

# Explainable AI für Deep Learning: Overview und Tutorial

Jörg Simon



**Huh? What?**





**Right....**

**But maybe we can explain  
the complex models a bit**

# Explainable AI für Deep Learning: Overview und Tutorial

Jörg Simon

# About me

- PhD on using DeepLearning 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.

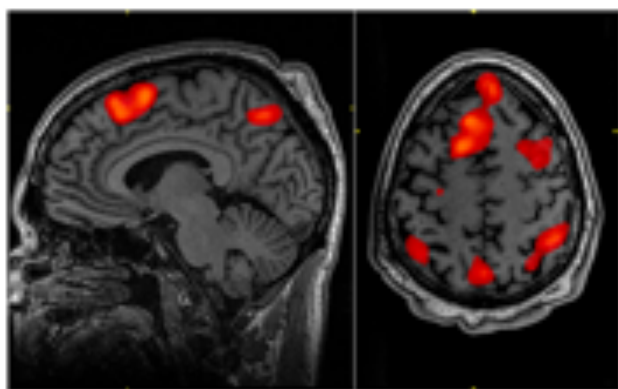


# Agenda

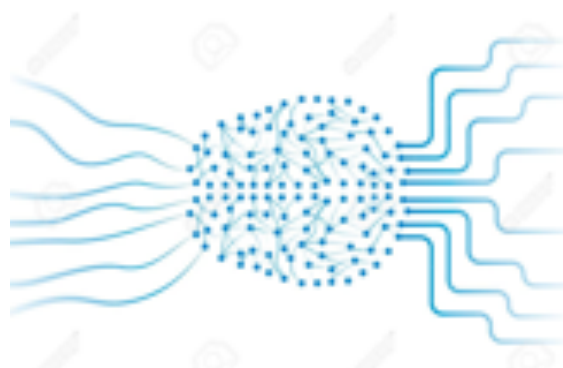
- Definitions and Stuff
- Hands On
- Discussion

# Definitions and Stuff

- Deep Learning



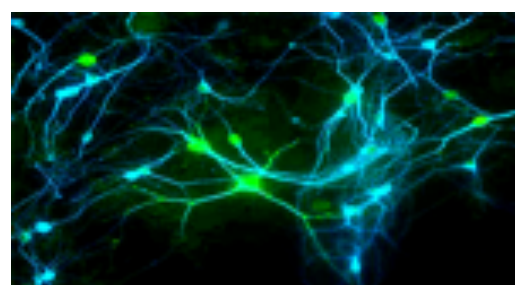
Distributed Representation



Super Simplified Model of Human Brain



Hinton



Spiking Frequency = weight



Yann LeCun

Simple Matrix Multiply + Non Linearity



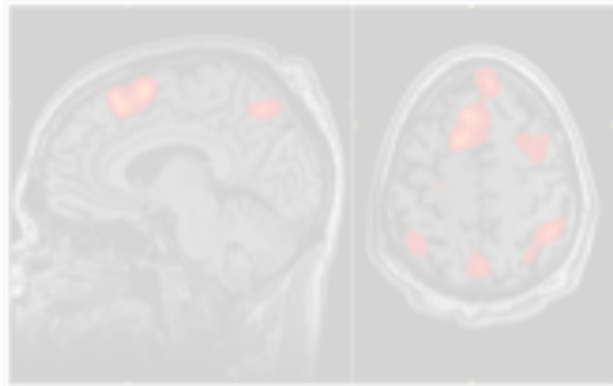
RNN



Bengio, Hochreiter, Schmidhuber

Deep Learning?





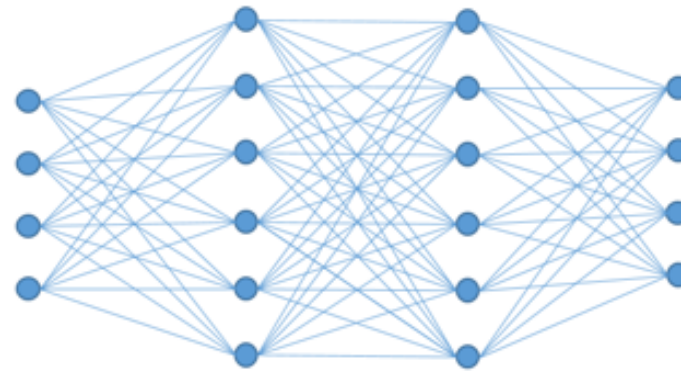
Distributed Representation



Spiking Frequency = weight



Super Simplified Model of Human Brain



Deep Learning?

Simple Matrix Multiply + Non Linearity



RNN

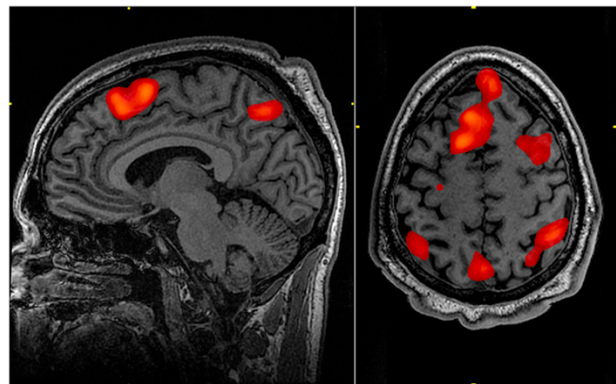


Bengio, Hochreiter, Schmidhuber

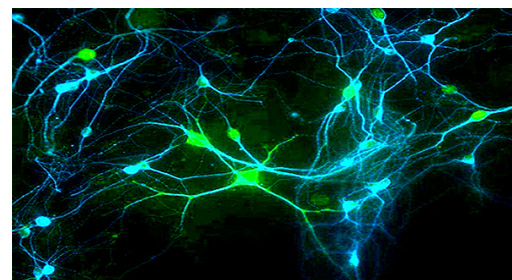


Hinton

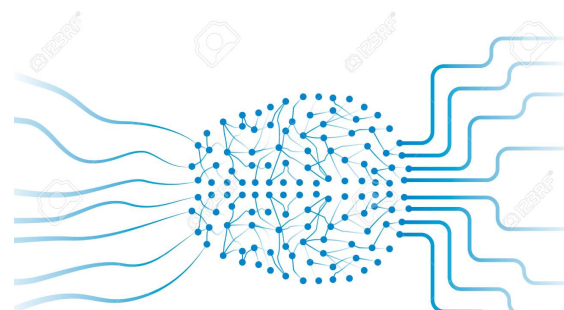
Yann LeCun



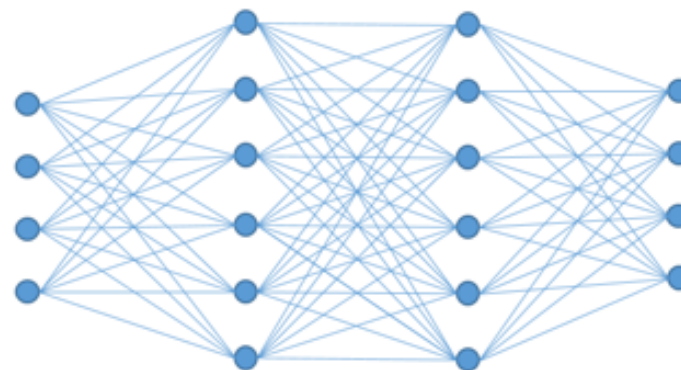
Distributed Representation



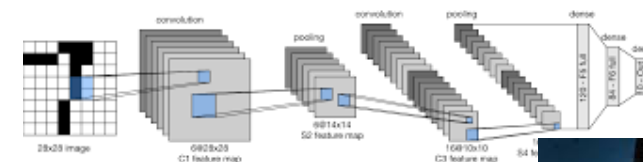
Spiking Frequency = weight



Super Simplified Model of Human Brain



Deep Learning?



CNN



Hinton



Yann LeCun

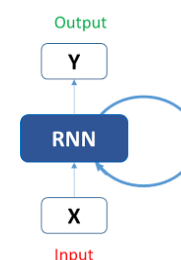


Bengio, Hochreiter, Schmidhuber

Simple Matrix Multiply + Non Linearity



RNN



# Definitions and Stuff

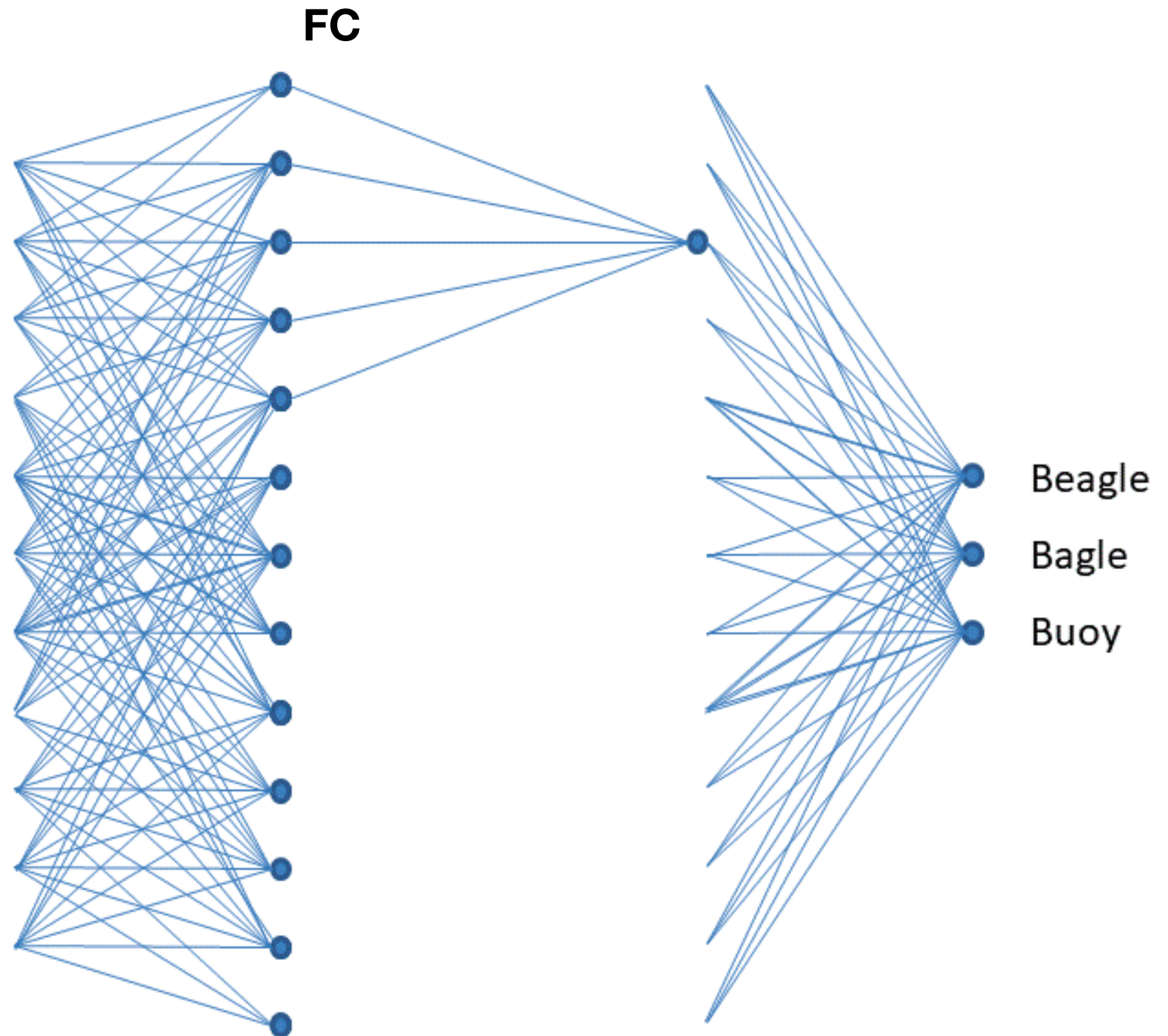
- Deep Learning
- Architectures

# Three main Classes of DL Architectures

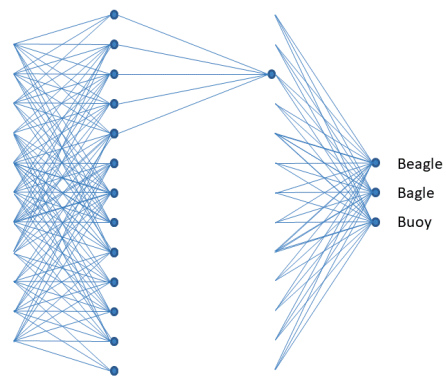
## Fully Connected / Feed Forward

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

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







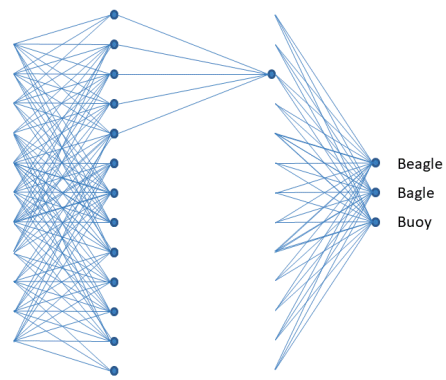
$$Z^i = W^i X + b^i 1$$

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

**Fully Connected / Feed Forward**

**FC**





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

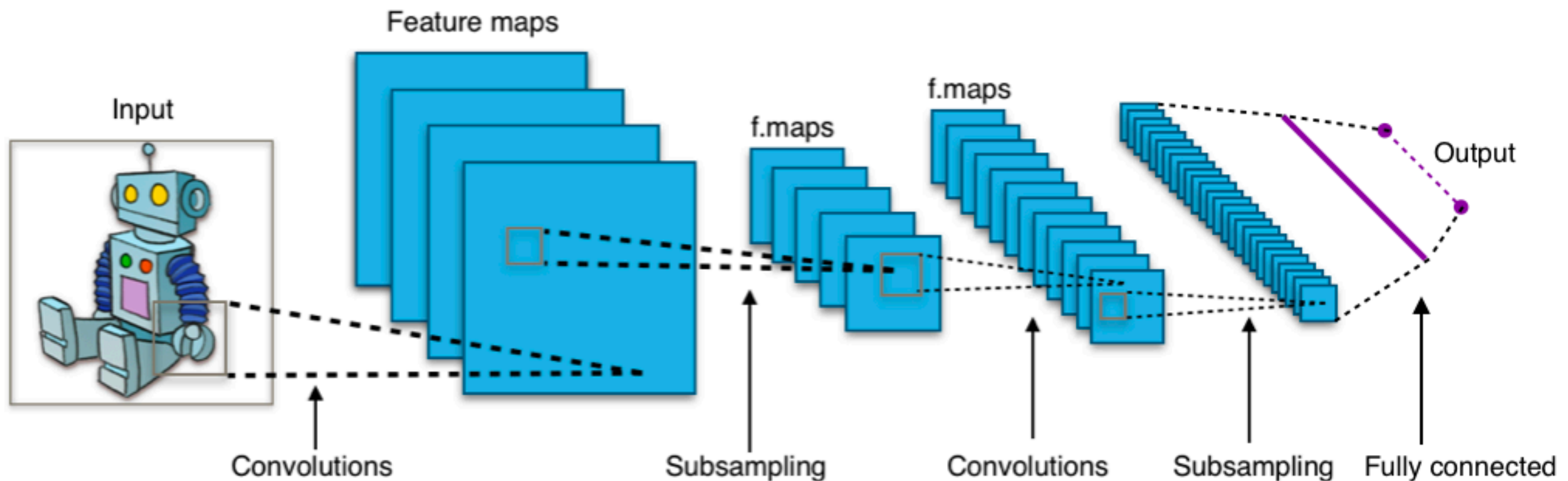
$$A^i = \mathbf{RELU}(Z^i)$$

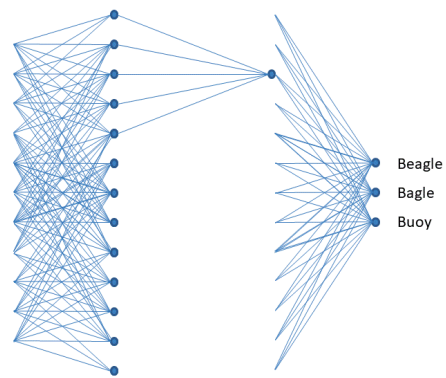
**Fully Connected / Feed Forward**

**FC**

**CNN**

**Convolutional Neural Networks**

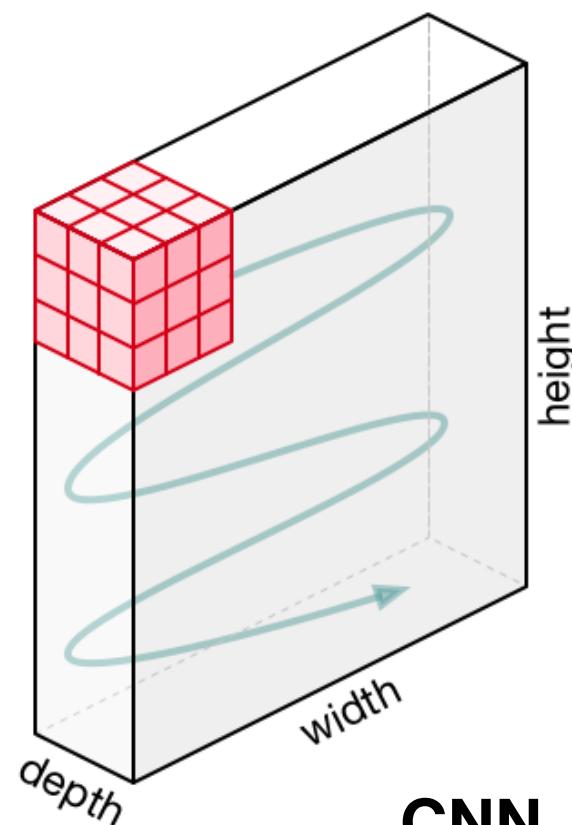




$$Z^i = W^i X + b^i 1$$

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

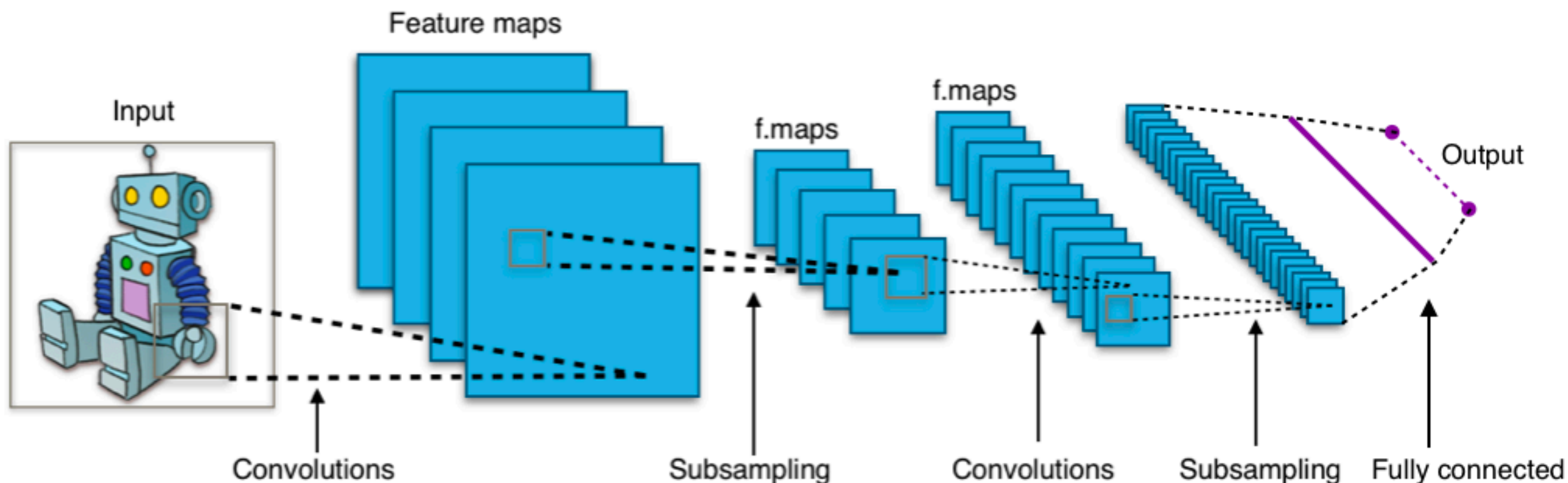
**Fully Connected / Feed Forward**  
**FC**

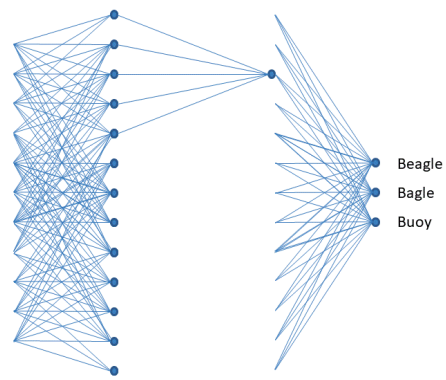


**Kernels**

**CNN**

**Convolutional Neural Networks**





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

$$A^i = \text{RELU}(Z^i)$$

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

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...	...	...	...	...	...	...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...	...	...	...	...	...	...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...	...	...	...	...	...	...

