# ISLab

**Universidade do Minho**
Escola de Engenharia
Departamento de Informática
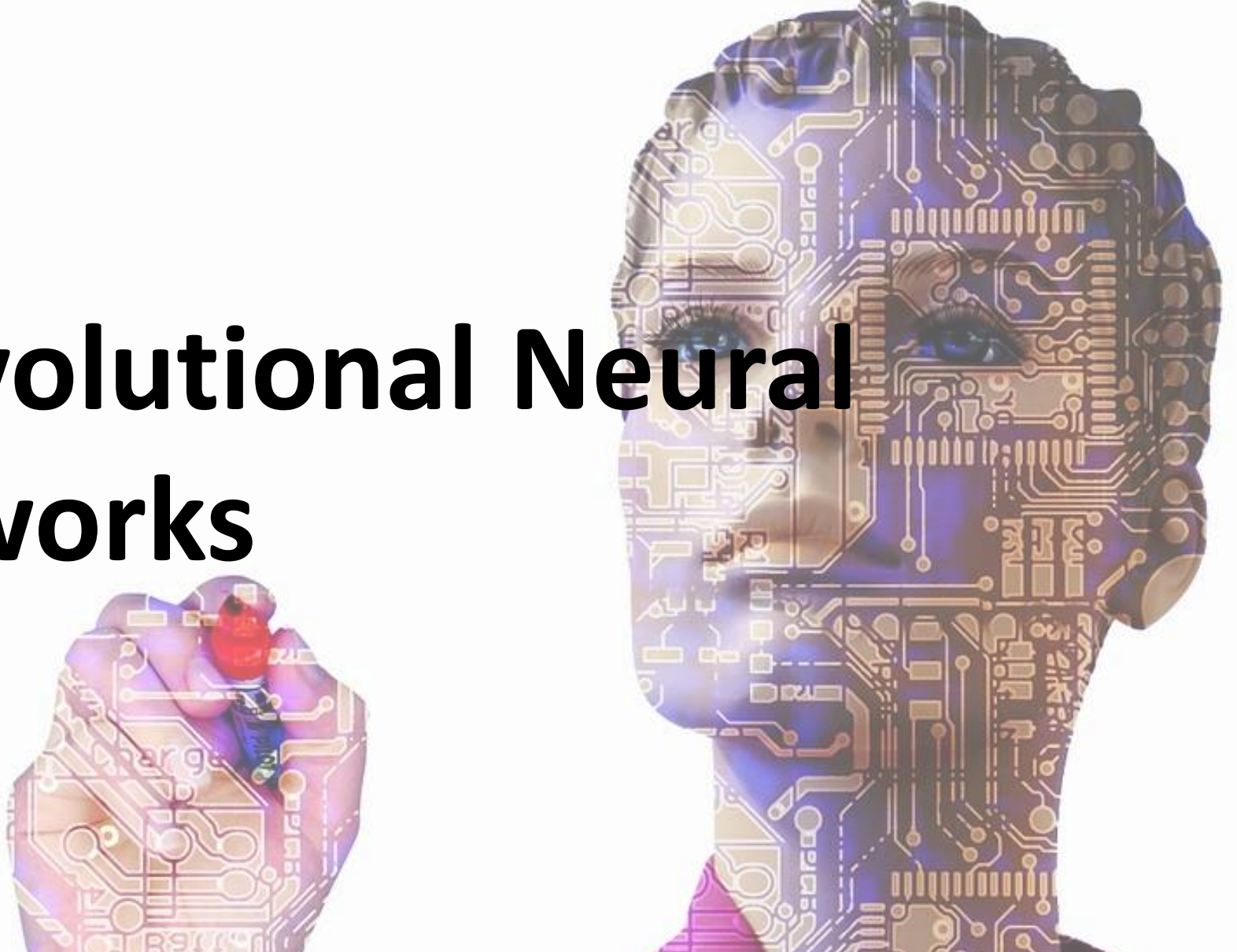
**Mestrado Integrado em Engenharia Informática**
**Mestrado em Engenharia Informática**
**Computação Natural**
**2020/2021**

Filipe Gonçalves, Paulo Novais

- Paulo Novais – pjon@di.uminho.pt

- Filipe Gonçalves – fgoncalves@algoritmi.uminho.pt


- Departamento de Informática
  Escola de Engenharia
  Universidade do Minho

- Grupo ISLab – (Synthetic Intelligence Lab)

