# Explainable AI für Deep Learning: Overview und Tutorial

Jörg Simon

# About me

- PhD on using deep learning to detect human factors from biosignals

- Prof. Eduardo Veas and Herbert Danzinger

- Sometimes very sparse data!

- Inspired to use interpretability results to change the training process itself
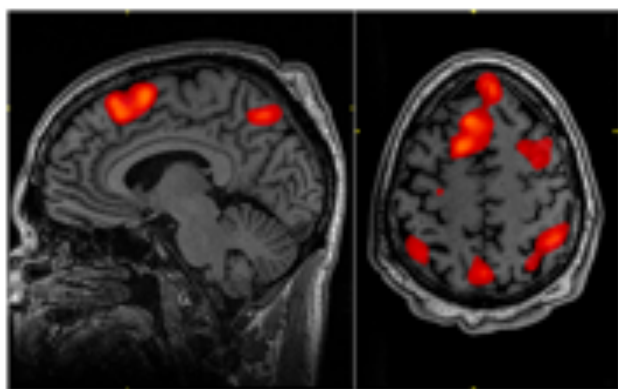
# Agenda

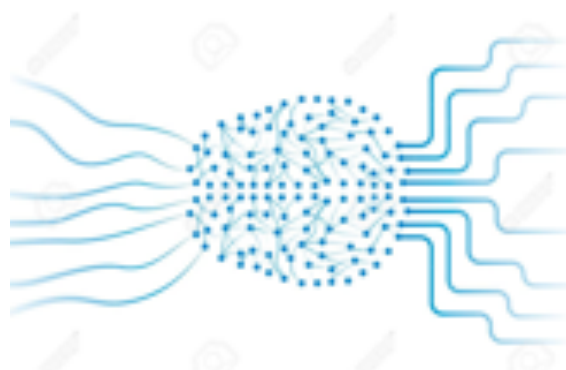- Definitions and Stuff

- Hands On

- Discussion

- Q&A on Discord

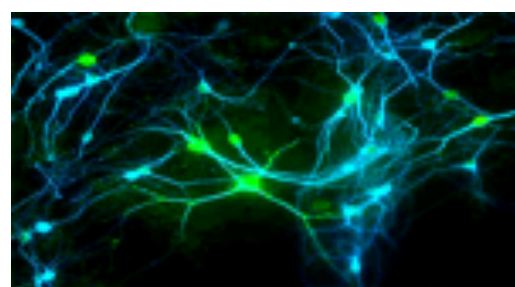# Definitions and Stuff

- Deep learning

Distributed Representation


Super Simplified Model of Human Brain


Hinton


Spiking Frequency = weight

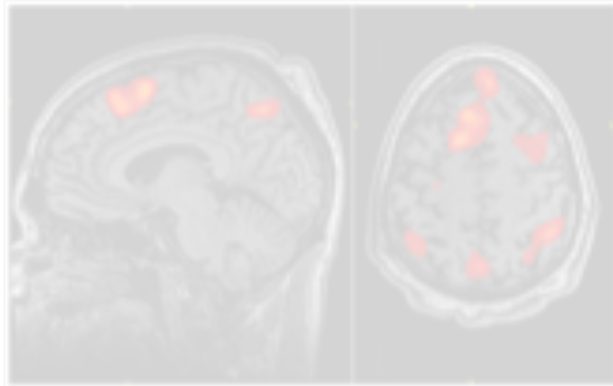Deep Learning?

CNN

Simple Matrix Multiply + Non Linearity

RNN

Yann LeCun

Bengio, Hochreiter, Schmidhuber

Distributed Representation

Spiking Frequency : weight

Simple Matrix Multiply + Non Linearity

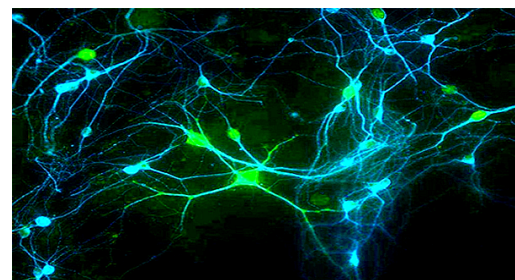Super Simplified Model of Human Brain

Deep Learning?

Hinton

CNN

RNN

Yann LeCun

Bengio, Hochreiter, Schmidhuber

Distributed Representation

Spiking Frequency = weight
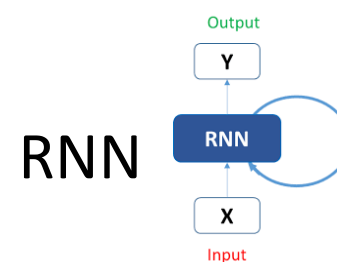
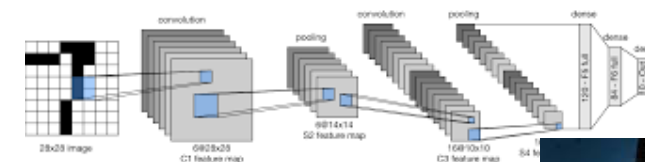Simple Matrix Multiply + Non Linearity

Super Simplified Model of Human Brain

Deep Learning?

CNN

RNN

Hinton

Yann LeCun

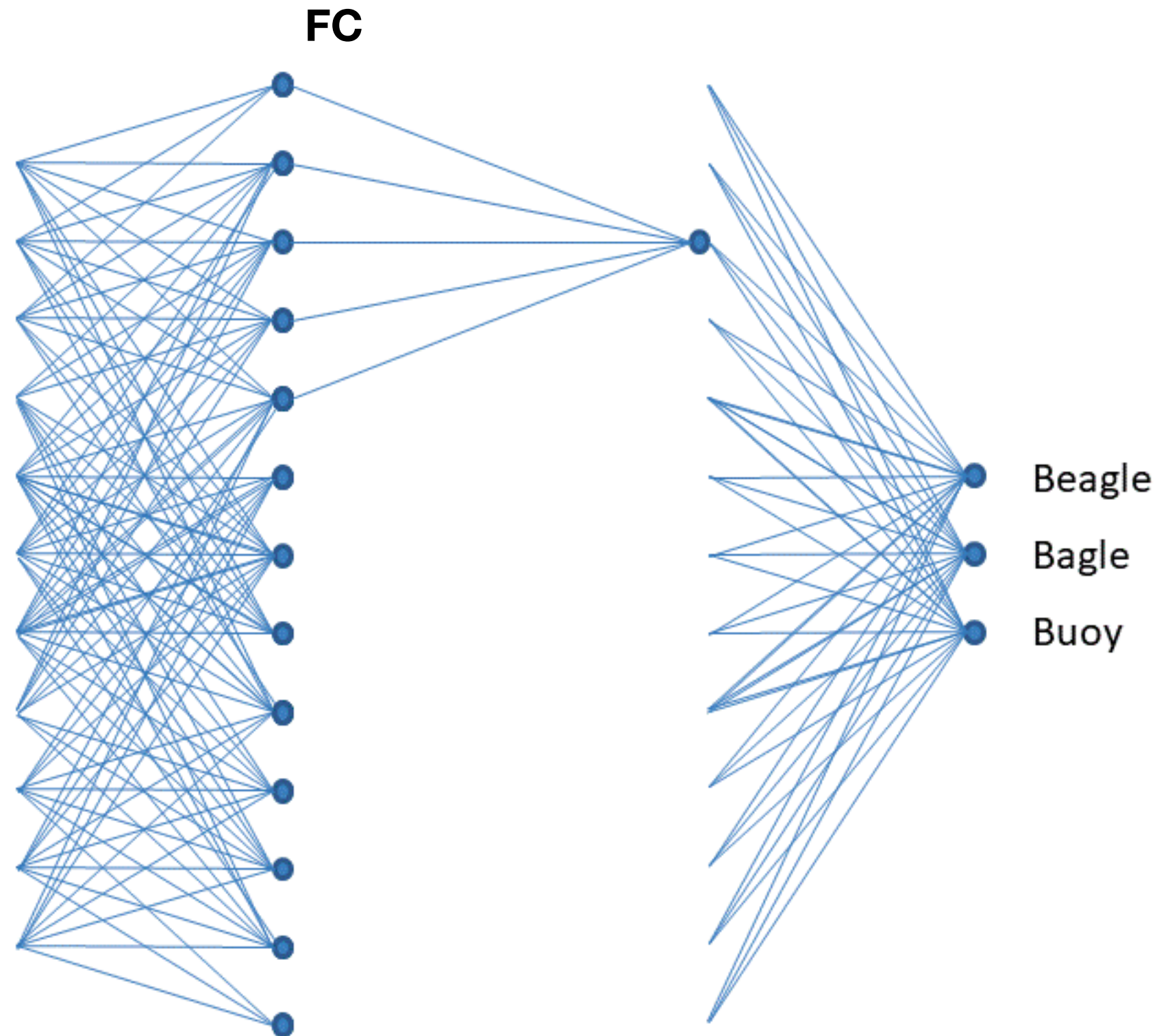Bengio, Hochreiter, Schmidhuber

# Definitions and Stuff

- Deep learning

- Architectures

# Three main classes of DL architectures
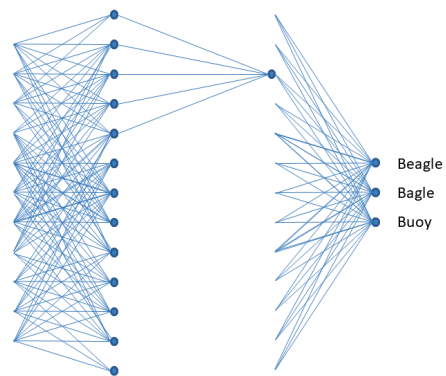
# Fully Connected / Feed Forward
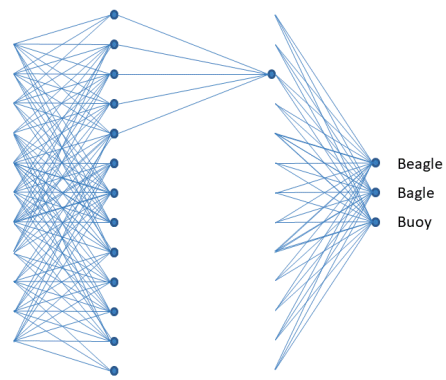
## FC

$$Z^i = W^i X + b^i 1$$
$$A^i = \textbf{RELU}\left(Z^i\right)$$

Beagle

Bagle

Buoy

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

Beagle
Bagle
Buoy

$$Z^i = W^i X + b^i 1$$

$$A^i = \mathbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**CNN**
**Convolutional Neural Networks**