Input Channel #3 (Blue)

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

Kernel Channel #1

↓

308

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

Kernel Channel #2

↓

-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

↓

164

+ 1 = -25

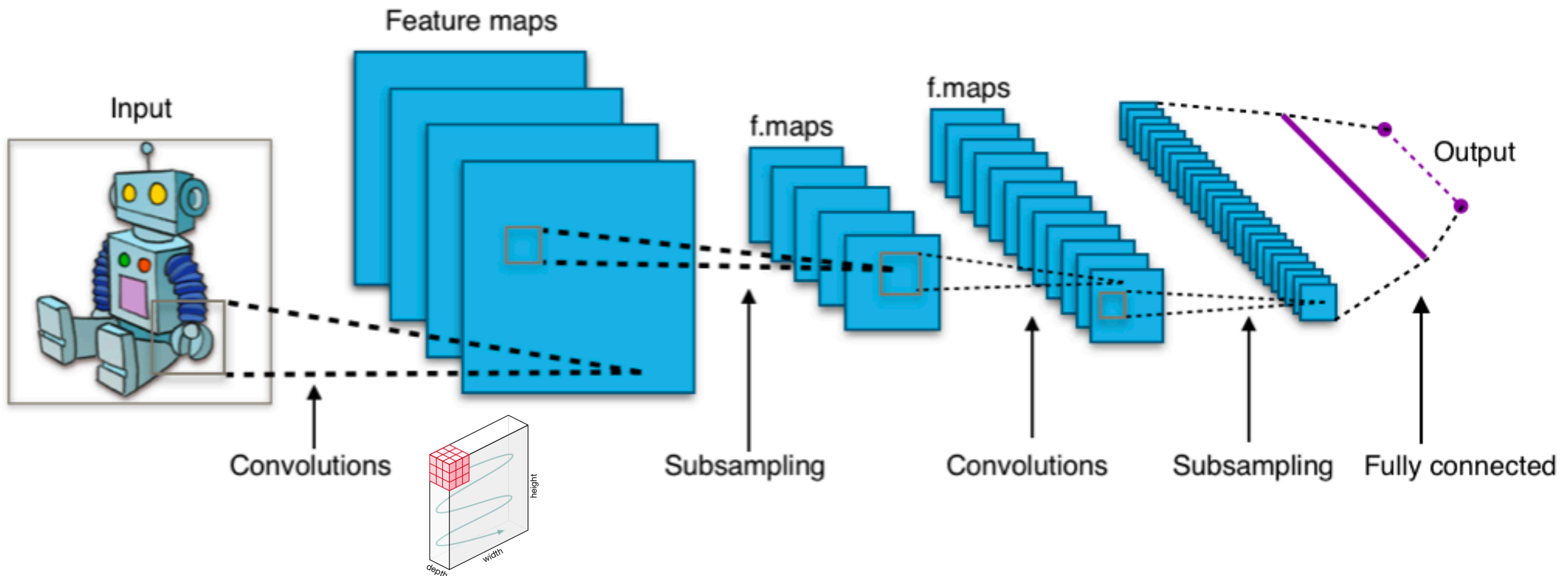
↑  
Bias = 1

-25				...
				...
				...
				...
...	...	...	...	...

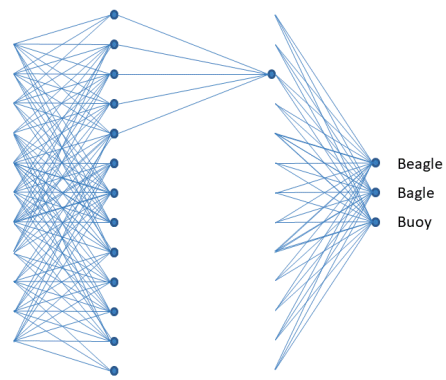
Output

**CNN**

**Convolutional Neural Networks**



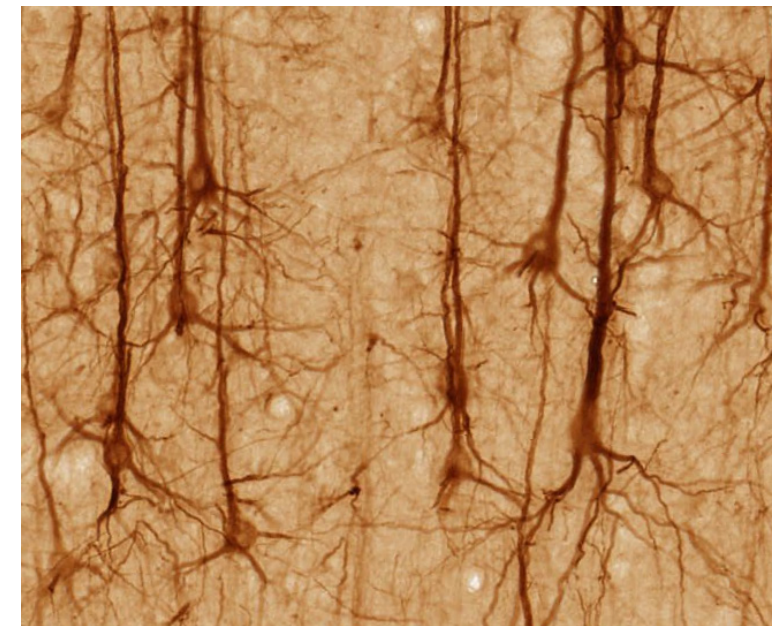




$$Z^i = W^i X + b^i 1$$

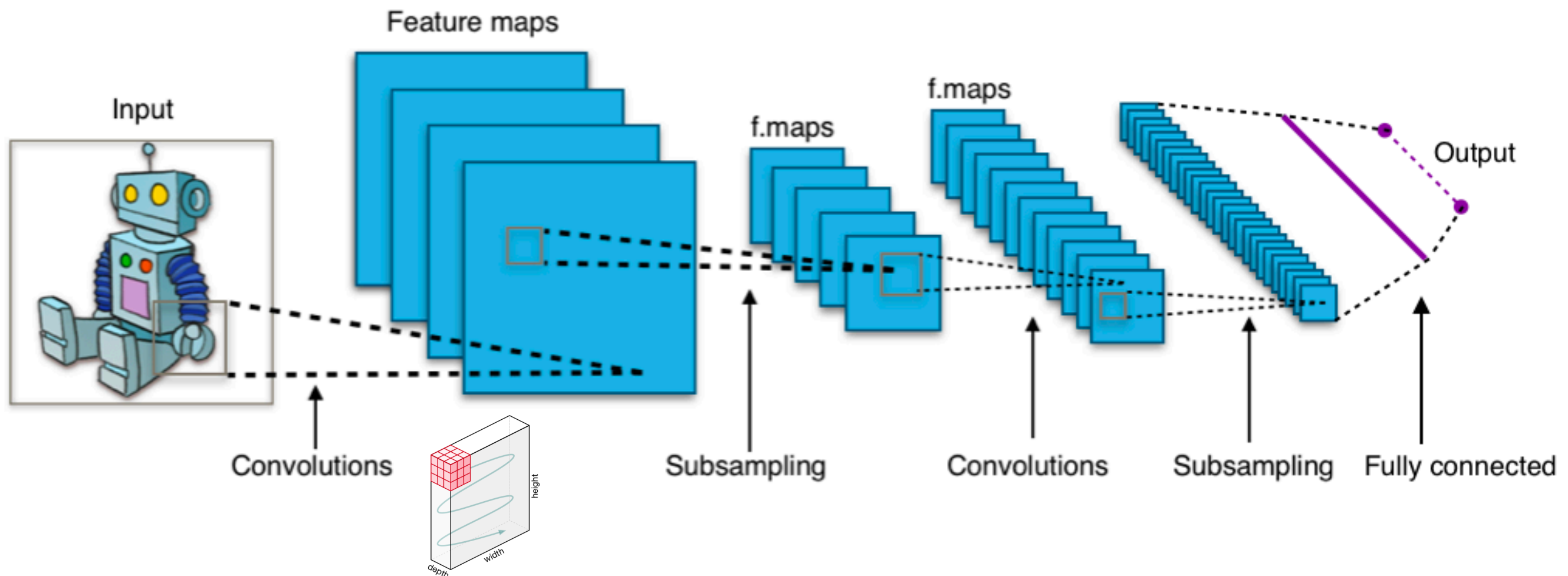
$$A^i = \mathbf{RELU}(Z^i)$$

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

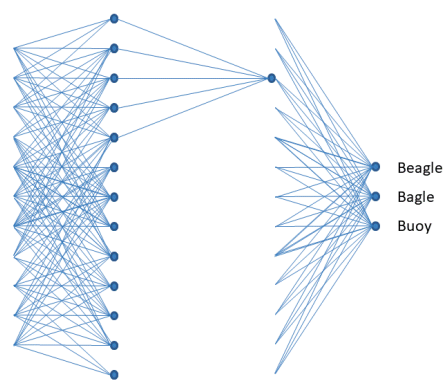


**CNN**

**Convolutional Neural Networks**





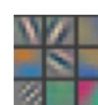
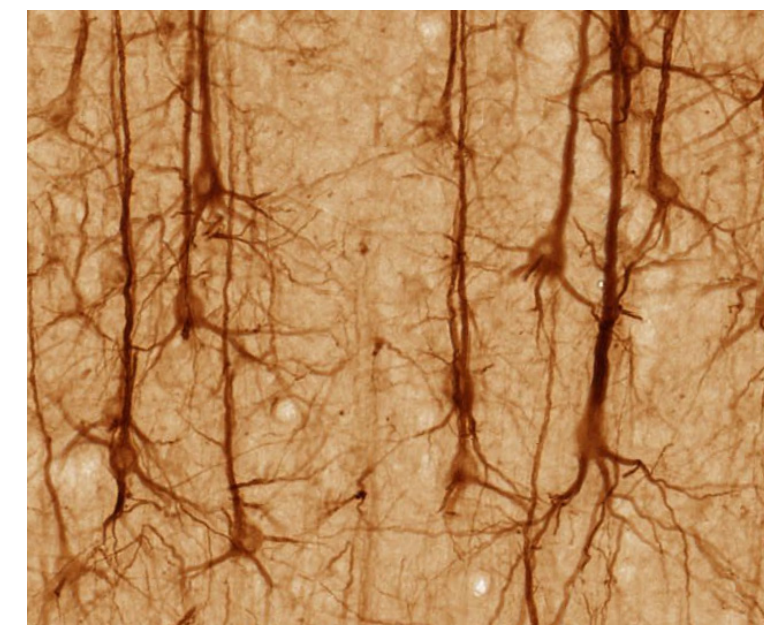


$$Z^i = W^i X + b^i 1$$

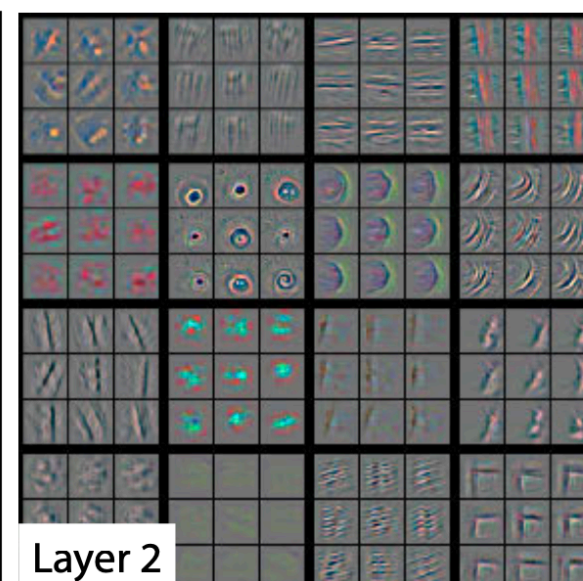
$$A^i = \text{RELU}(Z^i)$$

**Fully Connected / Feed Forward**

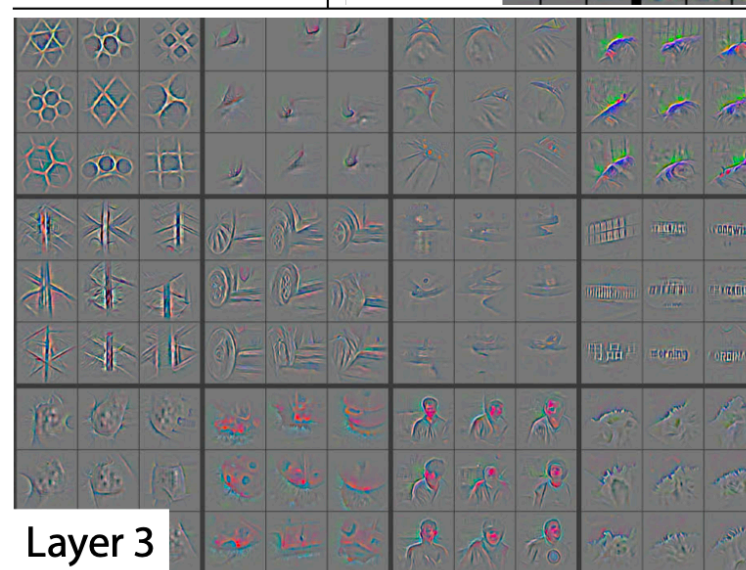
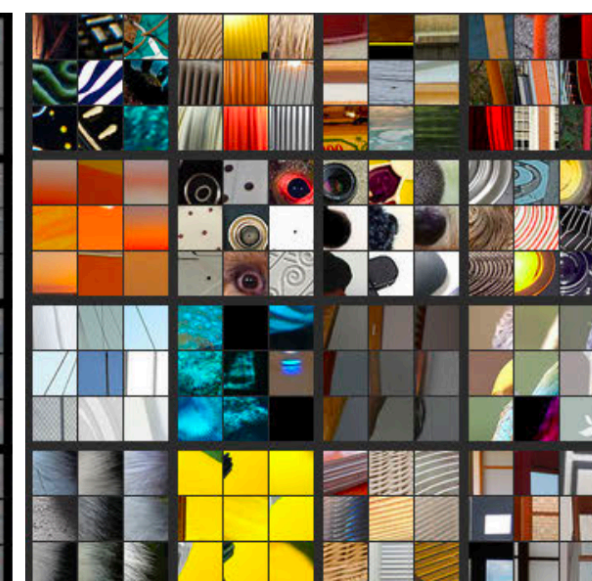
**FC**



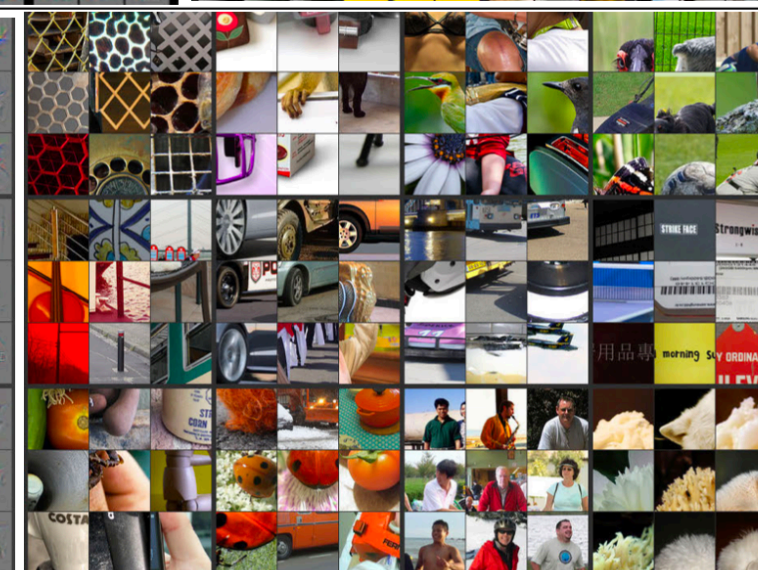
Layer 1

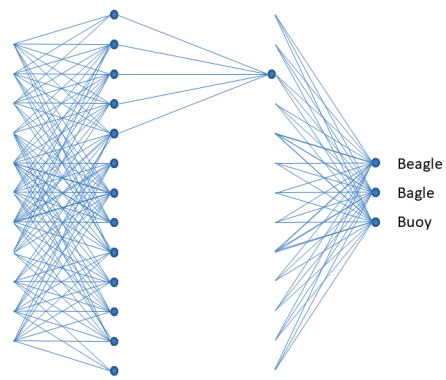


Layer 2



Layer 3





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

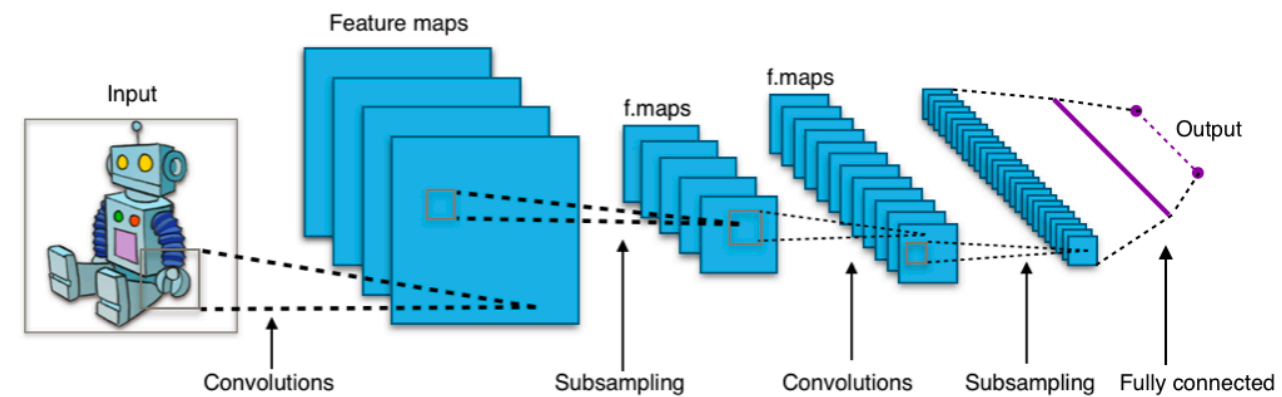
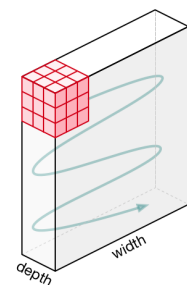
$$A^i = \mathbf{RELU}(Z^i)$$