- Centro ALGORITMI
  Universidade do Minho

# Convolutional Neural Networks

ISLab
Synthetic Intelligence Lab

**Convolutional Neural Networks (CNN's): what are they for?**
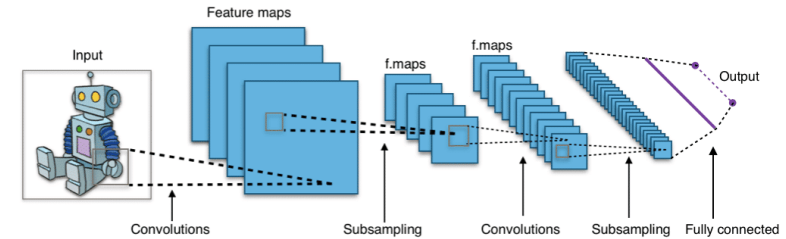
- Applied when you have data that doesn't neatly align into columns
  - o Images that you want to find features within
  - o Machine translation
  - o Sentence classification
  - o Sentiment analysis
- They can find features that aren't in a specific spot
  - o Like a stop sign in a picture
  - o Or words within a sentence
- They are "feature-location invariant"

ISLab
Synthetic Intelligence Lab

**Convolutional Neural Networks (CNN's): how do they work?**

▪ Inspired by the biology of the visual cortex

- o Local receptive fields are groups of neurons that only respond to a part of what your eyes see (sub-sampling)
- o They overlap each other to cover the entire visual field (convolutions)
- o They feed into higher layers that identify increasingly complex images
  - • Some receptive fields identify horizontal lines, lines at different angles, among other features (called feature maps or filters)
  - • These would feed into a layer that identifies shapes
  - • Which might feed into a layer that identifies objects
- o For color or RGB images, 3 layers are used to represent red, green and blue layers
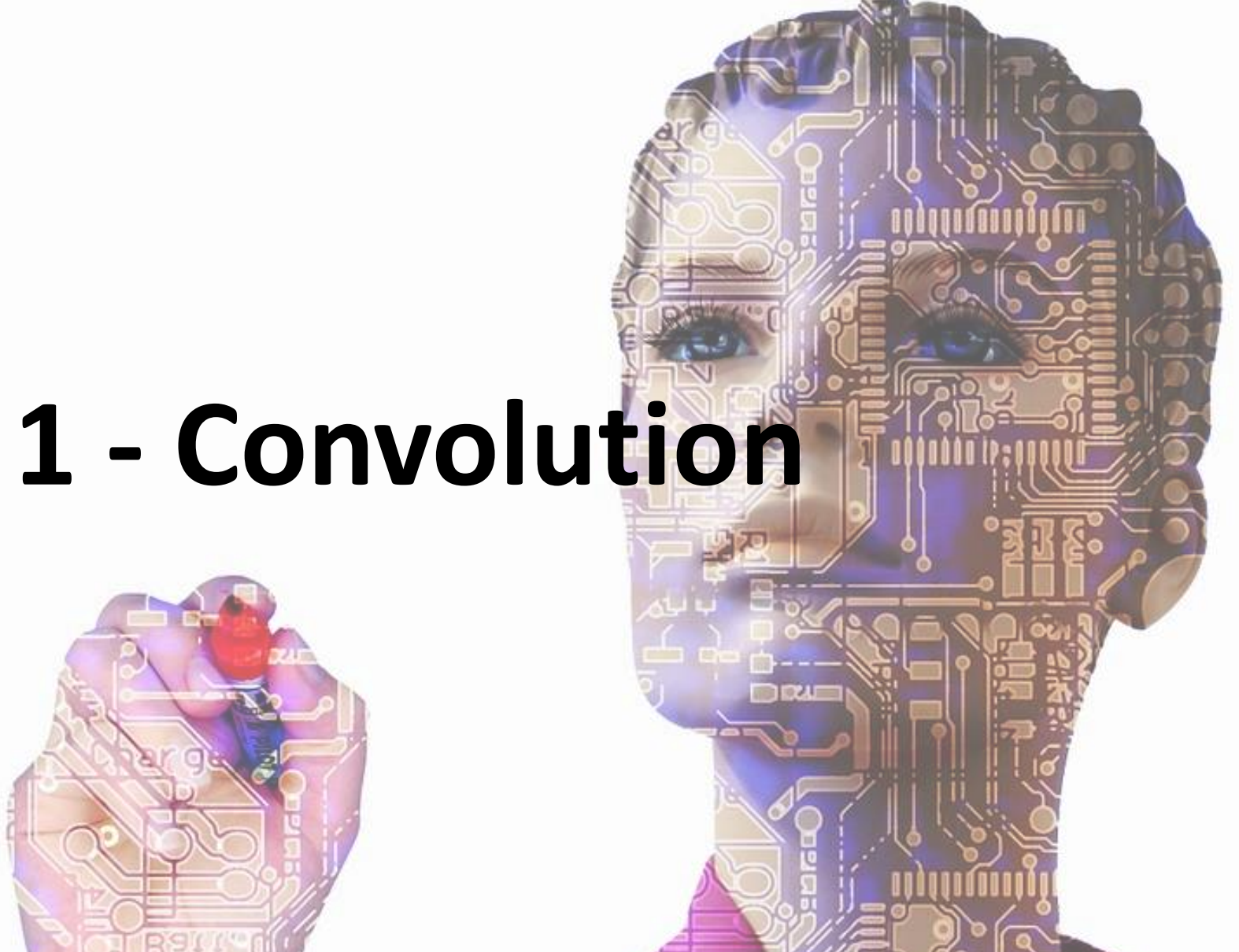
ISLab
Synthetic Intelligence Lab

