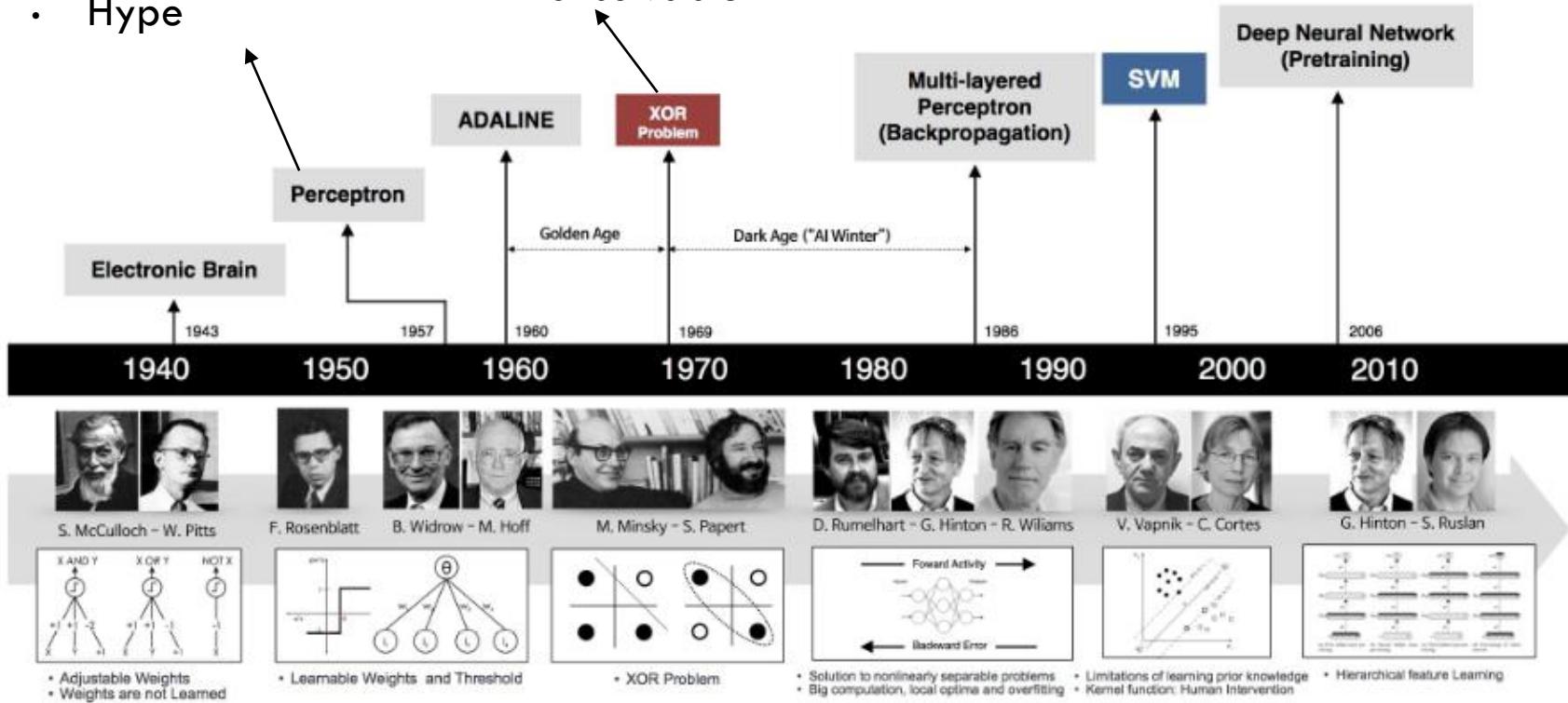
- Frank Rosenblatt, 1958

# History

- Linear
- No hidden layers
- Learnable weights
- Hype

## MILESTONES

- Simple non linear functions  
not solvable





# **Image Classification and other computer vision tasks**

# Image Classification

- A core task in computer vision



(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}



cat

# Image Classification

- Other task in computer vision

Image classification



Cat

Semantic Segmentation



Grass, Cat,  
Tree, Sky

Object  
Detection



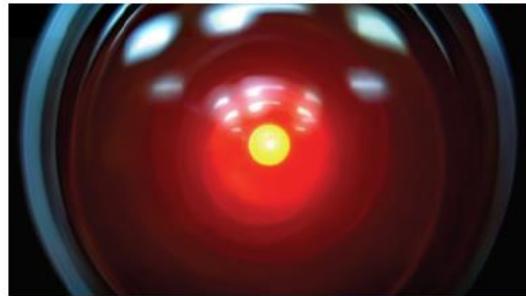
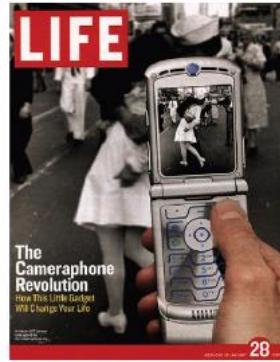
Cat

Object  
Detection

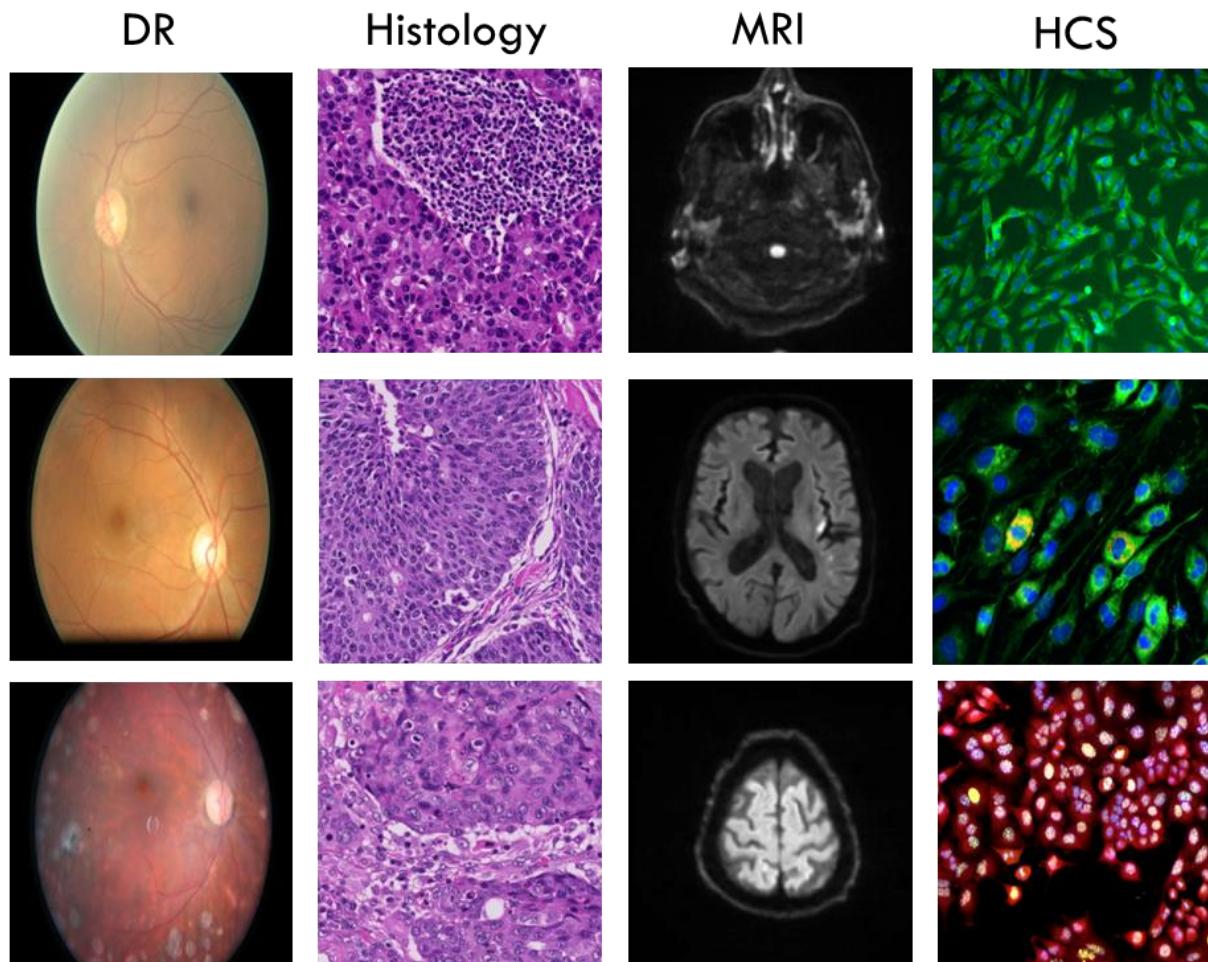


Dog, Dog, Cat

# Motivation

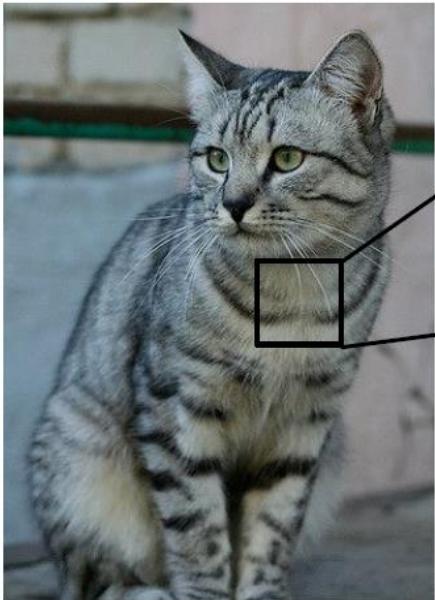


# Motivation



# Challenges

- Sematic gap



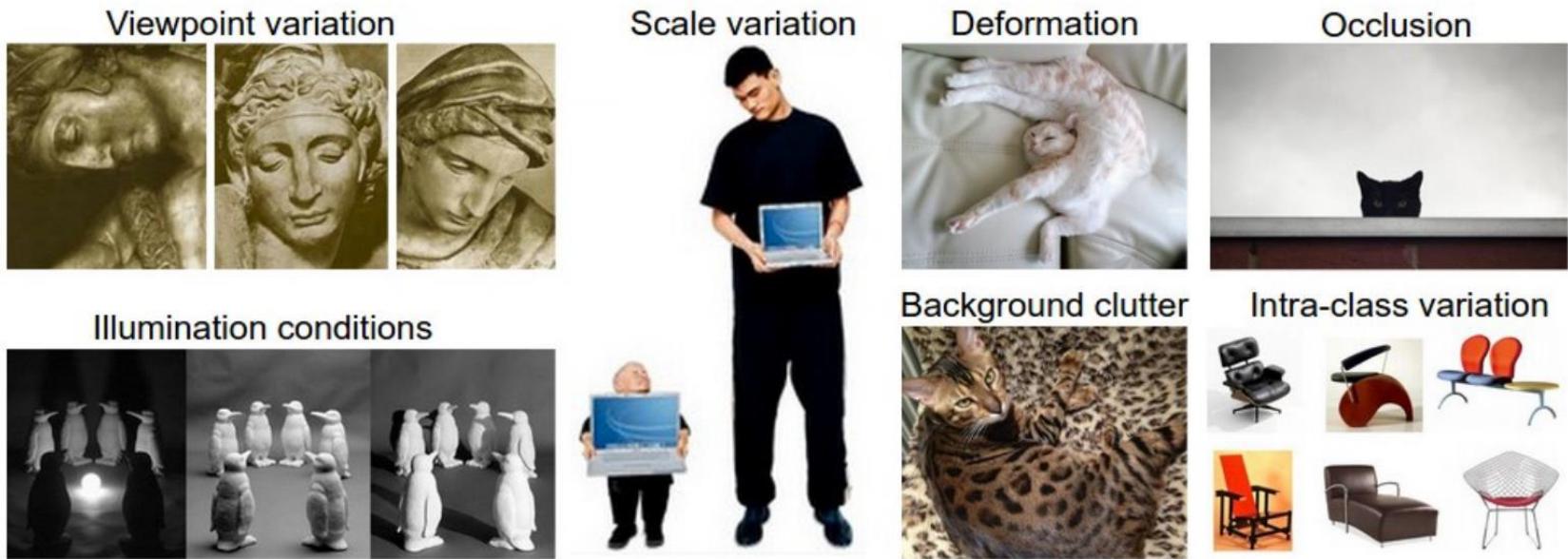
```
[1105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 99 85]
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[128 137 144 140 109 95 86 70 62 65 53 63 60 60 73 86 101]
[133 137 147 103 65 81 88 65 52 54 74 84 102 92 85 82]
[128 137 144 140 109 95 86 70 62 65 53 63 60 60 73 86 101]
[125 133 148 137 119 121 117 94 65 79 88 65 54 64 72 98]
[127 125 131 147 133 127 126 131 111 95 89 75 61 64 72 84]
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]
[ 89 93 90 97 108 147 131 118 113 114 113 109 106 95 77 80]
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]
[ 62 65 82 80 78 71 80 101 124 126 119 101 107 114 131 119]
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]
[ 87 65 71 87 106 95 68 45 76 130 126 107 92 94 105 112]
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]
[164 145 112 88 82 120 124 104 76 48 45 66 88 101 102 109]
[157 170 157 128 93 86 114 132 112 97 69 55 70 82 99 94]
[130 128 134 161 139 108 109 118 121 134 114 87 65 53 69 86]
[128 112 96 117 150 144 126 115 104 107 102 93 87 81 72 79]
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]
[122 184 148 103 71 55 78 83 93 103 119 130 102 61 69 84]]
```

What the computer sees

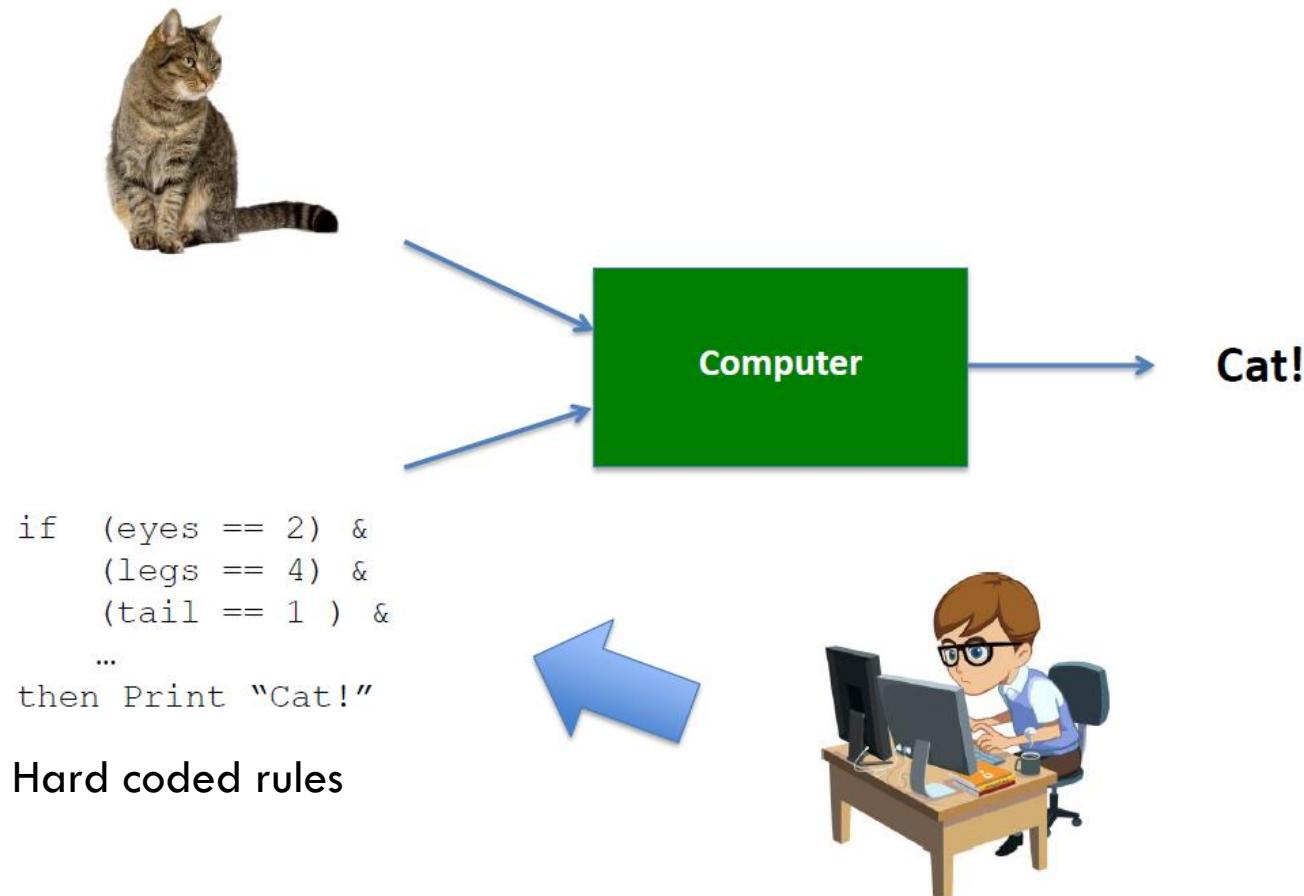
An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