**Fully Connected / Feed Forward**

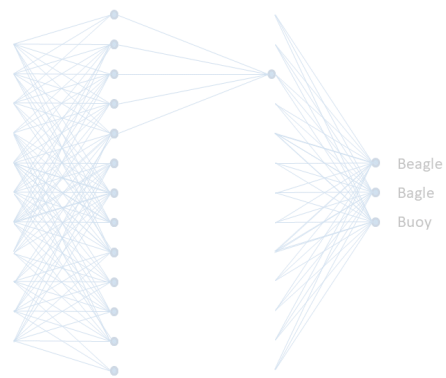
**FC**

**CNN**

**Convolutional Neural Networks**



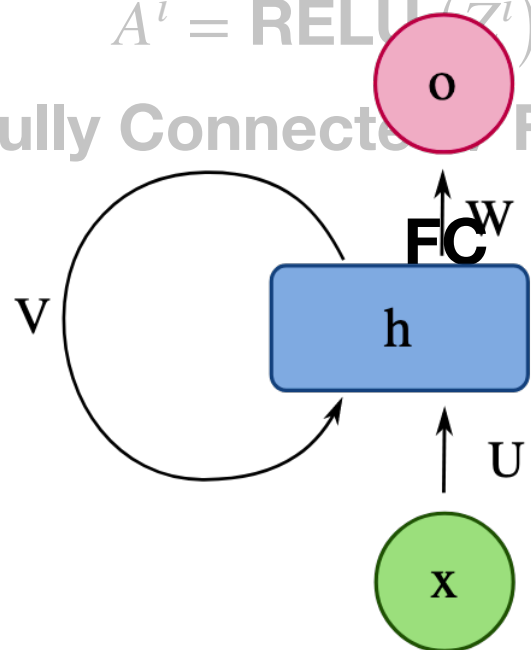




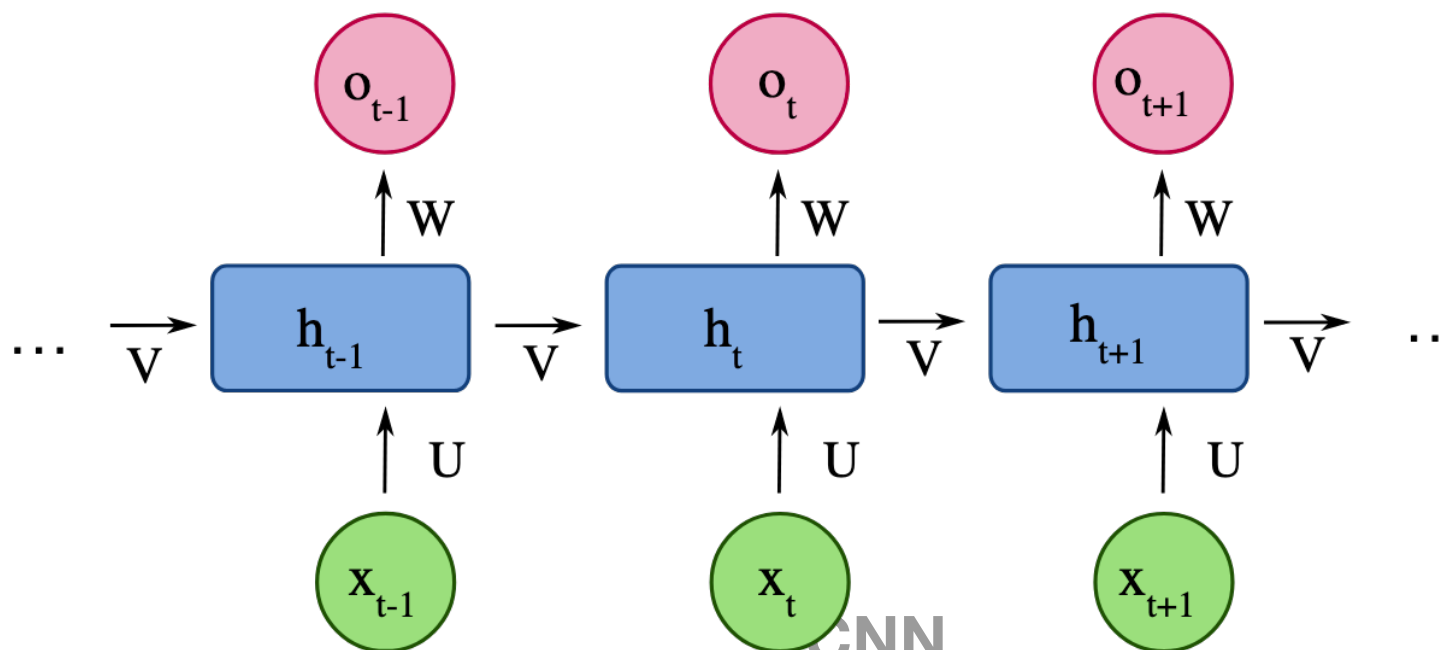
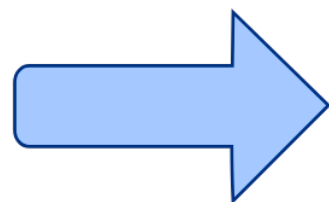
$$Z^i = W^i X + b^i 1$$

$$A^i = \text{RELU}(Z^i)$$

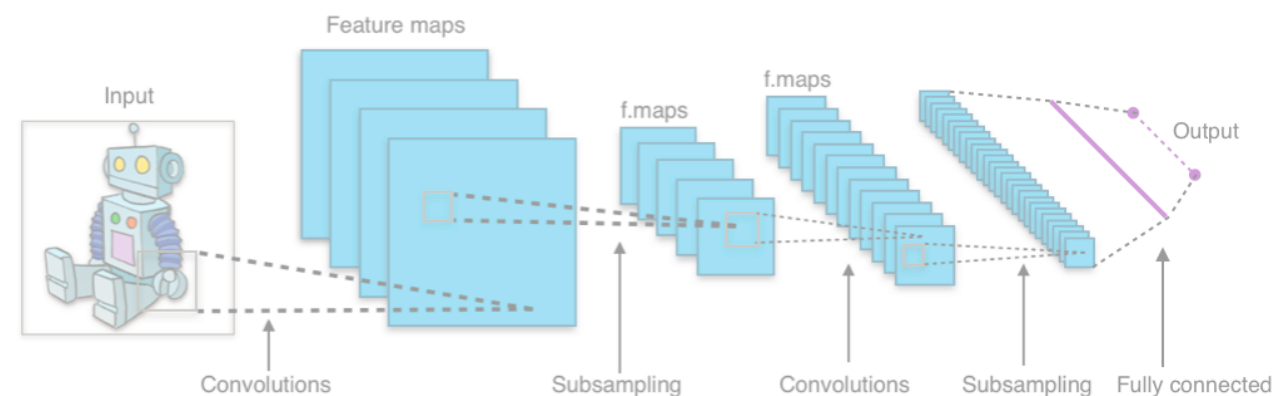
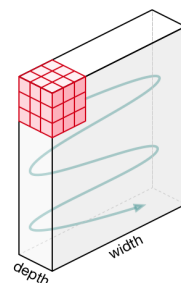
Fully Connected Feed Forward

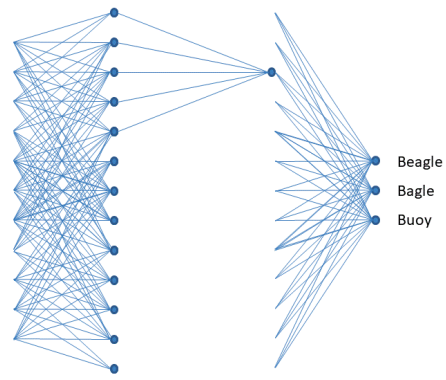


Unfold



Convolutional Neural Networks



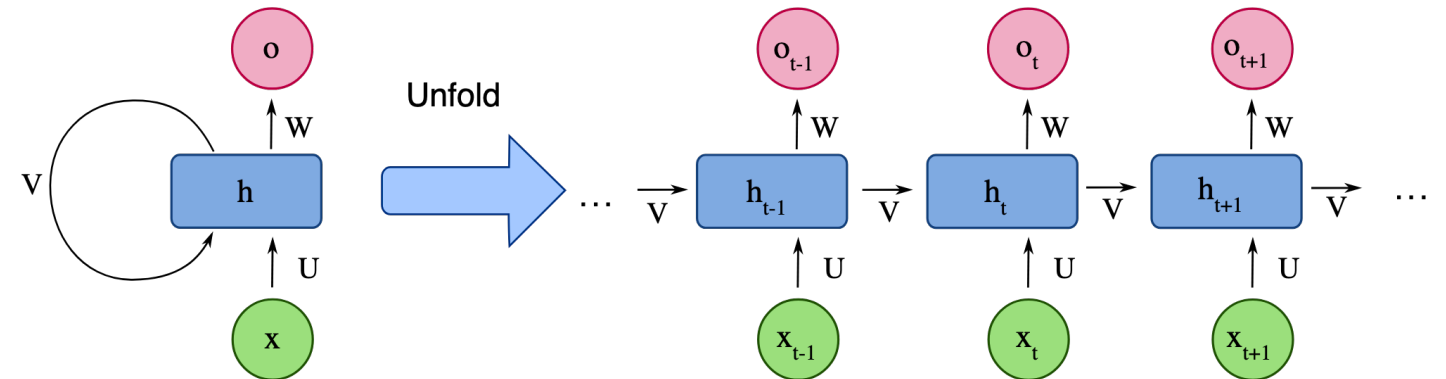


$$Z^i = W^i X + b^i 1$$

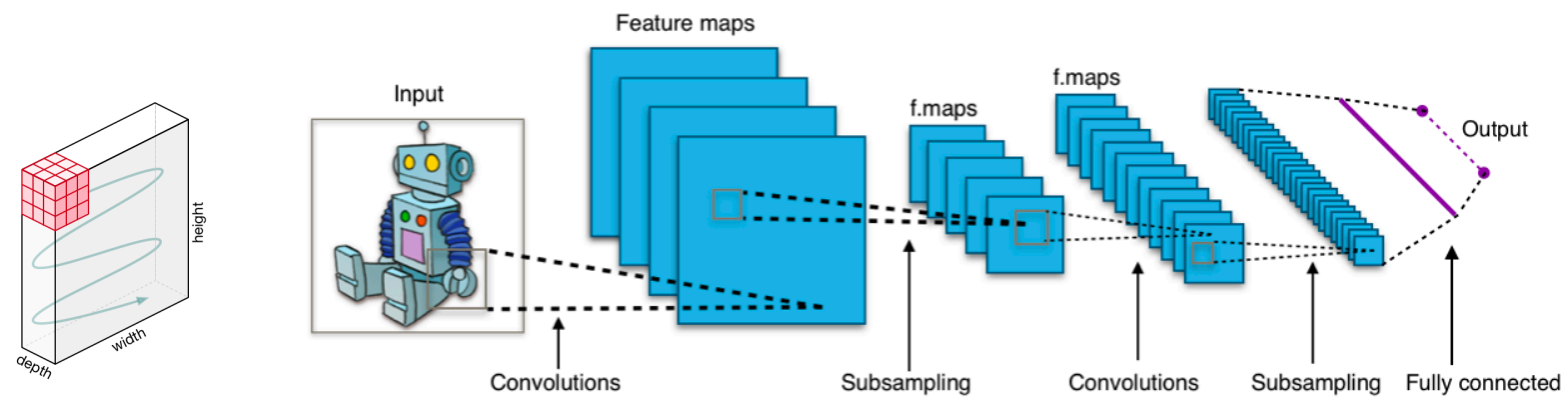
$$A^i = \text{RELU}(Z^i)$$

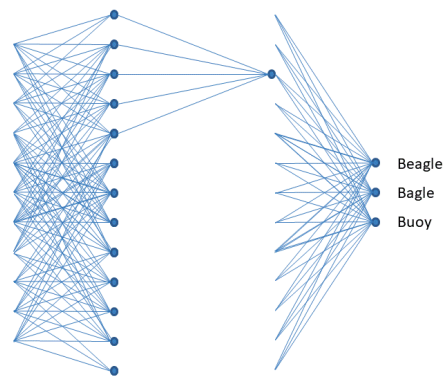
**Fully Connected / Feed Forward**  
**FC**

## RNN Recurrent Neural Network



## CNN Convolutional Neural Networks



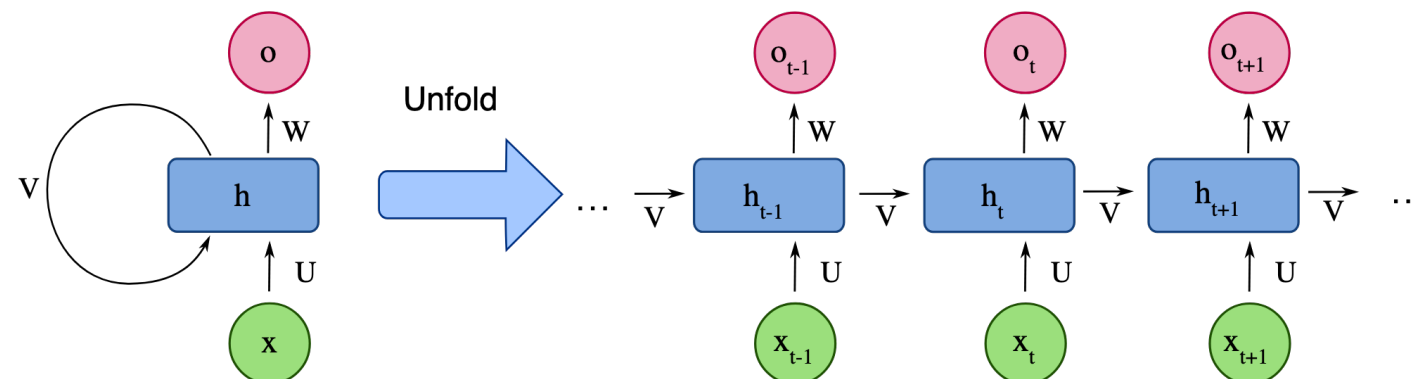


$$Z^i = W^i X + b^i 1$$

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

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

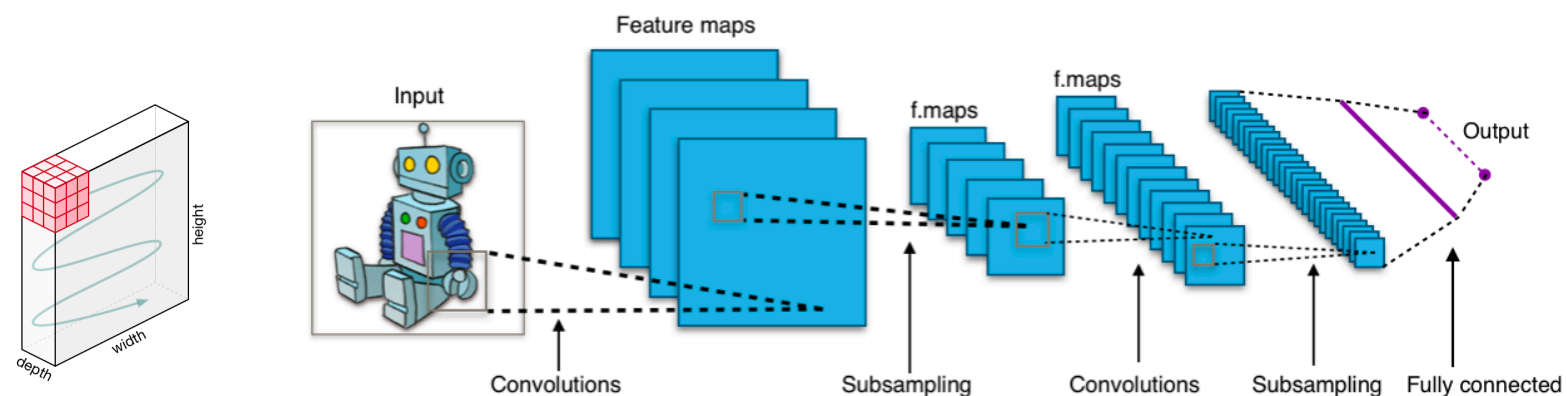
## RNN Recurrent Neural Network

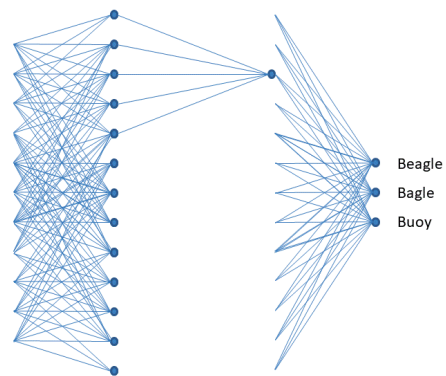


**GANs,**  
**Auto Encoders,**  
**ODE Networks,**  
**Invertible Flow Networks,**

.....

## CNN Convolutional Neural Networks





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

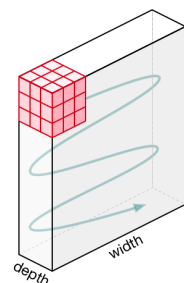
$$A^i = \text{RELU}(Z^i)$$

**Fully Connected / Feed Forward**

**FC**

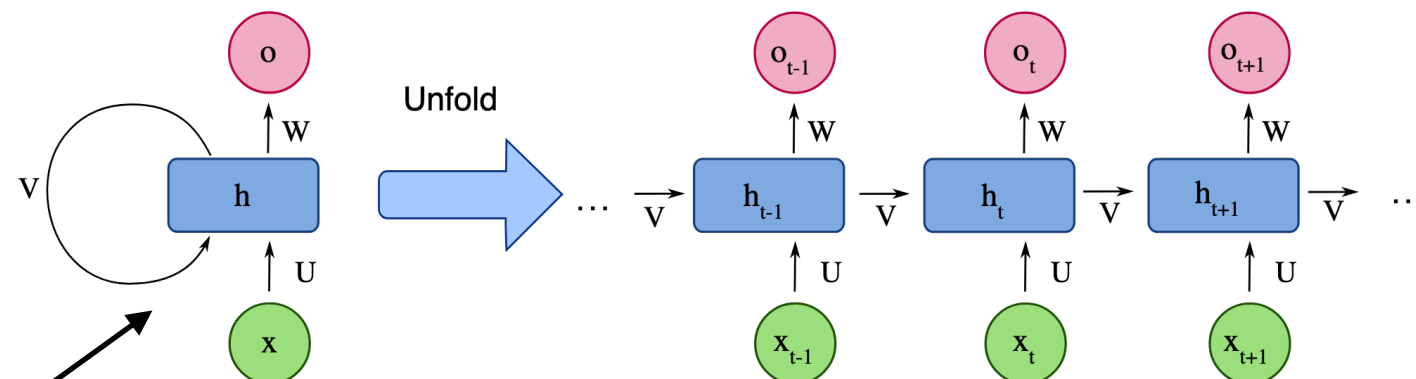
**GANs,  
Auto Encoders,  
ODE Networks,  
Invertible Flow Networks,**

.....



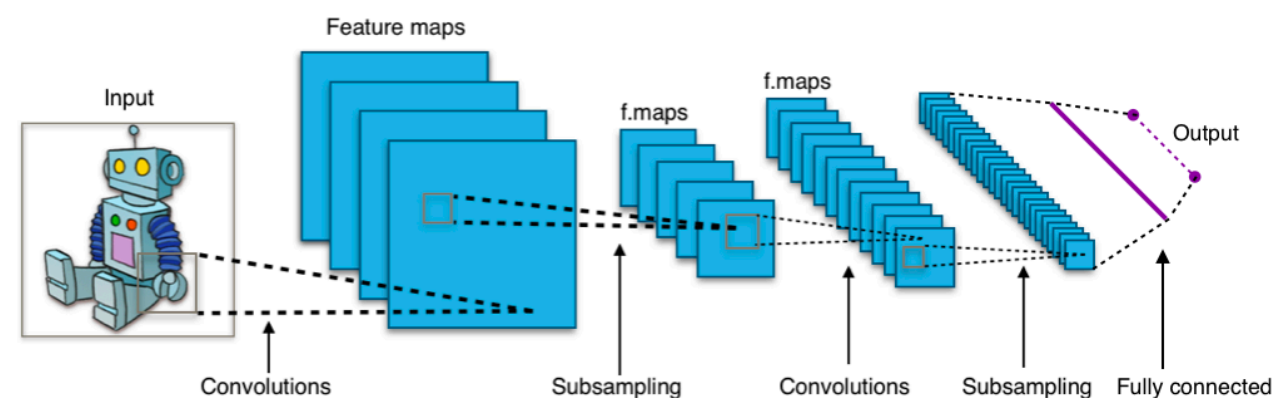
**RNN**

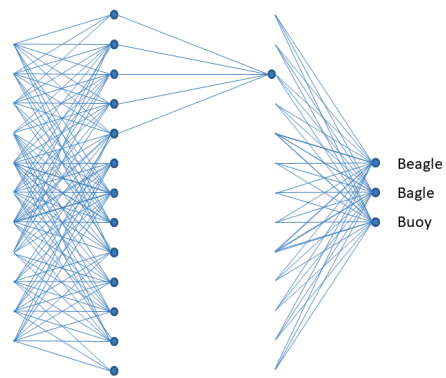
**Recurrent Neural Network**



**CNN**

**Convolutional Neural Networks**





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

$$A^i = \text{RELU}(Z^i)$$

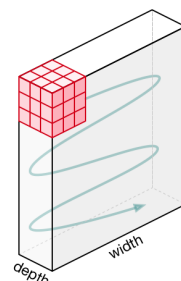
**Fully Connected / Feed Forward**

**FC**

**GANs,  
Auto Encoders,  
ODE Networks,**

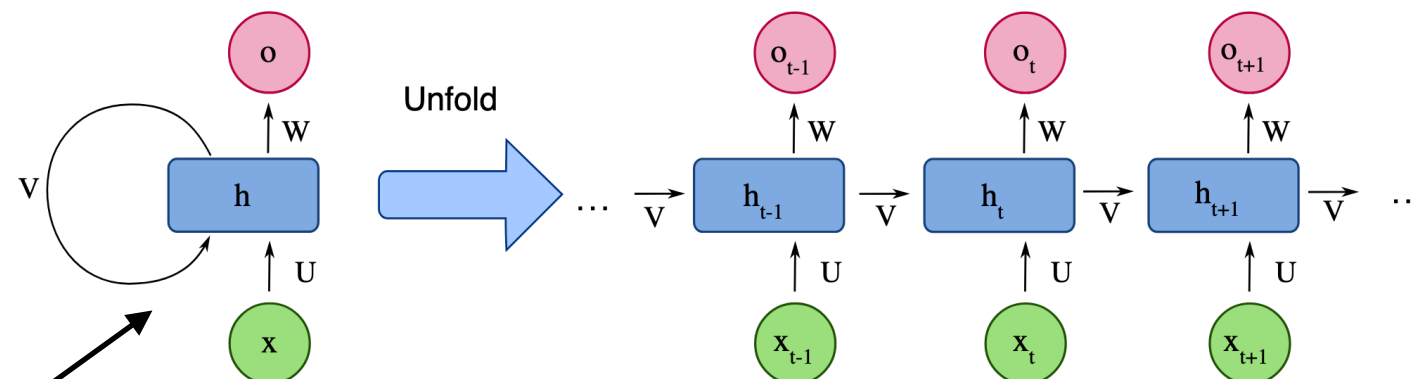
**Invertible Flow Networks,**

.....



