



Master in Computer Vision *Barcelona*

Module 3: Machine Learning for Computer Vision

Lecture: Understanding and visualizing CNNs

Lecturers: Maria Vanrell / Guillem Arias

Credits for Some Slides to: Ivet Rafegas

Motivation

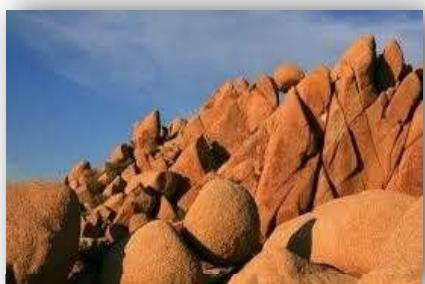
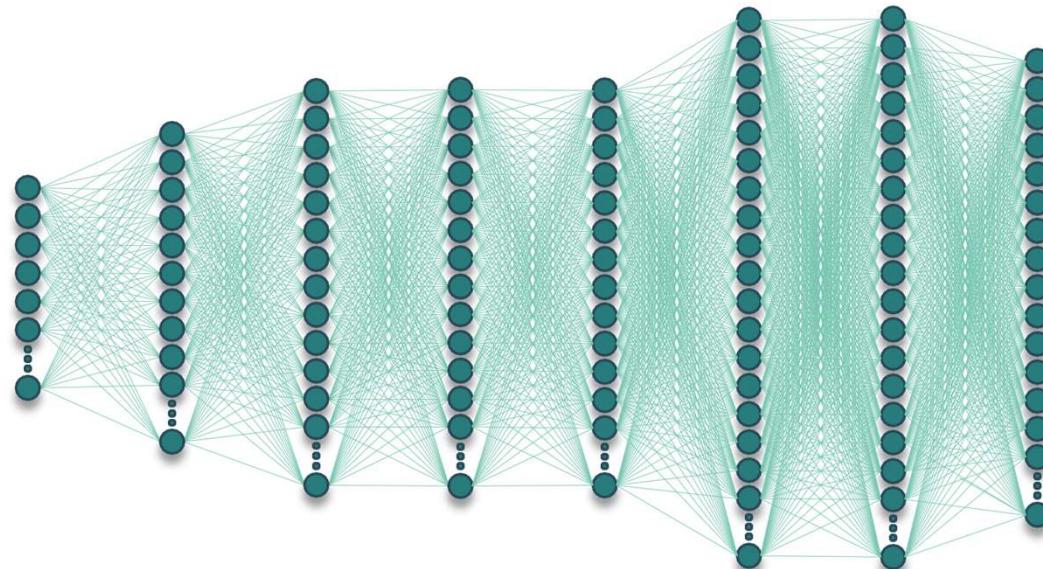
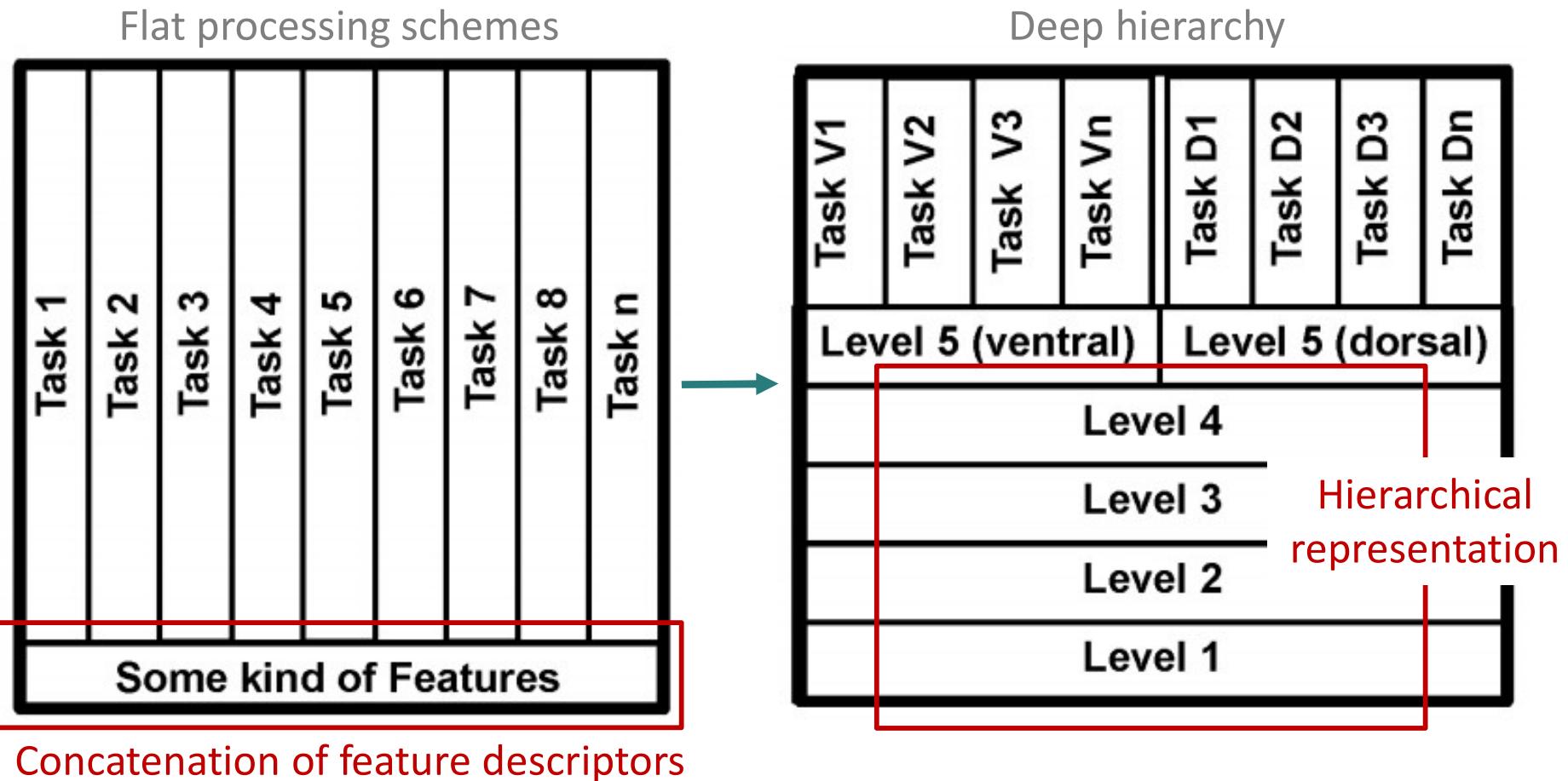