# Challenges

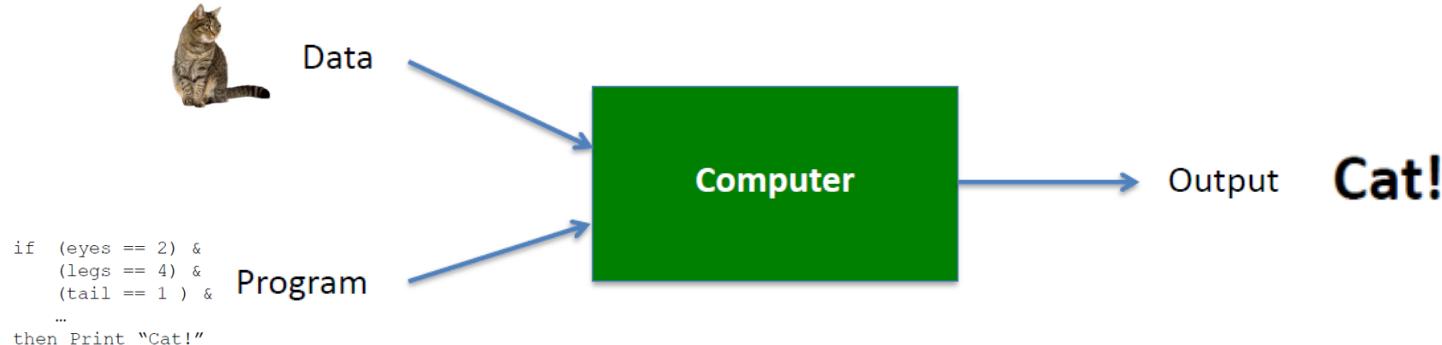


# Old image classification paradigm

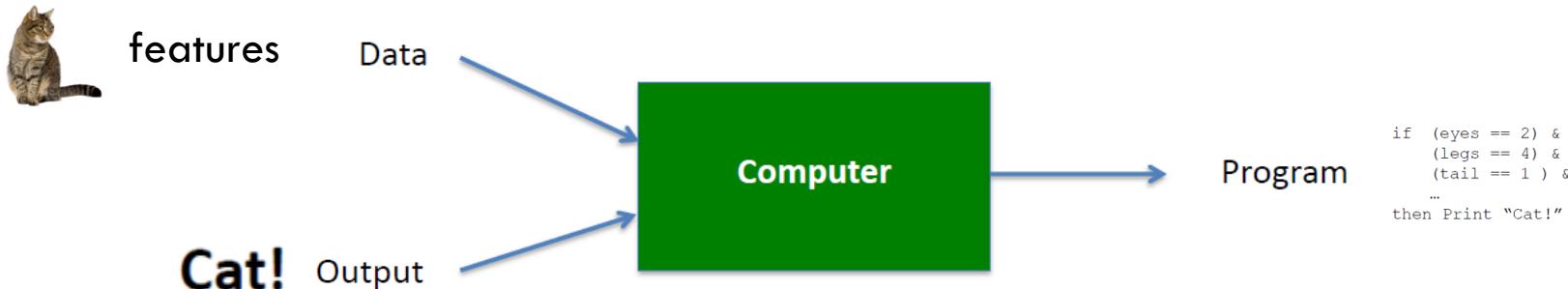


# New Machine learning paradigm

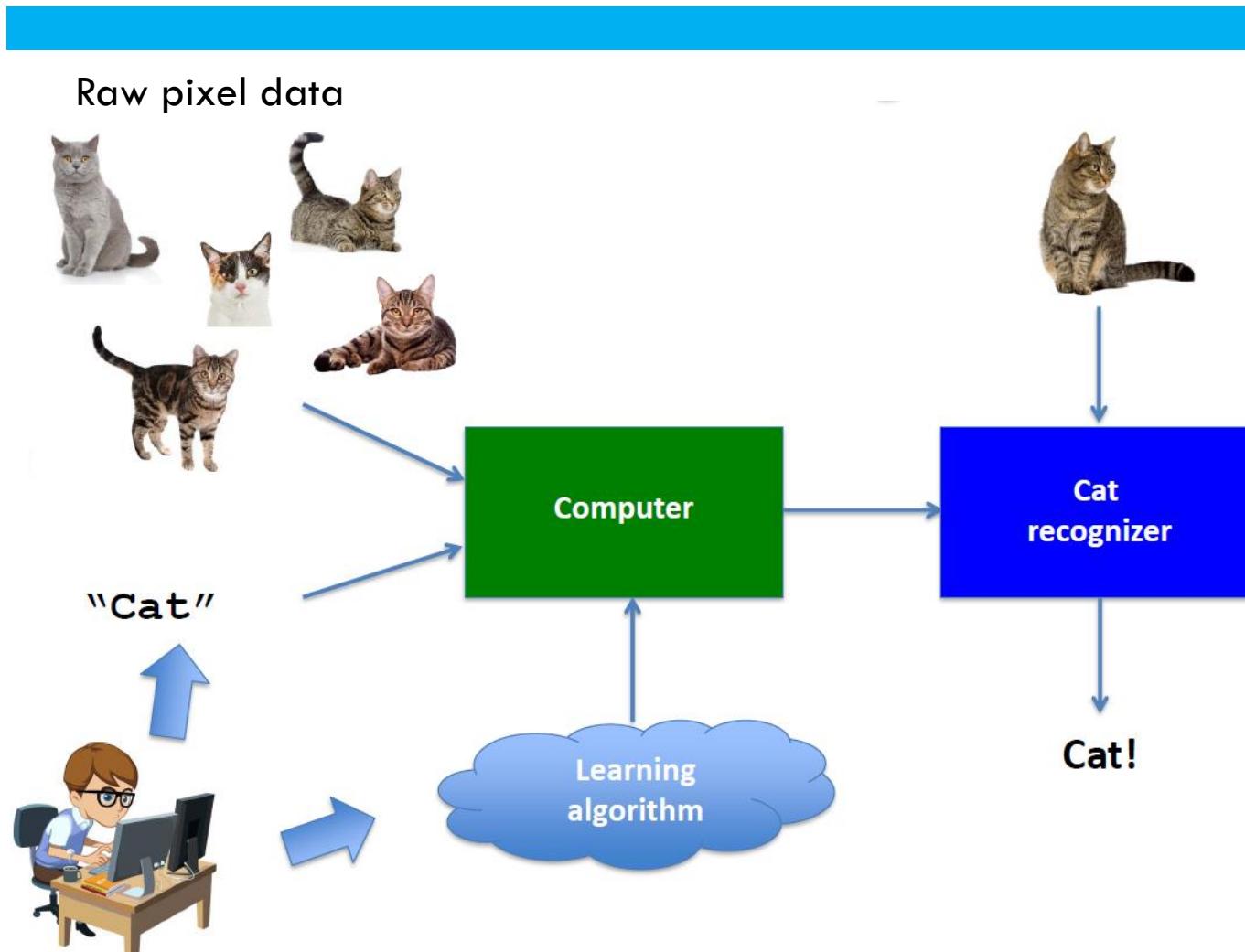
## Traditional programming



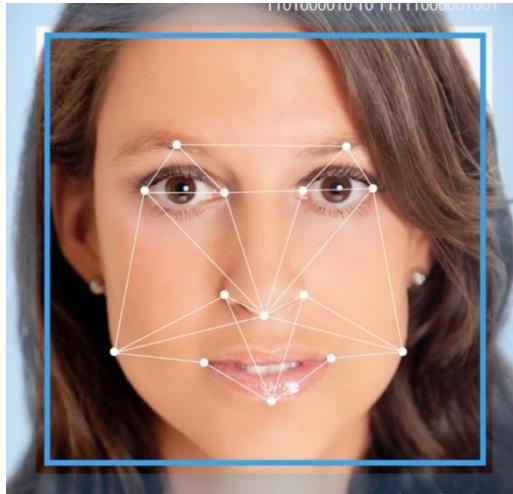
## Machine learning



# Deep Learning paradigm



# Traditional image analysis (ML) vs. Deep Learning

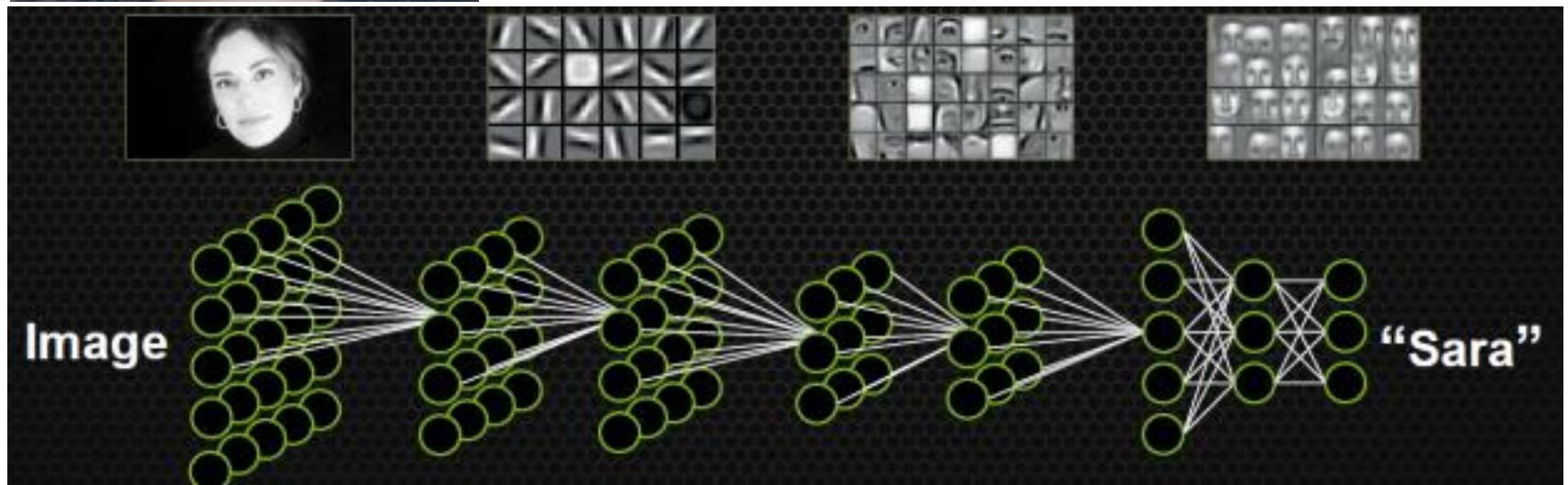


## Traditional:

Extract **handcrafted features** and use these features to **train/fit a model** (SVM, RF) and use fitted model to perform classification/prediction.

## Deep learning:

In **deep neural networks** start with raw data and **learn** during training/fitting to extract appropriate **hierarchical features** and to use them for classification/prediction.





# Imagenet

## The Deep Learning revolution

# ImageNet competition

1000 classes  
1 Mio samples



Human: 5% misclassification



A. Krizhevsky  
first CNN in 2012

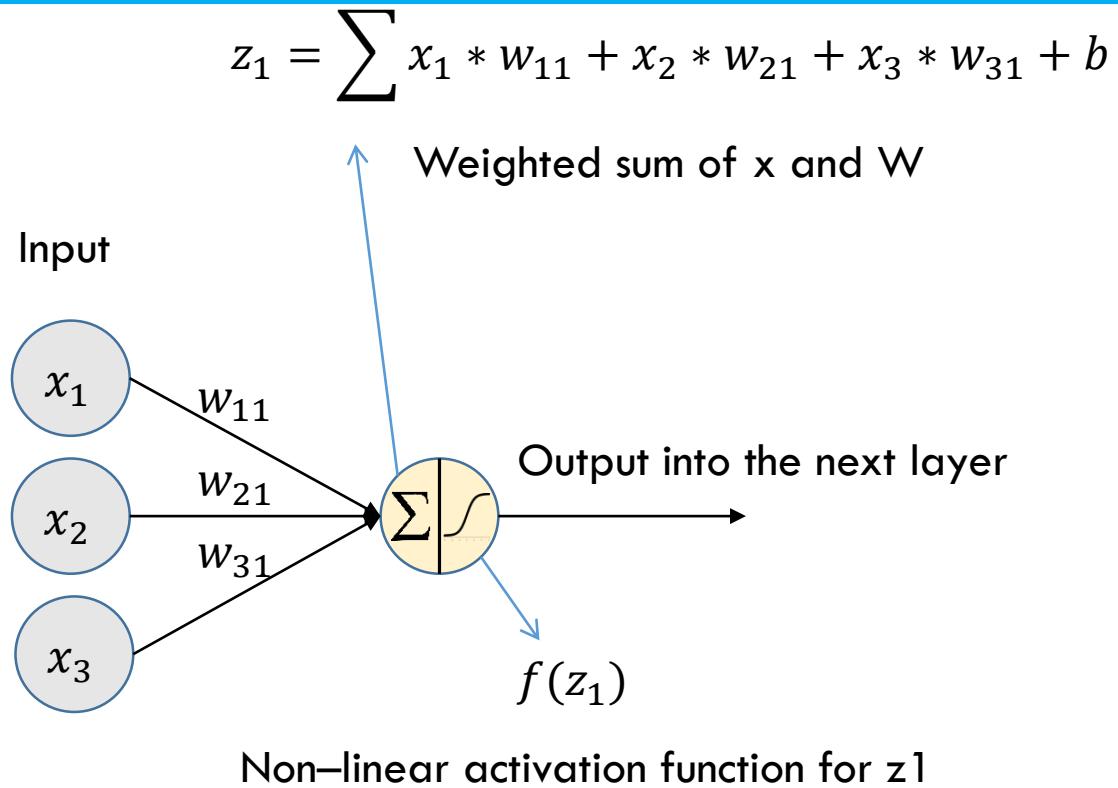
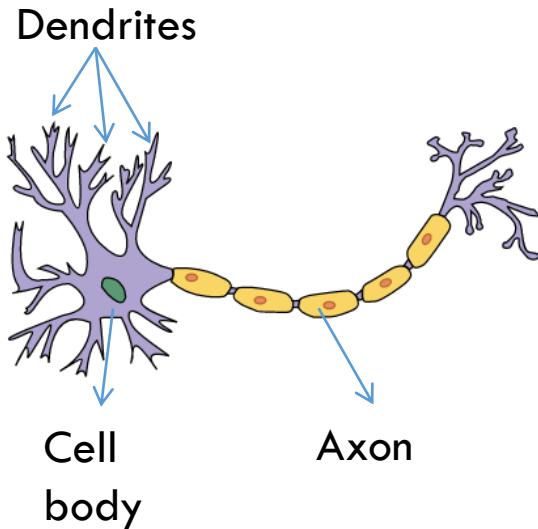
2015

4.95% Microsoft (surpassing human performance 5%)  
4.8% Google (further improved to 3.6%)  
3.57% Microsoft (Resnet winner 2015)

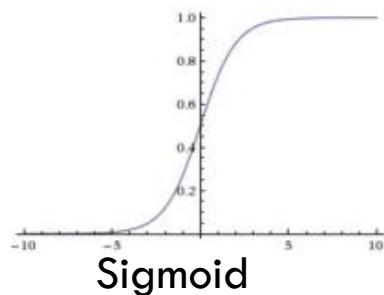


# Artificial neural networks

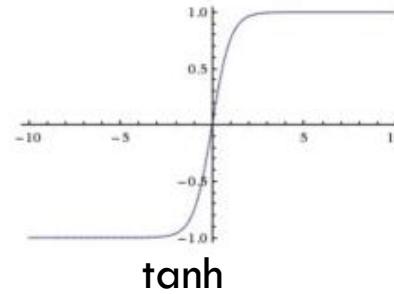
# Artificial neural networks



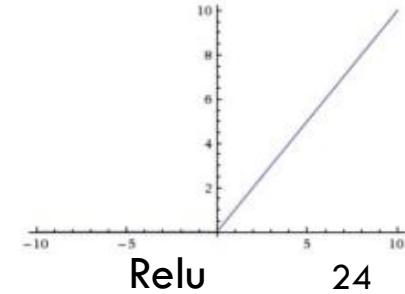
Non-linear activation function for  $z_1$



Sigmoid

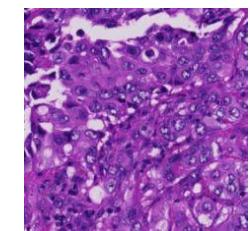


tanh

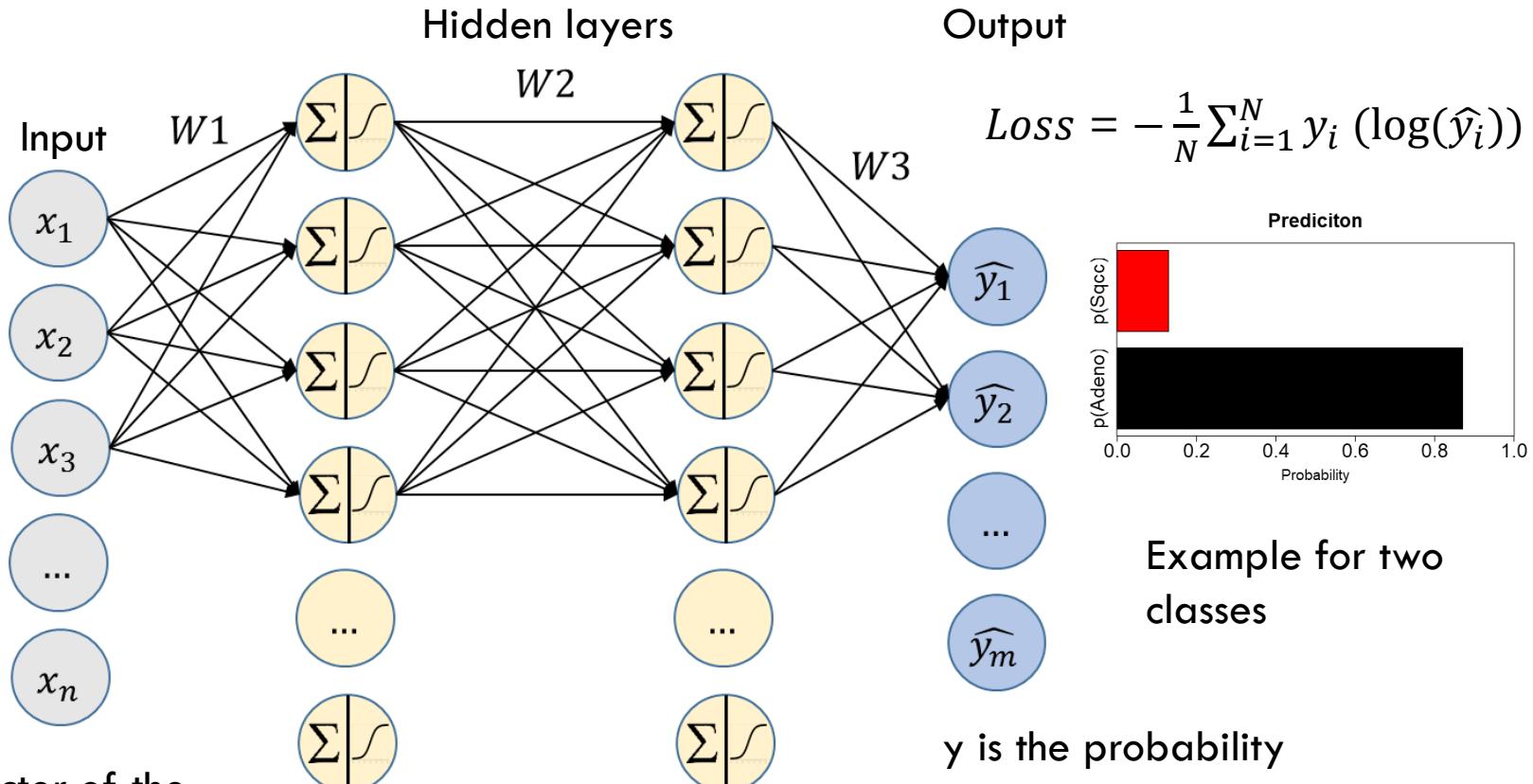


Relu

# Artificial neural networks



(100x100x3)