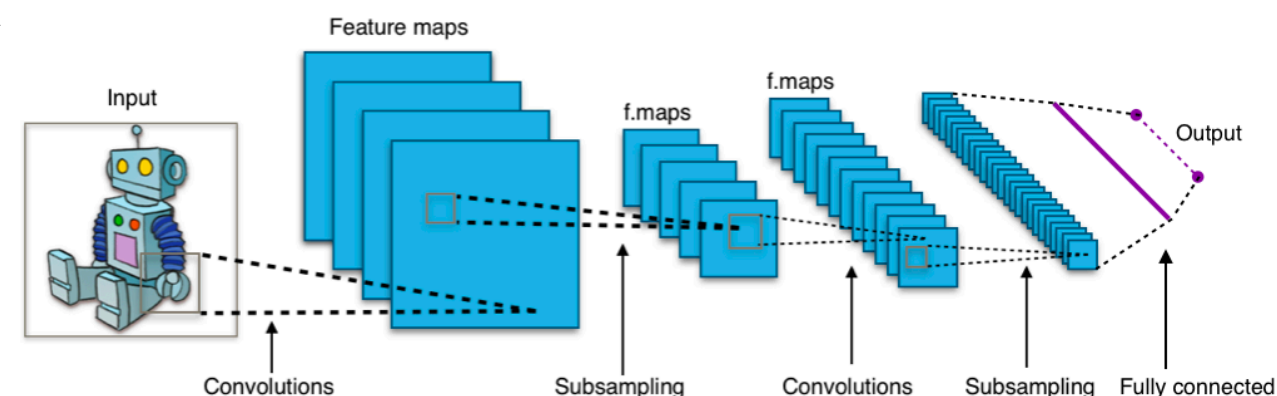
**RNN**

**Recurrent Neural Network**

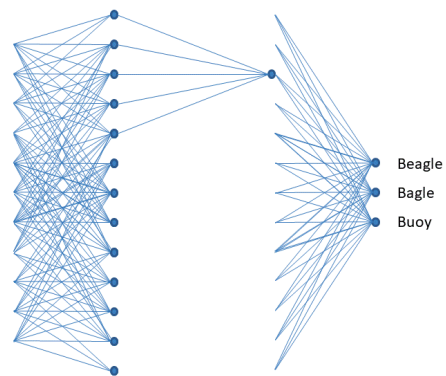


**CNN**

**Convolutional Neural Networks**



**All kind of Domains: Medical Imaging, Autonomous Driving, Emotion Recognition,  
Recommenders, Natural Language Processing**



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

$$A^i = \text{RELU}(Z^i)$$

**Fully Connected / Feed Forward**

**FC**

**GANs,  
Auto Encoders,  
ODE Networks,**

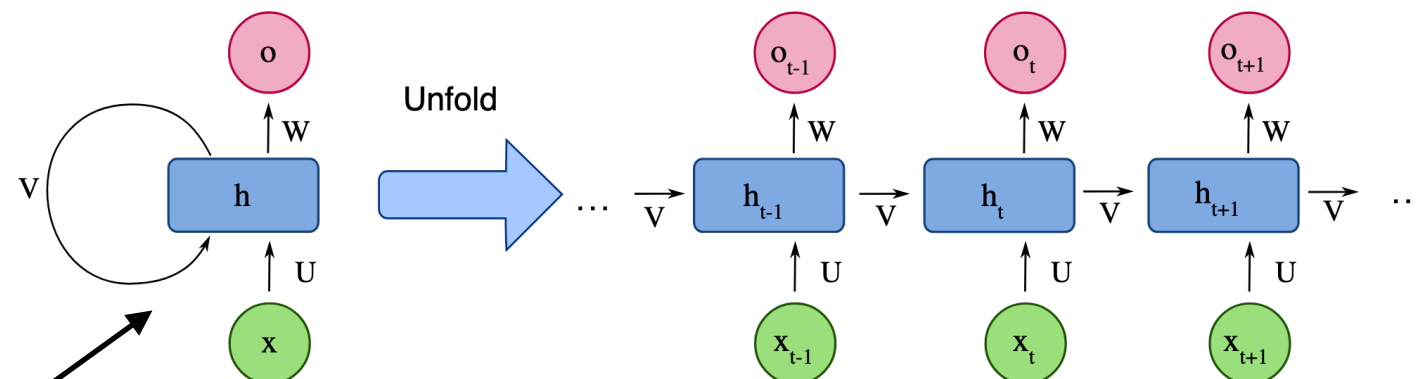
**Invertible Flow Networks,**

.....

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

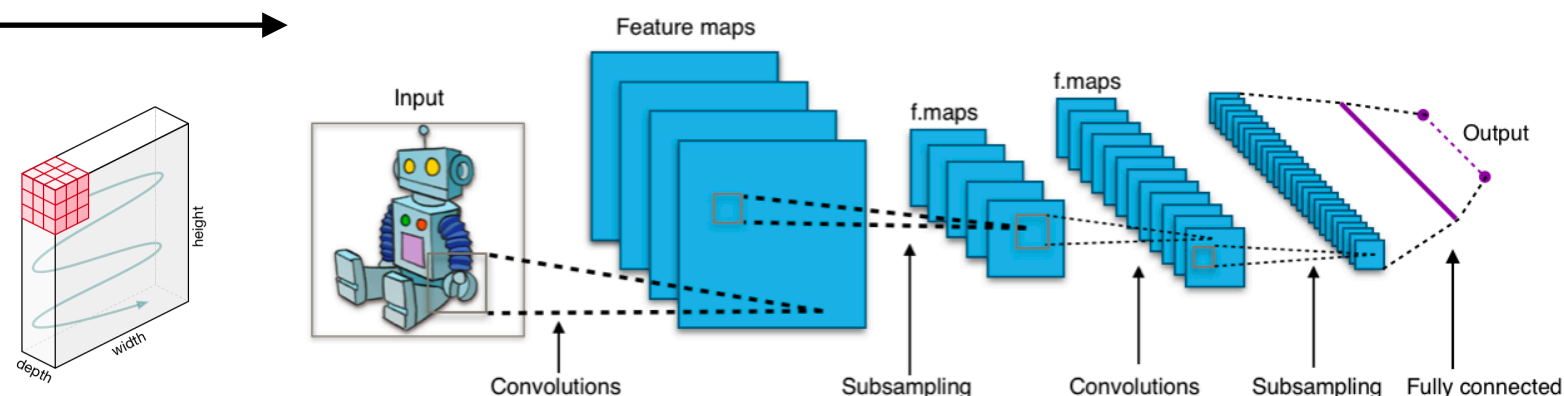
**RNN**

**Recurrent Neural Network**



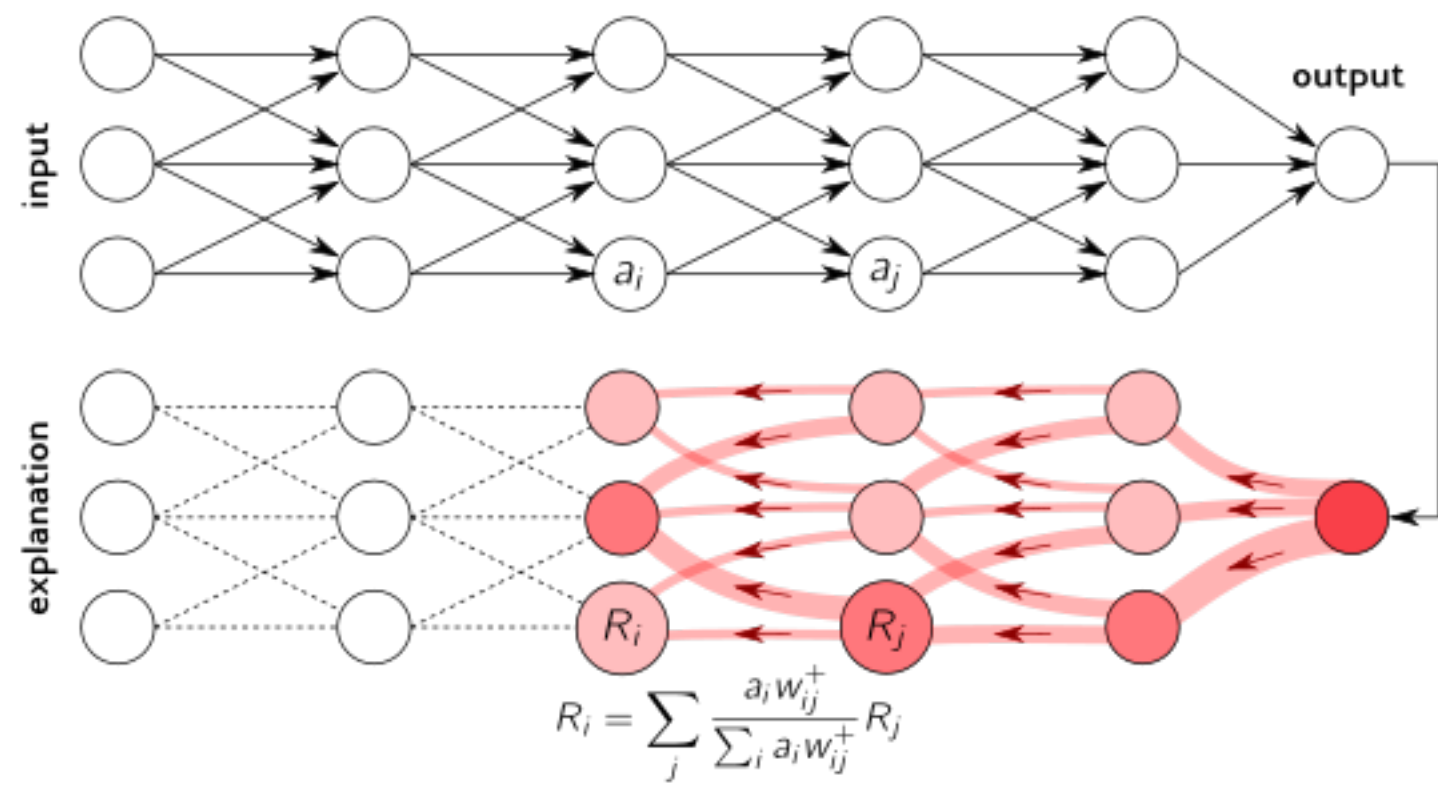
**CNN**

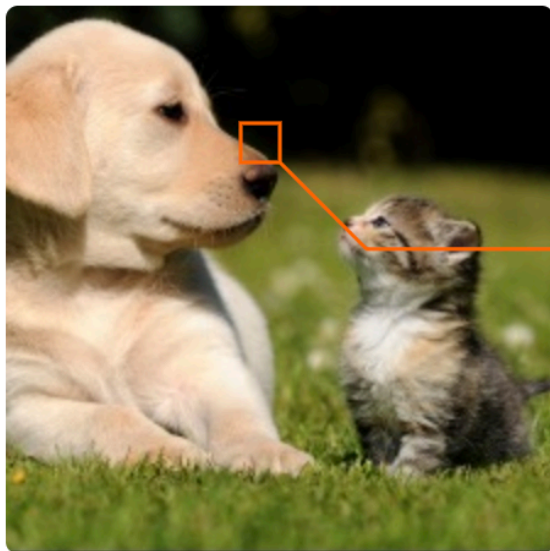
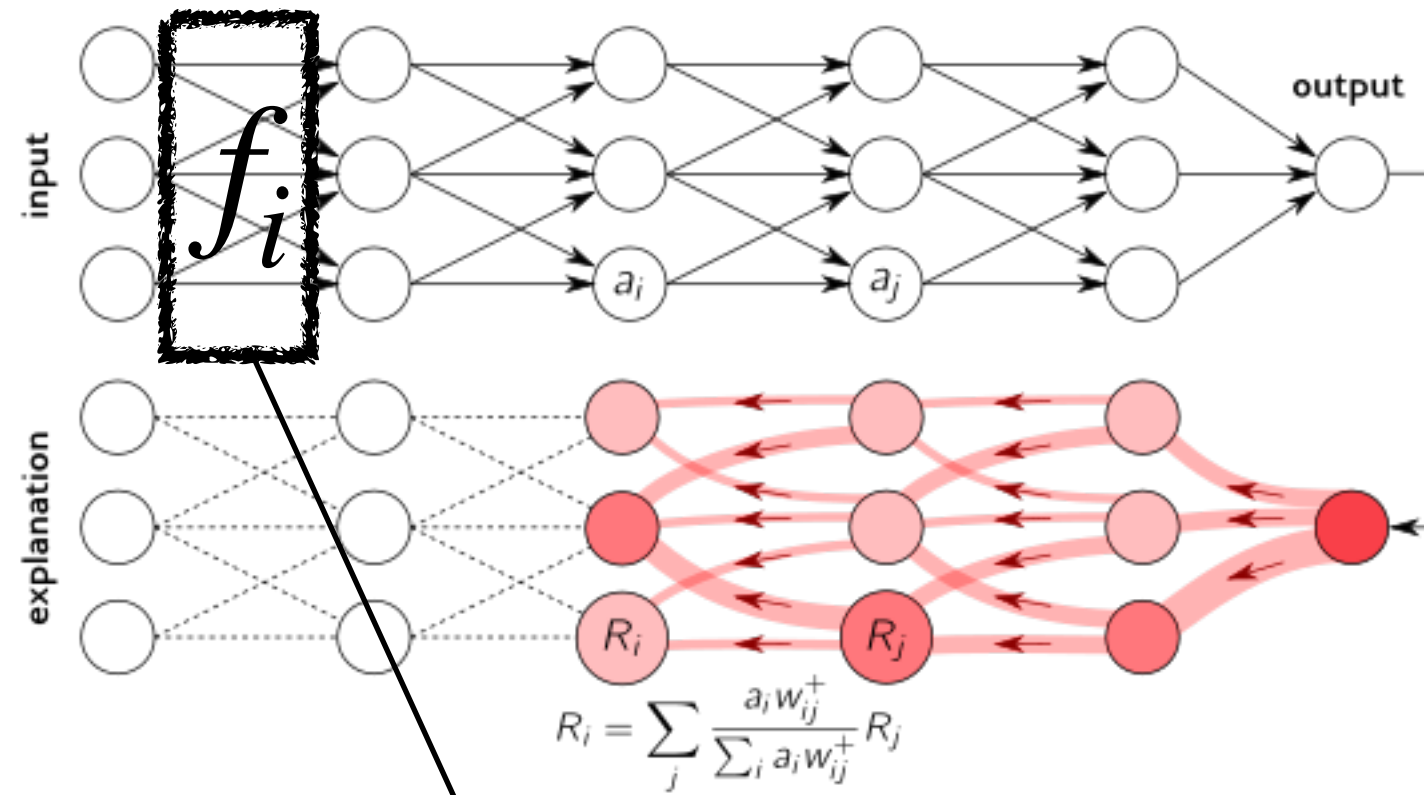
**Convolutional Neural Networks**



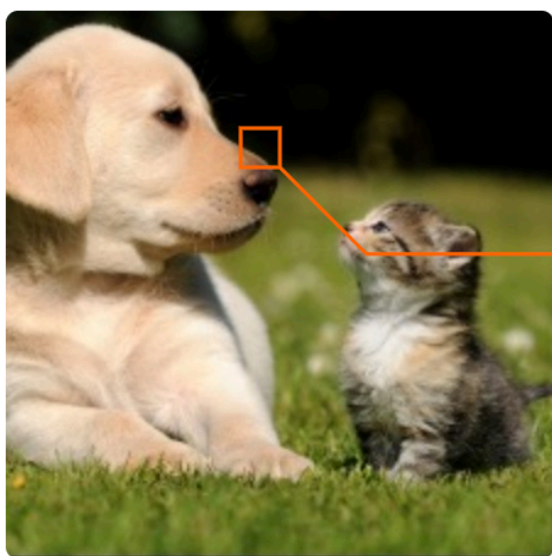
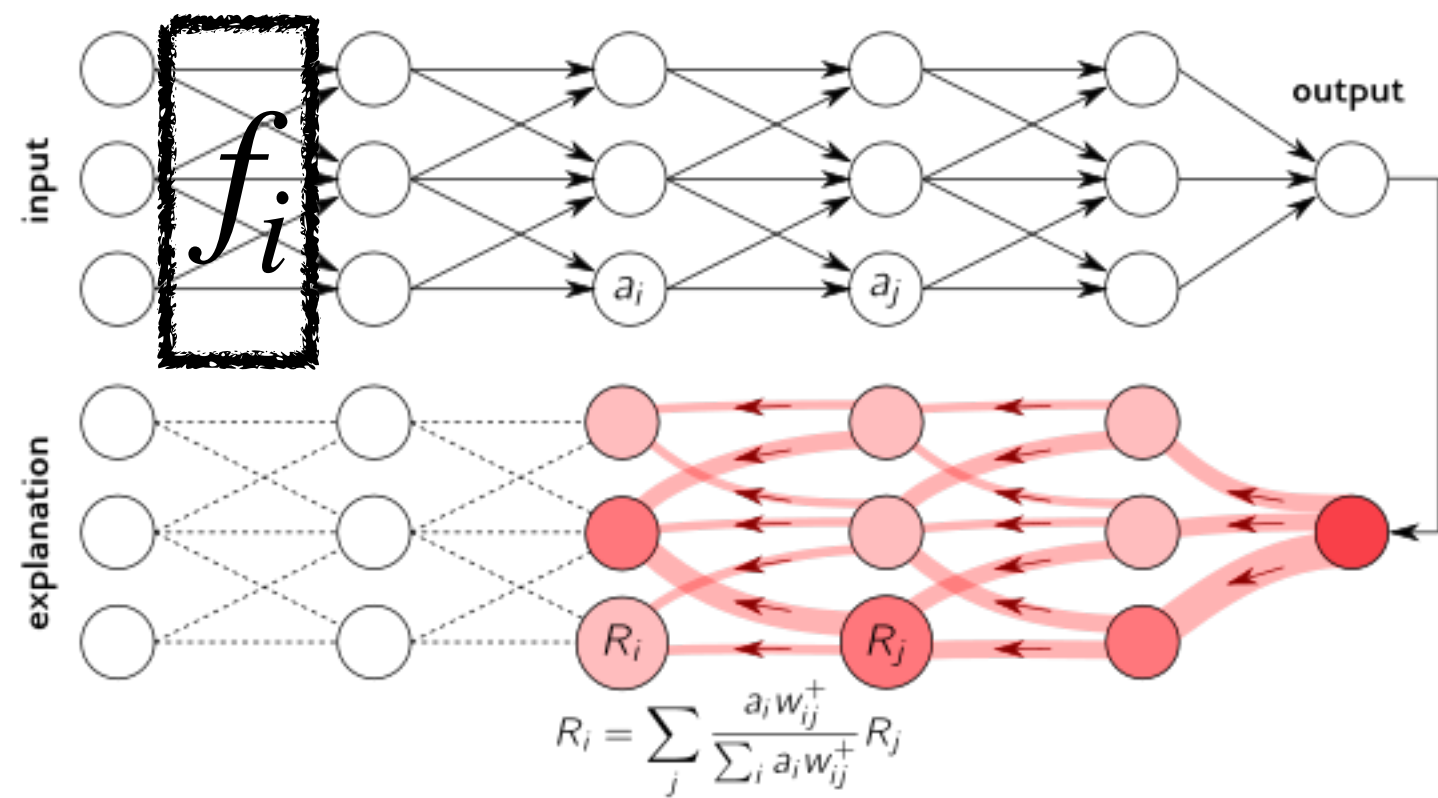
**All kind of Domains: Medical Imaging, Autonomous Driving, Emotion Recognition,  
Recommenders, Natural Language Processing**