Feature maps

Input

f.maps

f.maps

Output

Convolutions   Subsampling   Convolutions   Subsampling   Fully connected

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**Kernels**

height

width

depth

**CNN**

**Convolutional Neural Networks**

Feature maps

Input

f.maps

f.maps

Output

Convolutions

Subsampling

Convolutions

Subsampling

Fully connected

$$Z^i = W^i X + b^i 1$$
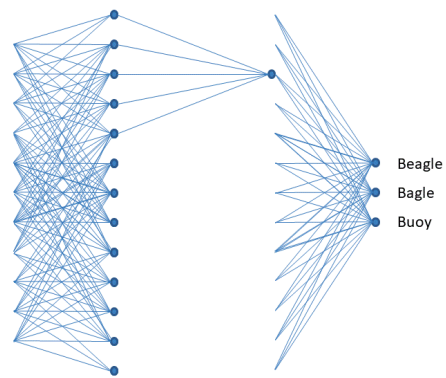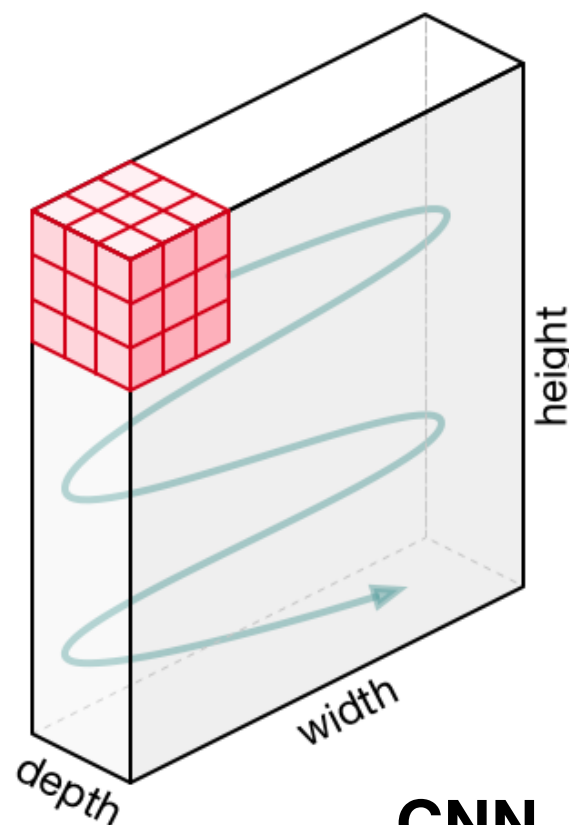$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**CNN**
**Convolutional Neural Networks**

Beagle
Bagle
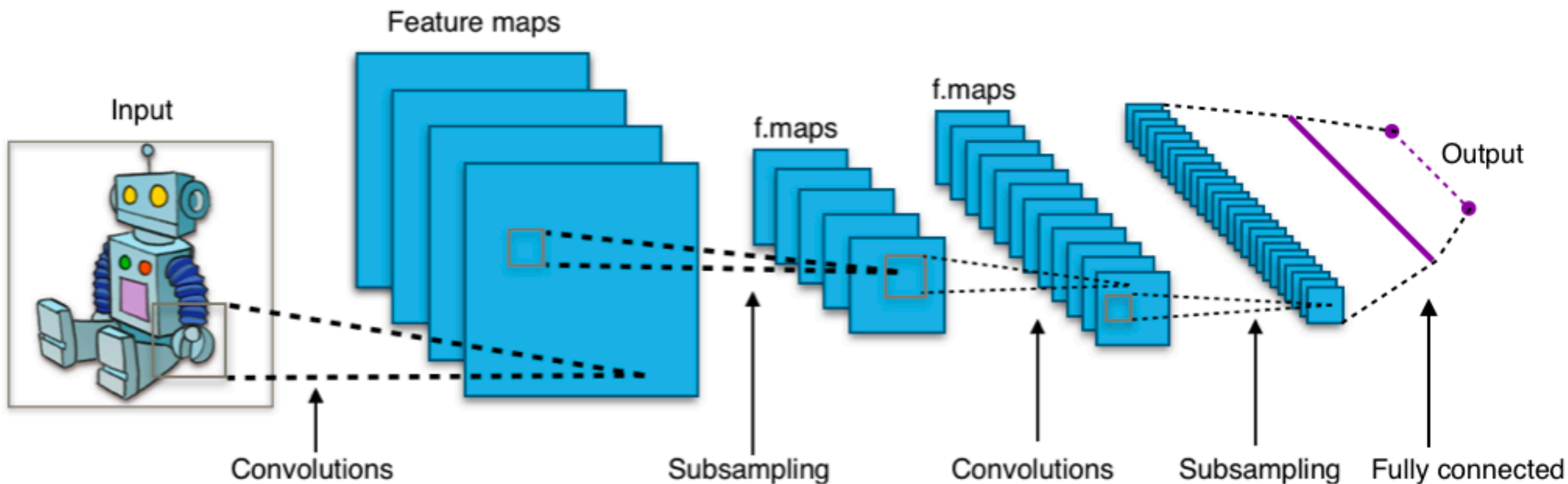Buoy

$$Z^i = W^i X + b^i 1$$
$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**
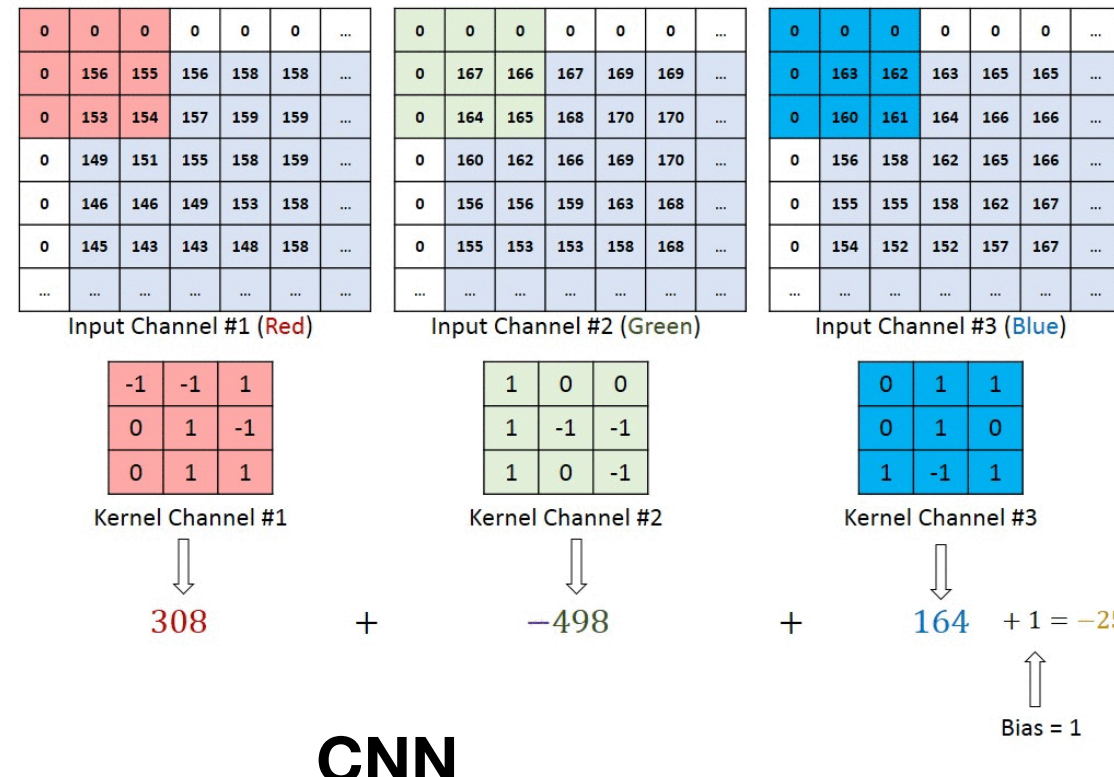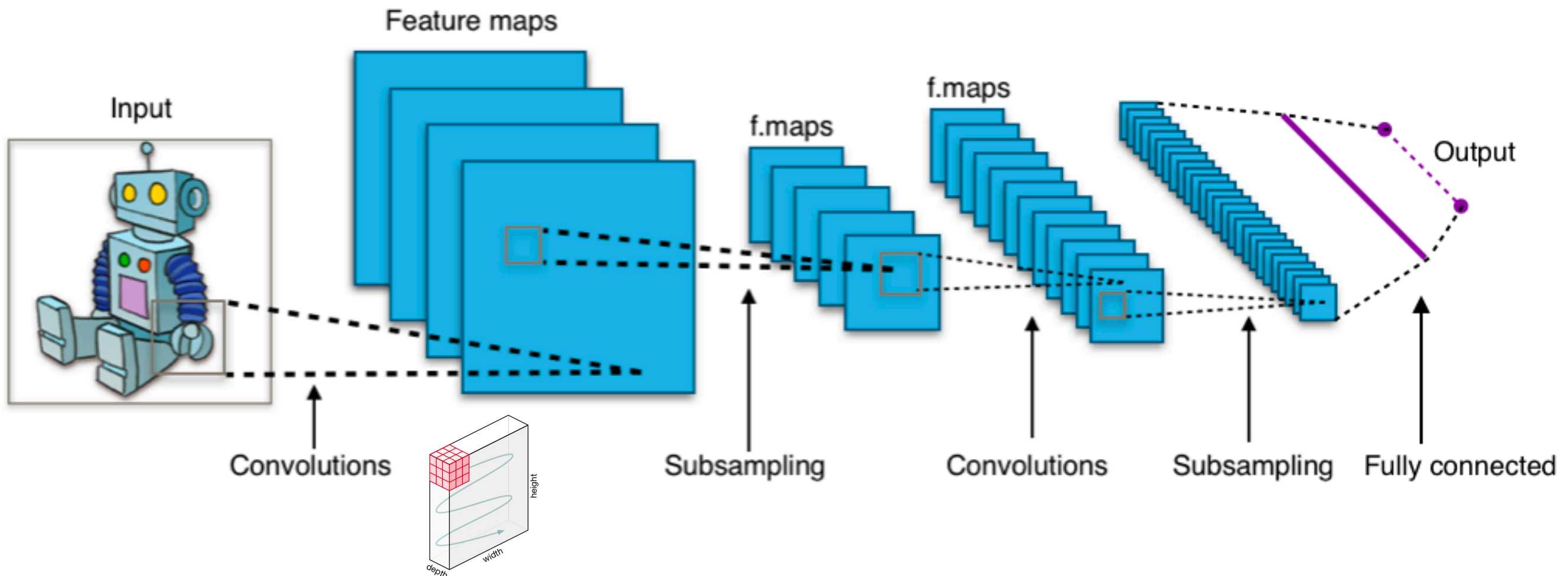
**FC**