$x$  is the vector of the flattened input image

$W$  are the weight matrices that are learned during the training phase of the network

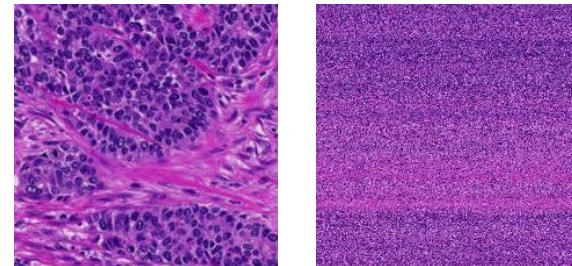
$y$  is the probability vector for each class

$$\text{Softmax} = \frac{e^i}{\sum_{k=1}^m e^k}$$

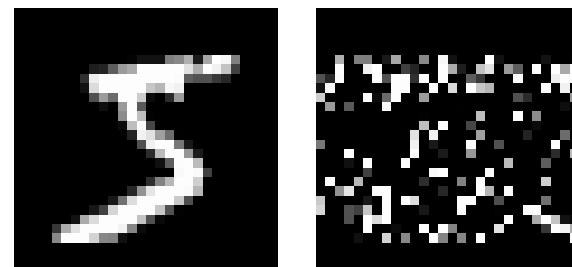
# Artificial neural networks

- Two main disadvantages
  - Images are often big and therefore you need a lot of weights (also for the hidden layers)
  - The spatial information of the image gets lost, because you flatten the input.
  - Images left and right are the same for a fully connected network

Adeno vs Sqcc example

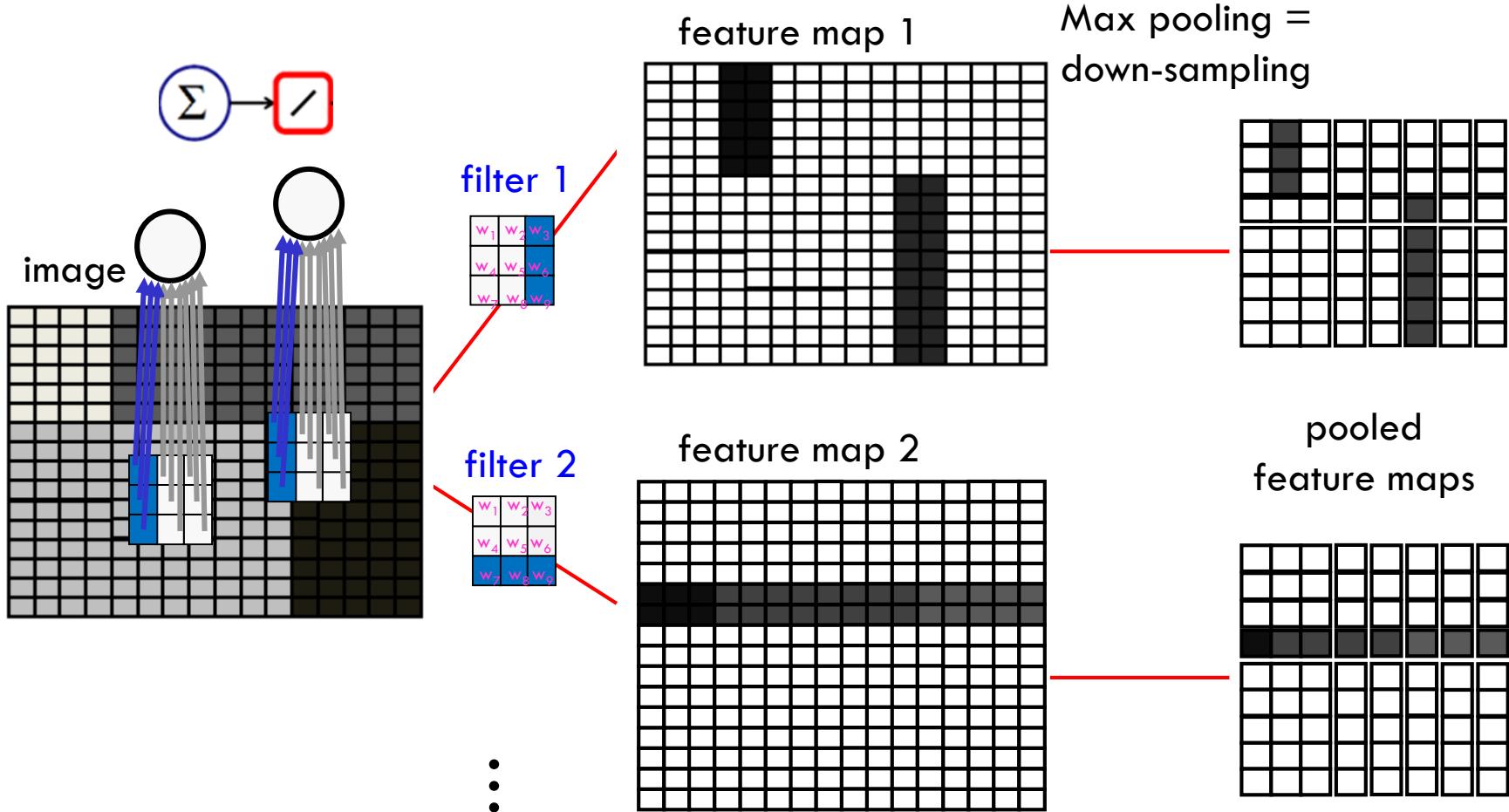


Mnist example



# Convolutional neural networks

# Convolutional neural networks



The weights of each filter are randomly initiated and then adapted during the training.

# Convolutional neural networks

Input image 6x6x1

|     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|
| 255 | 220 | 150 | 200 | 110 | 100 |
| 240 | 50  | 35  | 45  | 200 | 130 |
| 0   | 20  | 245 | 250 | 230 | 120 |
| 170 | 180 | 235 | 145 | 170 | 255 |
| 190 | 185 | 170 | 165 | 130 | 120 |
| 255 | 255 | 245 | 190 | 200 | 170 |

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

|      |      |     |     |     |     |
|------|------|-----|-----|-----|-----|
| -0.7 | 0.2  | 0.1 | 200 | 110 | 100 |
| 0.3  | 0.5  | 0.4 | 45  | 200 | 130 |
| -0.2 | -0.4 | 0.2 | 250 | 230 | 120 |
| 170  | 180  | 235 | 145 | 170 | 255 |
| 190  | 185  | 170 | 165 | 130 | 120 |
| 255  | 255  | 245 | 190 | 200 | 170 |

Feature map 4x4x1

|      |
|------|
| 32.5 |
|------|

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

|     |      |      |     |     |     |  |
|-----|------|------|-----|-----|-----|--|
| 255 | -0.7 | 0.2  | 0.1 | 110 | 100 |  |
| 240 | 0.3  | 0.5  | 0.4 | 200 | 130 |  |
| 0   | -0.2 | -0.4 | 0.2 | 230 | 120 |  |
| 170 | 180  | 235  | 145 | 170 | 255 |  |
| 190 | 185  | 170  | 165 | 130 | 120 |  |
| 255 | 255  | 245  | 190 | 200 | 170 |  |

Feature map 4x4x1

|      |        |
|------|--------|
| 32.5 | -105.5 |
|------|--------|

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

|     |     |      |      |     |     |
|-----|-----|------|------|-----|-----|
| 255 | 220 | -0.7 | 0.2  | 0.1 | 100 |
| 240 | 50  | 0.3  | 0.5  | 0.4 | 130 |
| 0   | 20  | -0.2 | -0.4 | 0.2 | 120 |
| 170 | 180 | 235  | 145  | 170 | 255 |
| 190 | 185 | 170  | 165  | 130 | 120 |
| 255 | 255 | 245  | 190  | 200 | 175 |

Feature map 4x4x1

|      |        |       |
|------|--------|-------|
| 32.5 | -105.5 | 185.5 |
|------|--------|-------|

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

A 6x6 grid representing an input image. The values are as follows:

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | -0.7 | 0.2  | 0.1 |
| 240 | 50  | 35  | 0.3  | 0.5  | 0.4 |
| 0   | 20  | 245 | -0.2 | -0.4 | 0.2 |
| 170 | 180 | 235 | 145  | 170  | 255 |
| 190 | 185 | 170 | 165  | 130  | 120 |
| 255 | 255 | 245 | 190  | 200  | 170 |

Feature map 4x4x1

|      |        |       |    |
|------|--------|-------|----|
| 32.5 | -105.5 | 185.5 | 54 |
|------|--------|-------|----|

A 3x3 grid representing a filter. The values are as follows:

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

A 6x6 grid representing an input image. The values are as follows:

|      |      |     |     |     |     |
|------|------|-----|-----|-----|-----|
| 255  | 220  | 150 | 200 | 110 | 100 |
| -0.7 | 0.2  | 0.1 | 45  | 200 | 130 |
| 0.3  | 0.5  | 0.4 | 250 | 230 | 120 |
| -0.2 | -0.4 | 0.2 | 145 | 170 | 255 |
| 190  | 185  | 170 | 165 | 130 | 120 |
| 255  | 255  | 245 | 190 | 200 | 170 |

Feature map 4x4x1

A 4x4 grid representing a feature map. The values are as follows:

|        |        |       |    |
|--------|--------|-------|----|
| 32.5   | -105.5 | 185.5 | 54 |
| -105.5 |        |       |    |

A 3x3 grid representing a filter. The values are as follows:

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

A 6x6 input image grid showing a cartoon dog's head. The image is labeled with pixel values ranging from 100 to 255. A 3x3 convolutional kernel is applied to the image, with weights shown in the bottom right corner of each 3x3 receptive field. The kernel has weights: -0.7, 0.2, 0.1; 0.3, 0.5, 0.4; and -0.2, -0.4, 0.2.

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | 200  | 110  | 100 |
| 240 | 50  | 35  | 45   | 200  | 130 |
| 0   | 20  | 245 | 250  | 230  | 120 |
| 170 | 180 | 235 | -0.7 | 0.2  | 0.1 |
| 190 | 185 | 170 | 0.3  | 0.5  | 0.4 |
| 255 | 255 | 245 | -0.2 | -0.4 | 0.2 |

Feature map 4x4x1

|        |        |       |      |
|--------|--------|-------|------|
| 32.5   | -105.5 | 185.5 | 54   |
| -105.5 | 104    | 217.5 | 31   |
| -44    | 224    | 38.5  | -18  |
| -60.5  | 213.5  | 52.5  | 37.5 |

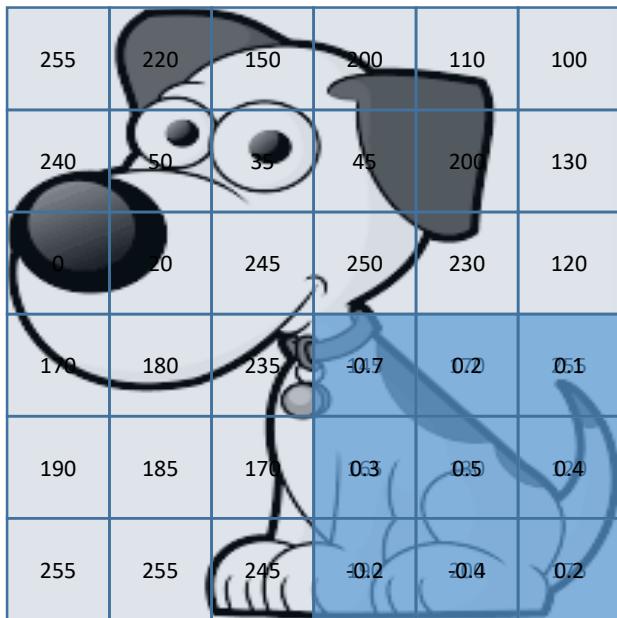
|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1

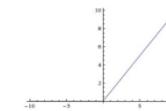


A 6x6 grid representing an input image of a cartoon dog's head. Numerical values are placed at various pixels, showing the intensity or color information. The values range from 100 to 255.

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | 200  | 110  | 100 |
| 240 | 50  | 35  | 45   | 200  | 130 |
| 0   | 20  | 245 | 250  | 230  | 120 |
| 170 | 180 | 235 | -0.7 | 0.2  | 0.1 |
| 190 | 185 | 170 | 0.3  | 0.5  | 0.4 |
| 255 | 255 | 245 | -0.2 | -0.4 | 0.2 |

Feature map 4x4x1

|        |        |       |      |
|--------|--------|-------|------|
| 32.5   | -105.5 | 185.5 | 54   |
| -105.5 | 104    | 217.5 | 31   |
| -44    | 224    | 38.5  | -18  |
| -60.5  | 213.5  | 52.5  | 37.5 |



Relu

|      |       |       |      |
|------|-------|-------|------|
| 32.5 | 0     | 185.5 | 54   |
| 0    | 104   | 217.5 | 31   |
| 0    | 224   | 38.5  | 0    |
| 0    | 213.5 | 52.5  | 37.5 |

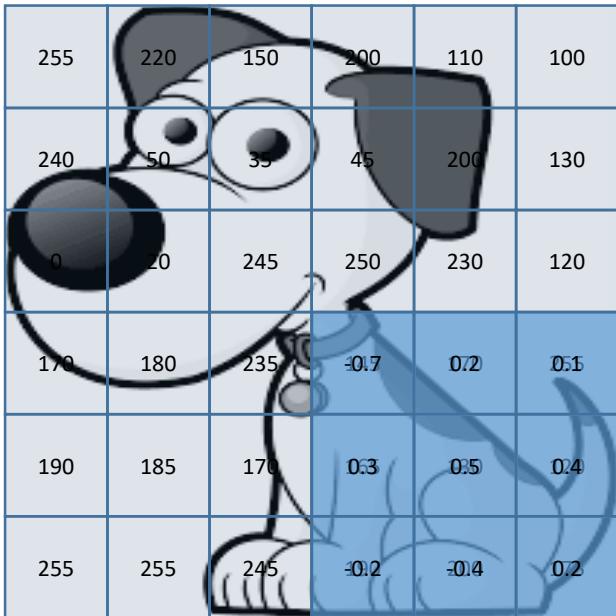
|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

Input image 6x6x1



A 6x6 input image grid showing a cartoon owl. A 3x3 kernel is applied to the central 3x3 area of the image. The values in the kernel are: -0.7, 0.2, 0.1; 0.3, 0.5, 0.4; -0.2, -0.4, 0.2. The resulting output values are: -0.7, 0.2, 0.1; 0.3, 0.5, 0.4; -0.2, -0.4, 0.2.

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | 200  | 110  | 100 |
| 240 | 50  | 35  | 45   | 200  | 130 |
| 0   | 20  | 245 | 250  | 230  | 120 |
| 170 | 180 | 235 | -0.7 | 0.2  | 0.1 |
| 190 | 185 | 170 | 0.3  | 0.5  | 0.4 |
| 255 | 255 | 245 | -0.2 | -0.4 | 0.2 |

Feature map 4x4x1

|        |        |       |      |
|--------|--------|-------|------|
| 32.5   | -105.5 | 185.5 | 54   |
| -105.5 | 104    | 217.5 | 31   |
| -44    | 224    | 38.5  | -18  |
| -60.5  | 213.5  | 52.5  | 37.5 |



Relu

|      |       |       |      |
|------|-------|-------|------|
| 32.5 | 0     | 185.5 | 54   |
| 0    | 104   | 217.5 | 31   |
| 0    | 224   | 38.5  | 0    |
| 0    | 213.5 | 52.5  | 37.5 |

Maxpool  
(2x2x1)

104

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

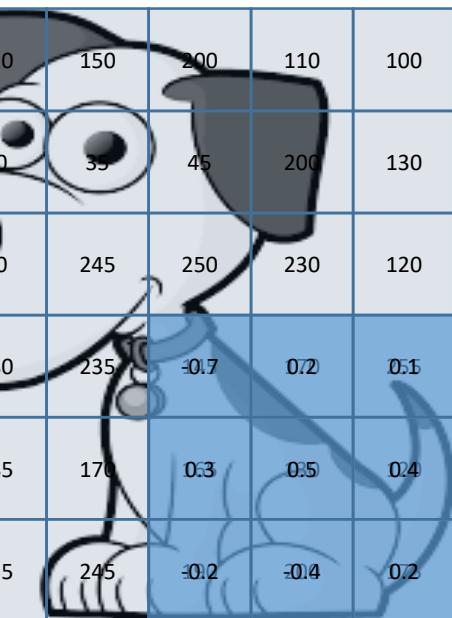
3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

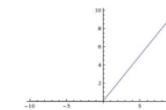
Input image 6x6x1

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | 200  | 110  | 100 |
| 240 | 50  | 35  | 45   | 200  | 130 |
| 0   | 20  | 245 | 250  | 230  | 120 |
| 170 | 180 | 235 | -0.7 | 0.2  | 0.1 |
| 190 | 185 | 170 | 0.3  | 0.5  | 0.4 |
| 255 | 255 | 245 | -0.2 | -0.4 | 0.2 |



Feature map 4x4x1

|        |        |       |      |
|--------|--------|-------|------|
| 32.5   | -105.5 | 185.5 | 54   |
| -105.5 | 104    | 217.5 | 31   |
| -44    | 224    | 38.5  | -18  |
| -60.5  | 213.5  | 52.5  | 37.5 |



Relu

|      |       |       |      |
|------|-------|-------|------|
| 32.5 | 0     | 185.5 | 54   |
| 0    | 104   | 217.5 | 31   |
| 0    | 224   | 38.5  | 0    |
| 0    | 213.5 | 52.5  | 37.5 |

Maxpool  
(2x2x1)

|     |       |
|-----|-------|
| 104 | 217.5 |
|-----|-------|

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

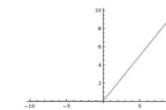
# Convolutional neural networks

Input image 6x6x1

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | 200  | 110  | 100 |
| 240 | 50  | 35  | 45   | 200  | 130 |
| 0   | 20  | 245 | 250  | 230  | 120 |
| 170 | 180 | 235 | -0.7 | 0.2  | 0.1 |
| 190 | 185 | 170 | 0.3  | 0.5  | 0.4 |
| 255 | 255 | 245 | -0.2 | -0.4 | 0.2 |