Image descriptors were designed to represent specific spatial features, such as edges at different directions, blobs, and combined with first order statistics



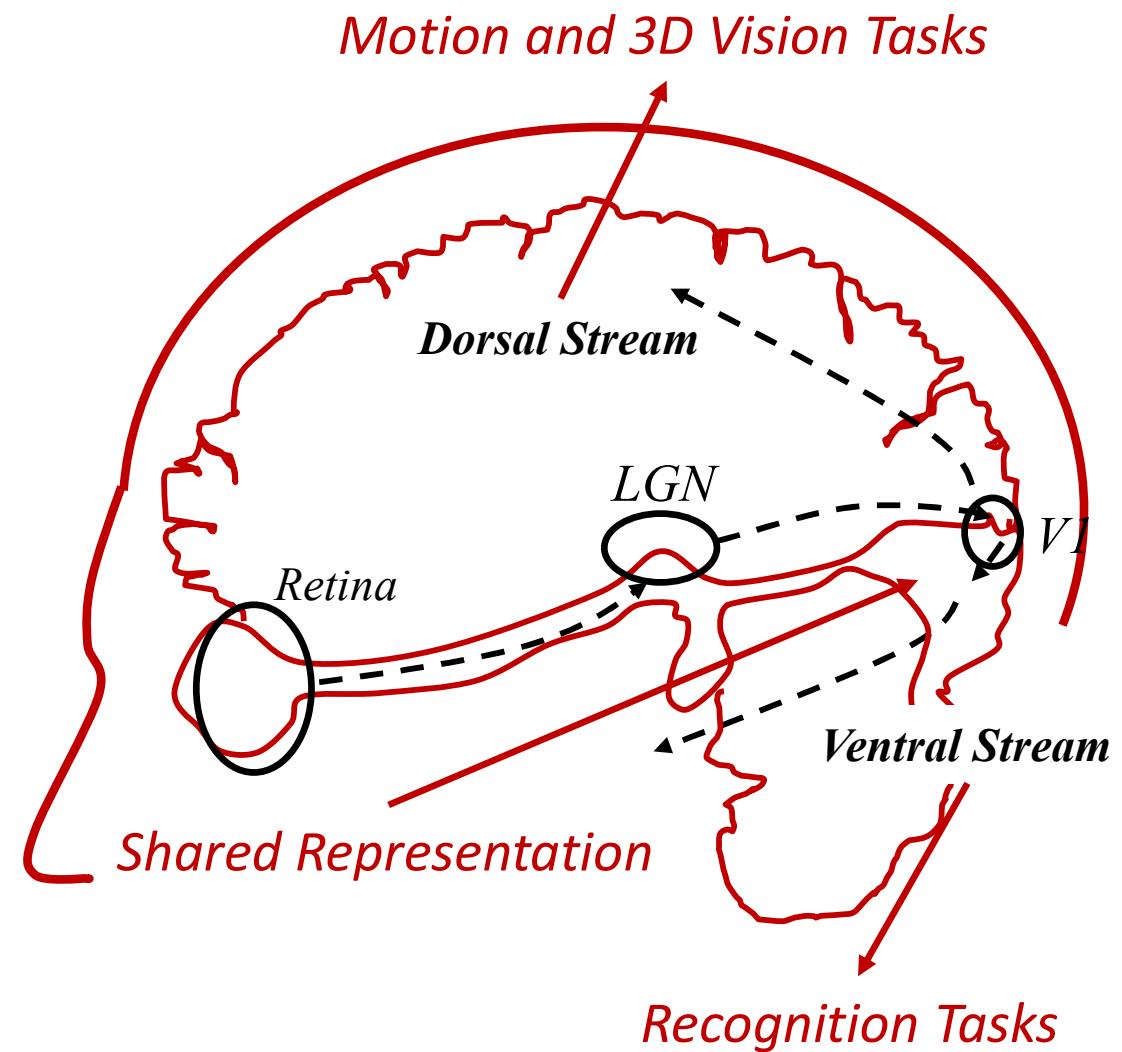
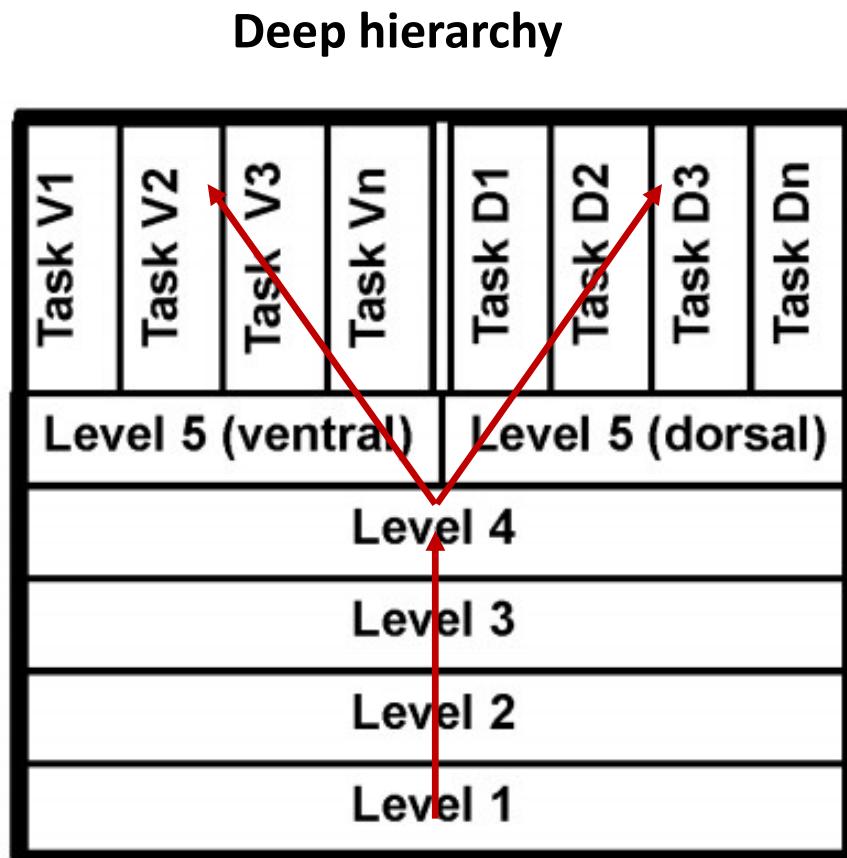
Since **2010** Convolutional Neural Networks have overcome all previous image descriptors by using a **distributed representation** of these descriptors in a **black-box of neurons**