**CNN**
**Convolutional Neural Networks**

Feature maps

Input

f.maps

f.maps

Output

Convolutions

height

depth  width

Subsampling

Convolutions

Subsampling

Fully connected

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**





Layer 1

Layer 2

Layer 3

Beagle
Bagle
Buoy

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**
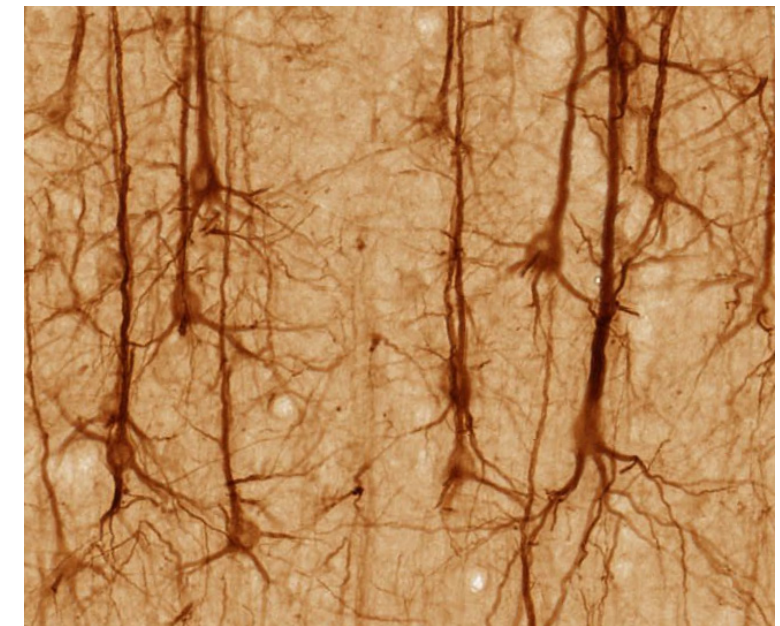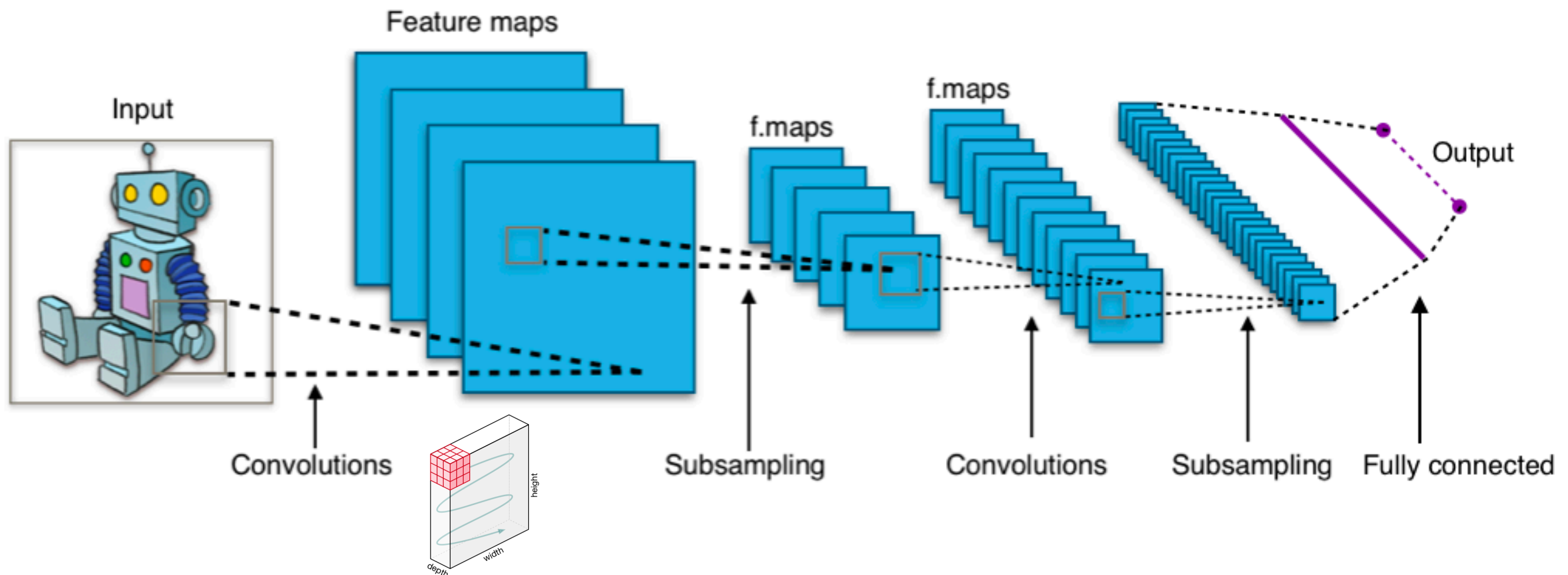
**FC**