Feature map 4x4x1

|        |        |       |      |
|--------|--------|-------|------|
| 32.5   | -105.5 | 185.5 | 54   |
| -105.5 | 104    | 217.5 | 31   |
| -44    | 224    | 38.5  | -18  |
| -60.5  | 213.5  | 52.5  | 37.5 |



Relu

|      |       |       |      |
|------|-------|-------|------|
| 32.5 | 0     | 185.5 | 54   |
| 0    | 104   | 217.5 | 31   |
| 0    | 224   | 38.5  | 0    |
| 0    | 213.5 | 52.5  | 37.5 |

Maxpool  
(2x2x1)

|     |       |
|-----|-------|
| 104 | 217.5 |
| 224 |       |

|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

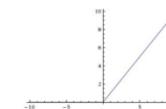
Input image 6x6x1

|     |     |     |      |      |     |
|-----|-----|-----|------|------|-----|
| 255 | 220 | 150 | 200  | 110  | 100 |
| 240 | 50  | 35  | 45   | 200  | 130 |
| 0   | 20  | 245 | 250  | 230  | 120 |
| 170 | 180 | 235 | -0.7 | 0.2  | 0.1 |
| 190 | 185 | 170 | 0.3  | 0.5  | 0.4 |
| 255 | 255 | 245 | -0.2 | -0.4 | 0.2 |



Feature map 4x4x1

|        |        |       |      |
|--------|--------|-------|------|
| 32.5   | -105.5 | 185.5 | 54   |
| -105.5 | 104    | 217.5 | 31   |
| -44    | 224    | 38.5  | -18  |
| -60.5  | 213.5  | 52.5  | 37.5 |



Relu

|      |       |       |      |
|------|-------|-------|------|
| 32.5 | 0     | 185.5 | 54   |
| 0    | 104   | 217.5 | 31   |
| 0    | 224   | 38.5  | 0    |
| 0    | 213.5 | 52.5  | 37.5 |

Maxpool  
(2x2x1)

|     |       |
|-----|-------|
| 104 | 217.5 |
| 224 | 52.5  |

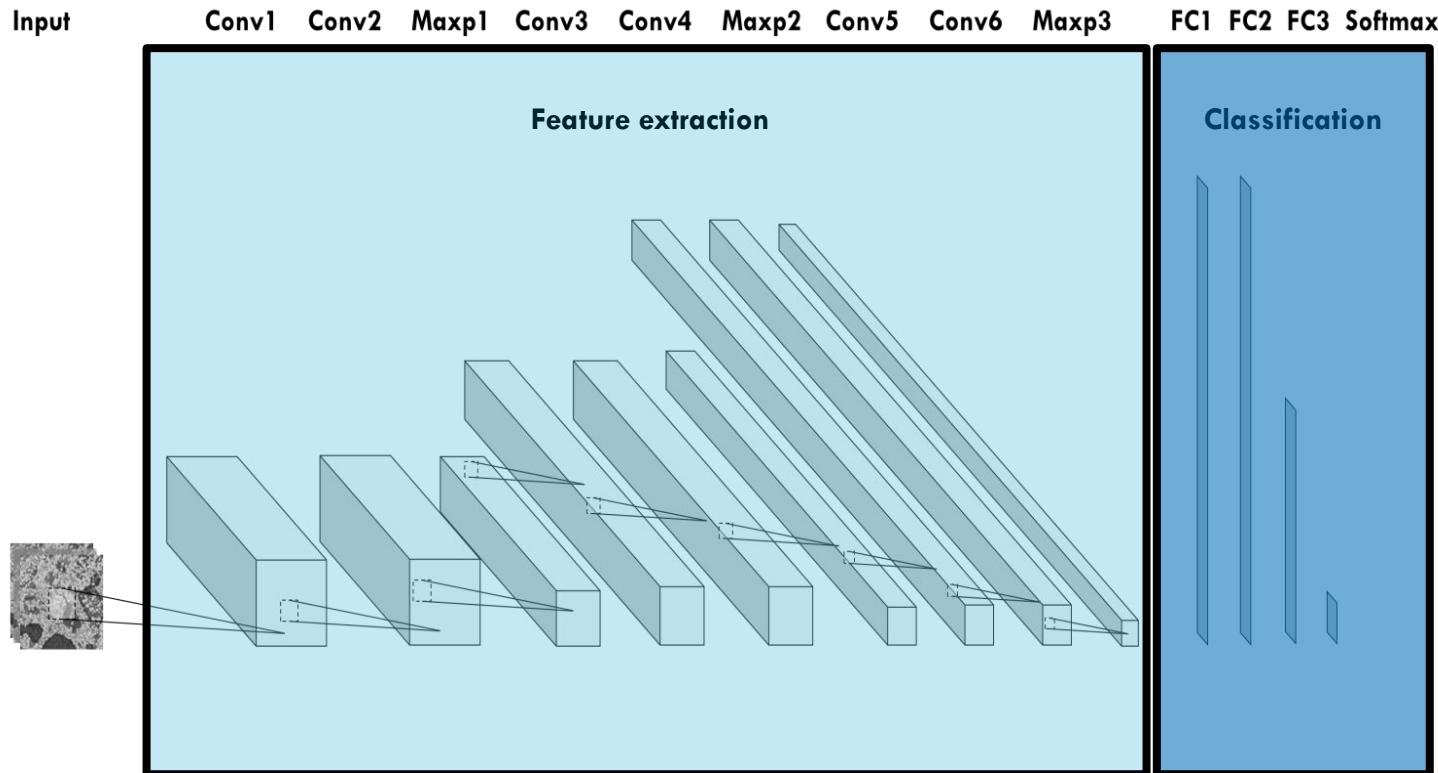
|      |      |     |
|------|------|-----|
| -0.7 | 0.2  | 0.1 |
| 0.3  | 0.5  | 0.4 |
| -0.2 | -0.4 | 0.2 |

3x3 filter

$$z = b + \sum_i x_i w_i$$

# Convolutional neural networks

## Typical architecture of a CNN

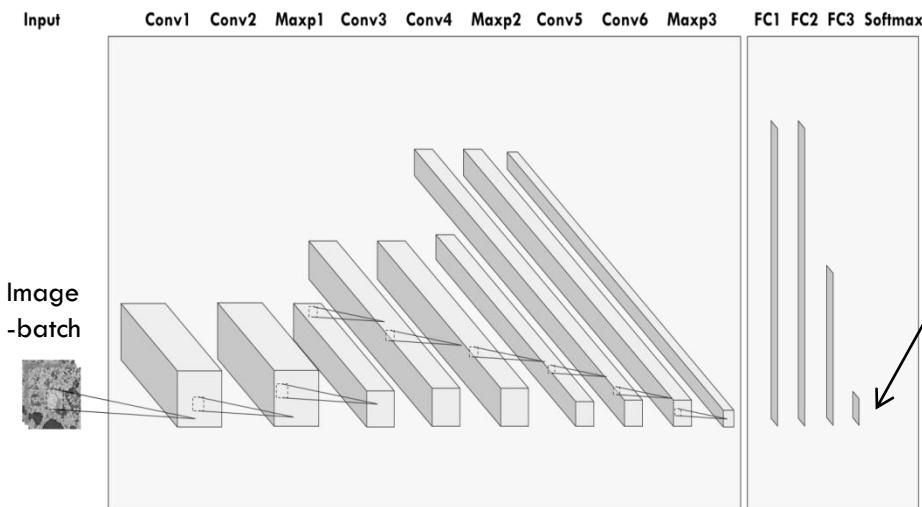


Spatial resolution is decreased with max-pooling while more abstract image features are detected in deeper layers.

# Convolutional neural networks



Calculate the output based on the current parameters

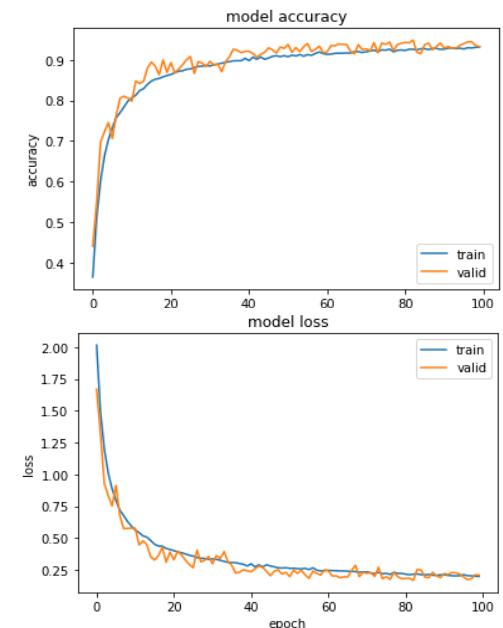


Update the parameters to minimize the loss



$$\text{Loss} = -\frac{1}{N} \sum_i^N y_i (\log(\hat{y}_i))$$

Categorical crossentropy



Weight update:

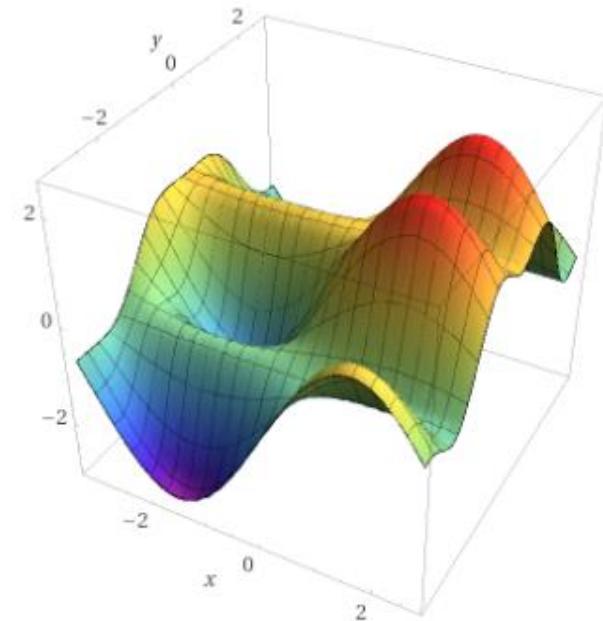
$$w_{i(t+1)} = w_{i(t)} - \eta * \frac{\partial L(w)}{\partial w_i}$$



**Gradient descent, stochastic gradient  
descent online gradient descent and  
mini-batch gradient descent to minimize  
the loss**

# Gradient descent in DL

Follow the slope!



2 dimensional loss function. In DL we have millions dimensions!  
We just know the current value (we are blind)

Image source: [http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture3.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture3.pdf) (cs231n)

Image source: <https://www.matroid.com/blog/post/the-hard-thing-about-deep-learning>

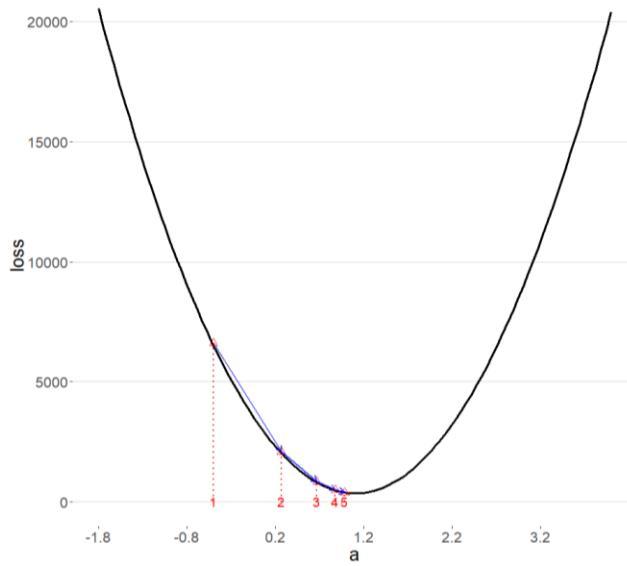
# Stochastic gradient in DL

- Calculate the gradient (vector of partial derivatives) of the Loss w.r.t to all weights, update weights and repeat
- Take a small step (learning rate) in the negative direction of the gradient which is the steepest descent
  - GD: calculate all the gradients with **all** training data and do one update step
    - Often not possible because (memory) of the size of the training data and impractical.
  - SGD: calculate gradient (approximation) of **one random sample** of the training data and do one update step.
  - MBSGD: calculate gradients (approximation) of **mini-batch random sample** of the training data and do one update step.
  - OGD: same as SGD

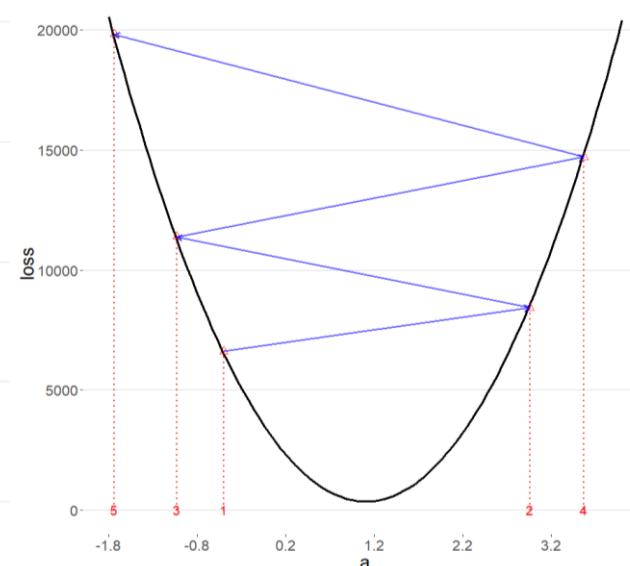
# Stochastic gradient in DL

- Learning rate is the most important hyper parameter in gradient descent!
- Toy example simple linear regression (convex Loss function):

small learning rate



large learning rate



**Keep clam and lower your learning rate!**

# Backpropagation and chain rule

# Chain rule recap

- If we have two functions

$$z = f(y) \text{ and}$$

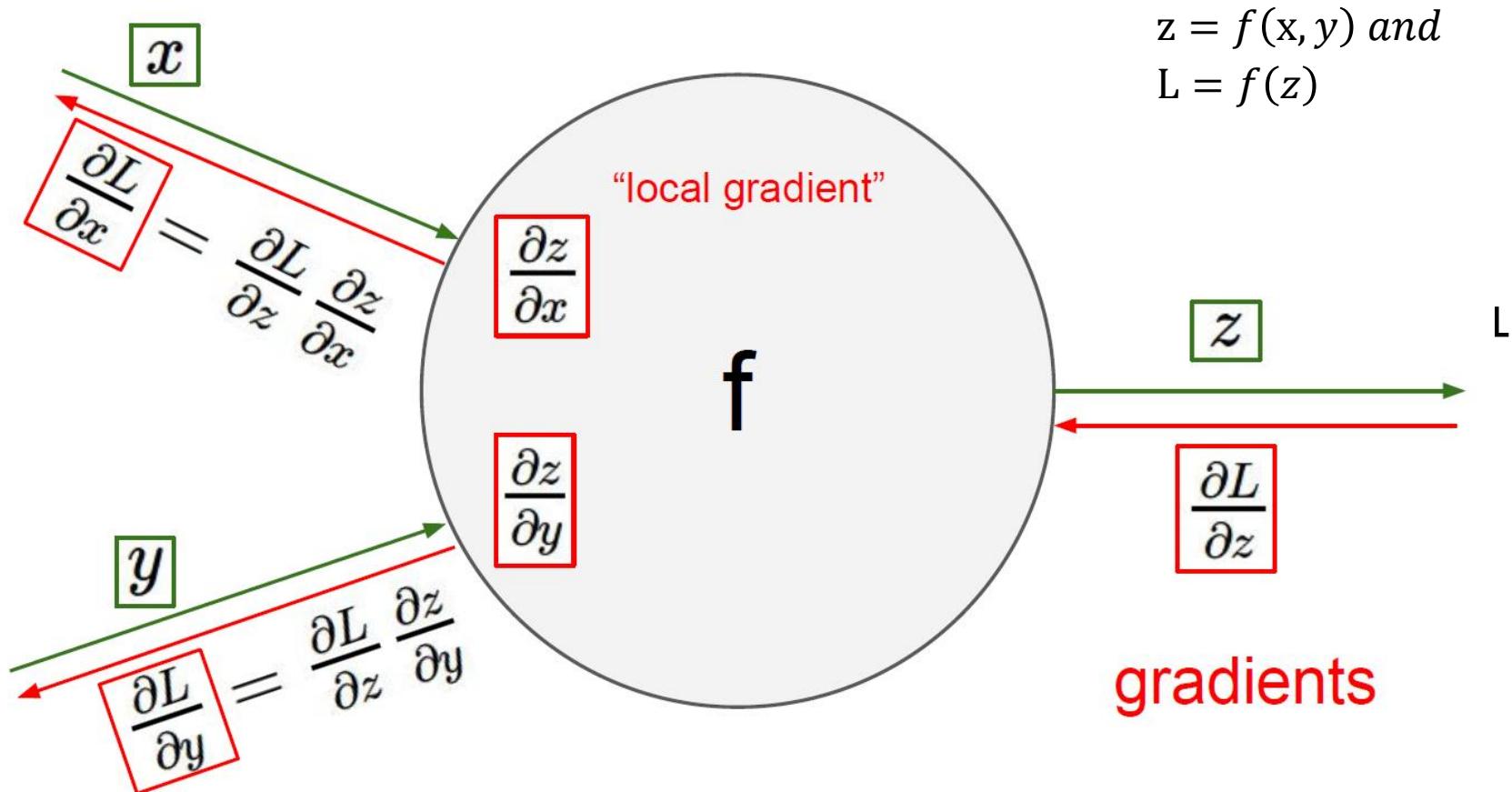
$$y = f(x)$$

then y and z are dependent variables.