The CNN paradigm Change According to Kruger et-al 2013



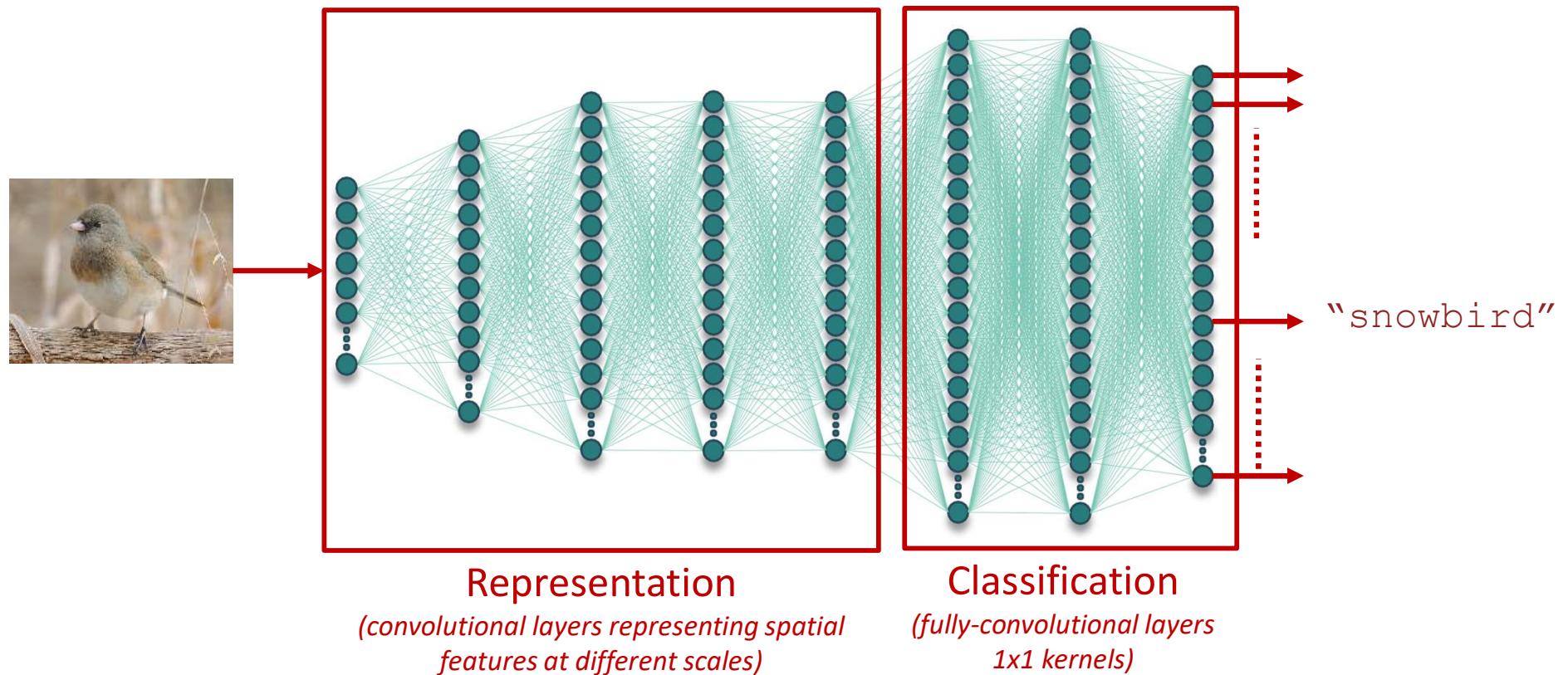
[Kruger,13] N. Kruger, P. Janseen, S. Kalkan, M. Lappe, A. Leonardis, J. Piater, A. J. Rodriguez-Sanchez, L. Wiskott.
Deep hierarchies in the primate visual cortex: What can we learn for computer vision?
IEEE Trans. Pattern Anal. Mach. Intell., 35 (8). 2013

The CNN paradigm has parallelisms with human brain



We can differentiate the Representation from the Task

Example: an object recognition network

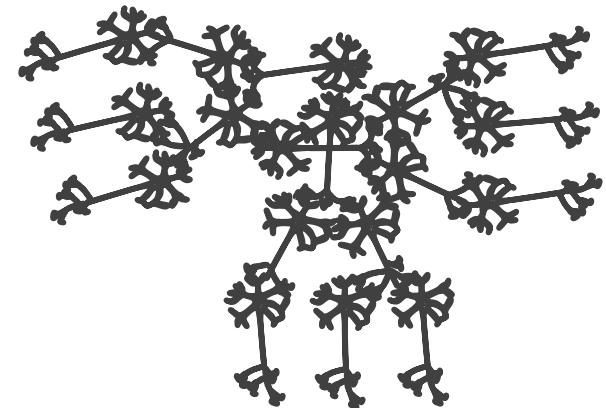
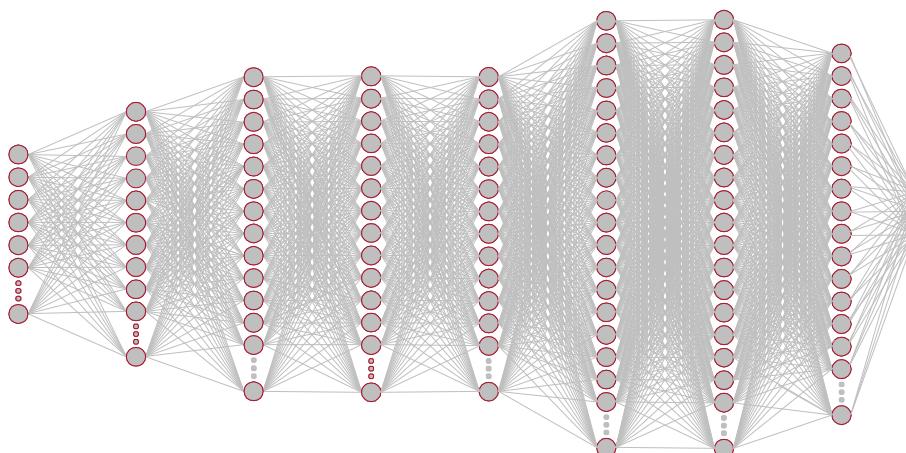


Although we have this global idea,
the network is a **black box**

Usual questions to be asked about the CNN based models:

- How a deep model concludes such a prediction ?
- Why are some features favored over others ?
- What changes are needed to improve model performance ?