**CNN**
**Convolutional Neural Networks**



Feature maps

Input

f.maps

f.maps

Output

depth
height
width

Convolutions    Subsampling    Convolutions    Subsampling    Fully connected

# RNN
# Recurrent Neural Network

$Z^i = W^i X + b^i 1$

$A^i = \mathbf{RELU}(z^i)$

**Fully Connected Feed Forward**

**FC**

**Unfold**

o

W

V

h

U

x

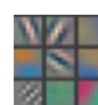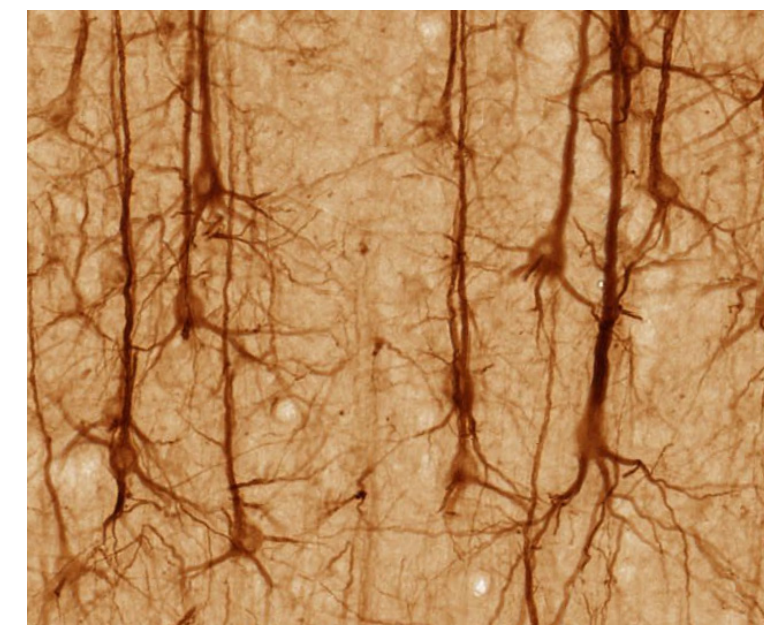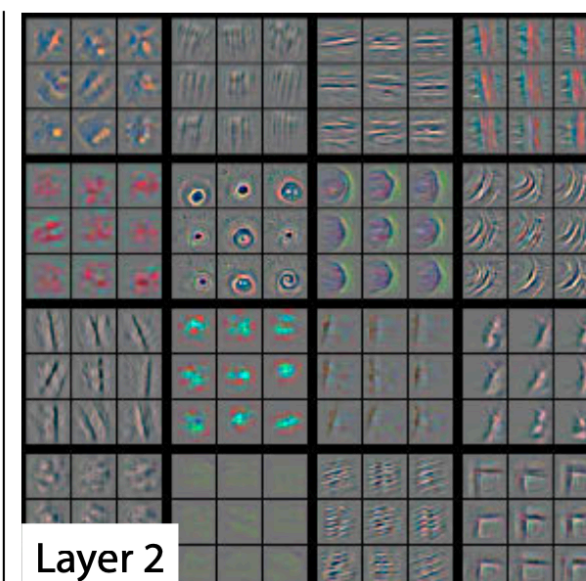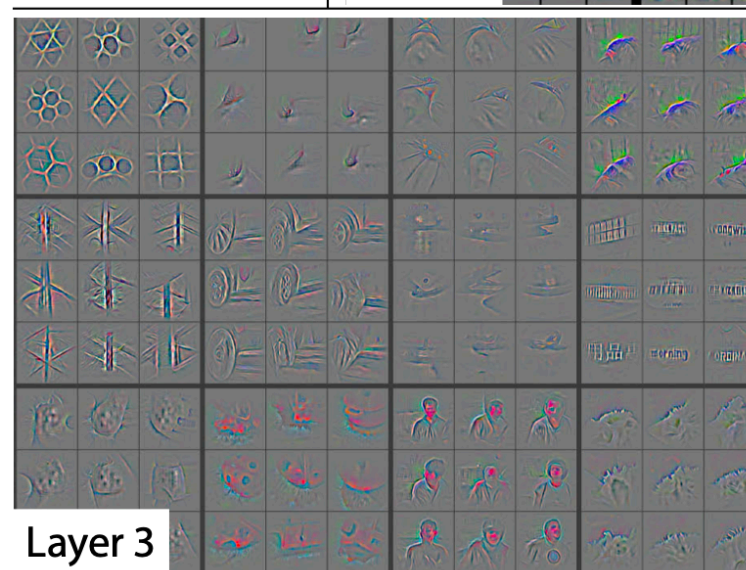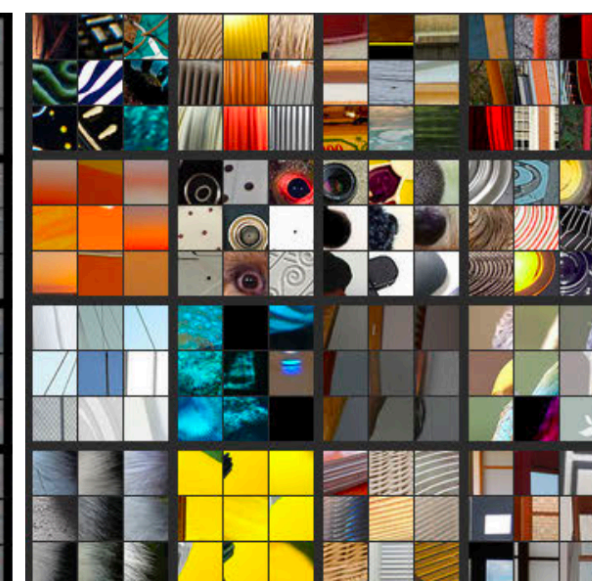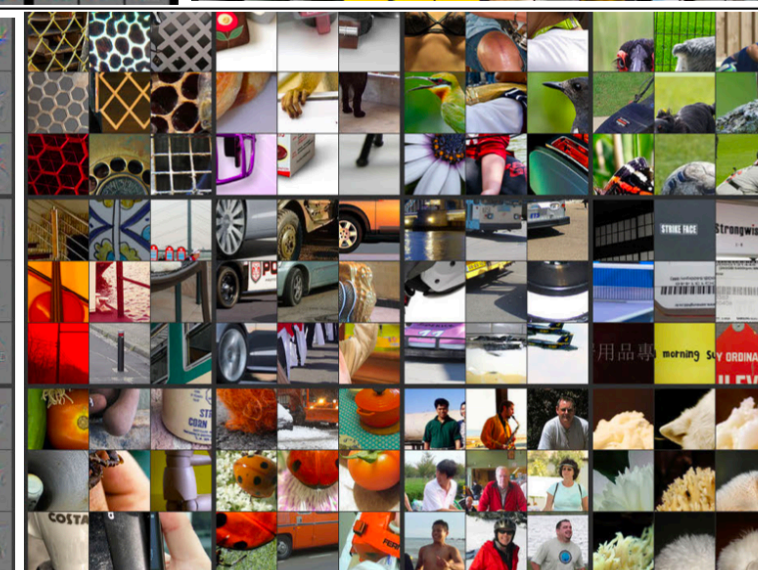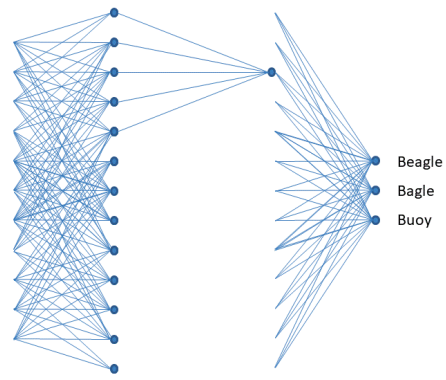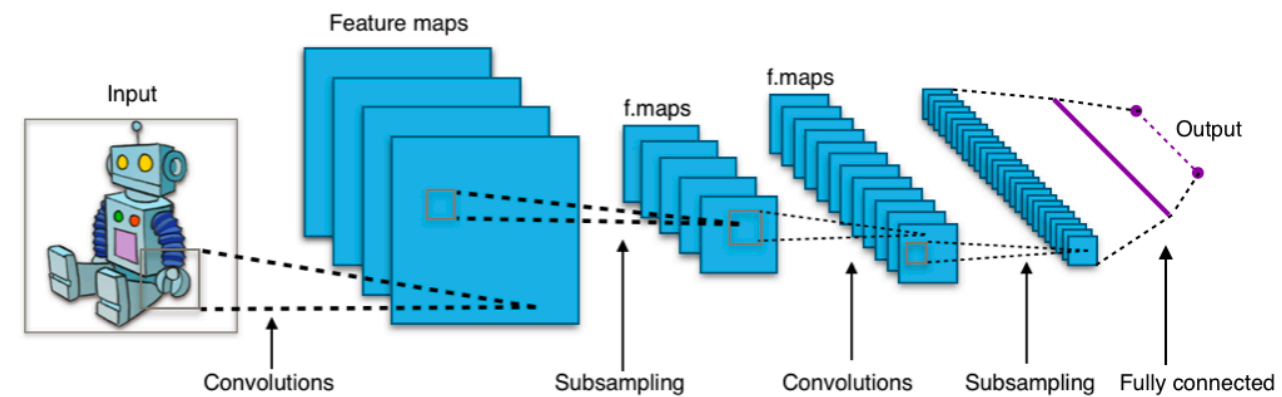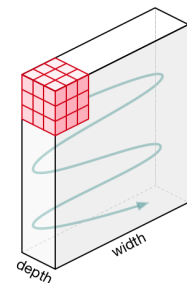$o_{t-1}$  $o_t$  $o_{t+1}$

W   W   W

$\cdots \rightarrow \overrightarrow{V} \rightarrow h_{t-1} \rightarrow \overrightarrow{V} \rightarrow h_t \rightarrow \overrightarrow{V} \rightarrow h_{t+1} \rightarrow \overrightarrow{V} \rightarrow \cdots$