- And by the chain rule:

$$\frac{\partial z}{\partial x} = \frac{\partial y}{\partial x} * \frac{\partial z}{\partial y}$$

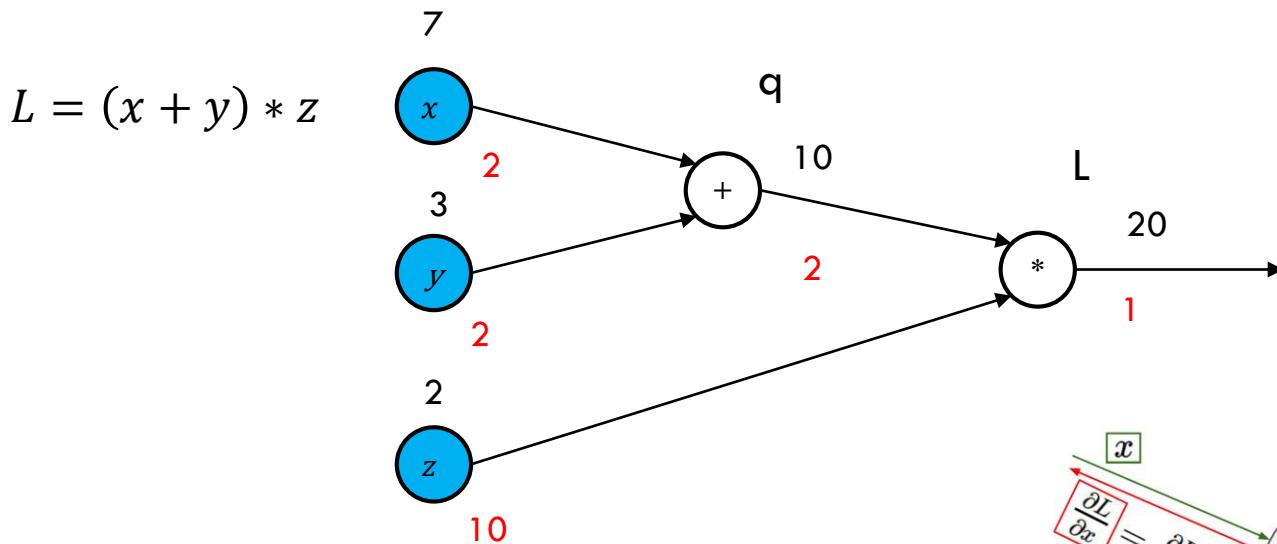
# Backpropagation and chain rule



# Simple example

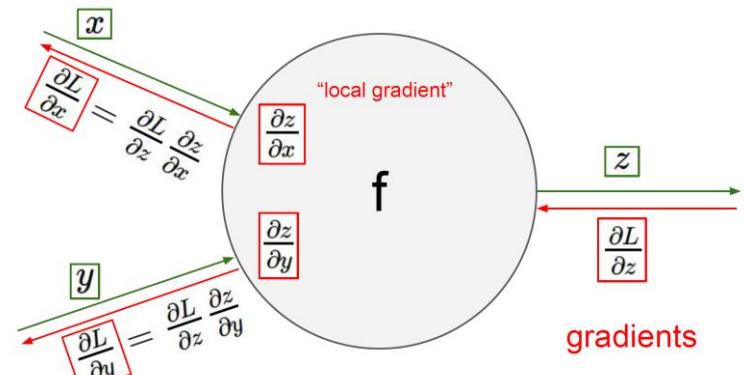
- Multiplications do switch gradients (times the incoming gradient)
- Additions are passing the gradient trough

$$\frac{\partial L}{\partial x} = ? ; \frac{\partial L}{\partial y} = ? ; \frac{\partial L}{\partial z} = ?$$



$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial q} * \frac{1}{\partial q / \partial y} = 1 * 20$$

|   |
|---|
| $f = a * b ; \frac{\partial f}{\partial a} = b ; \frac{\partial f}{\partial b} = a$ |
| $f = a + b ; \frac{\partial f}{\partial a} = 1 ; \frac{\partial f}{\partial b} = 1$ |



# Forward pass

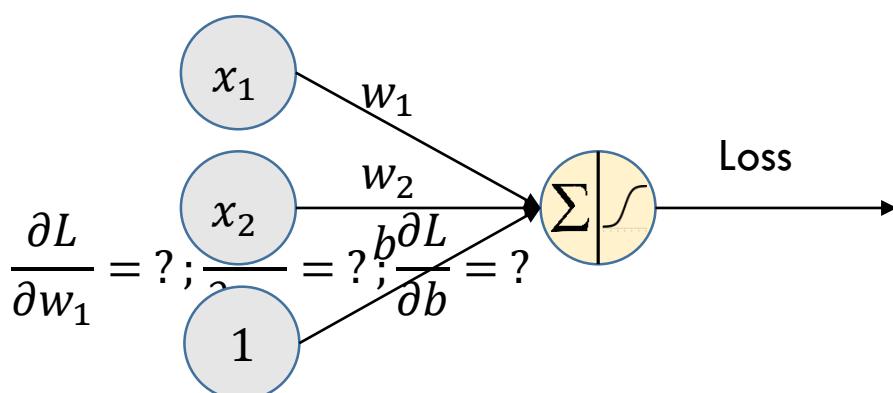
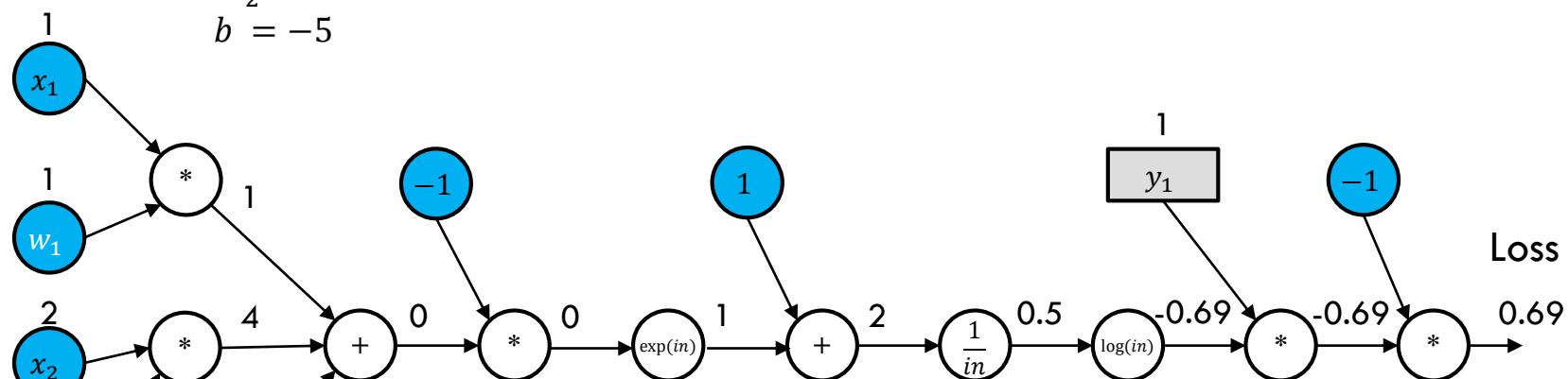
Training data:

$$\begin{aligned}x_1 &= 1 \\x_2 &= 2 \\y_1 &= 1\end{aligned}$$

Initial weights:

$$\begin{aligned}w_1 &= 1 \\w_2 &= 2 \\b &= -5\end{aligned}$$

$$p(y = 1|X) = \frac{1}{1 + e^{-(x_1 \cdot w_1 + x_2 \cdot w_2 + b)}}$$



# Backward pass

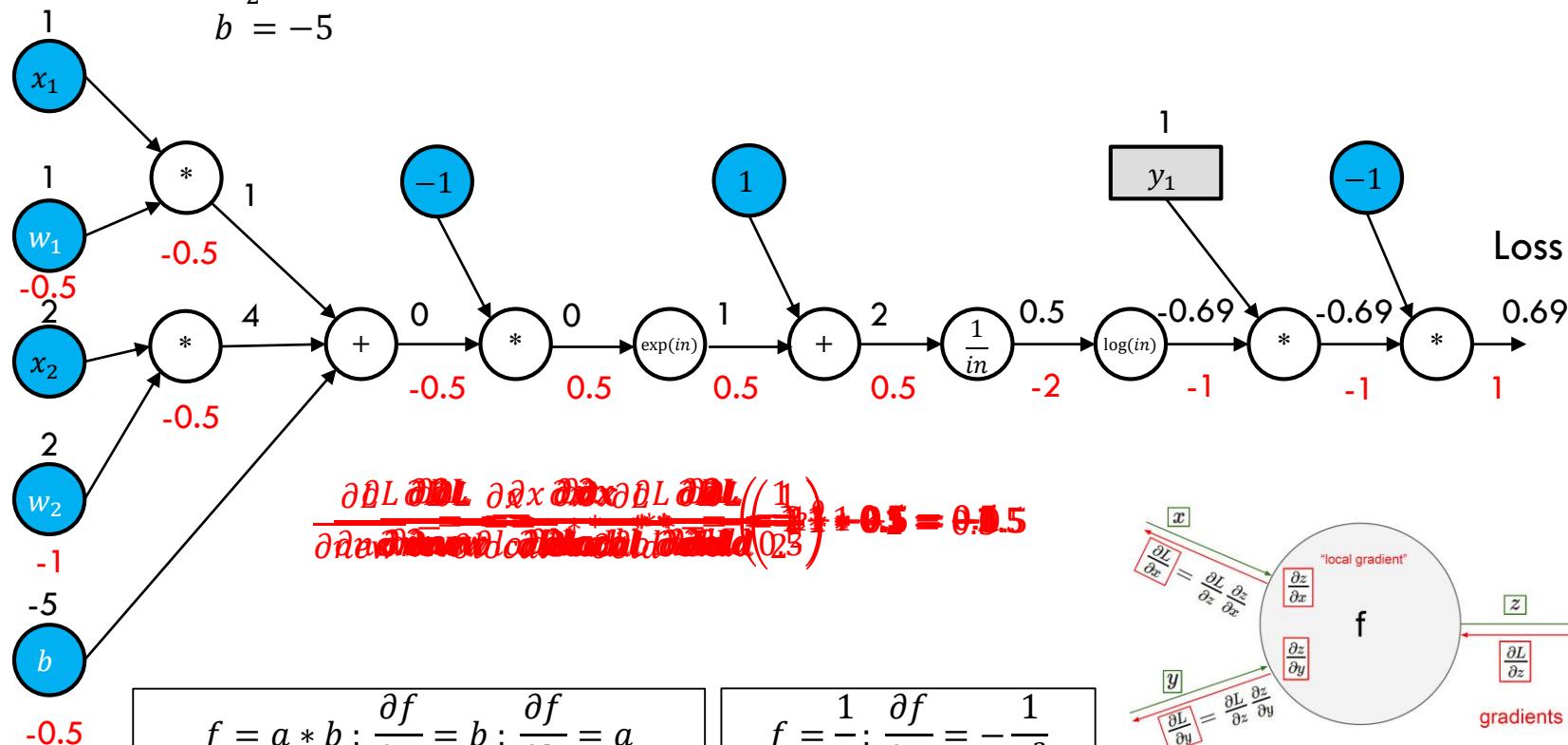
Training data:

$$\begin{aligned}x_1 &= 1 \\x_2 &= 2 \\y_1 &= 1\end{aligned}$$

Initial weights:

$$\begin{aligned}w_1 &= 1 \\w_2 &= 2 \\b &= -5\end{aligned}$$

$$p(y = 1|X) = \frac{1}{1 + e^{-(x_1 * w_1 + x_2 * w_2 + b)}}$$



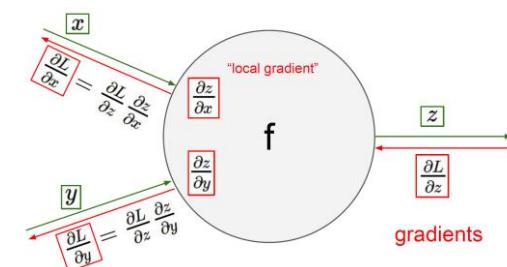
$$f = a * b ; \frac{\partial f}{\partial a} = b ; \frac{\partial f}{\partial b} = a$$

$$f = \frac{1}{a} ; \frac{\partial f}{\partial a} = -\frac{1}{a^2}$$

$$f = a + b ; \frac{\partial f}{\partial a} = 1 ; \frac{\partial f}{\partial b} = 1$$

$$f = e^a ; \frac{\partial f}{\partial a} = e^a$$

$$f = \log(a) ; \frac{\partial f}{\partial a} = \frac{1}{a}$$



# Forward pass

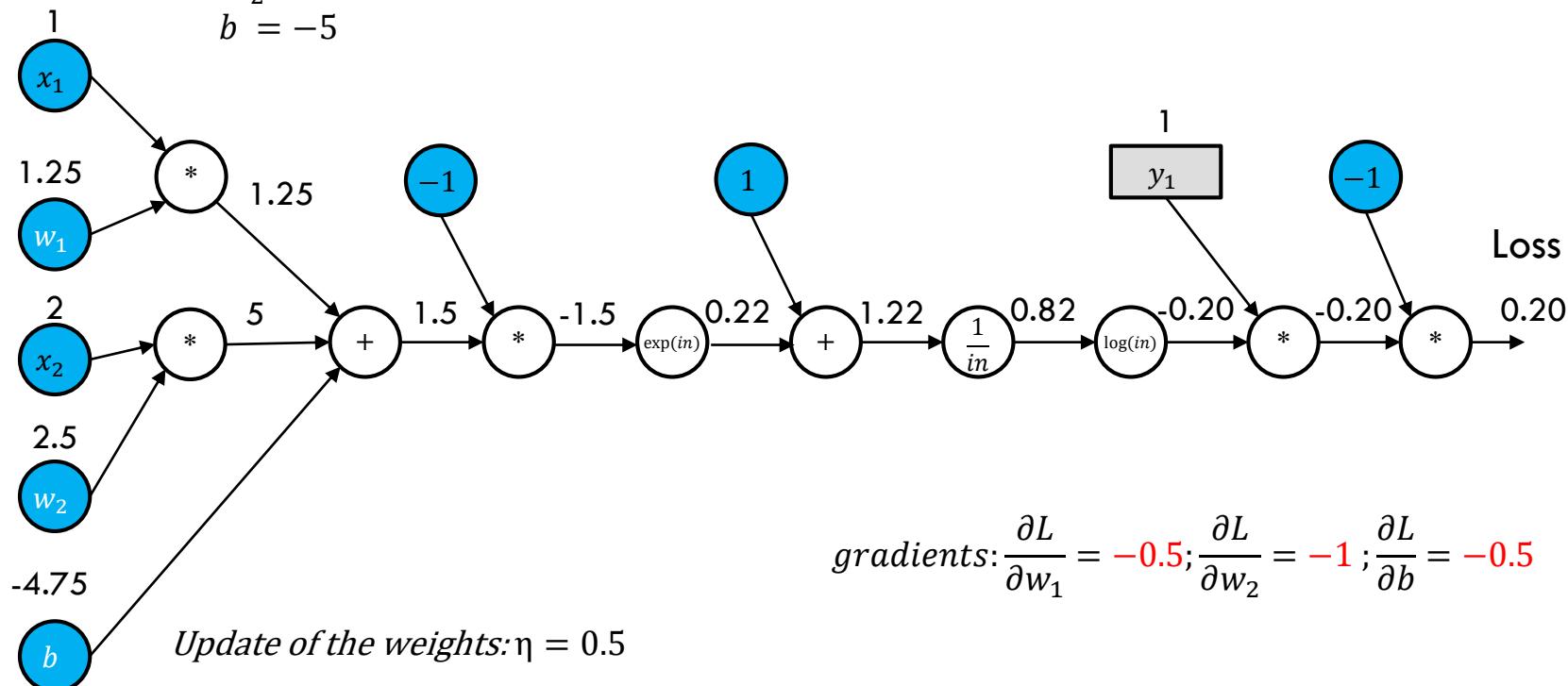
Training data:

$$\begin{aligned}x_1 &= 1 \\x_2 &= 2 \\y_1 &= 1\end{aligned}$$

Initial weights:

$$\begin{aligned}w_1 &= 1 \\w_2 &= 2 \\b &= -5\end{aligned}$$

$$p(y = 1|X) = \frac{1}{1 + e^{-(x_1 \cdot w_1 + x_2 \cdot w_2 + b)}}$$



$$\text{gradients: } \frac{\partial L}{\partial w_1} = -0.5; \frac{\partial L}{\partial w_2} = -1; \frac{\partial L}{\partial b} = -0.5$$

$$w_{1(t+1)} = w_{1(t)} - \eta * \frac{\partial L}{\partial w_1} = 1 - 0.5 * (-0.5) = 1.25$$

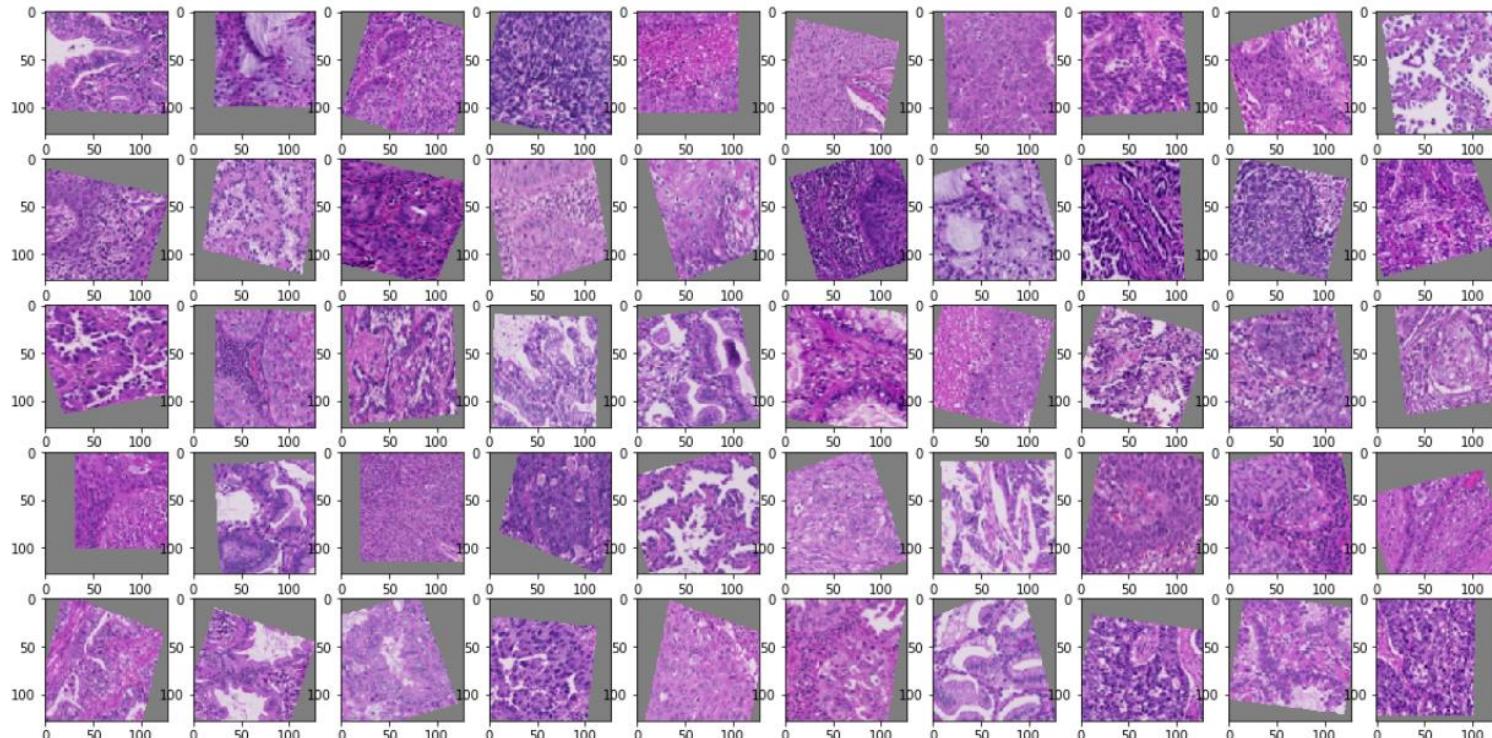
$$w_{2(t+1)} = w_{2(t)} - \eta * \frac{\partial L}{\partial w_2} = 2 - 0.5 * (-1) = 2.5$$

$$b_{(t+1)} = b_{(t)} - \eta * \frac{\partial L}{\partial b} = -5 - 0.5 * (-0.5) = -4.75$$