If we open the box of our CNN we just see Neurons:



Like in the brain

A new field emerged in parallel with the application of DNN

EXplainable Artificial Intelligence (XAI), the objective of the XAI method is to provide an explanation for the deep learning model that is understandable by humans

Explainable artificial intelligence (xAI)

D. Gunning 2017. Technical Report, Defense Advanced Research Projects Agency (DARPA) (2017)

A Survey Of Methods For Explaining Black Box Models

Guidotti et-al 2018. ACM Computing Surveys, Vol. 51 (5)

Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI

Arrieta et-al 2020. Information Fusion. Vol. 58.

On Interpretability of Artificial Neural Networks: A Survey

Feng-Lei Fan et-al 2021. IEEE Trans. On Radiation and Plasma medical sciences, vol. 5 (6).

Explainability of artificial intelligence methods, applications and challenges: A comprehensive survey

Weiping Ding et-al. 2022. Information Sciences 615.

[...] the findings show that there is no common consensus on how an explanation must be expressed, and how its quality and dependability should be evaluated [...].

We will use this survey as a basis for this lecture:

On Interpretability of Artificial Neural Networks: A Survey

Feng-Lei Fan, *Student Member, IEEE*, Jinjun Xiong, *Senior Member, IEEE*, Mengzhou Li, *Student Member, IEEE*,
and Ge Wang, *Fellow, IEEE*

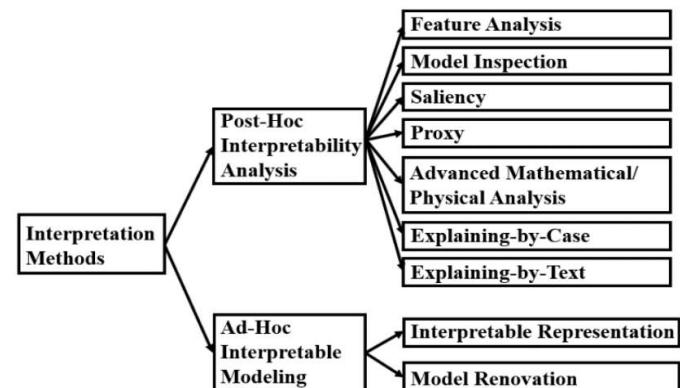
Abstract— Deep learning as represented by the artificial deep neural networks (DNNs) has achieved great success in many important areas that deal with text, images, videos, graphs, and so on. However, the black-box nature of DNNs has become one of the primary obstacles for their wide acceptance in mission-critical applications such as medical diagnosis and therapy. Due to the huge potential of deep learning, interpreting neural networks has recently attracted much research attention. In this paper, based on our comprehensive taxonomy, we systematically review recent studies in understanding the mechanism of neural networks, describe applications of interpretability especially in medicine, and discuss future directions of interpretability

provides in-depth perspectives but is limited in scope. For example, only 49 references are cited there. The review of M. Du *et al.* (2018) has a similar weakness, only covering 40 papers which are divided into post-hoc and ad-hoc explanations, as well as global and local interpretations. Their taxonomy is coarse-grained and neglects a number of important publications, such as *explaining-by-text*, *explaining-by-case*, etc. In contrast, our review is more detailed and comprehensive, which includes the latest results. While publications in L. H. Gilpin *et al.* (2018) are classified into understanding the workflow of a neural network.

IEEE Trans. On Radiation and Plasma medical sciences, vol. 5 (6).

<https://arxiv.org/vc/arxiv/papers/2001/2001.02522v2.pdf>

Proposed Taxonomy:



Proposed Taxonomy:

Post-hoc analysis
To explain existing models

Interpretation Methods

Ad-hoc modelling
To build explainable models

Post-Hoc Interpretability Analysis

Ad-Hoc Interpretable Modeling

Feature Analysis

Model Inspection

Saliency

Proxy

Advanced Mathematical/
Physical Analysis

Explaining-by-Case

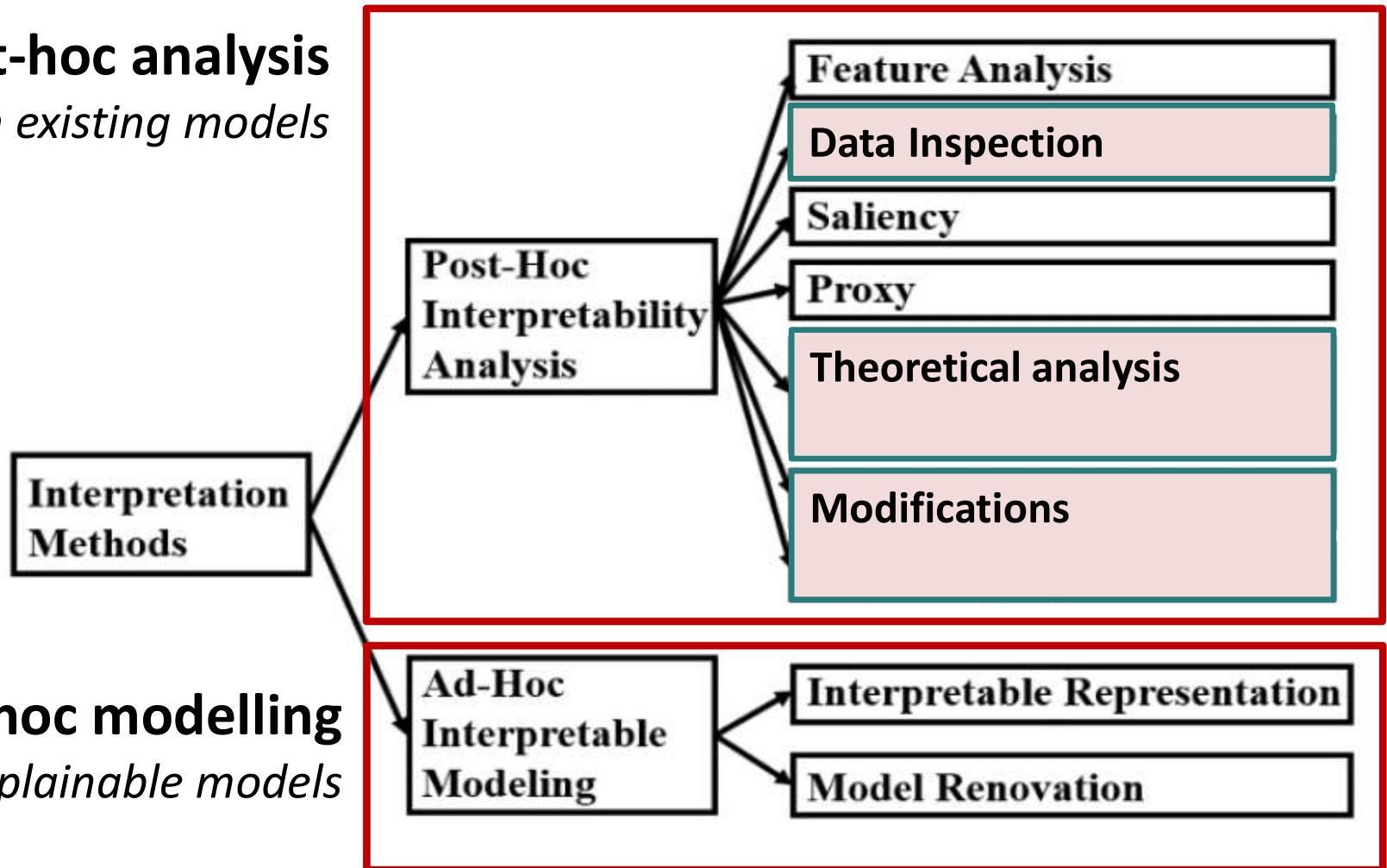
Explaining-by-Text

Interpretable Representation

Model Renovation

Proposed Taxonomy:

Post-hoc analysis
To explain existing models



Ad-hoc modelling
To build explainable models

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis
- Interpretable representation
- Model Renovation

Ad-hoc modelling

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Preliminary considerations

1. What is a neuron?
2. What information can we use to characterize neurons?

What is a neuron?

We already know some concepts related to neurons,

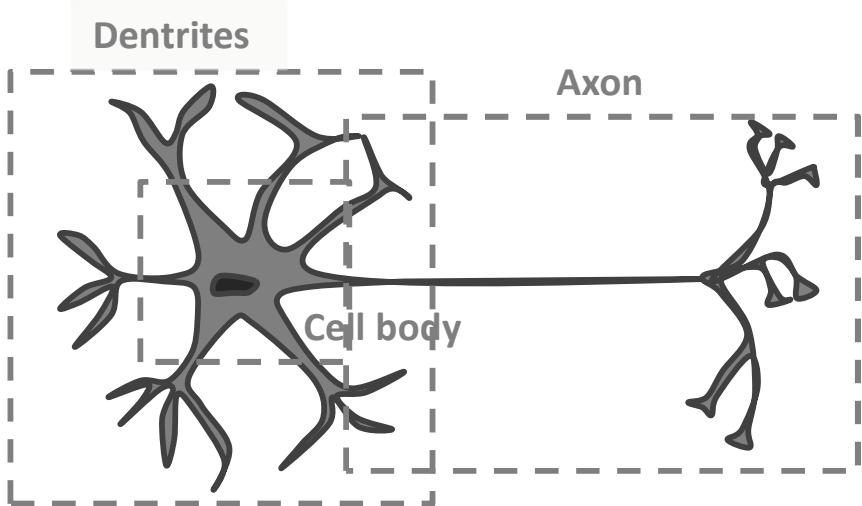
- **Neurons or Units** are the basis of the layers that a CNN relies on
- They are the main responsible for the ability of representing the input image in **high-dimensional feature spaces**
- Each neuron is **associated to some weights** that has been learnt during the CNN training step
- Neurons **give activation responses** depending on their inputs and their weights

So,

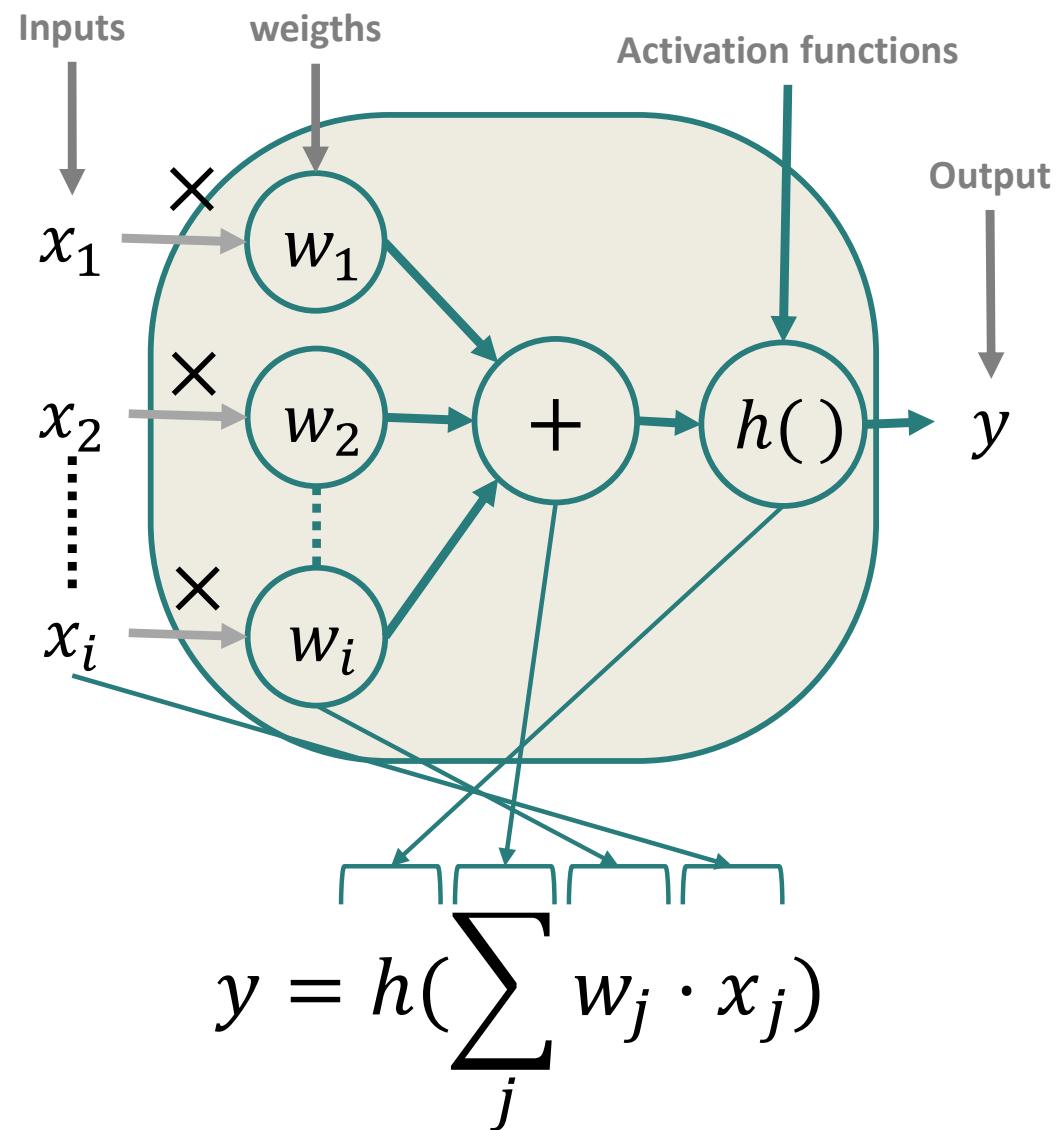
what these weights or activations are explaining about the neuron?