U   U   U

$x_{t-1}$  $x_t$  $x_{t+1}$

**CNN**
**Convolutional Neural Networks**

Beagle
Bagle
Buoy

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**RNN**

**Recurrent Neural Network**

Unfold

**CNN**

**Convolutional Neural Networks**

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**RNN**

**Recurrent Neural Network**

Unfold

**GANs,**
**Auto Encoders,**
**ODE Networks,**
**Invertible Flow Networks,**

**.....**

**CNN**

**Convolutional Neural Networks**
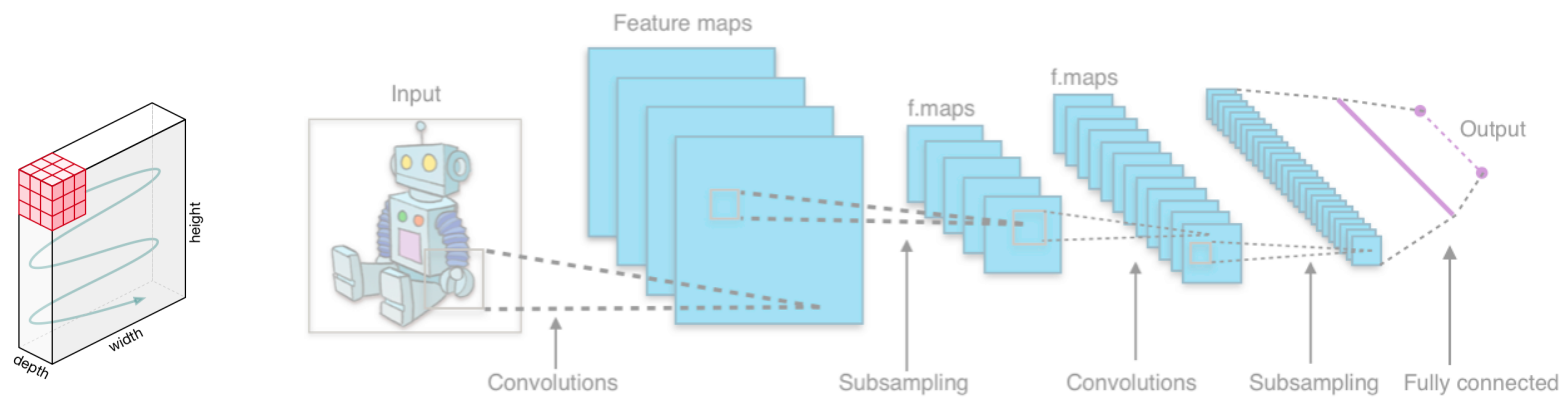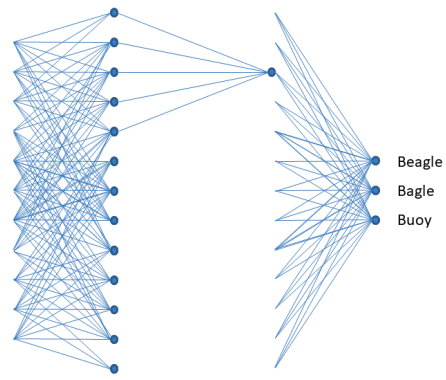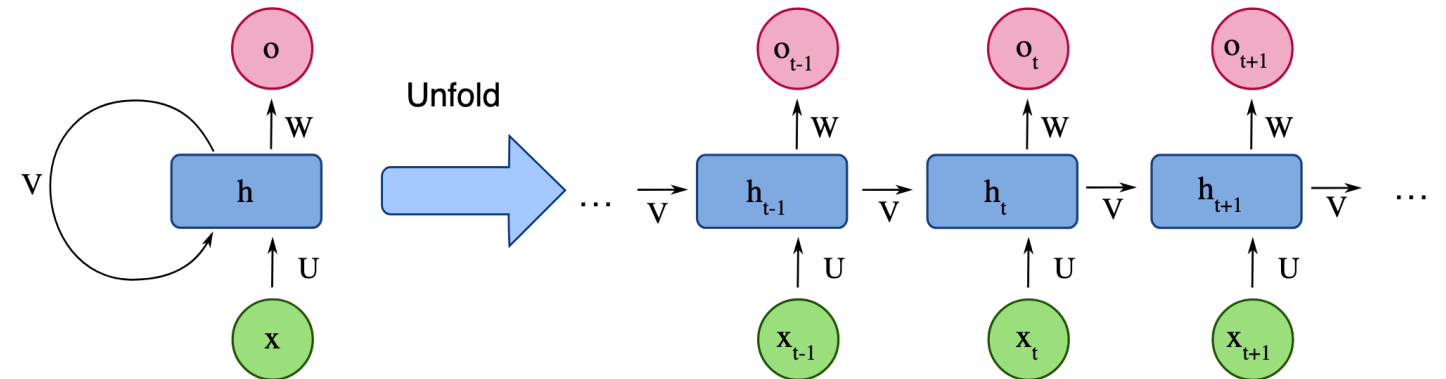
$$Z^i = W^i X + b^i 1$$

$$A^i = \mathbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**GANs,**
**Auto Encoders,**
**ODE Networks,**
**Invertible Flow Networks,**

**.....**

**RNN**

**Recurrent Neural Network**

**CNN**
**Convolutional Neural Networks**
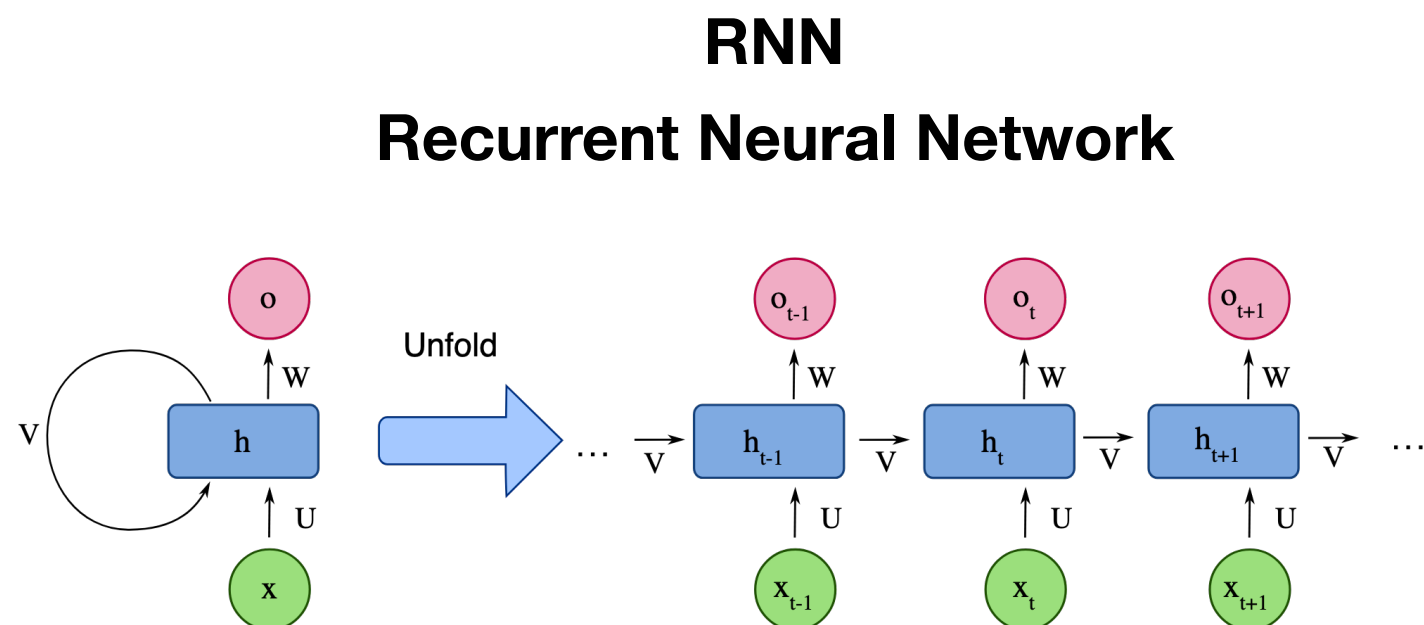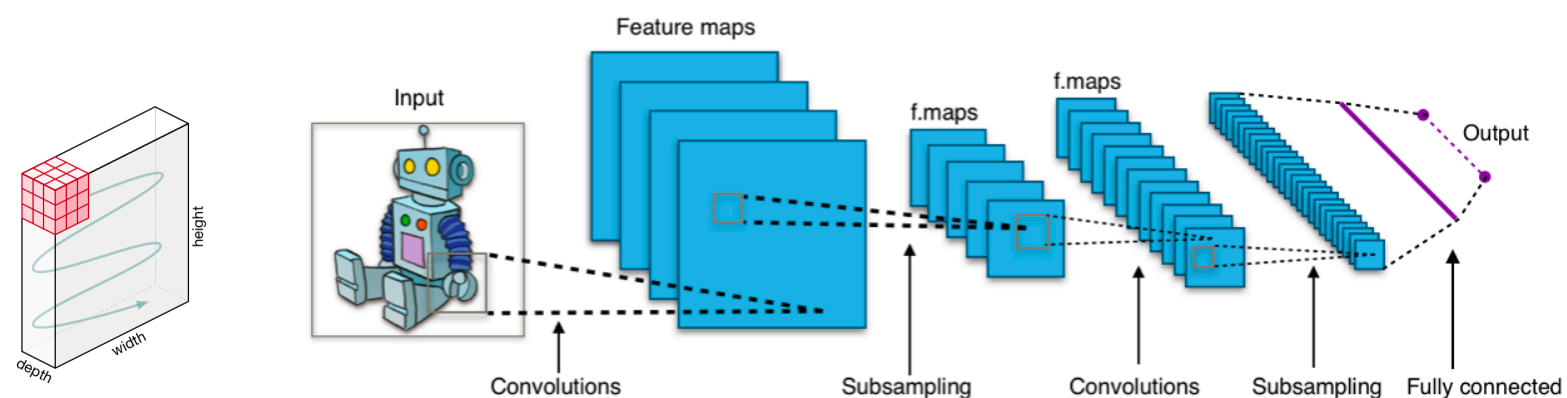
Beagle
Bagle
Buoy

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**GANs,**
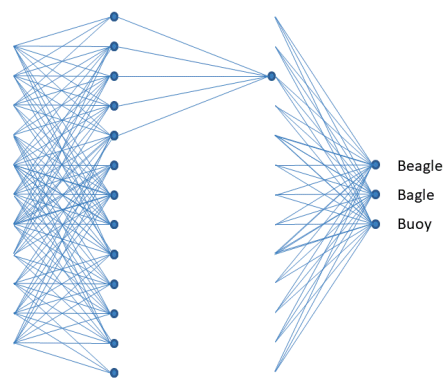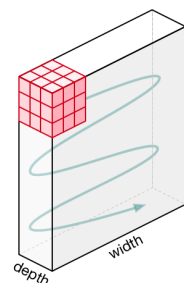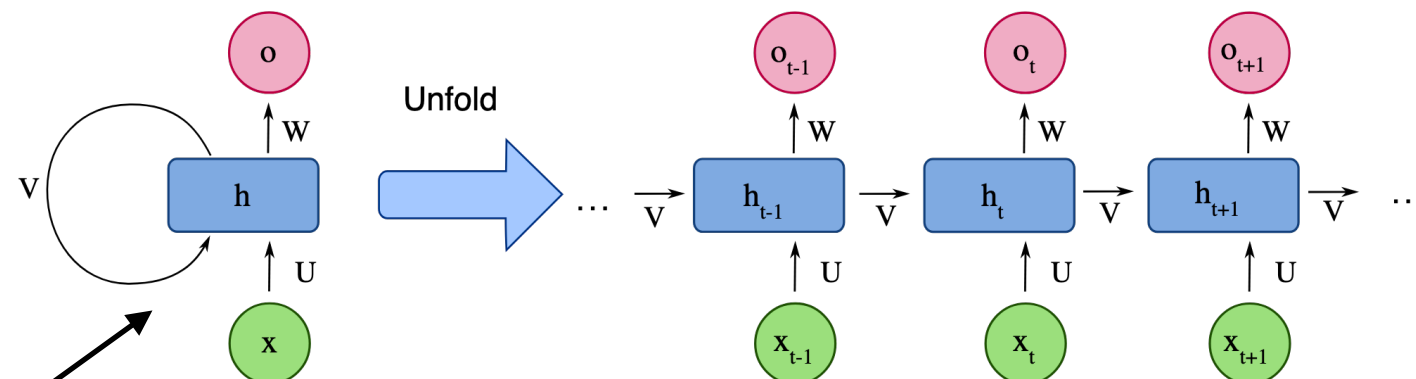**Auto Encoders,**
**ODE Networks,**
**Invertible Flow Networks,**
**.....**