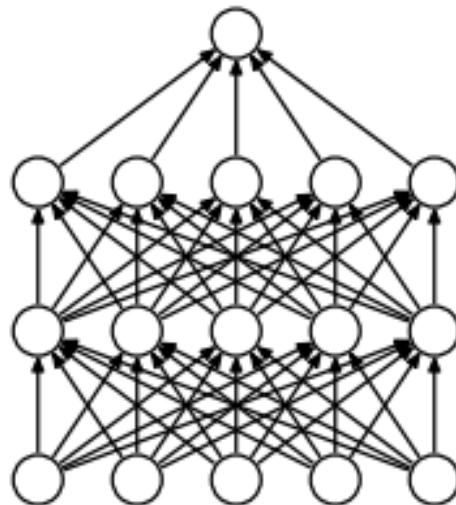
# Overfitting

# Fighting Overfitting

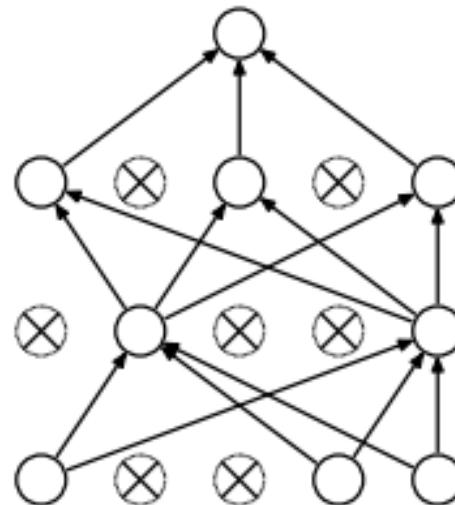
- Data Augmentation
  - Generate additional training data by rotating, zooming, shearing, flipping...



# Fighting Overfitting



(a) Standard Neural Net



(b) After applying dropout.



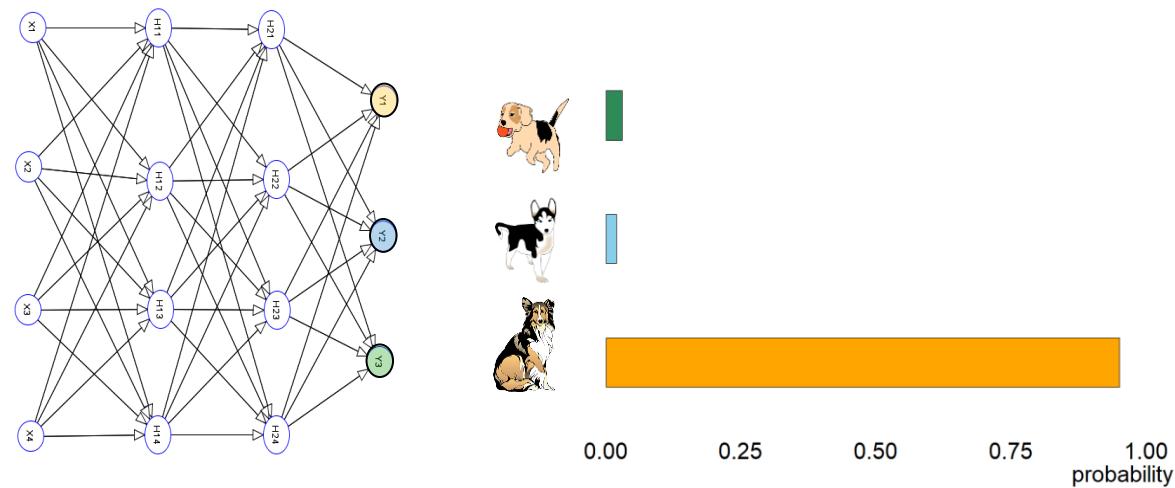
- In each training step we train another sparse NN
- Dropout prevents co-adaptation and overfitting



# Uncertainty

# Uncertainty: A first thought experiment

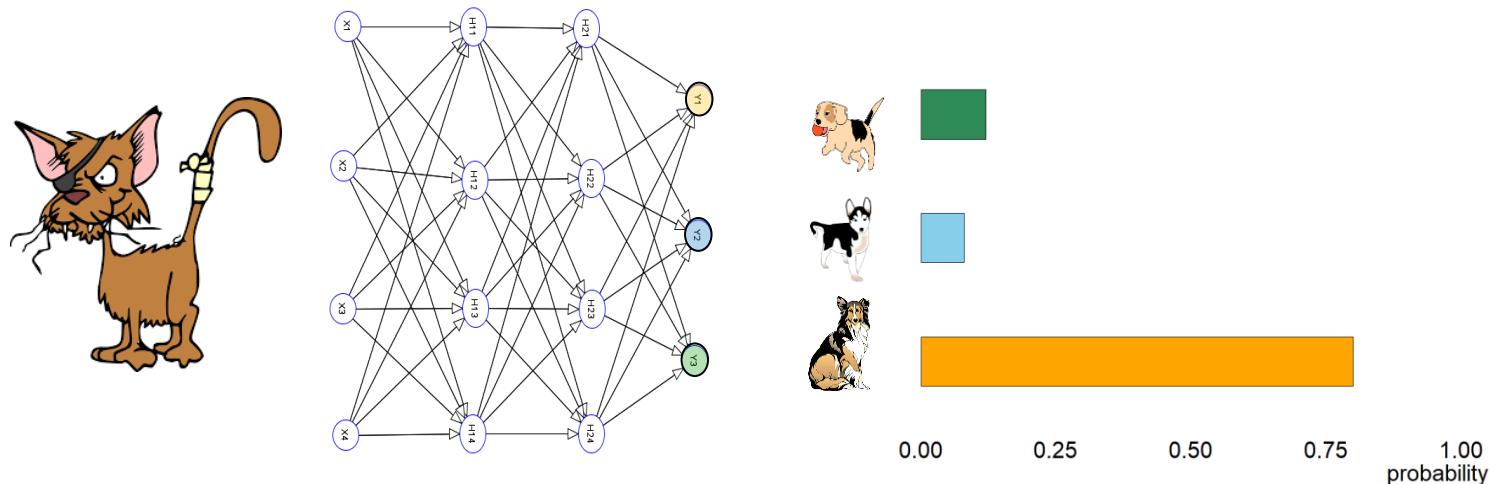
Suppose you train a classifier on dogs and show it a new dog that was in the training set.



# Uncertainty: A first thought experiment

Suppose you show the same classifier a cat

What will be the result?



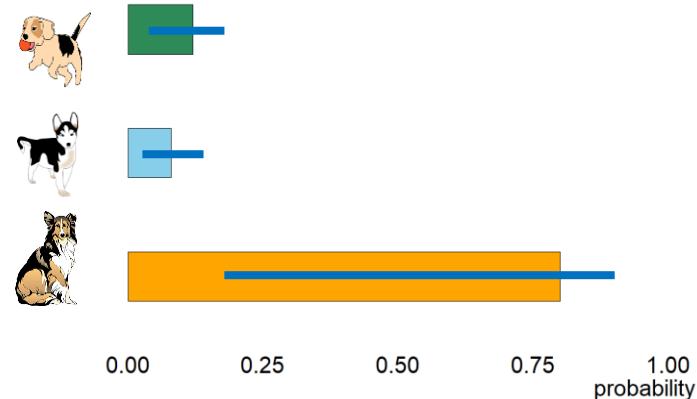
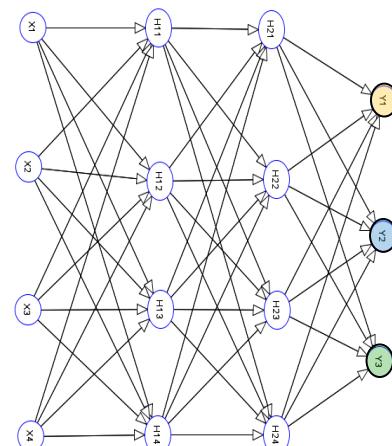
How can that be?

- Forced to classify it into one of the dogs.
- If it's a dog, than most probably a collie
- No confidence of the prediction given

# Uncertainty: A first thought experiment

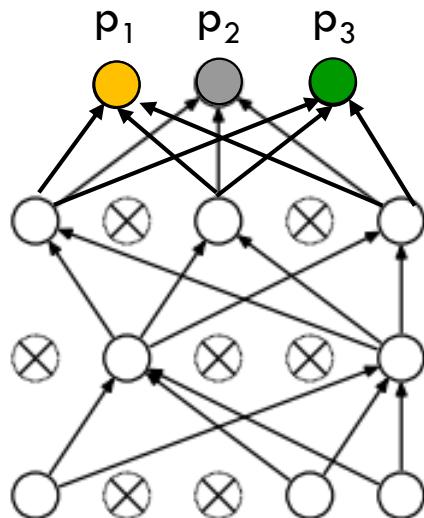
Suppose you show the same classifier a cat

What will be the result?



We need uncertainty measures to our predictions!

# MC Dropout and Bayesian Neural Networks

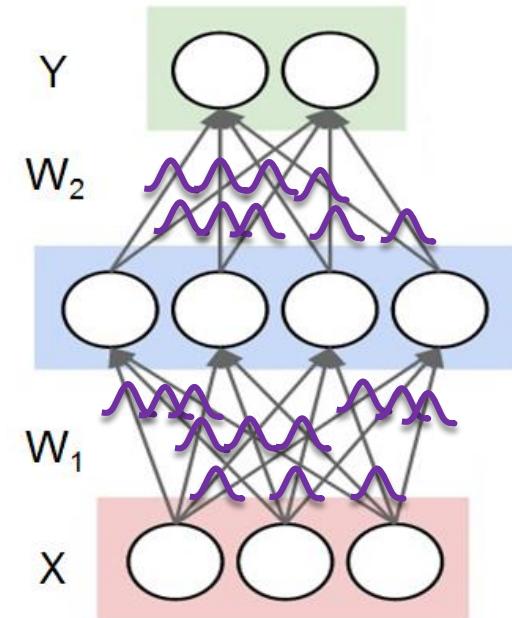
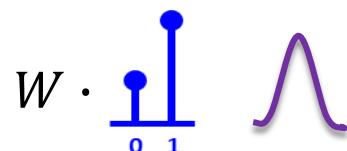