Usual parallelisms between biological and computational neurons

Real Neuron



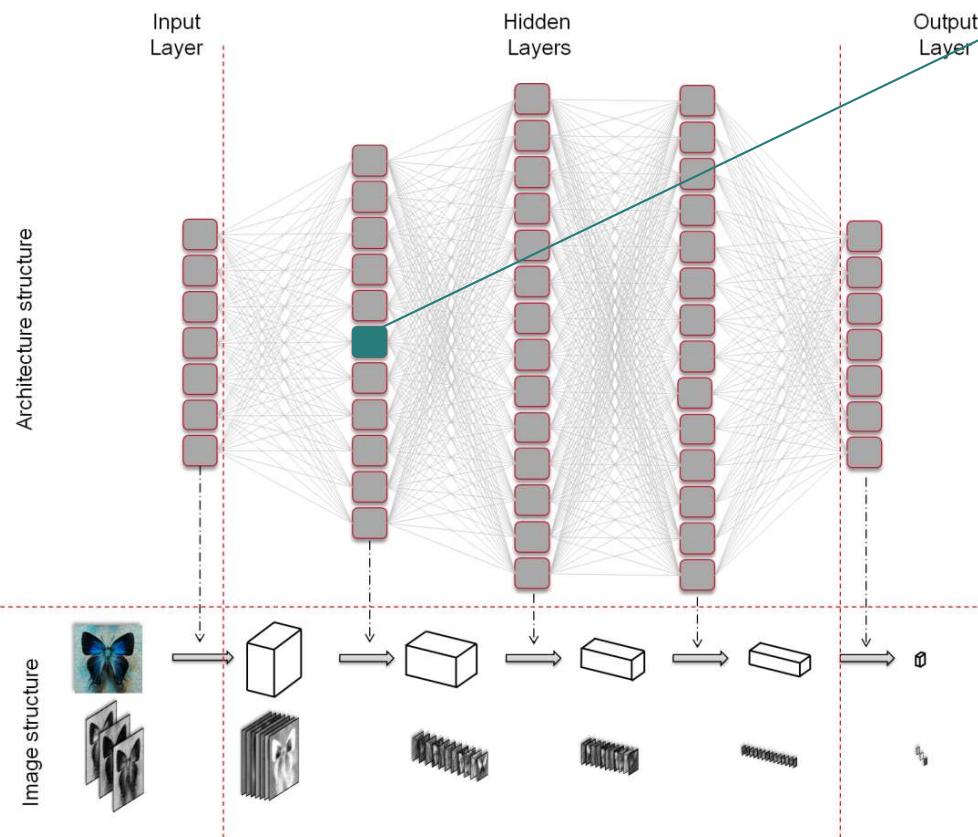
Artificial Neuron



When neurons are grouped in layers ...

All **neuron outputs** are grouped in large tensors

Usual questions about individual neurons:



Which feature is this neuron selecting from the input image?

Which is the task of this neuron within the global CNN task?

In summary, what characterizes this neurons?

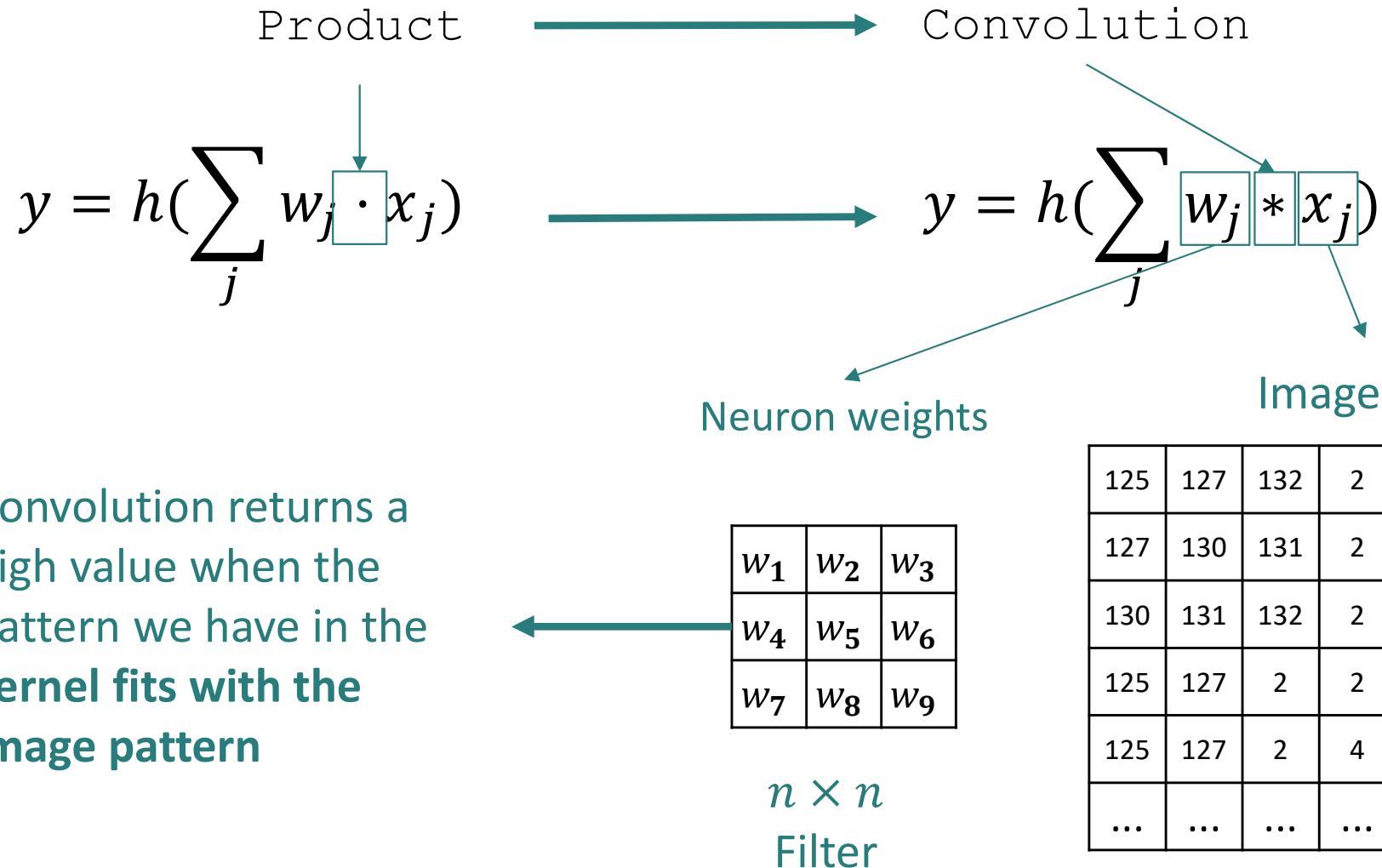
Preliminary considerations :