**How do we know a traffic signal is a stop sign?**

- Individual local receptive fields scan the image looking for edges, and pick up the edges of the stop sign in a layer

- Those edges are used by a higher-level convolution that identifies the stop sign's shape (among other features, e.g., letters)

- The shape then gets matches against the pattern of what a stop sign looks like, also using the strong red signal coming from the red layers

- The information keeps getting processed upward until a decision is made (i.e., classification)

- A CNN works the same way

Step 1 - Convolution

Input Image ⊗ Feature Detector = Feature Map

Input Image ⊗ Feature Detector = Feature Map

Input Image ⊗ Feature Detector = Feature Map

Input Image ⊗ Feature Detector = Feature Map

Input Image ⊗ Feature Detector = Feature Map

We create many feature maps to obtain our first convolution layer

Feature Maps

Input Image

Convolutional Layer

ISLab
Synthetic Intelligence Lab

Sharpen:

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 0 | -1 | 0 | 0 |
| 0 | -1 | 5 | -1 | 0 |
| 0 | 0 | -1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 |



Image Source: docs.gimp.org/en/plug-in-convmatrix.html

Blur:

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



*Image Source: docs.gimp.org/en/plug-in-convmatrix.html*

Emboss:

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



*Image Source: docs.gimp.org/en/plug-in-convmatrix.html*

ISLab
Synthetic Intelligence Lab

Edge Enhance:

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



*Image Source: docs.gimp.org/en/plug-in-convmatrix.html*

ISLab
Synthetic Intelligence Lab

Edge Detect:

| 0 | 1 | 0 |
|---|---|---|
| 1 | -4 | 1 |
| 0 | 1 | 0 |



*Image Source: docs.gimp.org/en/plug-in-convmatrix.html*

Step 2 – ReLU Layer

We create many feature maps to obtain our first convolution layer

Feature Maps

Input Image

Convolutional Layer

ISLab
Synthetic Intelligence Lab



Input Image

Convolutional Layer

Feature Maps

Rectifier

$$\phi(x) = \max(x, 0)$$

$$\sum_{i=1}^{m} w_i x_i$$

Sigmoid

ReLU

Leaky ReLU

Image Source: http://mlss.tuebingen.mpg.de/2015/slides/fergus/Fergus_1.pdf

Black = negative; white = positive values

Image Source: http://mlss.tuebingen.mpg.de/2015/slides/fergus/Fergus_1.pdf

Only non-negative values

Image Source: http://mlss.tuebingen.mpg.de/2015/slides/fergus/Fergus_1.pdf

Step 3 – Pooling

Feature Map

Max Pooling

Pooled Feature Map

Feature Map

Max Pooling

Pooled Feature Map

Max Pooling

Feature Map

Pooled Feature Map

Max Pooling

Feature Map

Pooled Feature Map

Feature Map

Max Pooling

Pooled Feature Map

Input Image

Convolution

Pooling

Convolutional Layer

Pooling Layer

# Step 4 – Flattening

Pooled Feature Map

Flattening

# Step 5 – Fully Connected Layer
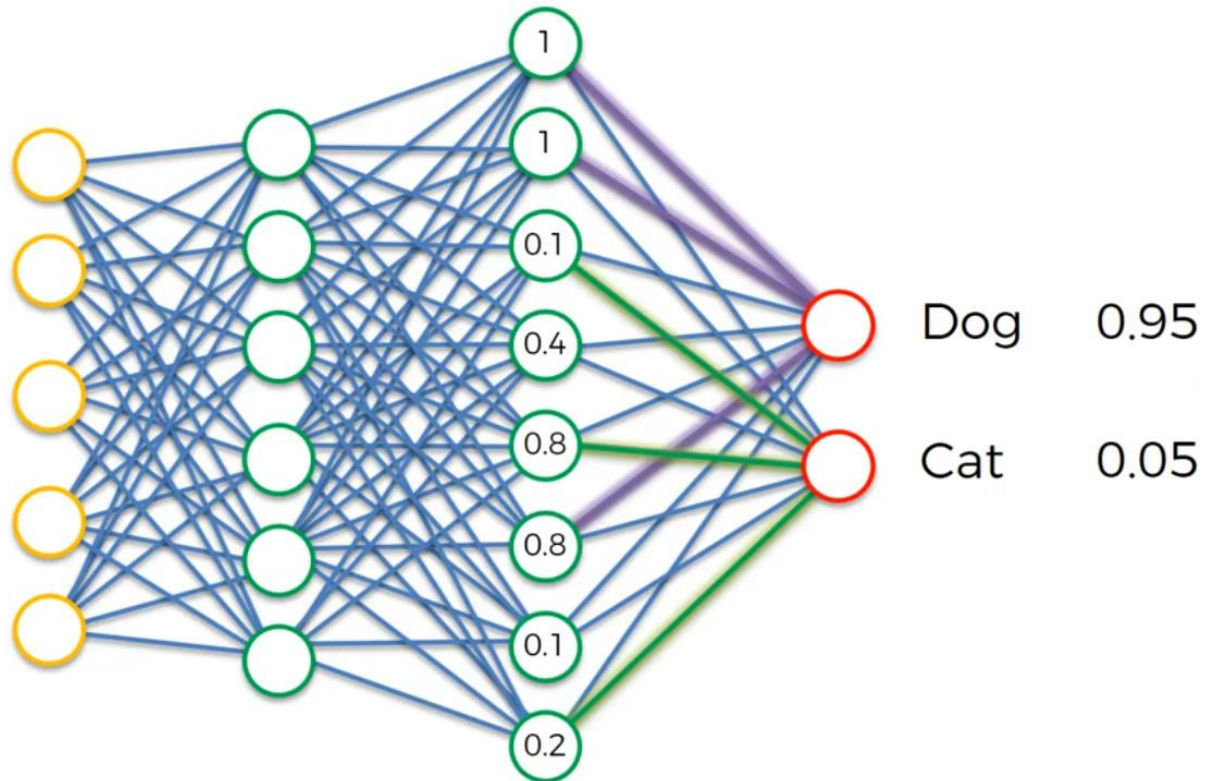
Flattening
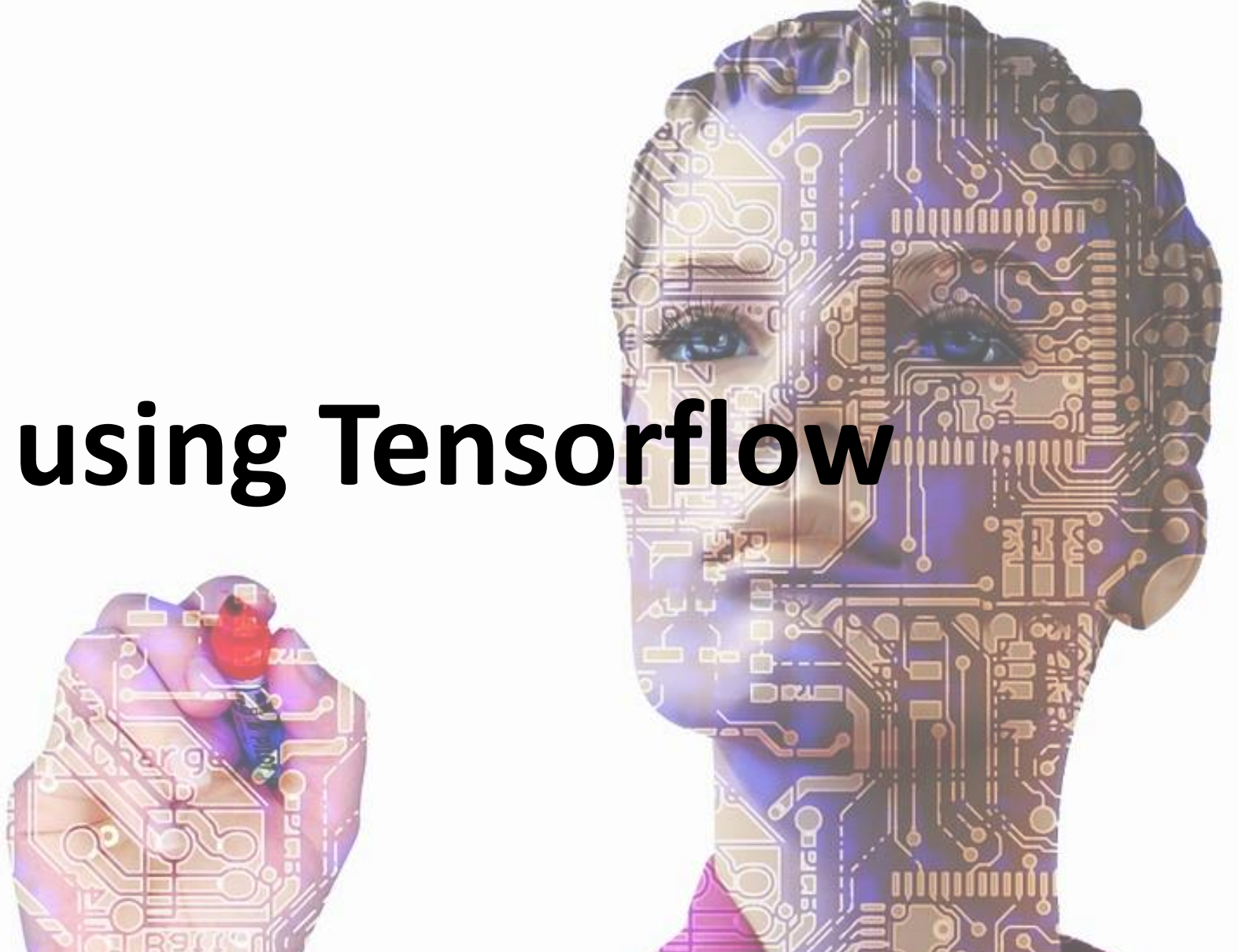
Dog     0.95