## MC Dropout

Randomly drop nodes  
in each run  
→ Usually done  
during training

Dropout in  
test time

Yarin Gal\* (2015):  
we learn a whole  
weight distribution

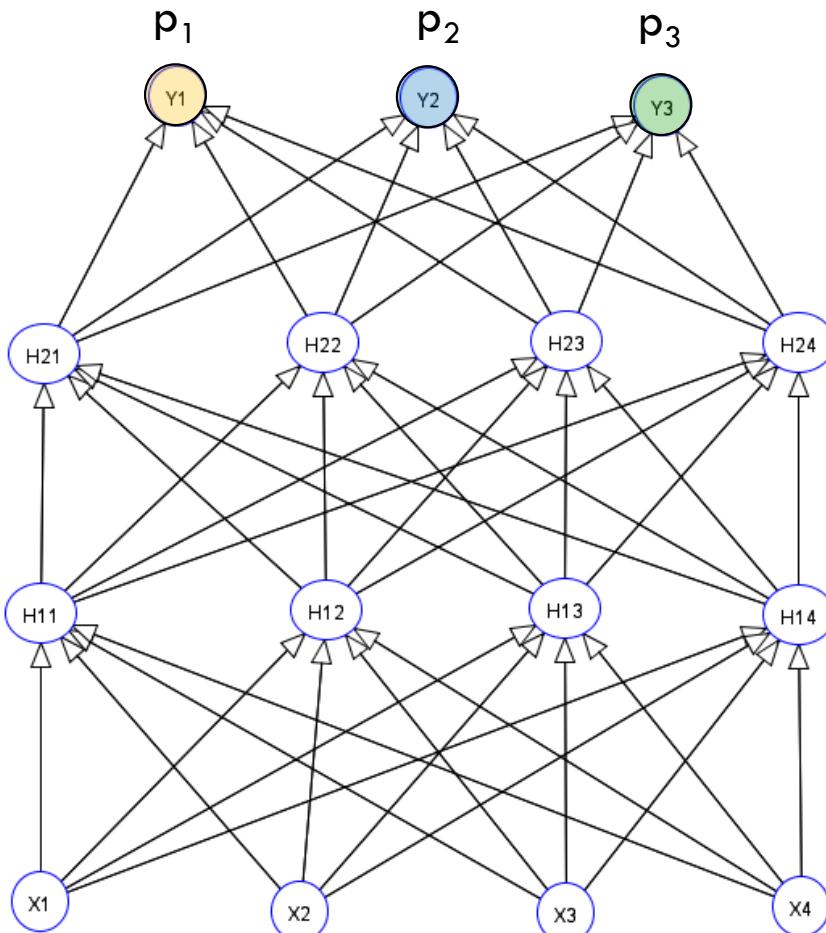


## Bayesian NN

→ We should sample from  
weight distribution during  
test time

# No MC Dropout

Output: probabilities for each class



Done in training anyway  
Why not use it also at test time

Input: image pixel values

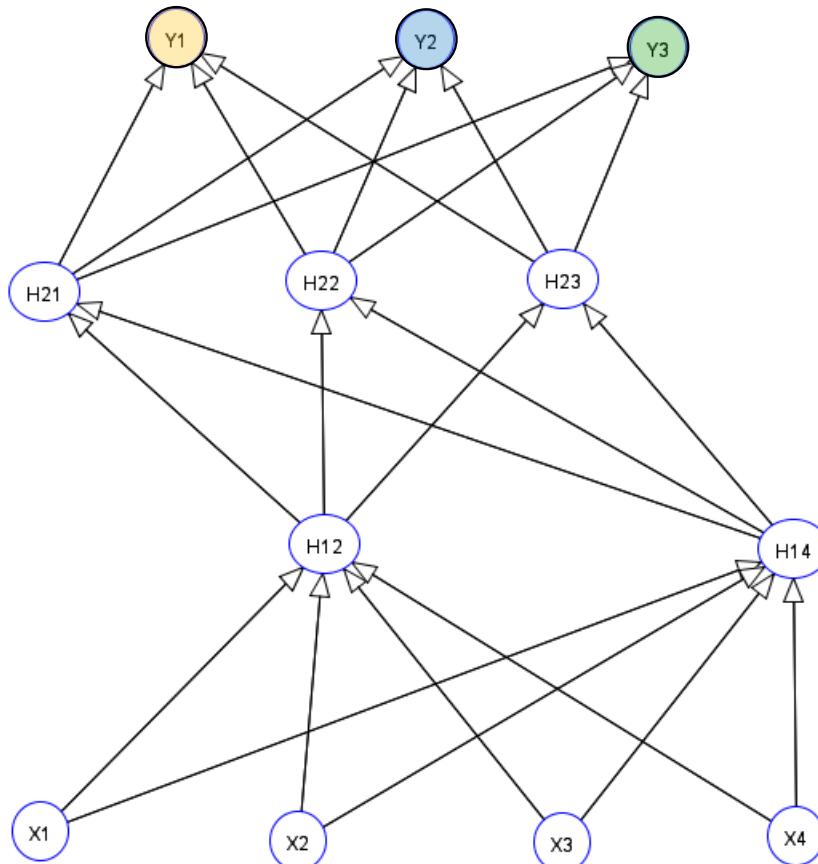


# Use dropout at test time: Run 1

$$p_1 = 0.08$$

$$p_2 = 0.89$$

$$p_3 = 0.03$$



Output depends on dropout

Stochastic dropout of units

Same input image

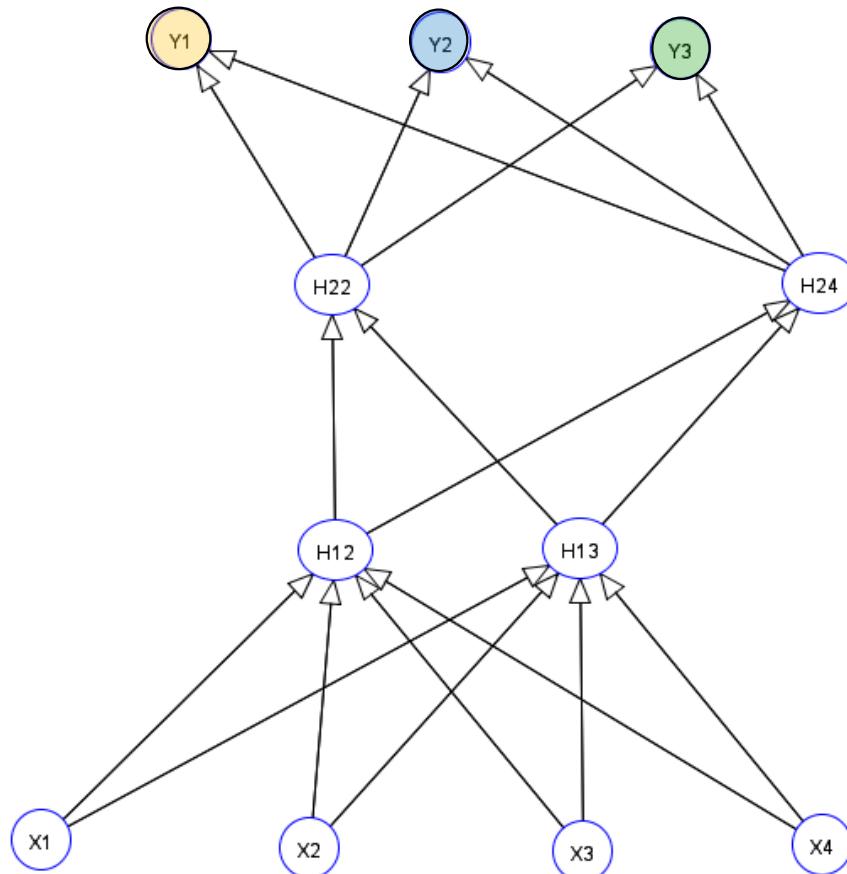


# Use dropout at test time: Run 2

$$p_1=0.11$$

$$p_2=0.81$$

$$p_3=0.08$$



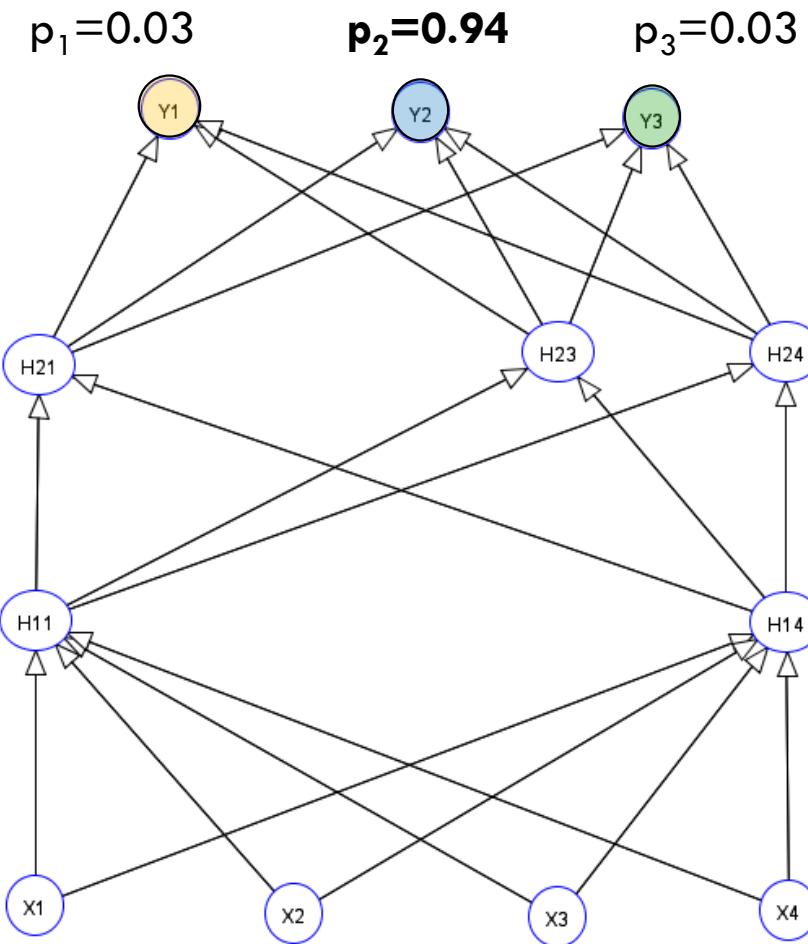
Output depends on dropout

Stochastic dropout of units

Same input image



# Use dropout at test time: Run 3



Output depends on dropout

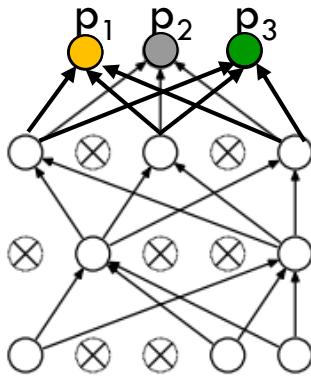
Stochastic dropout of units

Same input image



**...Repeat 500 times**

# Distributions of predicted probabilities



From the predicted distributions we can derive different measures

## Probability estimates

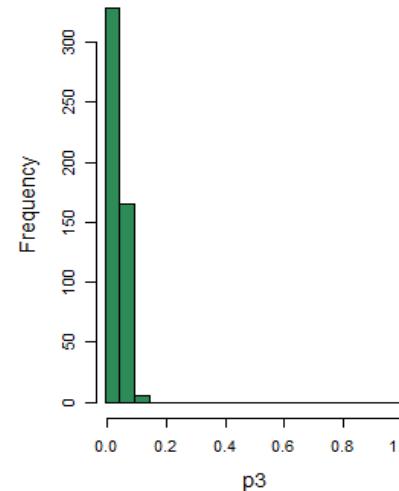
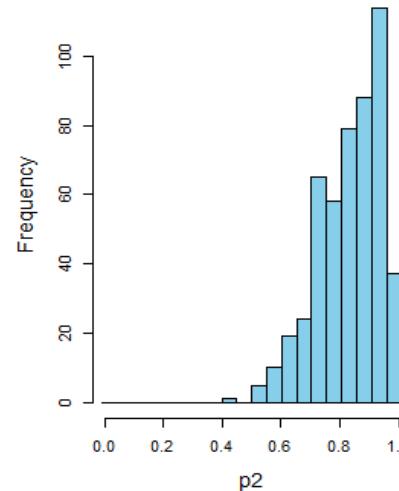
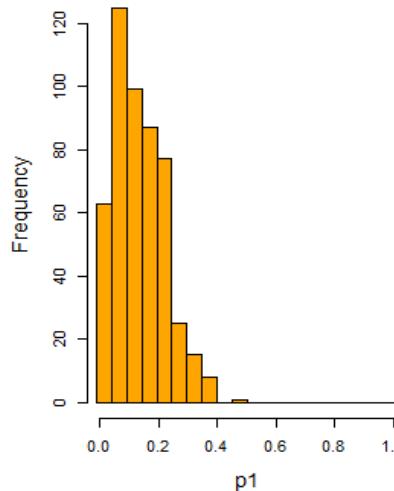
- $p_{\max}$  (when **not** using mc dropout)
- $p_{\max}^*$  (when using mc dropout)

## Uncertainty estimates

- $\sigma^*$  total standard deviation
- PE\* entropy

## Predicted Entropy:

$$PE = - \sum_{l=1}^L (\bar{p}_c l \cdot \log(\bar{p}_c l))$$



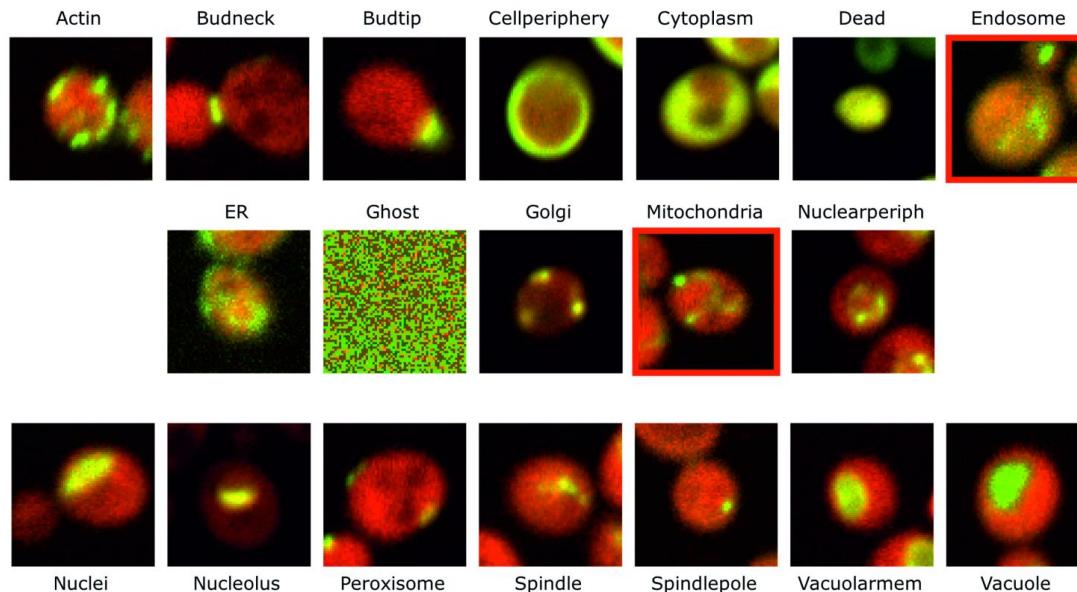
# HCS application

# HCS yeast protein localization data set

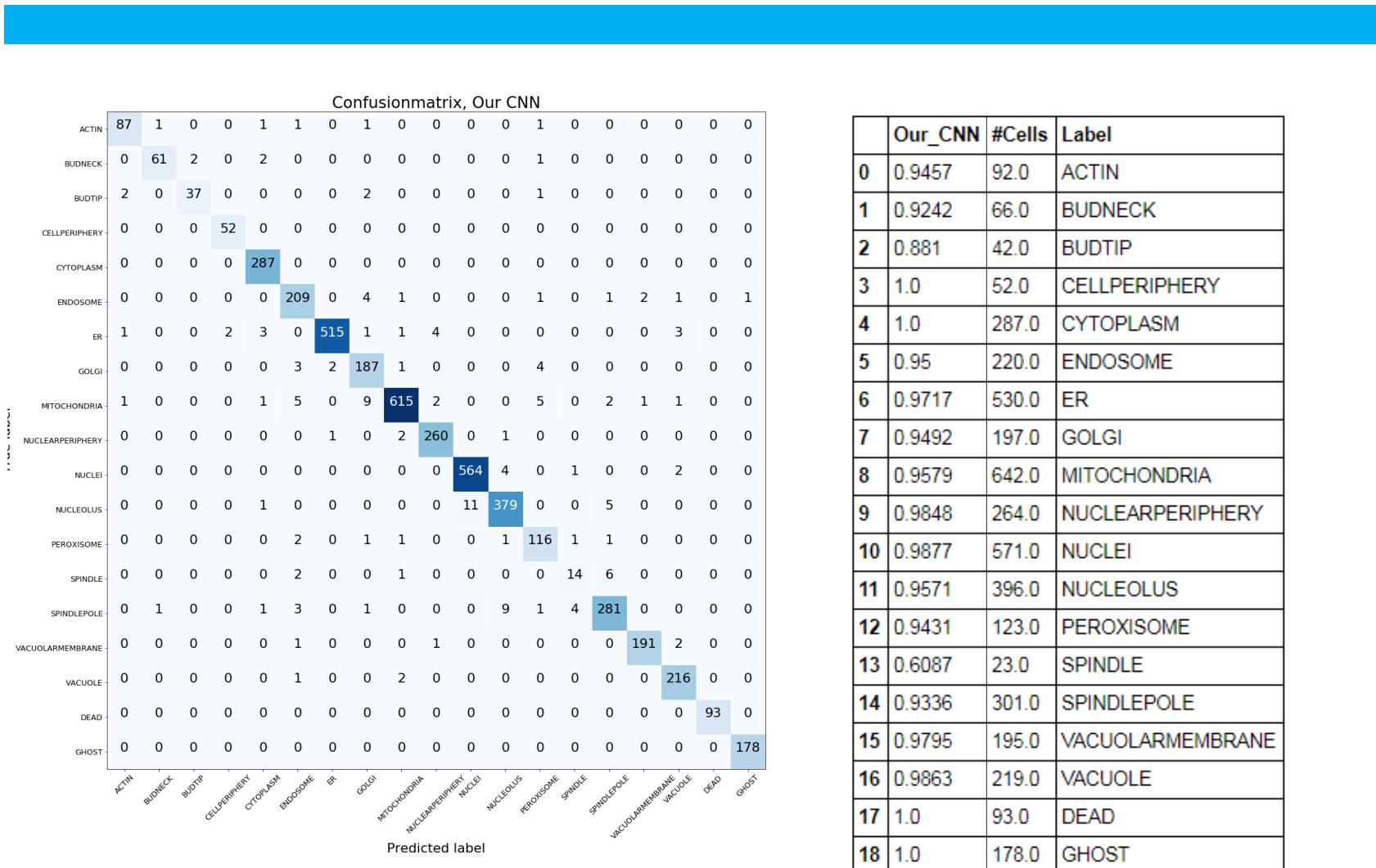
Data from budding yeast protein localizations

(all available from <https://github.com/okraus/DeepLoc>)

- 19 Classes
- 2 Channels
- 21882 64x64x2 segmented images in training set
- 4491 validation and 4516 for testing



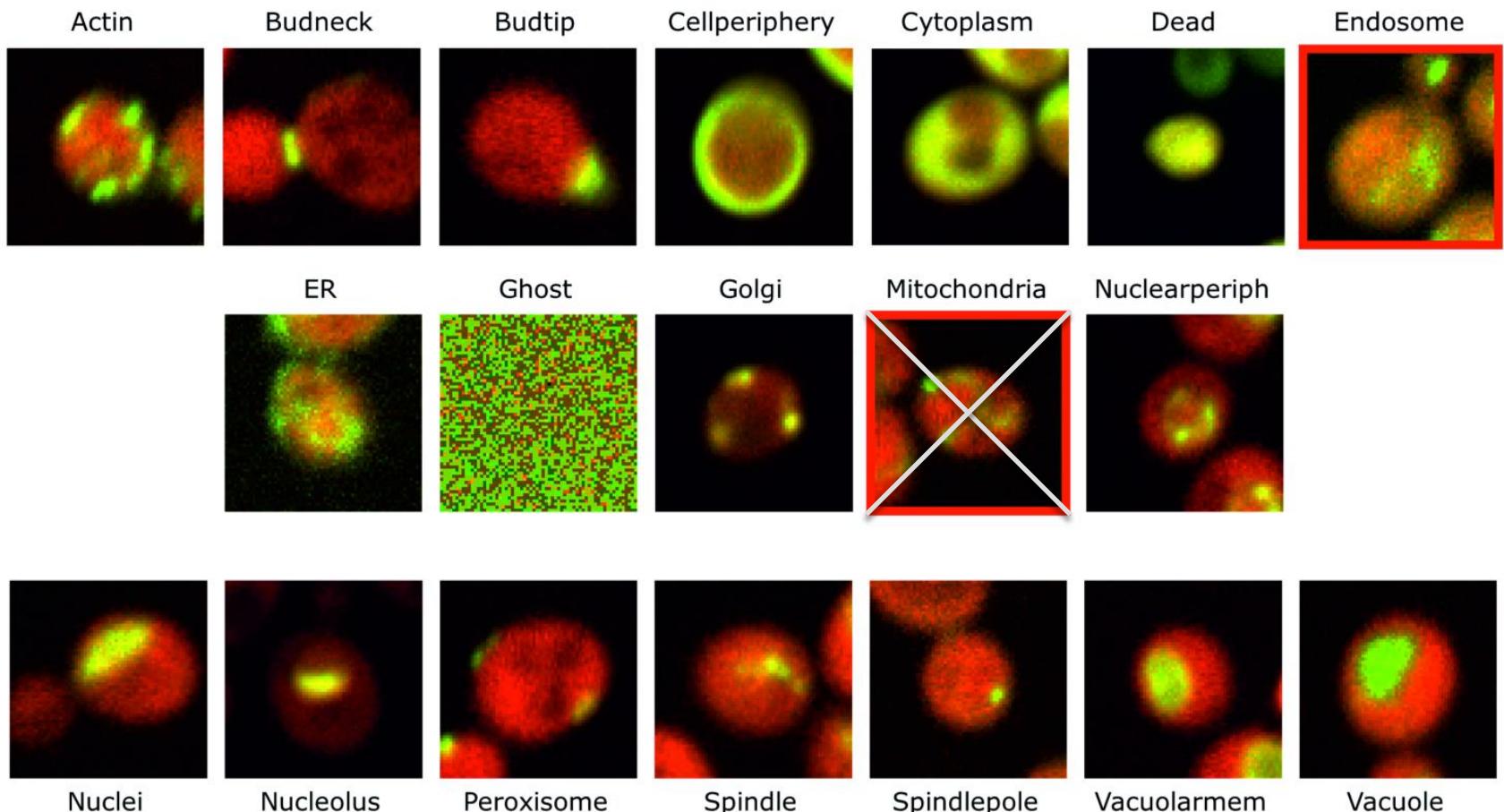
# Classification performance



|    | Our_CNN | #Cells | Label            |
|----|---------|--------|------------------|
| 0  | 0.9457  | 92.0   | ACTIN            |
| 1  | 0.9242  | 66.0   | BUDNECK          |
| 2  | 0.881   | 42.0   | BUDTIP           |
| 3  | 1.0     | 52.0   | CELLPERIPHERY    |
| 4  | 1.0     | 287.0  | CYTOPLASM        |
| 5  | 0.95    | 220.0  | ENDOSOME         |
| 6  | 0.9717  | 530.0  | ER               |
| 7  | 0.9492  | 197.0  | GOLGI            |
| 8  | 0.9579  | 642.0  | MITOCHONDRIA     |
| 9  | 0.9848  | 264.0  | NUCLEARPERIPHERY |
| 10 | 0.9877  | 571.0  | NUCLEI           |
| 11 | 0.9571  | 396.0  | NUCLEOLUS        |
| 12 | 0.9431  | 123.0  | PEROXISOME       |
| 13 | 0.6087  | 23.0   | SPINDLE          |
| 14 | 0.9336  | 301.0  | SPINDLEPOLE      |
| 15 | 0.9795  | 195.0  | VACUOLARMEMBRANE |
| 16 | 0.9863  | 219.0  | VACUOLE          |
| 17 | 1.0     | 93.0   | DEAD             |
| 18 | 1.0     | 178.0  | HOST             |