**RNN**

**Recurrent Neural Network**

o

W

V    h

U

x

Unfold

$o_{t-1}$    $o_t$    $o_{t+1}$

W    W    W

$h_{t-1}$    $h_t$    $h_{t+1}$

V    V    V

U    U    U

$x_{t-1}$    $x_t$    $x_{t+1}$

**CNN**

**Convolutional Neural Networks**

Feature maps

Input

f.maps

f.maps

Output

depth

width

height

Convolutions    Subsampling    Convolutions    Subsampling    Fully connected

**All kind of domains: medical imaging, autonomous driving, emotion recognition,**
**recommenders, natural language processing**

Beagle
Bagle
Buoy

$$Z^i = W^i X + b^i 1$$

$$A^i = \textbf{RELU}\left(Z^i\right)$$

**Fully Connected / Feed Forward**

**FC**

**GANs,
Auto Encoders,
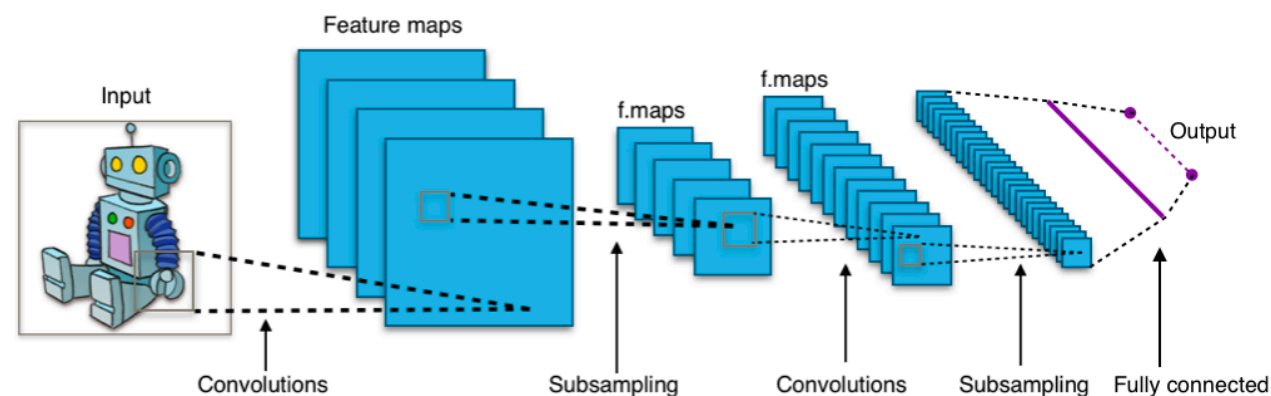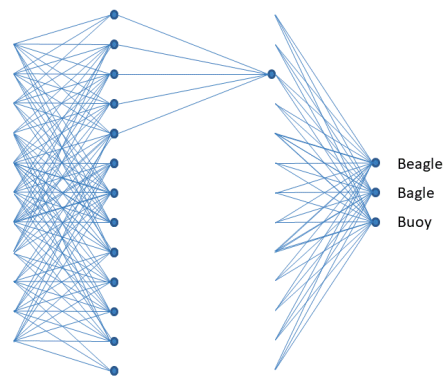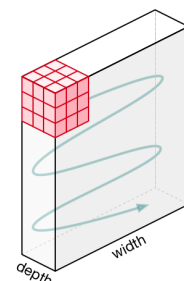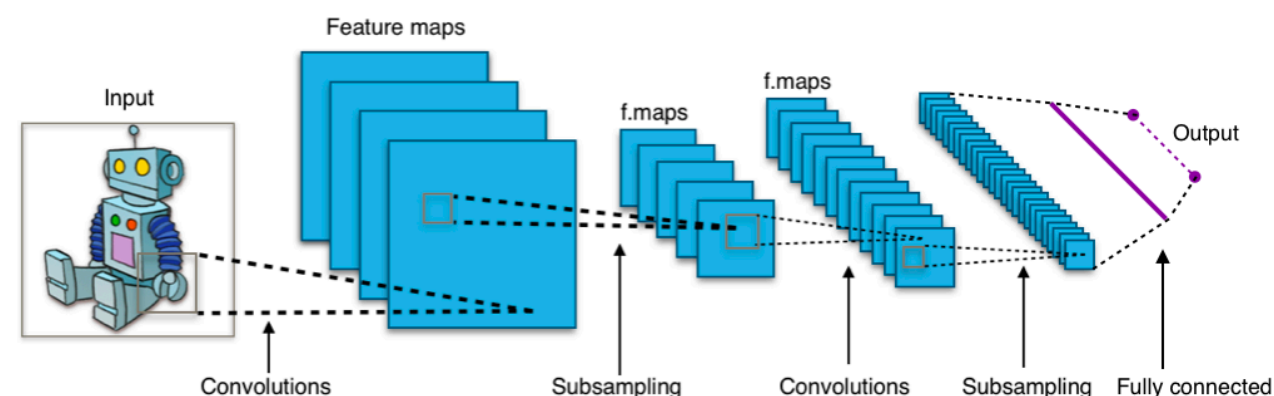ODE Networks,
Invertible Flow Networks,**