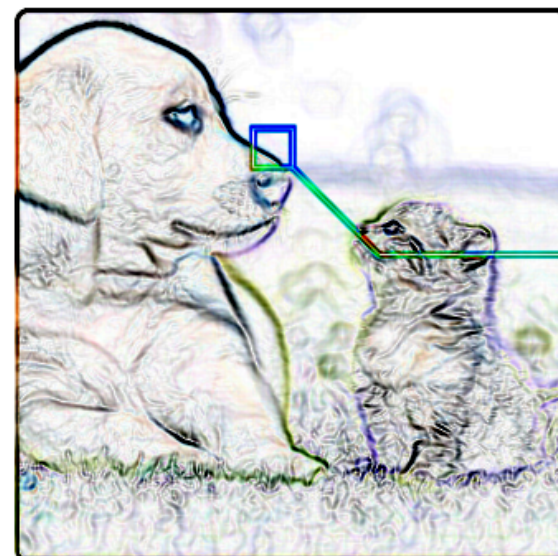


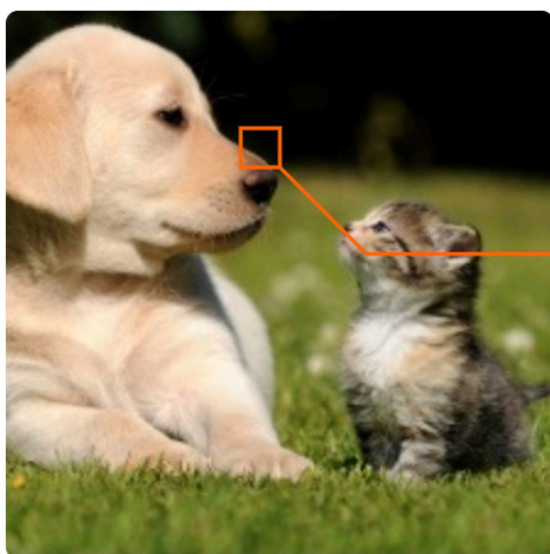
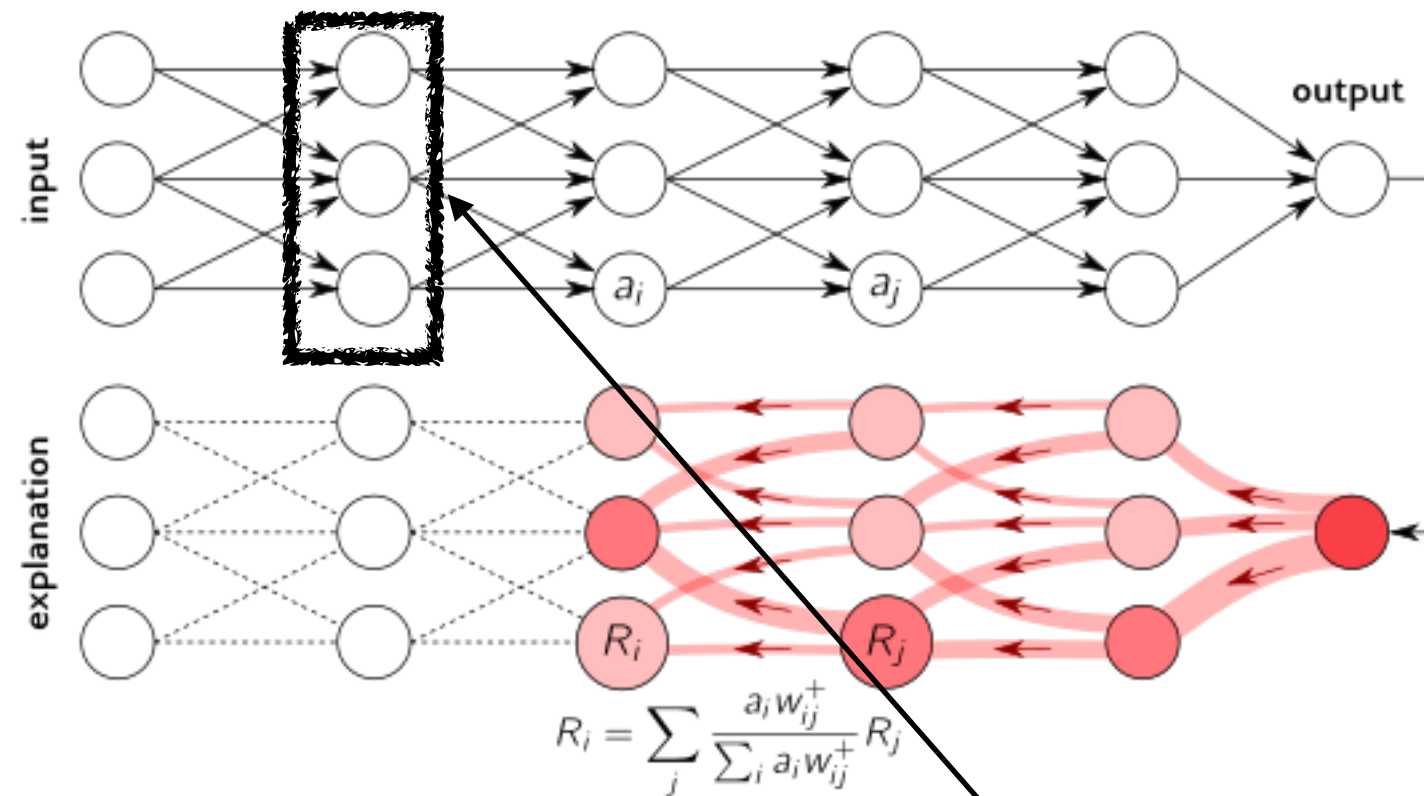


$*f_i$

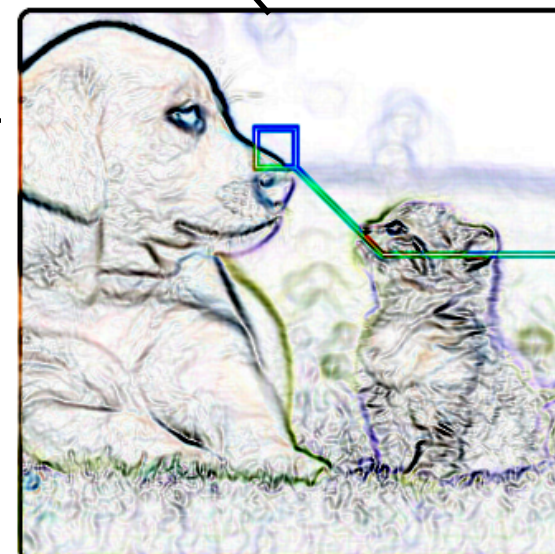


$* f_i$



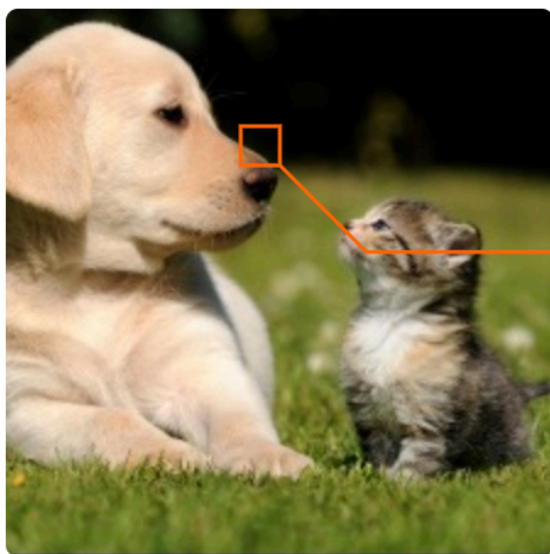
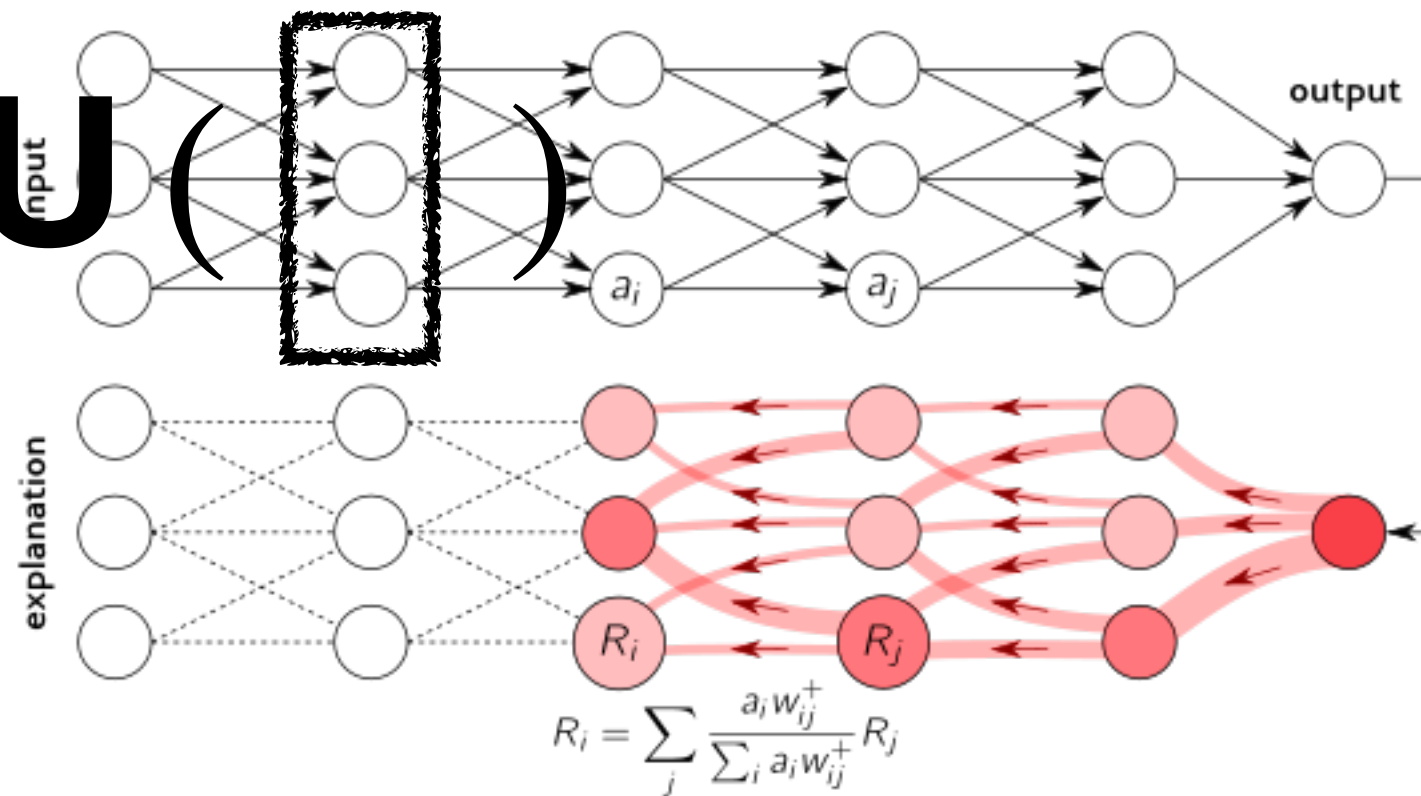


$\rightarrow$   
 $* f_i$

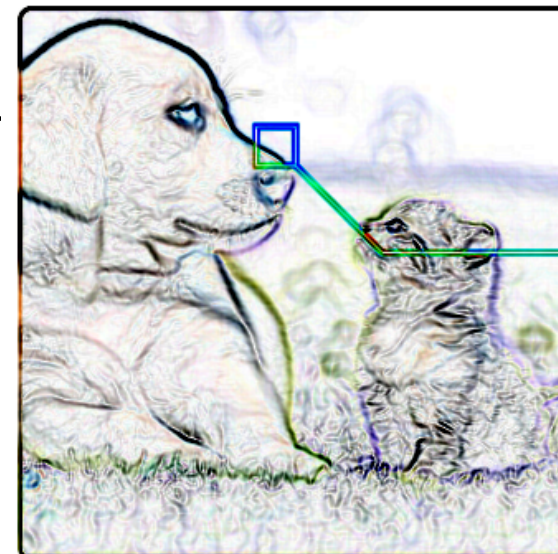




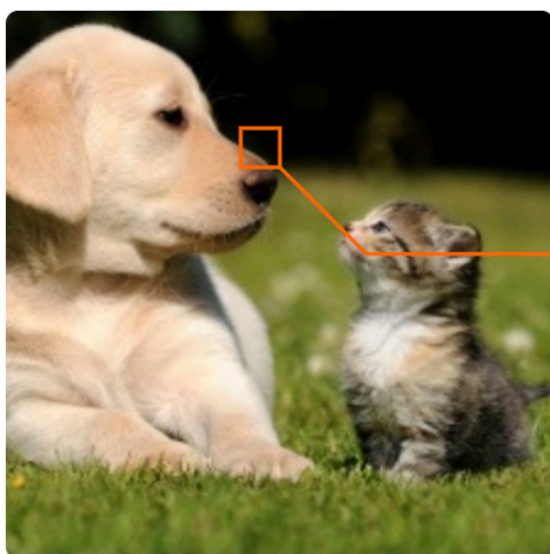
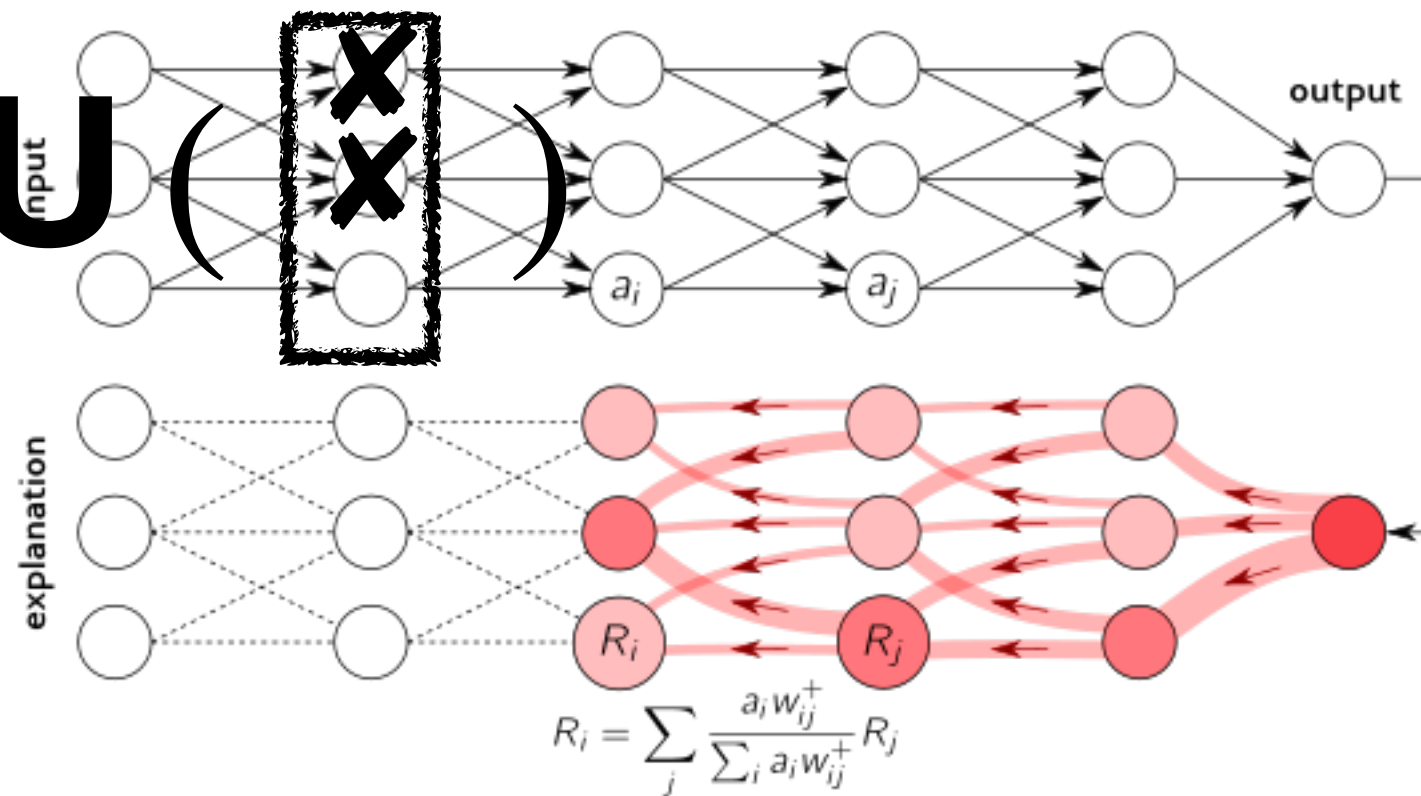
# RELU



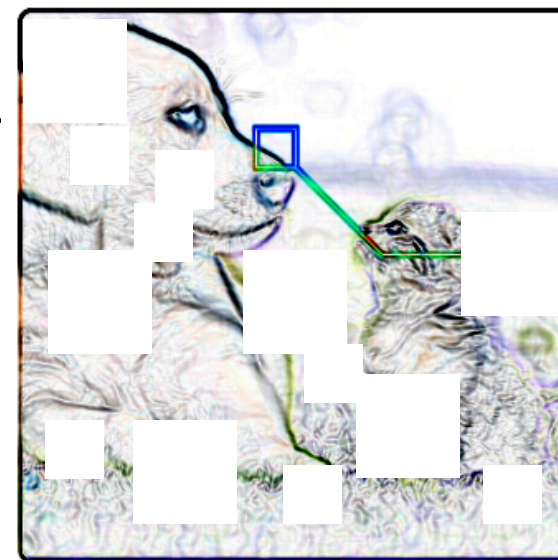
$\rightarrow$   
 $* f_i$



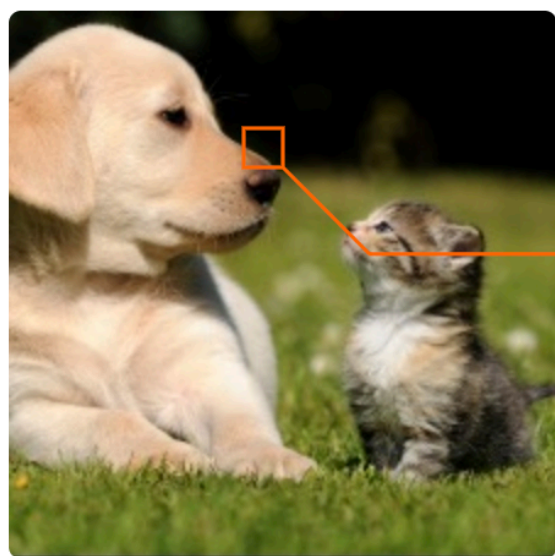
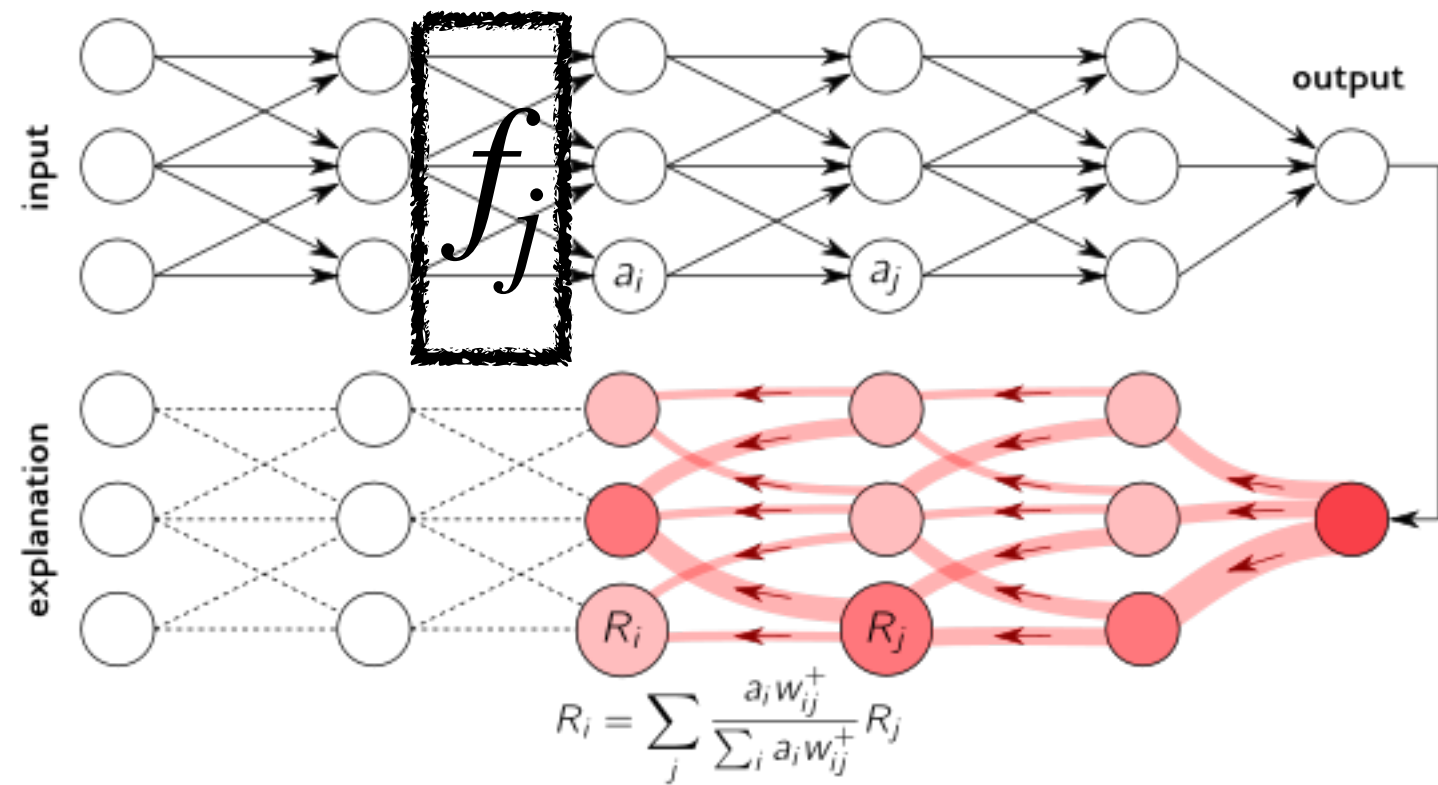
# RELU