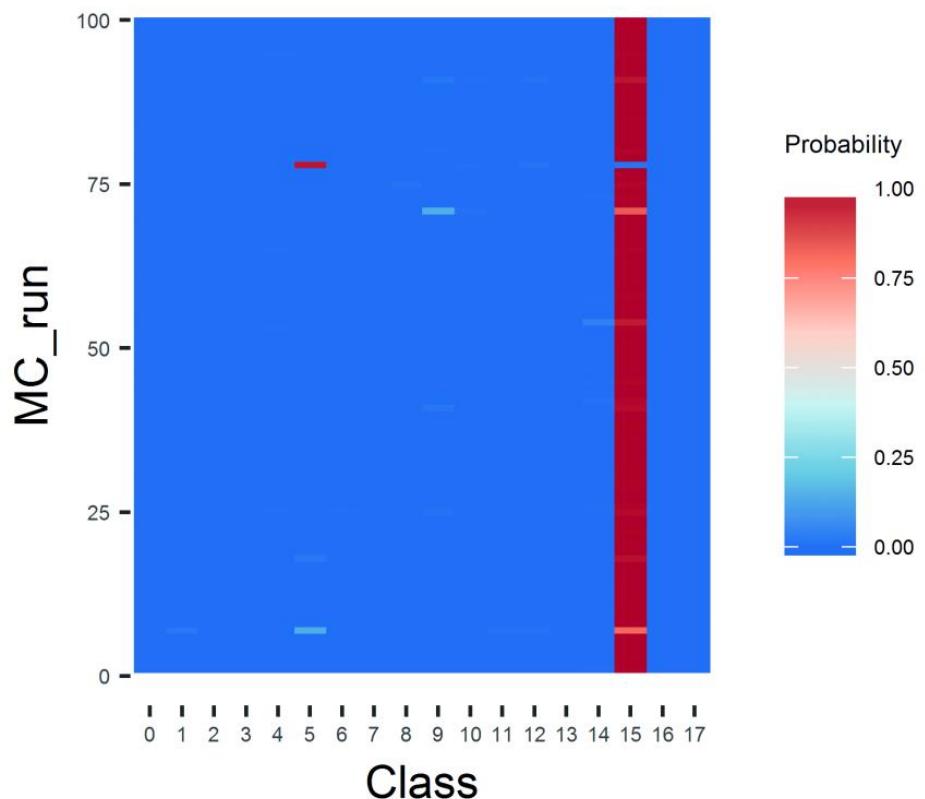
Overall test acc: 96.3% [95.7%,96.8%]

# Experiment with unknown phenotype



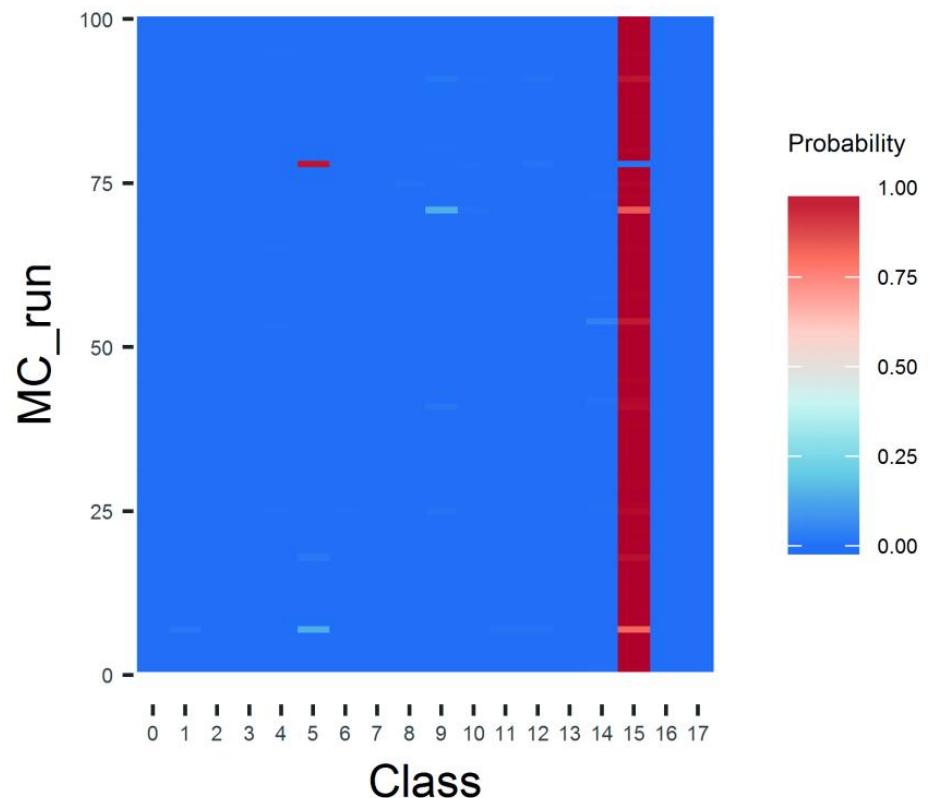
# Experiment with unknown phenotype

100 MC predictions for an image with known phenotype 15

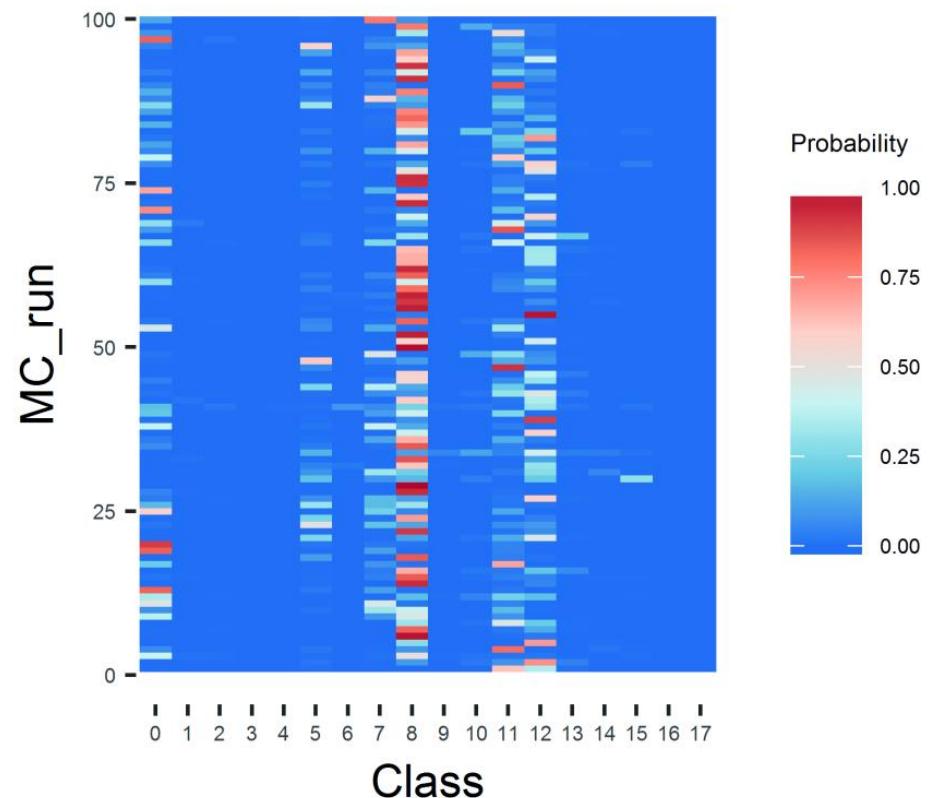


# Experiment with unknown phenotype

100 MC predictions for an image with known phenotype 15



100 MC predictions for an image with an unknown phenotype

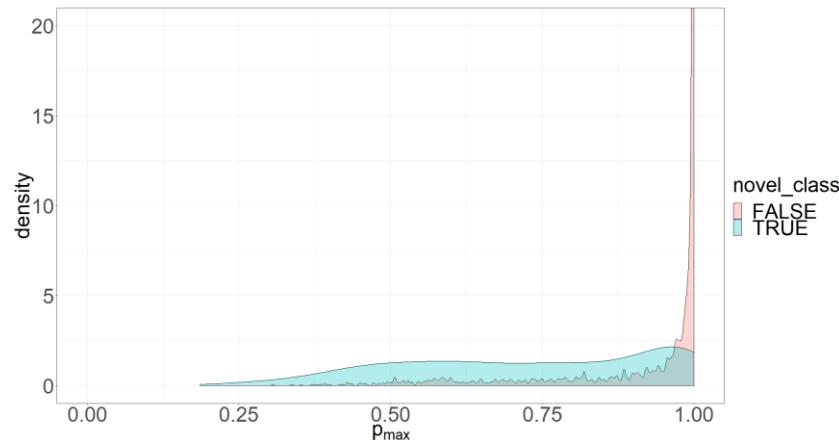


# Experiment with unknown phenotype

## No MC Dropout (classical)

### Probability estimate $p_{\max}$

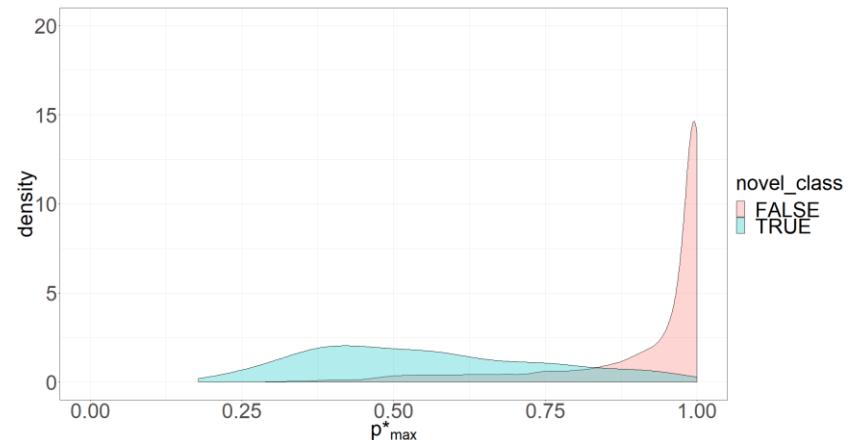
- accuracy (only not novel classes considered):
- 0.9367 [0.9286, 0.9442]**



## MC Dropout

### Probability estimate $p^*_{\max}$

- accuracy (only not novel classes considered):
- 0.9543 [0.9472, 0.9607]**

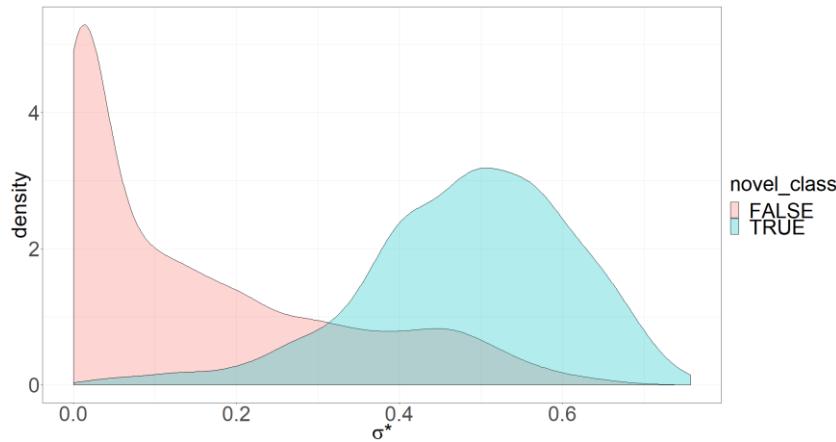


# Experiment with unknown phenotype

## MC Dropout

### Uncertainty estimate $\sigma^*$

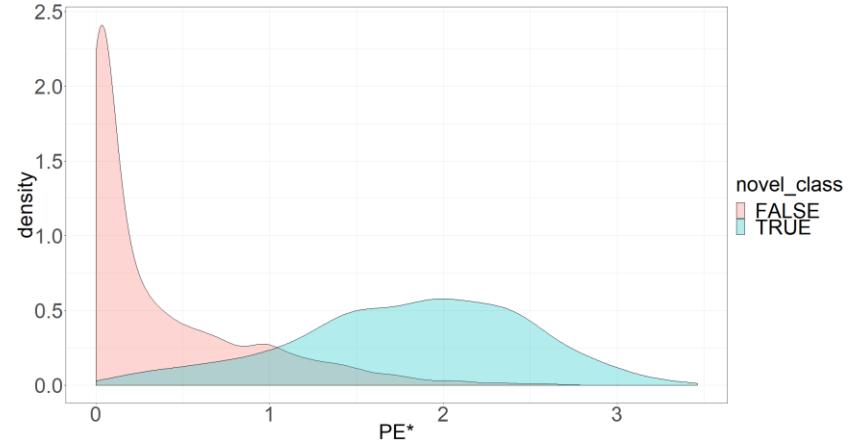
- High  $\sigma^*$  indicates a class not seen in during training
- Low  $\sigma^*$  indicates a confident prediction



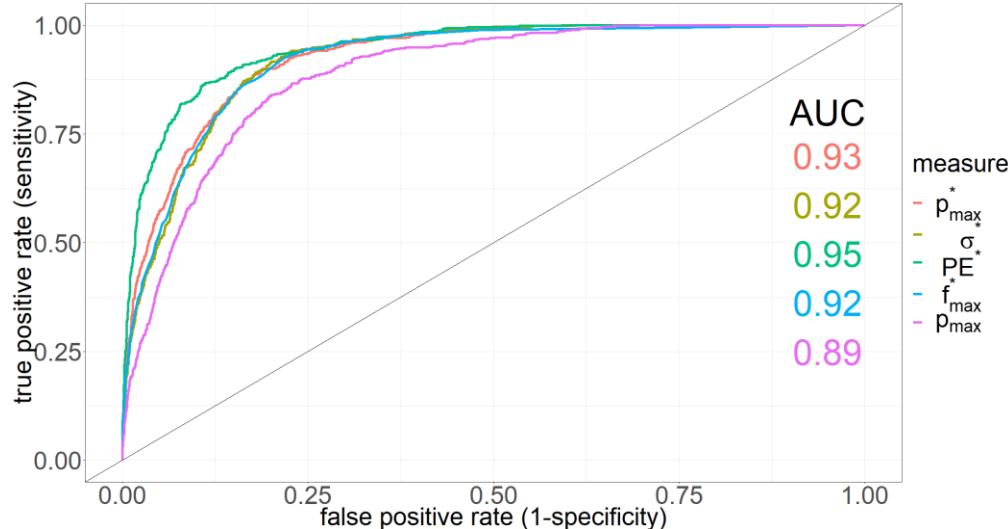
## MC Dropout

### Uncertainty estimate $PE^*$

- High  $PE^*$  indicates a class not seen in during training
- Low  $PE^*$  indicates a confident prediction



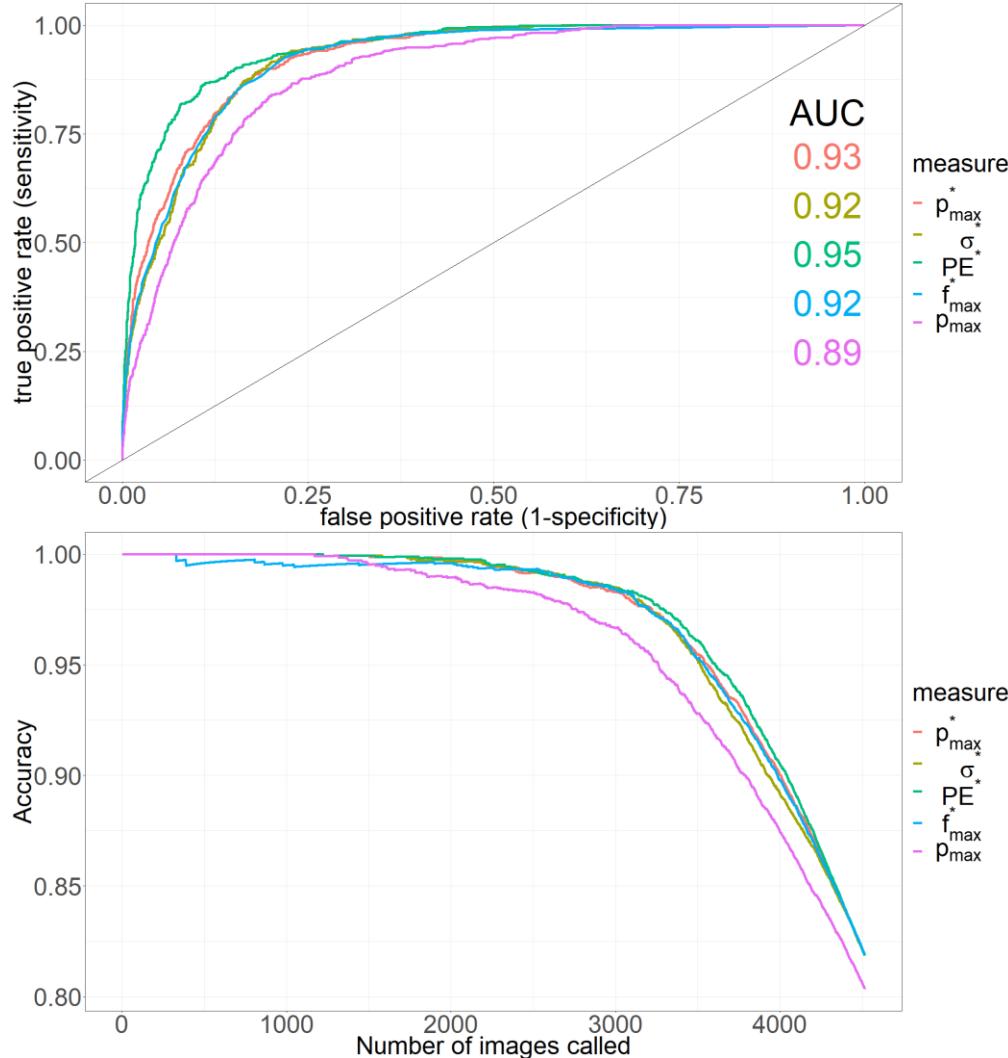
# Experiment with unknown phenotypes



**Discriminative power between known and unkown class**

- All **MC Dropout based approaches are superior to the non-MC based approaches**

# Experiment with unknown phenotypes



## Discriminative power between known and unkown class

- All MC Dropout based approaches are superior to the non-MC based approaches

## Filter uncertain predictions

- In the lift chart, we order the classified images according to the certainty of their call
- $p_{\max}$  which corresponds to the classical probability without MC dropout, is clearly worse compared with the MC dropout approaches

# Main Resources

- Stanford university, CS231: Convolutional Neural Networks for Visual Recognition
  - Li Fei Fei, Andrej Karpathy, Justin Johnson, Serena Yeung
- ZHAW, CAS Machine Intelligence, Modul B: Deep Learning
  - Beate Sick, Oliver Dürr
- MSE, Machine Learning FS2018
  - Oliver Dürr, Thilo Stadelmann
- Deep Learning with Python
  - François Chollet (Google)
- Uncertainty in Deep Learning
  - Yarin Gal (University of Cambridge)

# Thank you, Questions?