1. What is a neuron?
2. What information can we use to characterize neurons?

Let's take a look to:

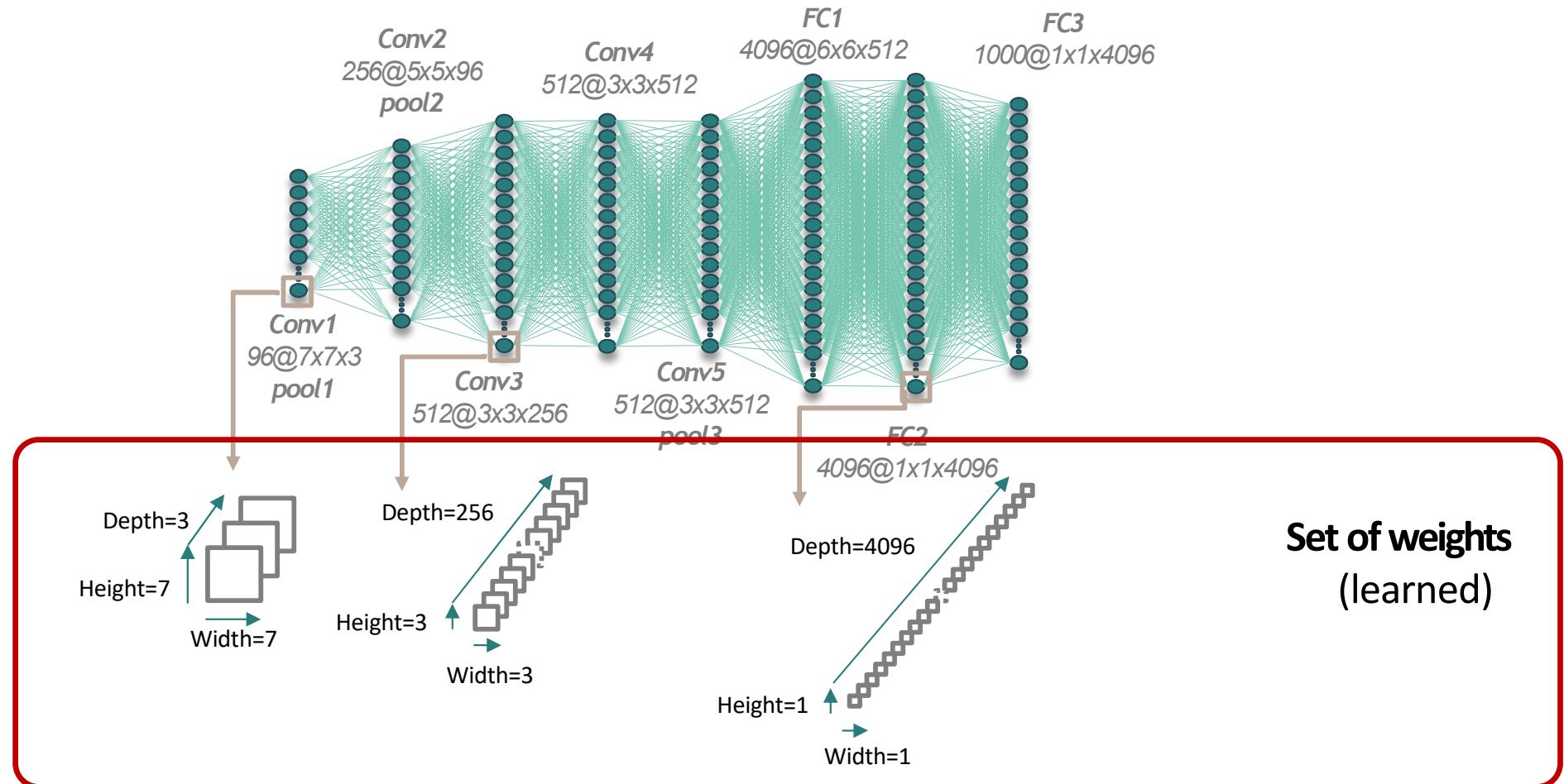
- Neuron weights
- Neuron outputs

In Convolutional Neural Networks



How many weights for each Neuron?

we have a **3D tensor of weights** in the 1st layer neurons



Size (height,width) encode the spatial dimension of the convolution kernels

Depth of a neuron is fixed by the number of neurons (channels) in the previous convolutional layer.

Preliminary definitions :