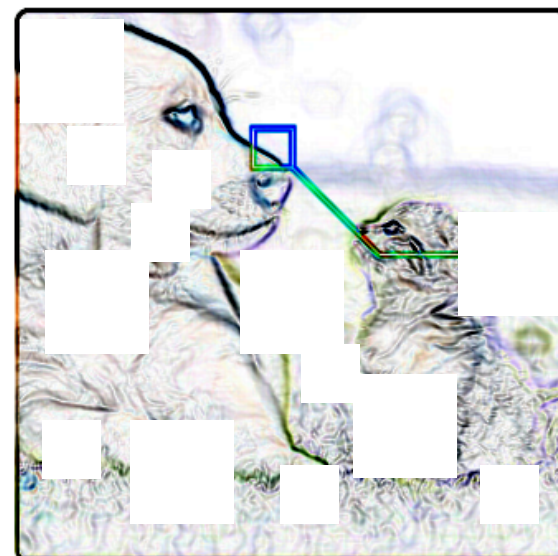
$\xrightarrow{*f_i}$



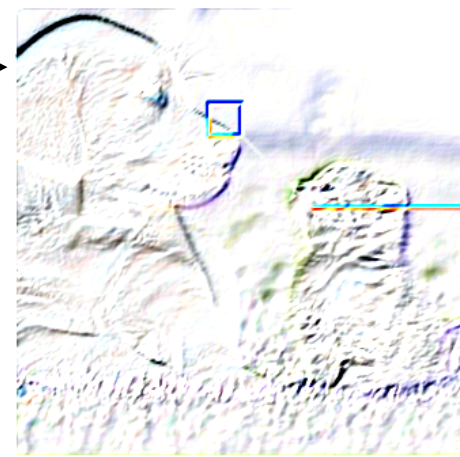


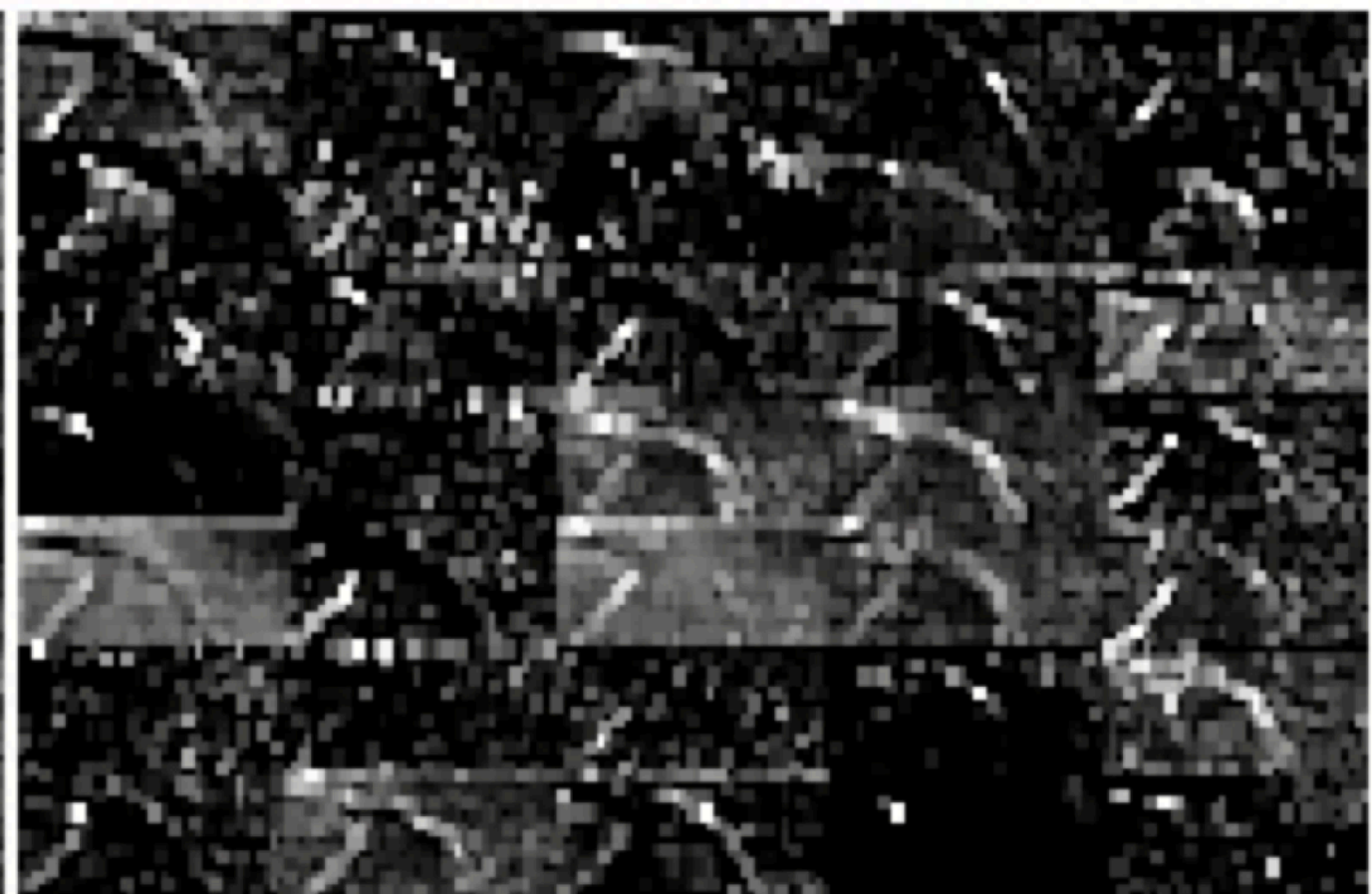
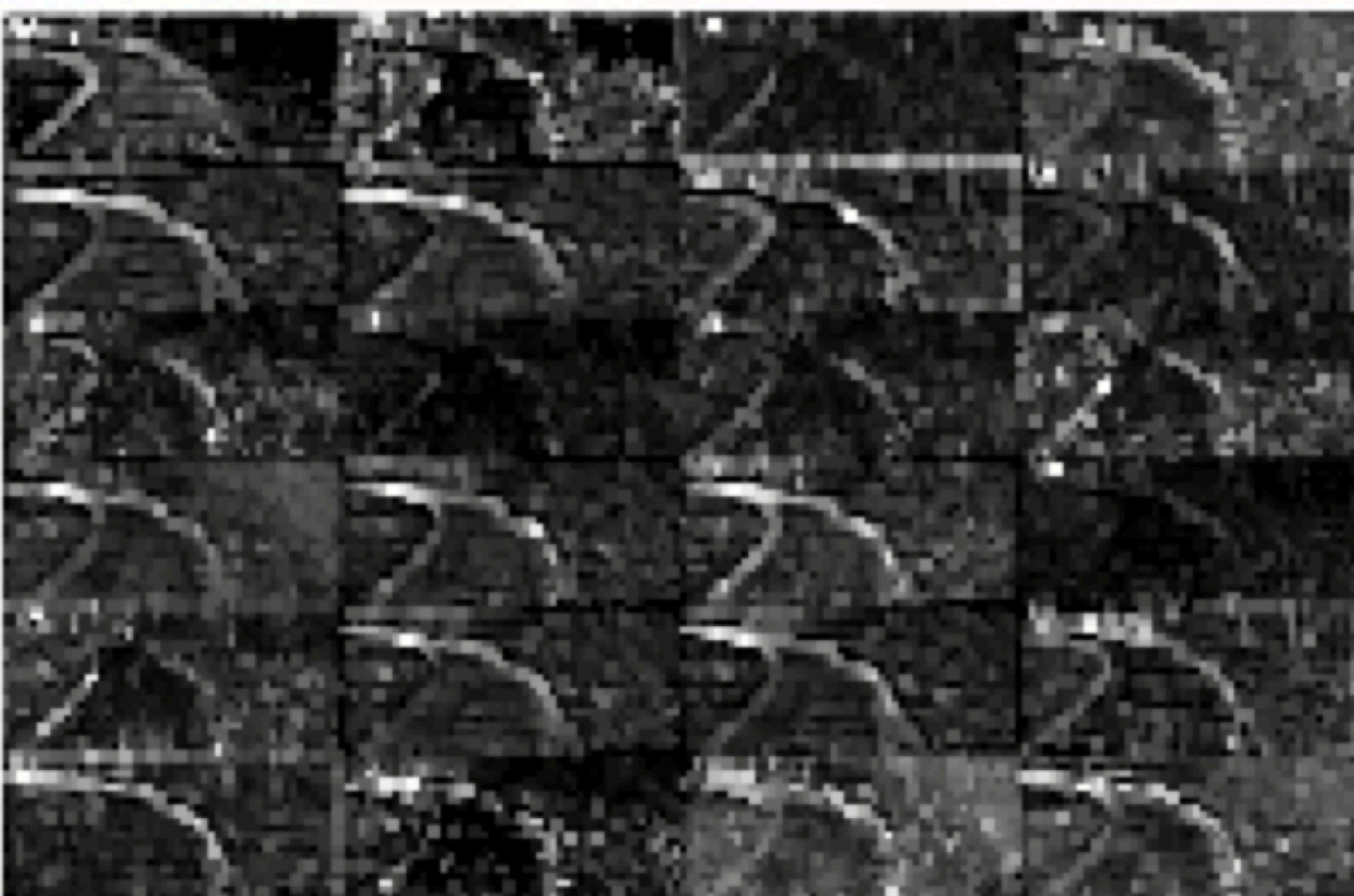


$\rightarrow * f_i$



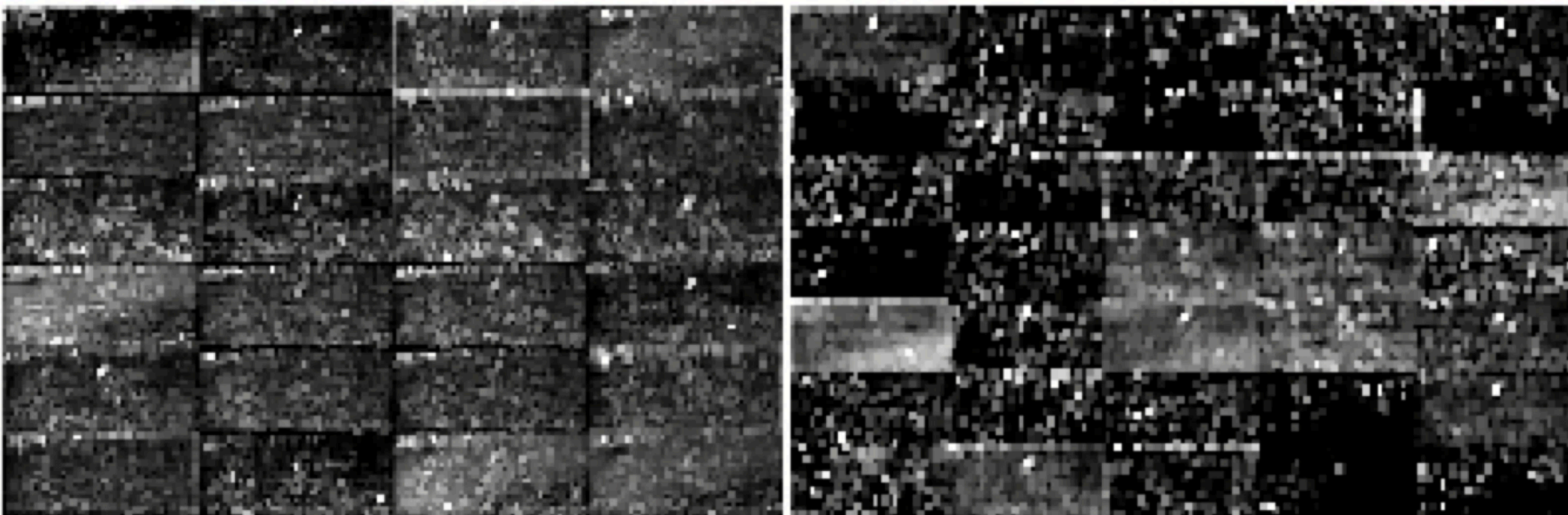
$\rightarrow * f_j$



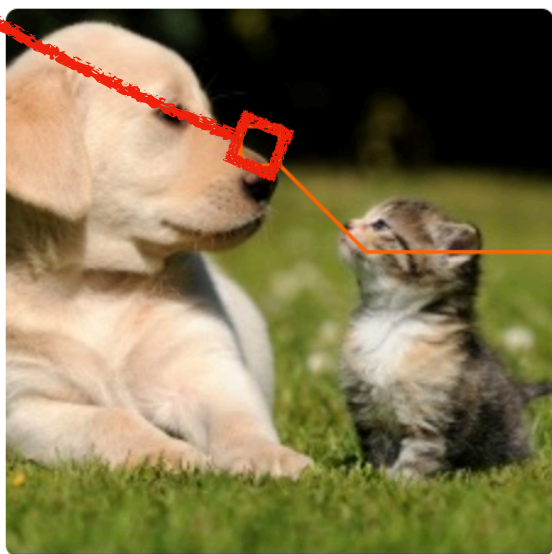
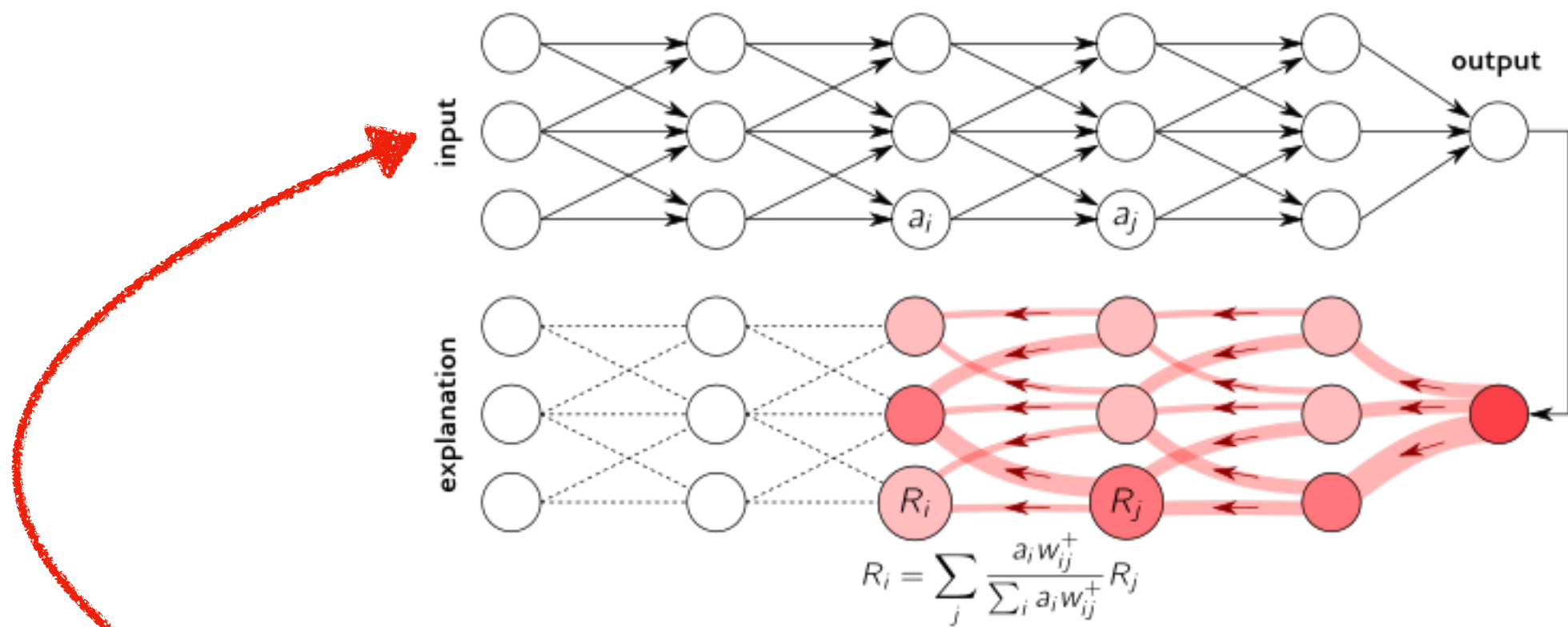


2016 Bojarski et al., End to End Learning for Self-Driving Cars

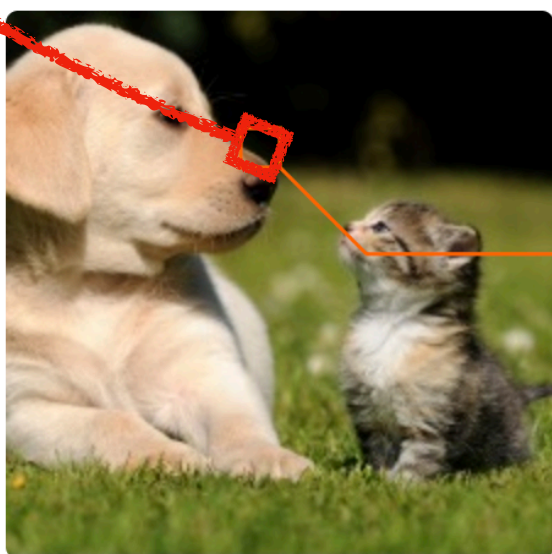
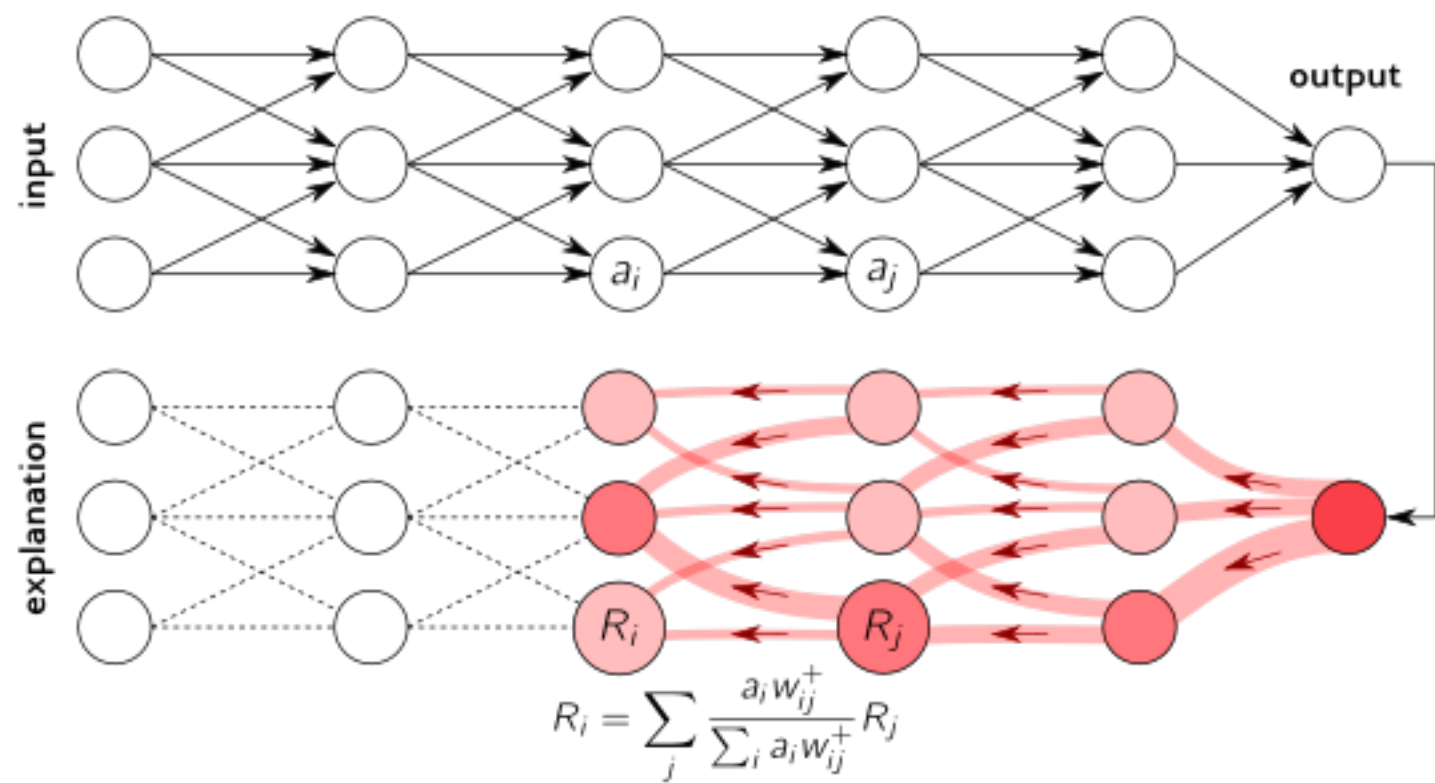
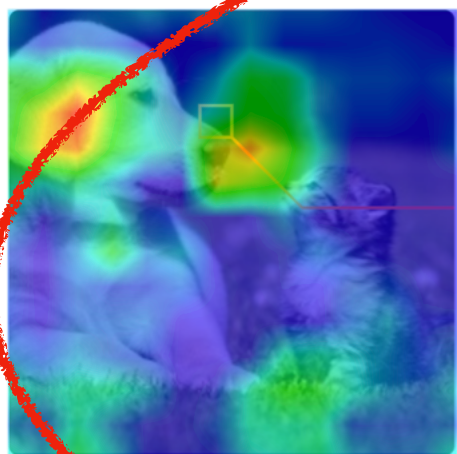