**.....**

**Supervised,
Unsupervised,
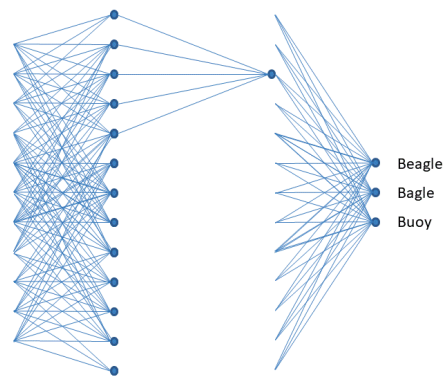Self-Supervised,
Reinforcement Learning**

**RNN**

**Recurrent Neural Network**

o

W

V

h
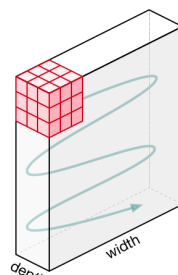
U

x

Unfold

$\cdots$

V

$o_{t-1}$

W

$h_{t-1}$

V

U

$x_{t-1}$

$o_t$

W

$h_t$

V

U

$x_t$

$o_{t+1}$

W

$h_{t+1}$

V

U

$x_{t+1}$

$\cdots$

**CNN**

**Convolutional Neural Networks**

Feature maps

Input

f.maps

f.maps

height

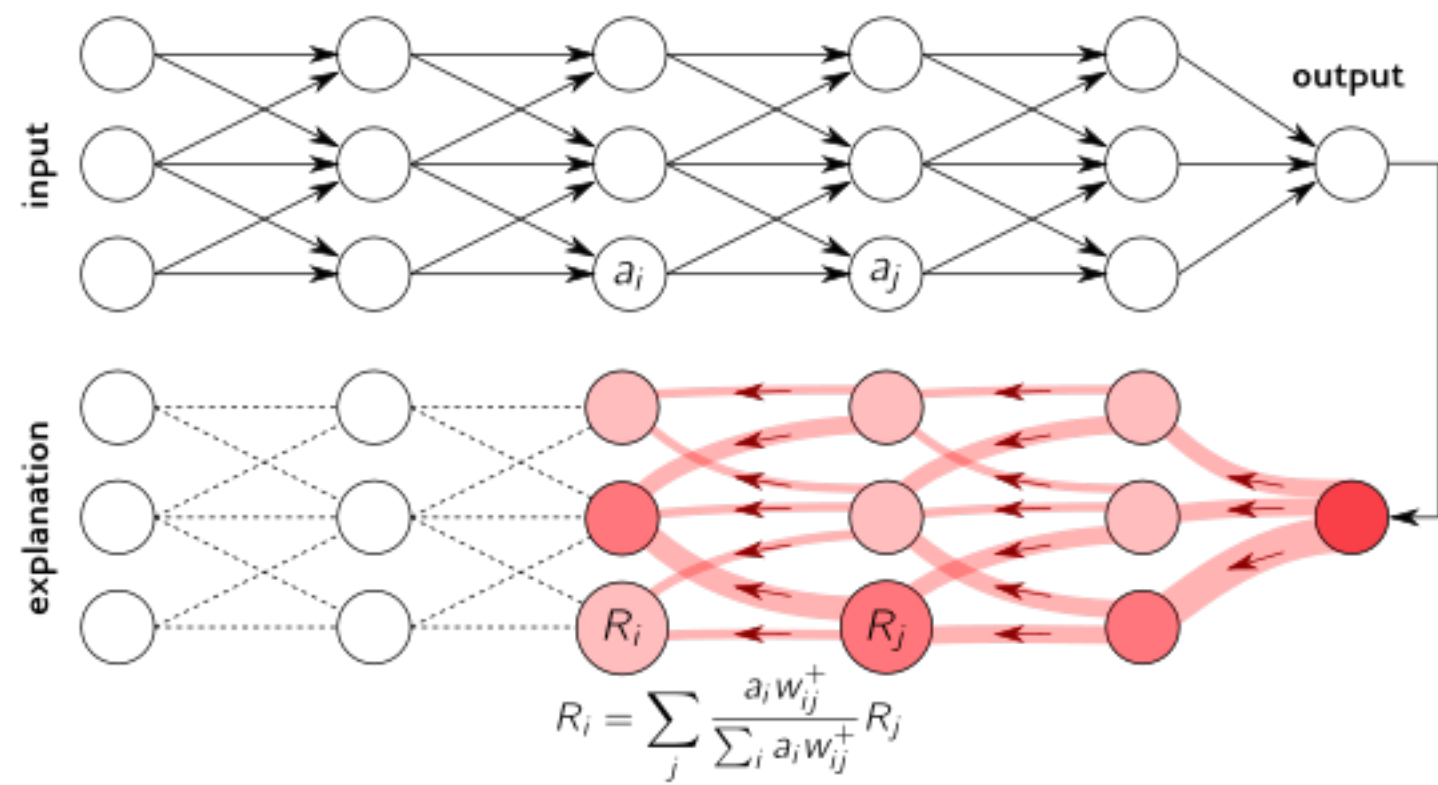depth

width

Convolutions

Subsampling

Convolutions

Subsampling
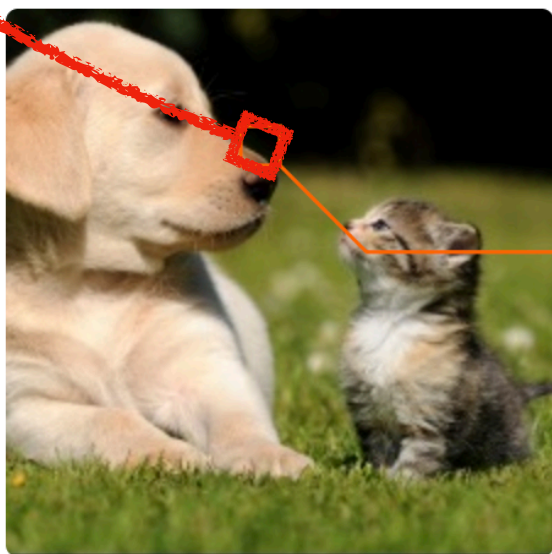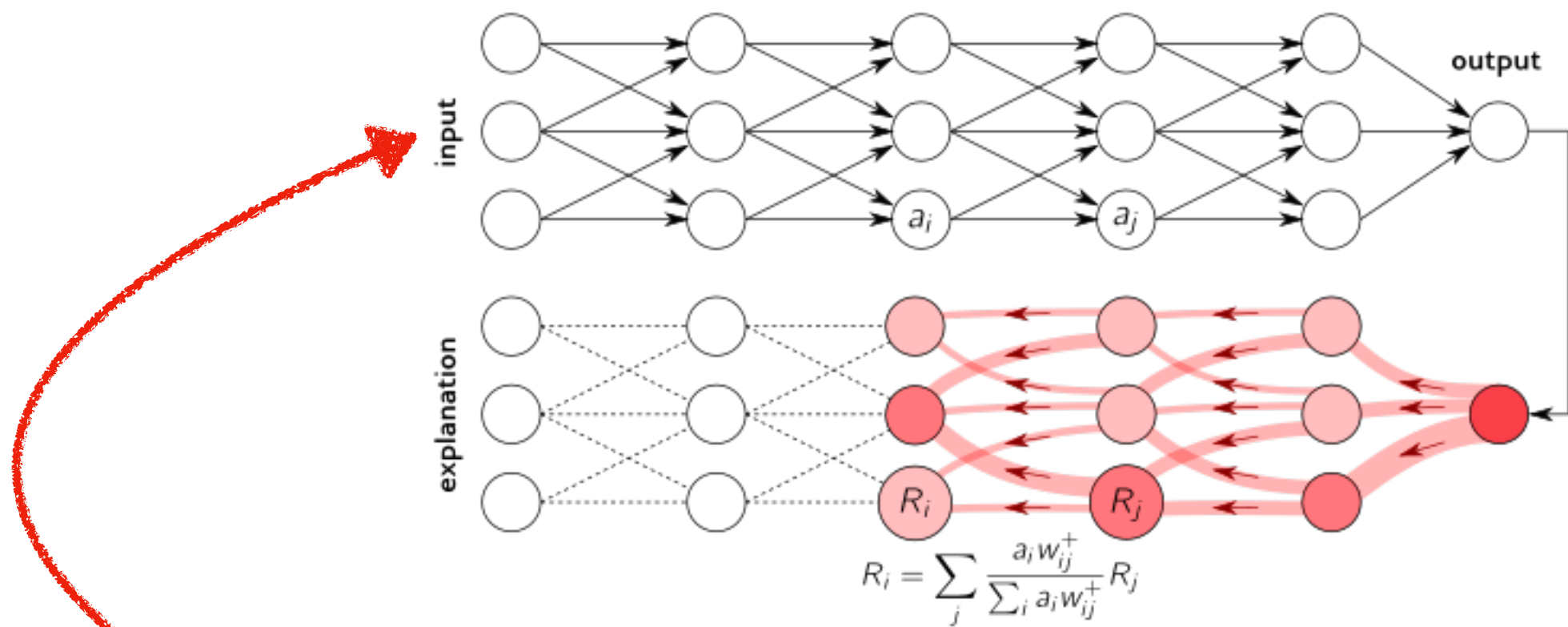
Output

Fully connected

**All kind of domains: medical imaging, autonomous driving, emotion recognition,
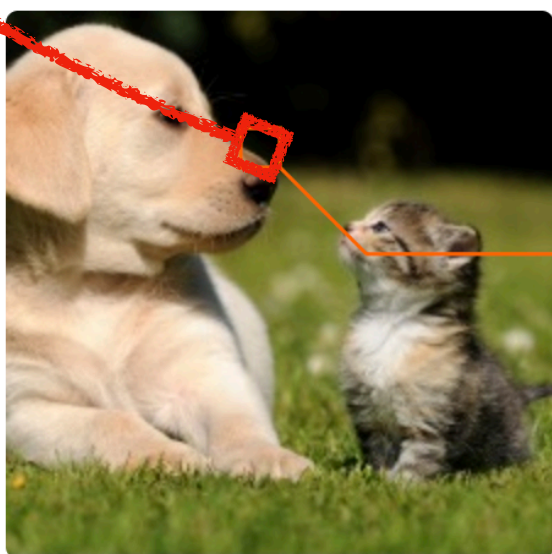recommenders, natural language processing**

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

output

input

explanation

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$\{$ : 714. , : 293. , : 288. , : 250. , ... $\}$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

{ 714. : , 293. : , 288. : , 250. : ... }

output

input

explanation

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$\{$  : 714. ,  : 293. ,  : 288. ,  : 250. , ... $\}$

output

input

explanation

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$a_i$ $a_j$

$R_i$ $R_j$

{ : 714. , : 293. , : 288. , : 250. , ... }

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$*f_i$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$* f_i$

**RELU**( )

output

input

explanation

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$*f_i$

**RELU** ( )

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$* f_i$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

# What questions can we currently answer?

- Given **one** manually selected input:

  - On **which parts** of the input the does the model **focus**? (f.e. *LRP*)

- Given **one** selected output:

  - What different **strategies (clusters) exist** for focussing on images? (f.e. *SpRAy*)

  - What **kind of template** does it look for? (f.e. *Max Activation*)

- Given a **representative set** of inputs for a **latent factor**:

  - Are there any **geometric properties** of the features? (f.e. *de-biasing*)

# Hands On

[https://github.com/grazai/xai-tutorial-march-2020](https://github.com/grazai/xai-tutorial-march-2020)

# Side Step: Data

- We use **MNIST** here

  - Super **simple**, super **fast to train**, good for a demo

- *Better*: For images, datasets for segmentation like **COCO** provide perfect ground truth for the attribution.

- *simply-clevr-dataset* https://github.com/ahmedmagdiosman/simply-clevr-dataset

- Don't know a similar dataset for TimeSeries (if anyone knows, please tell me!)

VGG-16

conv1

conv2

conv3

conv4

conv5

fc6    fc7    fc8

$1 \times 1 \times 4096$    $1 \times 1 \times 1000$

$14 \times 14 \times 512$

$7 \times 7 \times 512$

$28 \times 28 \times 512$

$56 \times 56 \times 256$

$112 \times 112 \times 128$

$224 \times 224 \times 64$

convolution+ReLU

max pooling

fully connected+ReLU

**We use something VGG like**
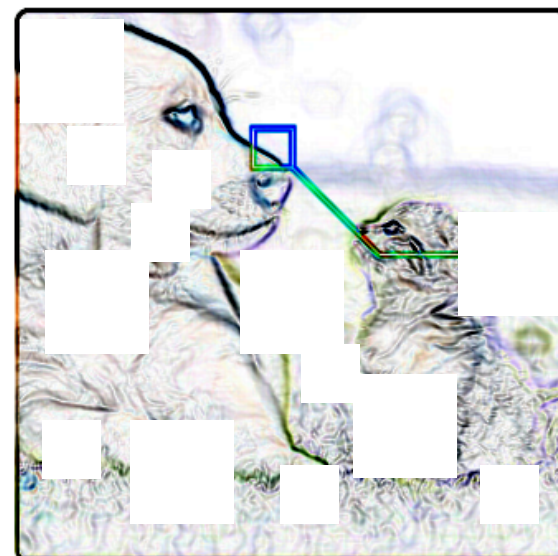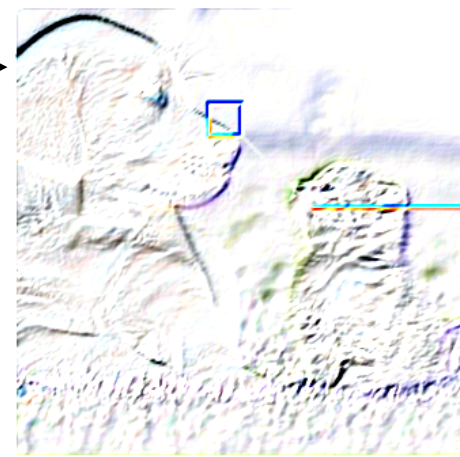
# What questions can we currently answer?