Cat     0.05

Examples from the test set
(with the network's guesses)

Image Source: a talk by Geoffrey Hinton

Synthetic Intelligence Lab

**CNN's with Tensorflow**

- Source data must be of appropriate dimensions
  - i.e., width x height x color channels
- Conv2D layer type does the actual convolution on a 2D image
  - Conv1D and Conv3D also available – doesn't have to be image data, e.g., Signal Data and Video Data
- MaxPooling2D layers can be used to reduce a 2D layer down by taking the maximum value in a given kernel
- Flatten layers will convert the 2D layer to a 1D layer for passing into a flat hidden layer of neurons
- Typical architecture use:
  - Conv2D -> MaxPooling2D -> Dropout -> Flatten -> Dense -> Dropout -> Softmax

ISLab
Synthetic Intelligence Lab

**CNN's are resource-intensive**

▪ Uses a lot of computational resources (CPU, GPU and RAM)

▪ Lots of hyper-parameters

   o Kernel sizes, multiple layers with different number of units, amount of pooling, number of layers, choice of optimizers, etc.

▪ Getting the training data is often the hardest part (as well as storing and accessing it)

ISLab
Synthetic Intelligence Lab

**Specialized CNN architectures**

▪ Defines specific arrangement of layers, padding, and hyper-parameters

▪ LeNet-5

   o Good for handwriting recognition

▪ AlexNet

   o Image Classification, deeper than LeNet

▪ VGG

   o Upgrade version of AlexNet
   o Used in multiple contexts with good overall performance

▪ GoogLeNet

   o Even deeper, but with better performance
   o Introduces inception modules (groups of convolution layers)

▪ ResNet (Residual Network)

   o Even deeper – maintains performance via skip connections

**Universidade do Minho**
Escola de Engenharia
Departamento de Informática

**Mestrado Integrado em Engenharia Informática**
**Mestrado em Engenharia Informática**
**Computação Natural**
**2020/2021**

Filipe Gonçalves, Paulo Novais