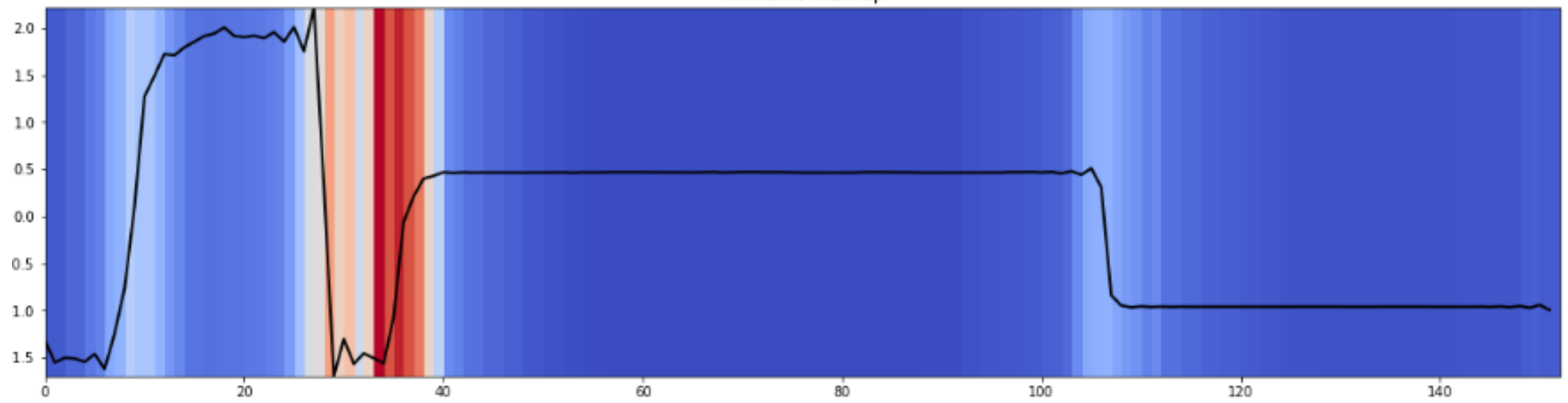
2016 Bojarski et al., End to End Learning for Self-Driving Cars



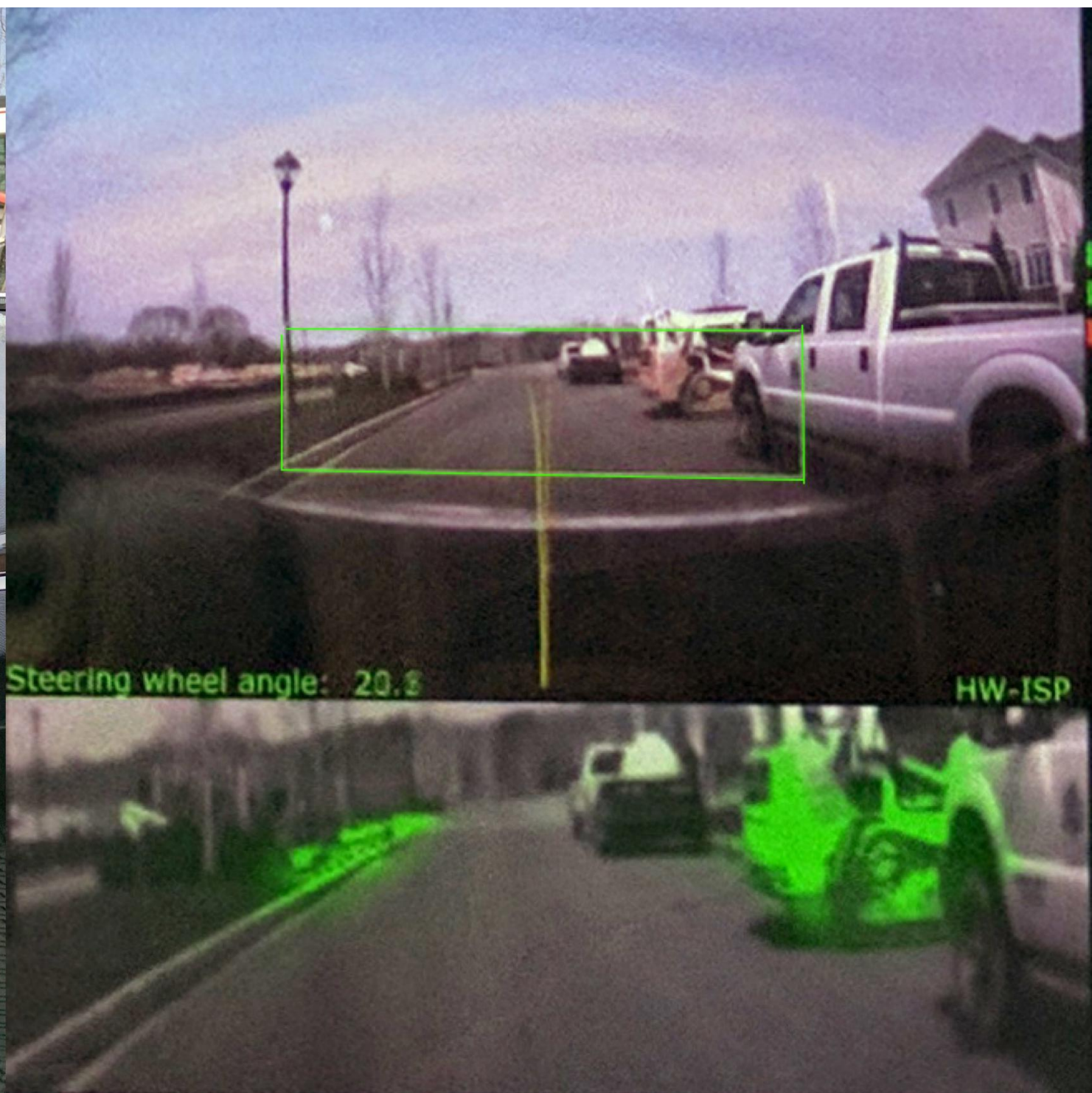




Relevance heatmap





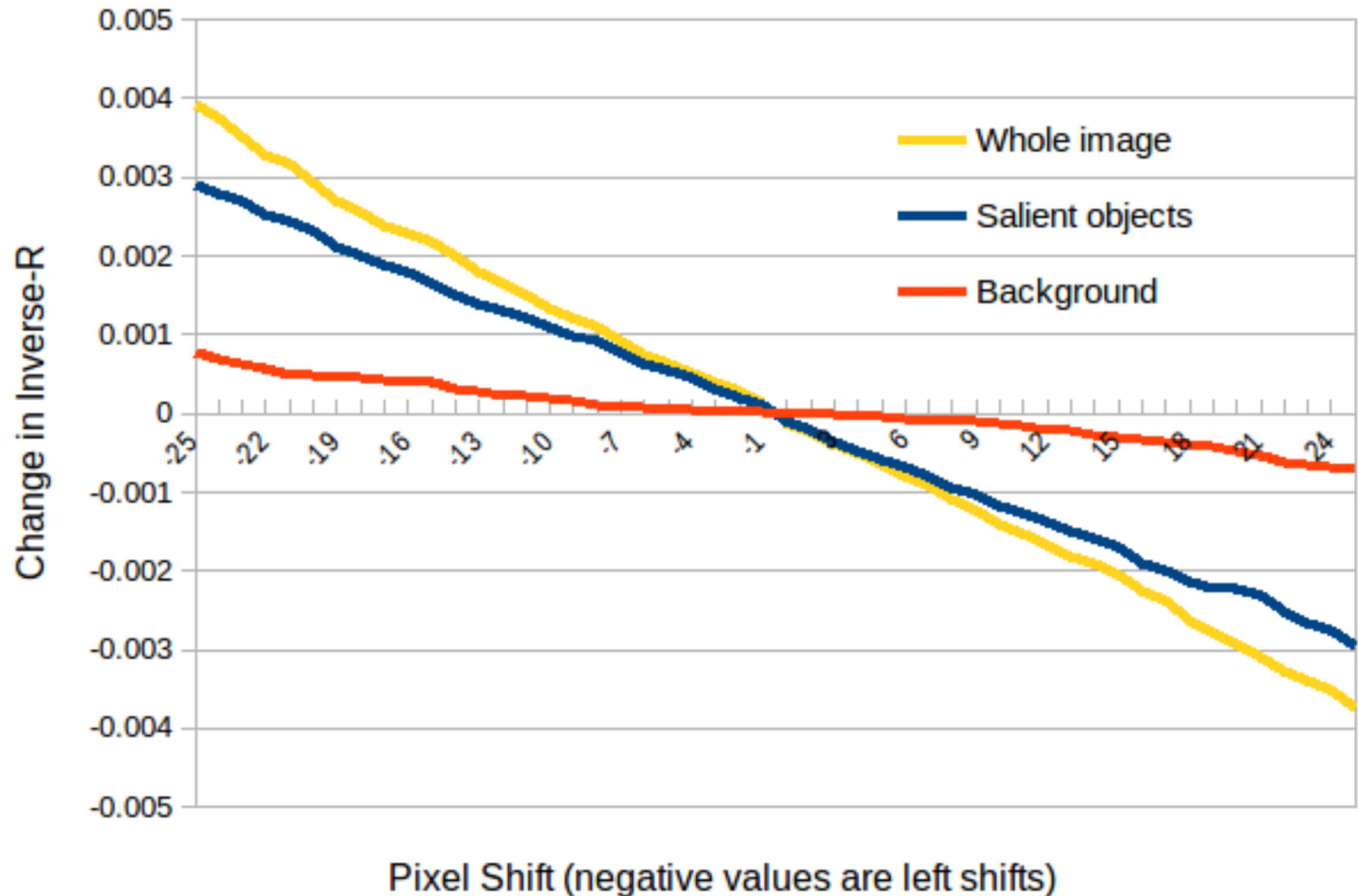


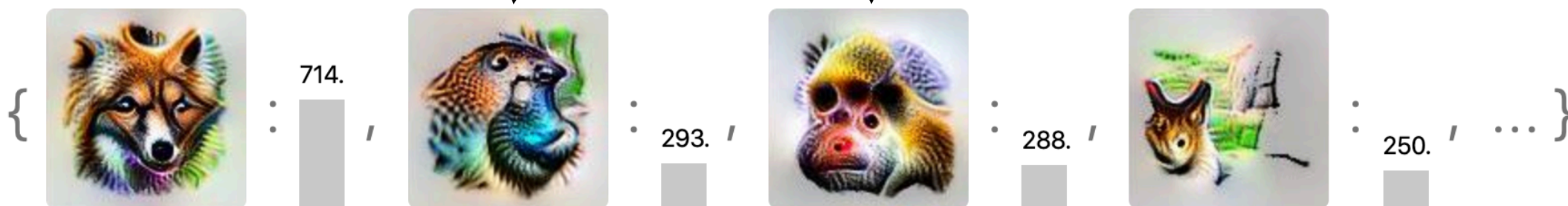
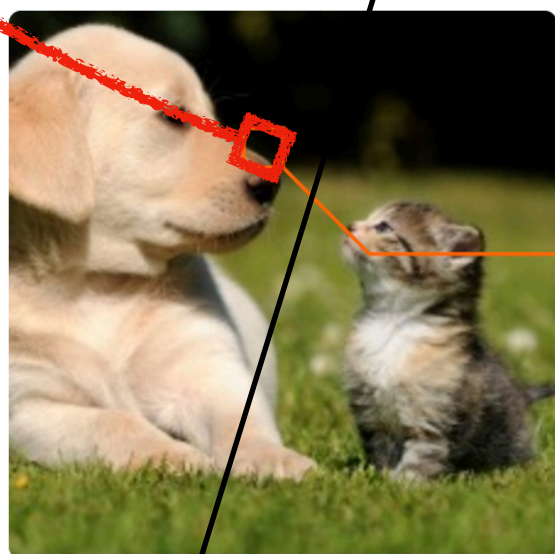
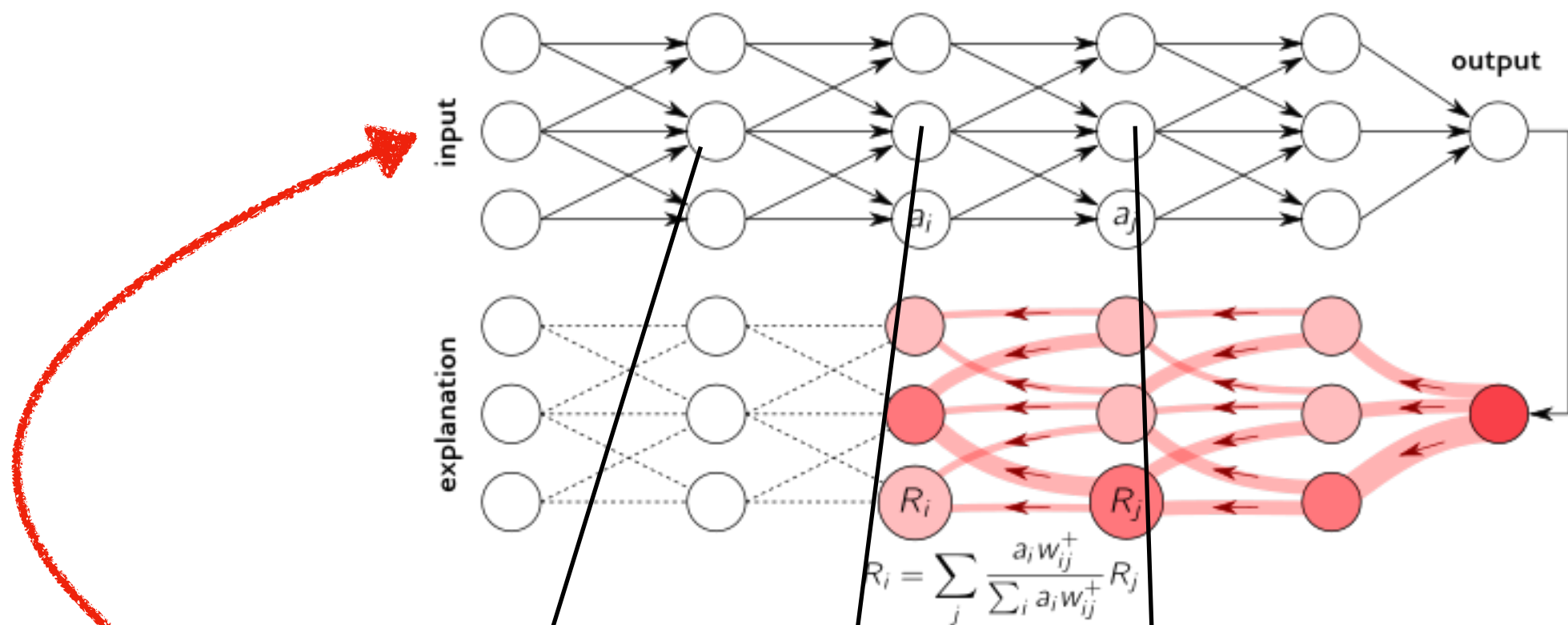
2017 Bojarski et al., Explaining How End-to-End Deep Learning Steers a Self-Driving Car



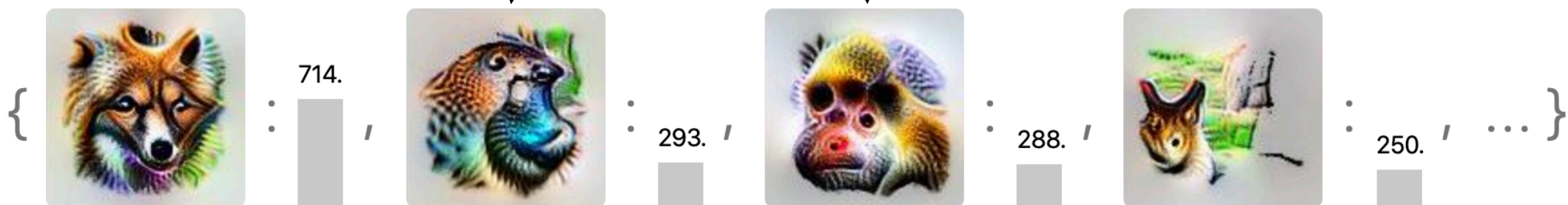
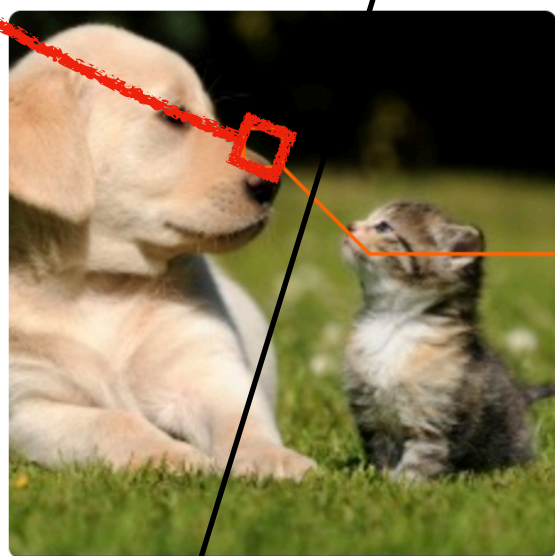
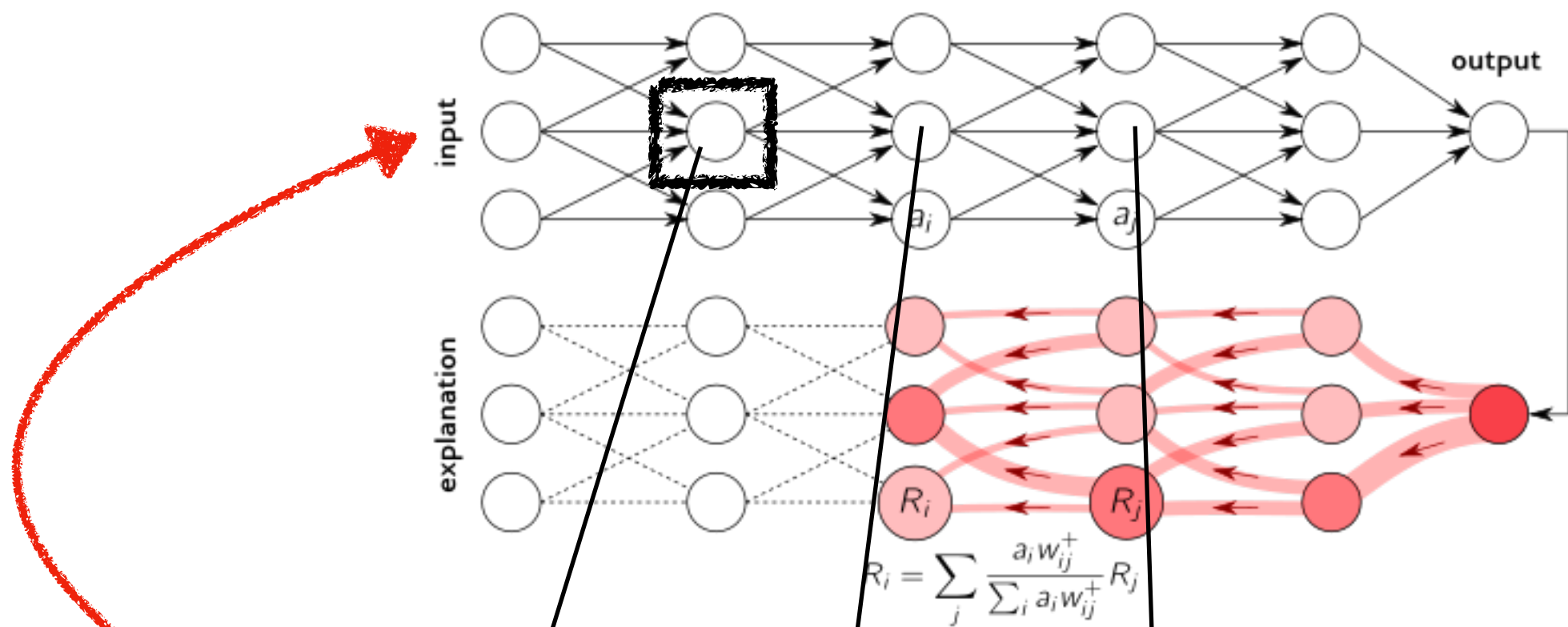