- Given **one** manually selected input:

  - On **which parts** of the input does the model **focus**?

    - Attention mechanisms, LRP, GradCAM, IntegratedGradients, ….

    - https://human-centered.ai/wordpress/wp-content/uploads/2020/03/706.046-AK-explainable-AI-Introduction-MiniProjects-Class-of-2020.pdf for more (Prof. Holzinger)

# What questions can we currently answer?

- Given **one** selected output:

  - Are there **clusters** on the parts the model focuses?

    - SpRAy, Sampling, …

    - https://human-centered.ai/wordpress/wp-content/uploads/2020/03/706.046-AK-explainable-AI-Introduction-MiniProjects-Class-of-2020.pdf for more (Prof. Holzinger)

# What questions can we currently answer?

- Given **one** selected output:

  - What **kind of template** does it look for?

    - **Max Activation**, Project Lucid, Activation Atlas

  - distill.pup

# What questions can we currently answer?

- Given a **representative set** of inputs for a **latent factor**:

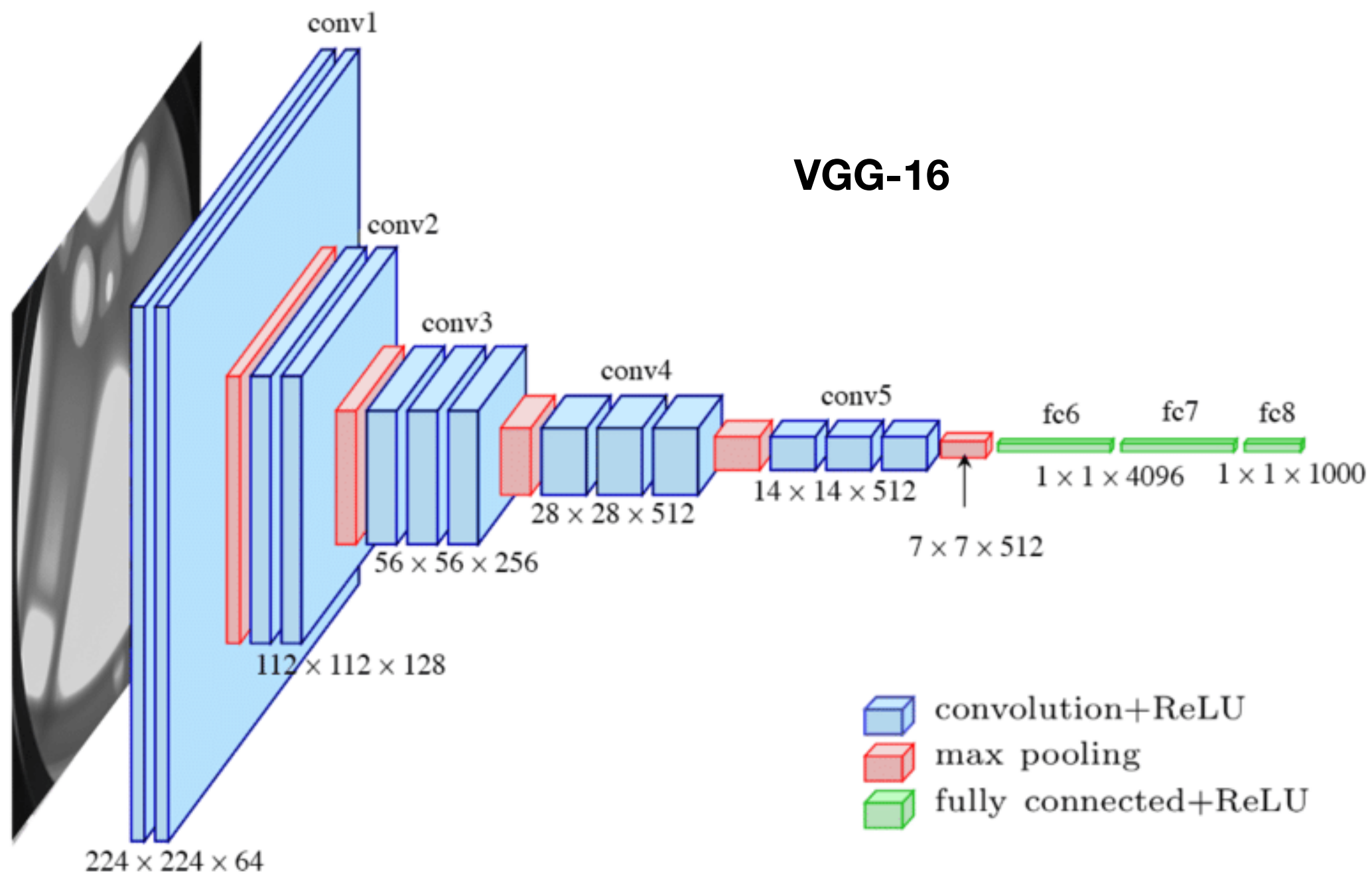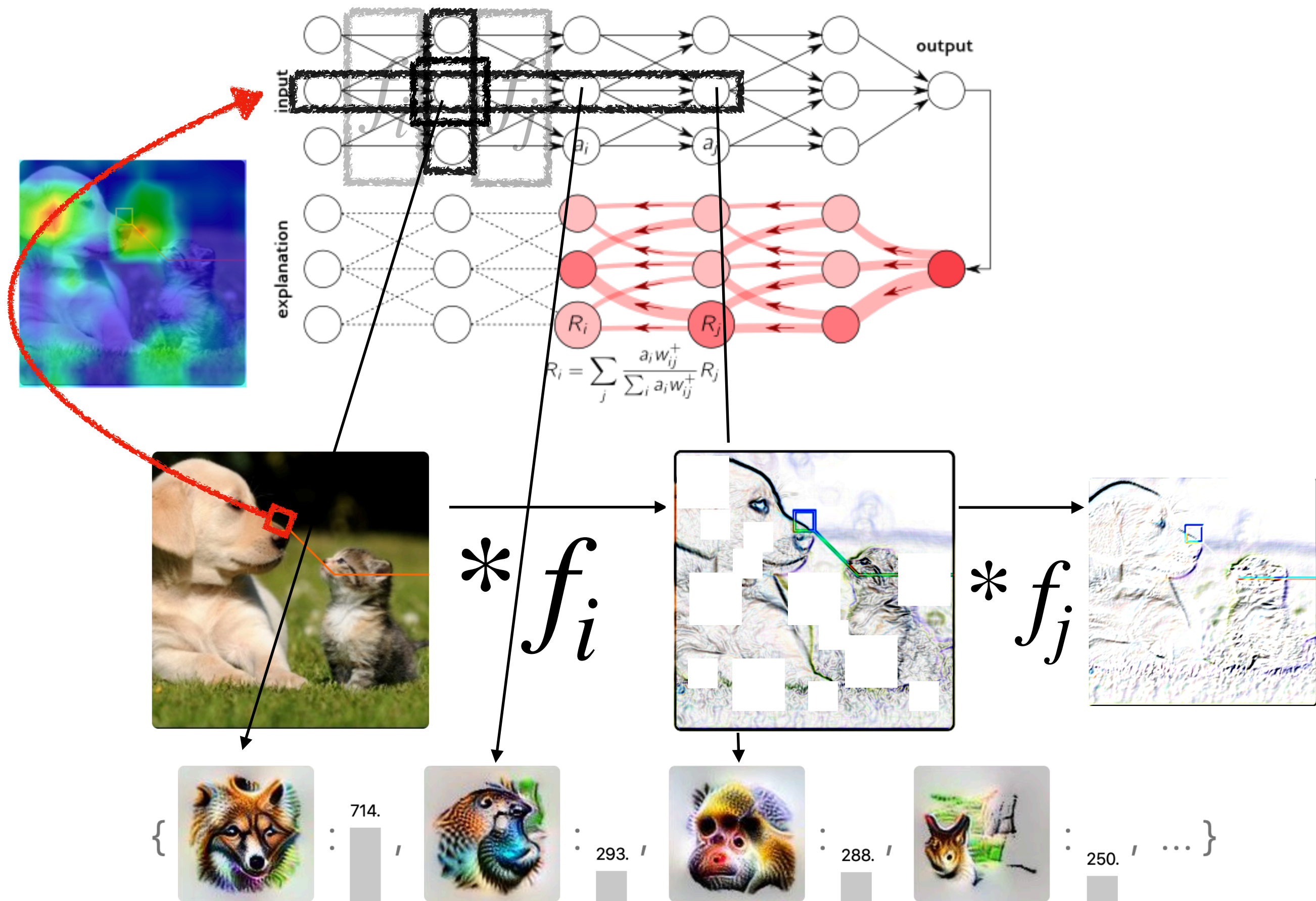  - Are there any **geometric properties** of the features?

    - Embeddings and De-Biasing

# I did lie to you!

- **Adversarial** images

- **Sensitivity** instead of importance

- Not the **complete** picture

- Not completely **mature** in case of frameworks

- But already **ok** for the *knowledgeable* and a **great promise**

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$* f_i$

$* f_j$

{ : 714. , : 293. , : 288. , : 250. , ... }

# Thanks for listening

I hope there was something of value for you?

We can have some Q&A in the DeepLearning Discord chat

https://discord.gg/nvdxH7