# Applying Displacement to Salient Objects, Background, and Whole Image And Measuring the Median Change in Predicted Inverse-R Across a Sample of 200 Images

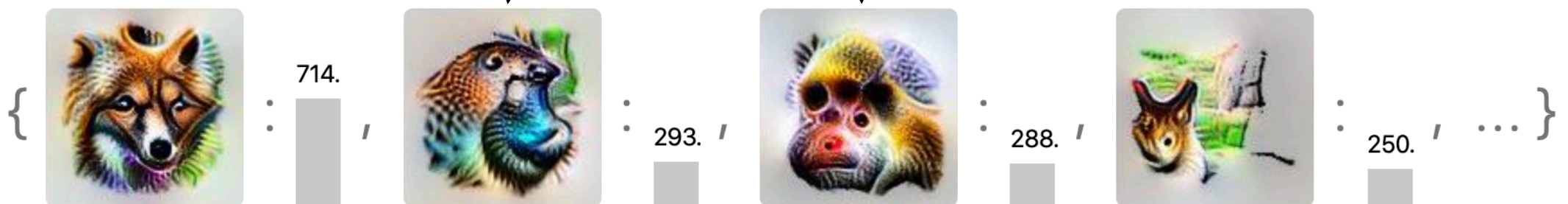
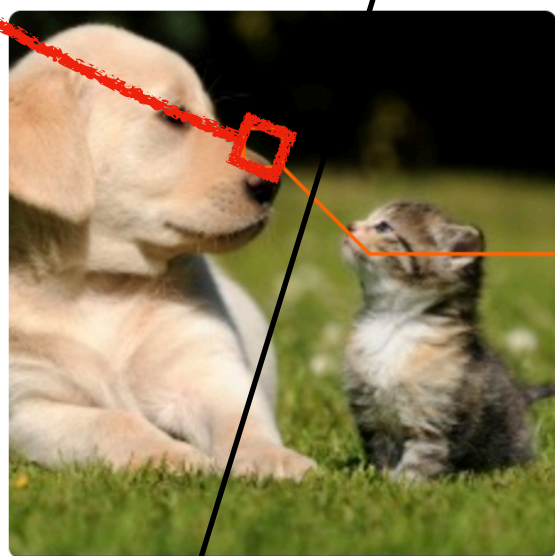
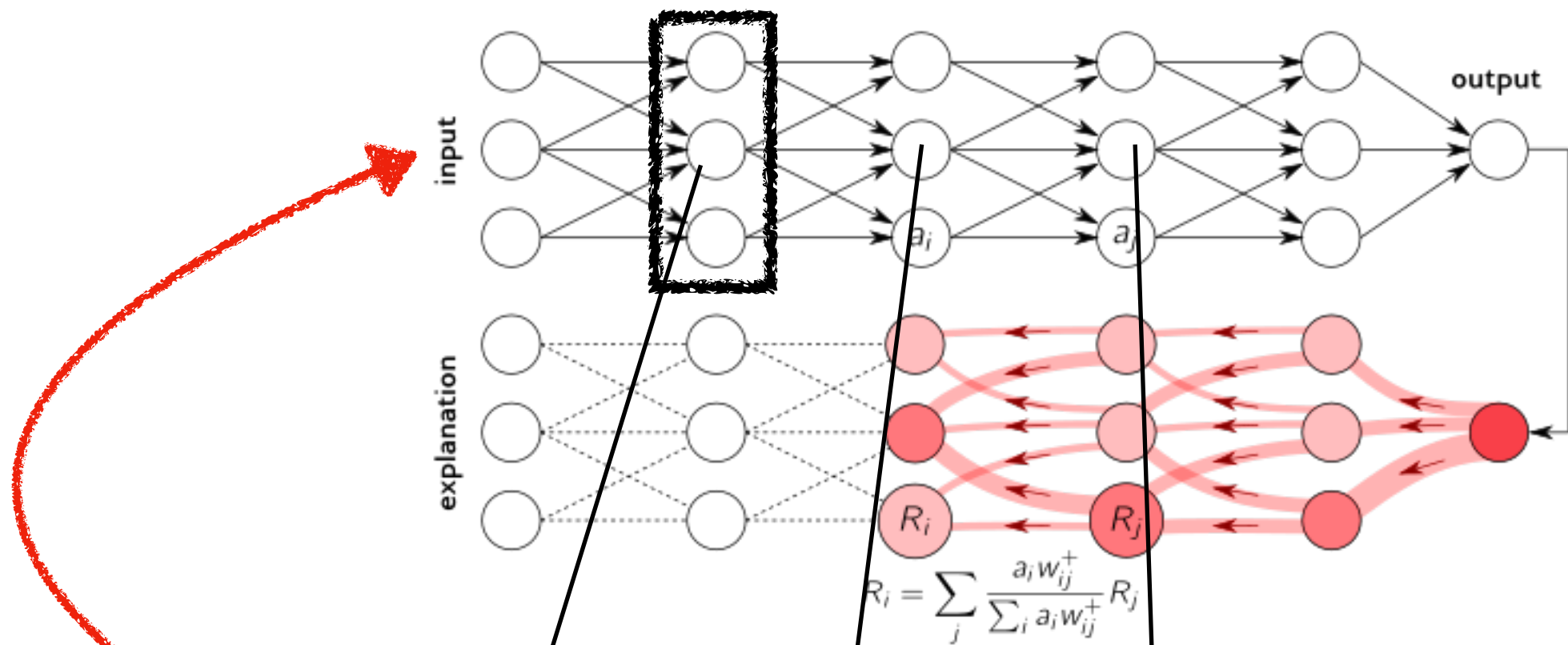


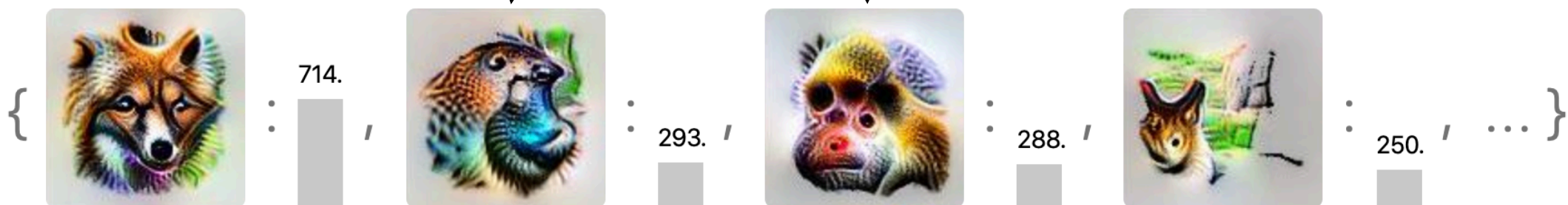
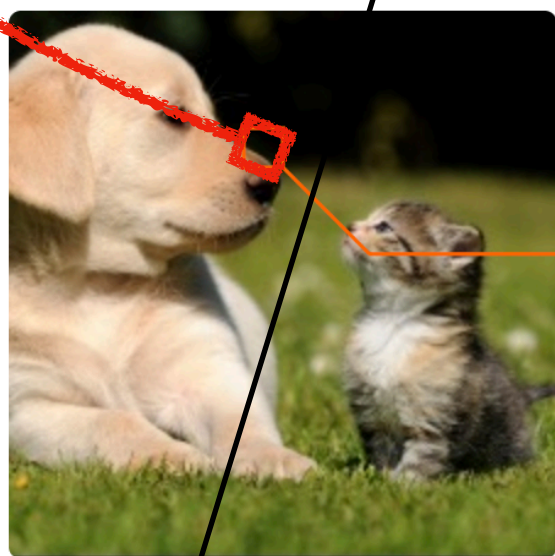
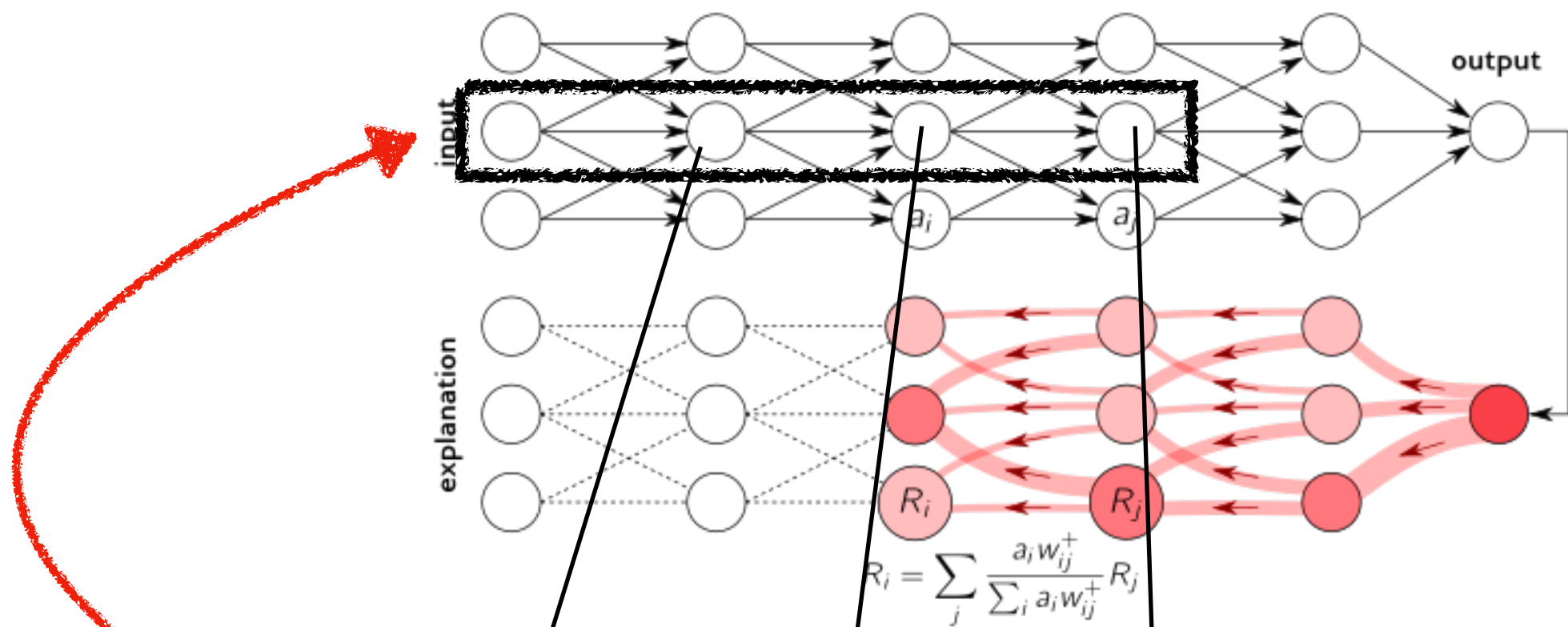


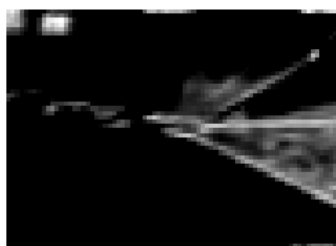
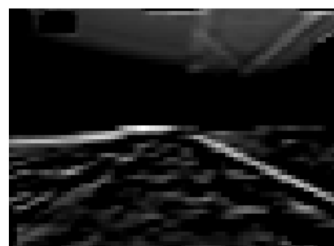












# What questions can we currently answer?

- Given **one** manually selected input:
  - On **which parts** of the input the model **focusses**? (f.e. *LRP*)
- Given **one** selected output:
  - What different **strategies (clusters) exist** for the 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*)

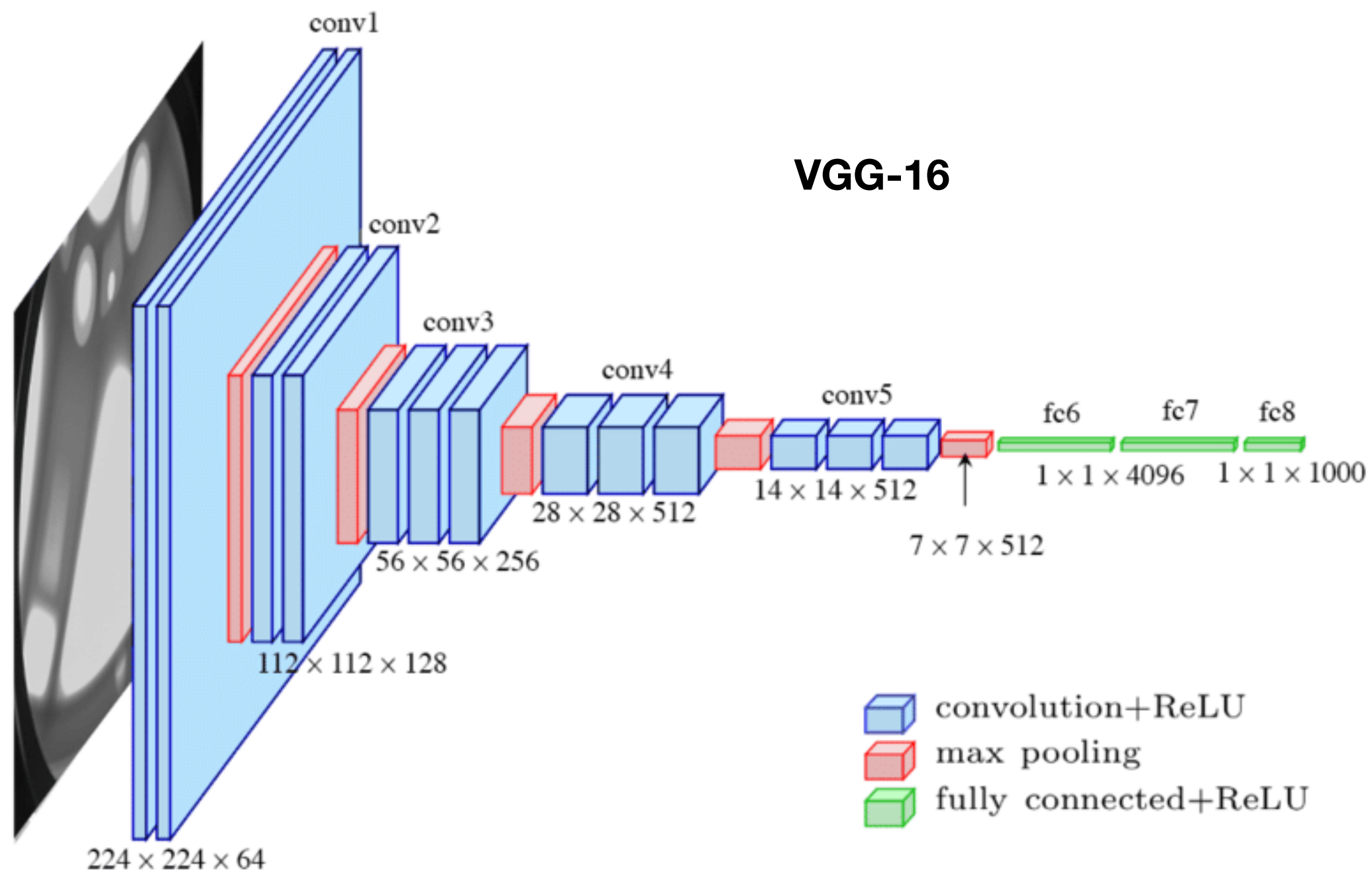
# Hands On



**[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>
- Diverse Automotive Related Datasets to play around



**We use something VGG like**

# What questions can we currently answer?

- Given **one** manually selected input:
  - On **which parts** of the input the model **focusses**?
    - 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?

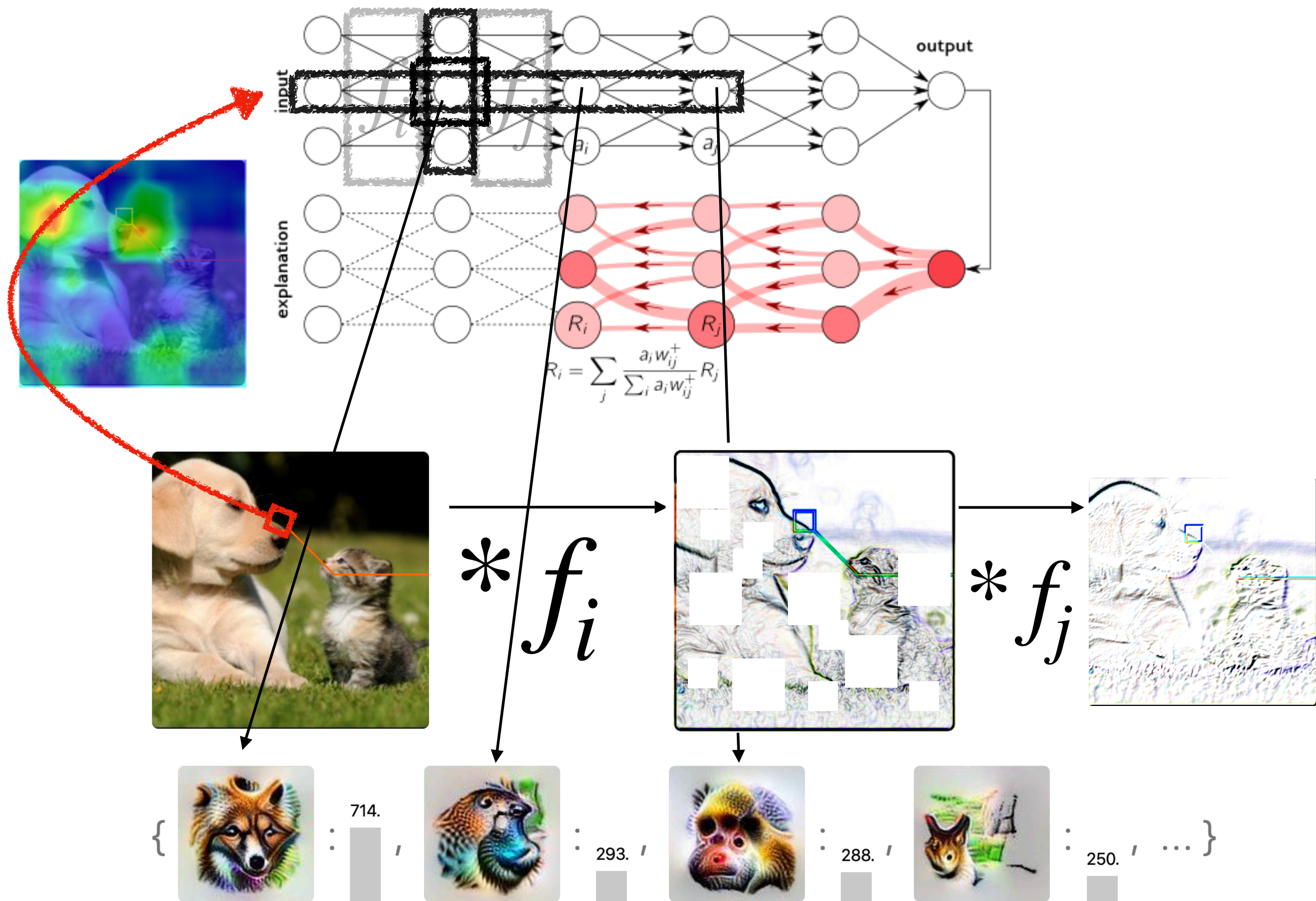
- Given **one** selected output:
  - What **kind of template** does it look for?
    - **Max Activation**, Project Lucid, Activation Atlas
    - `destill.pup`

# 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**



# Thanks for listening

I hope there was something of value for you?

We can have some Q&A in the teams channel