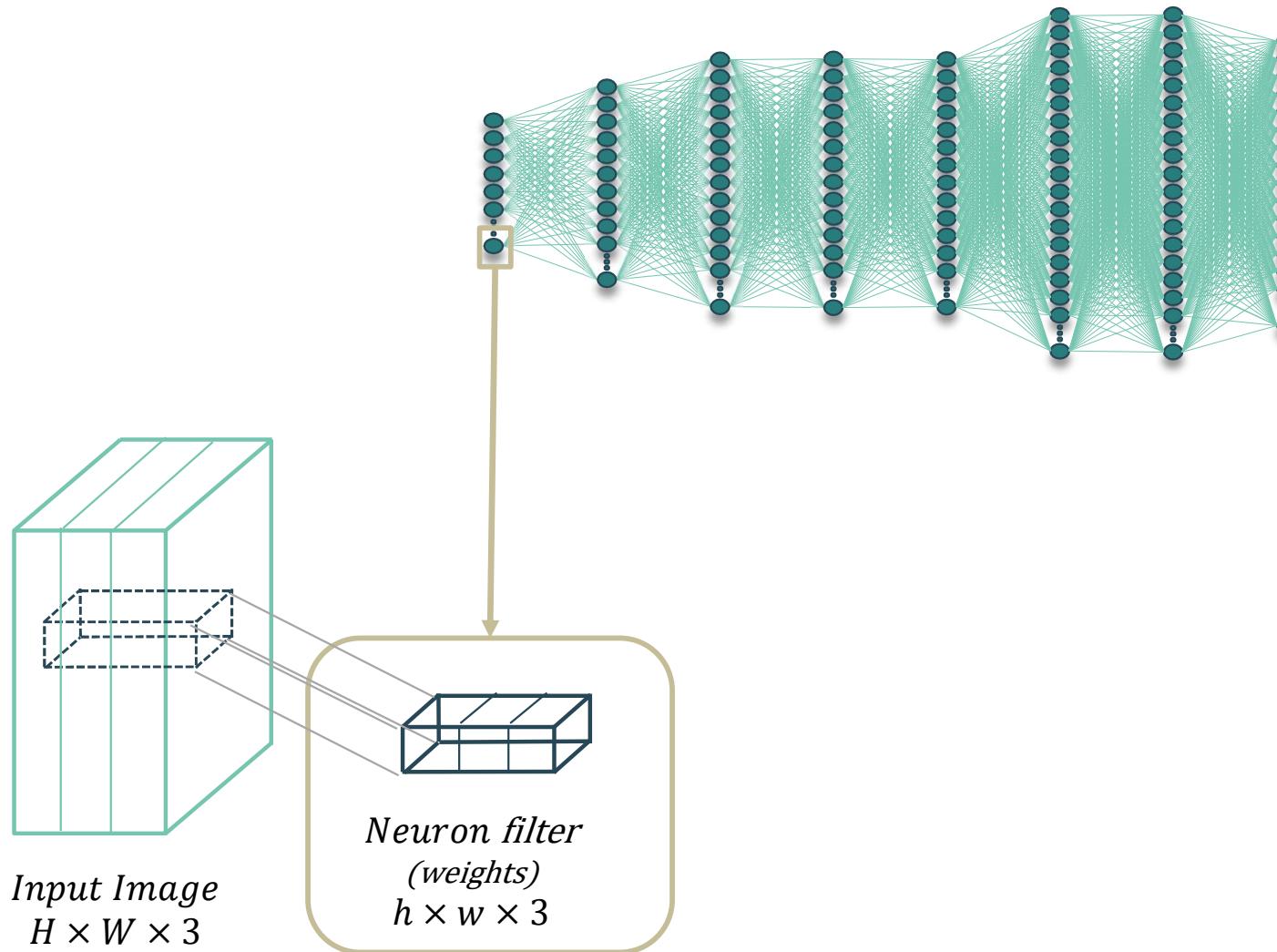
1. What is a neuron?
2. What information can we use to characterize neurons?

Let's take a look to:

- Neuron weights
- Neuron outputs

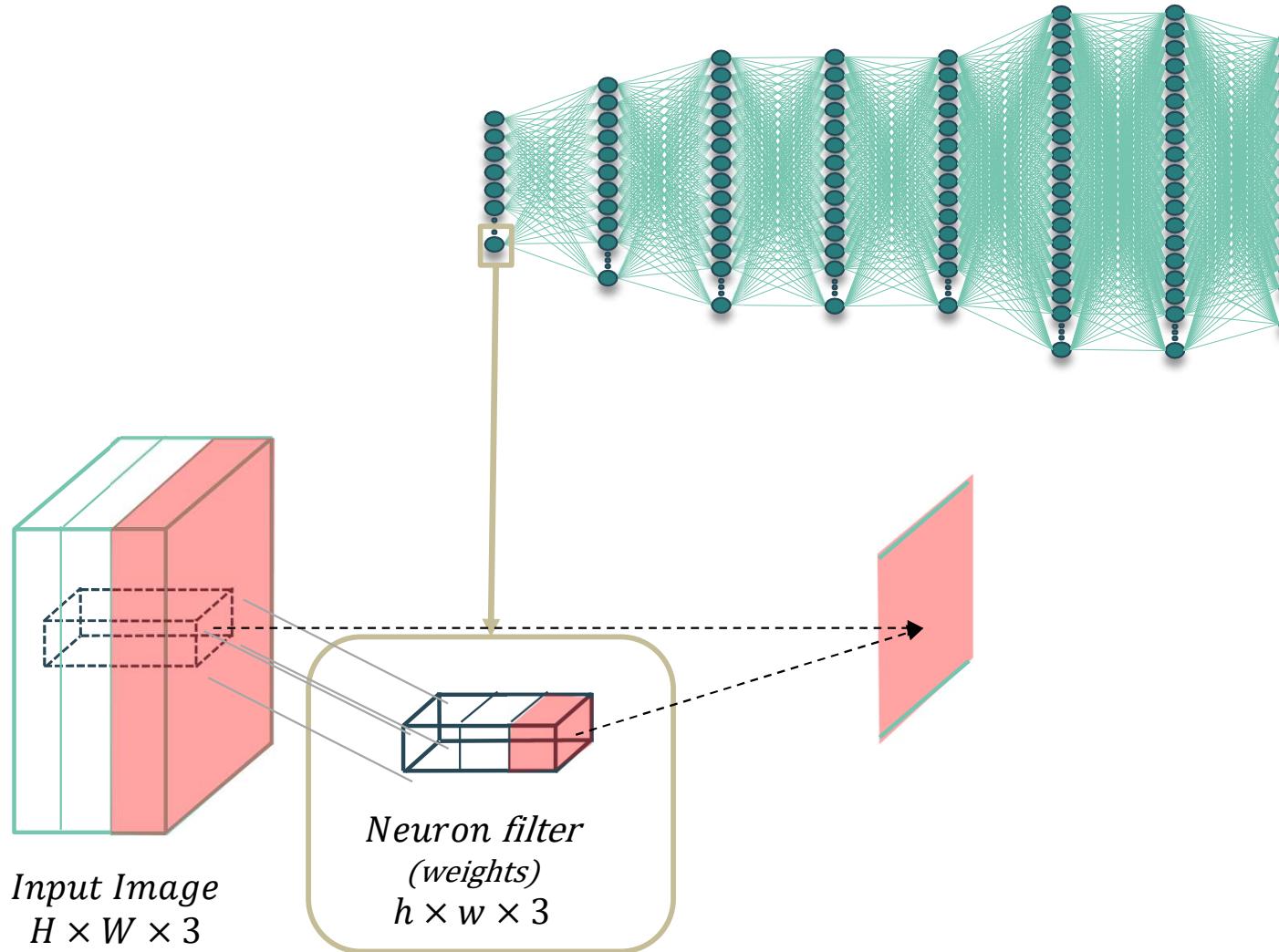
How is the Neuron Output?

each neuron has an activation for every image pixel



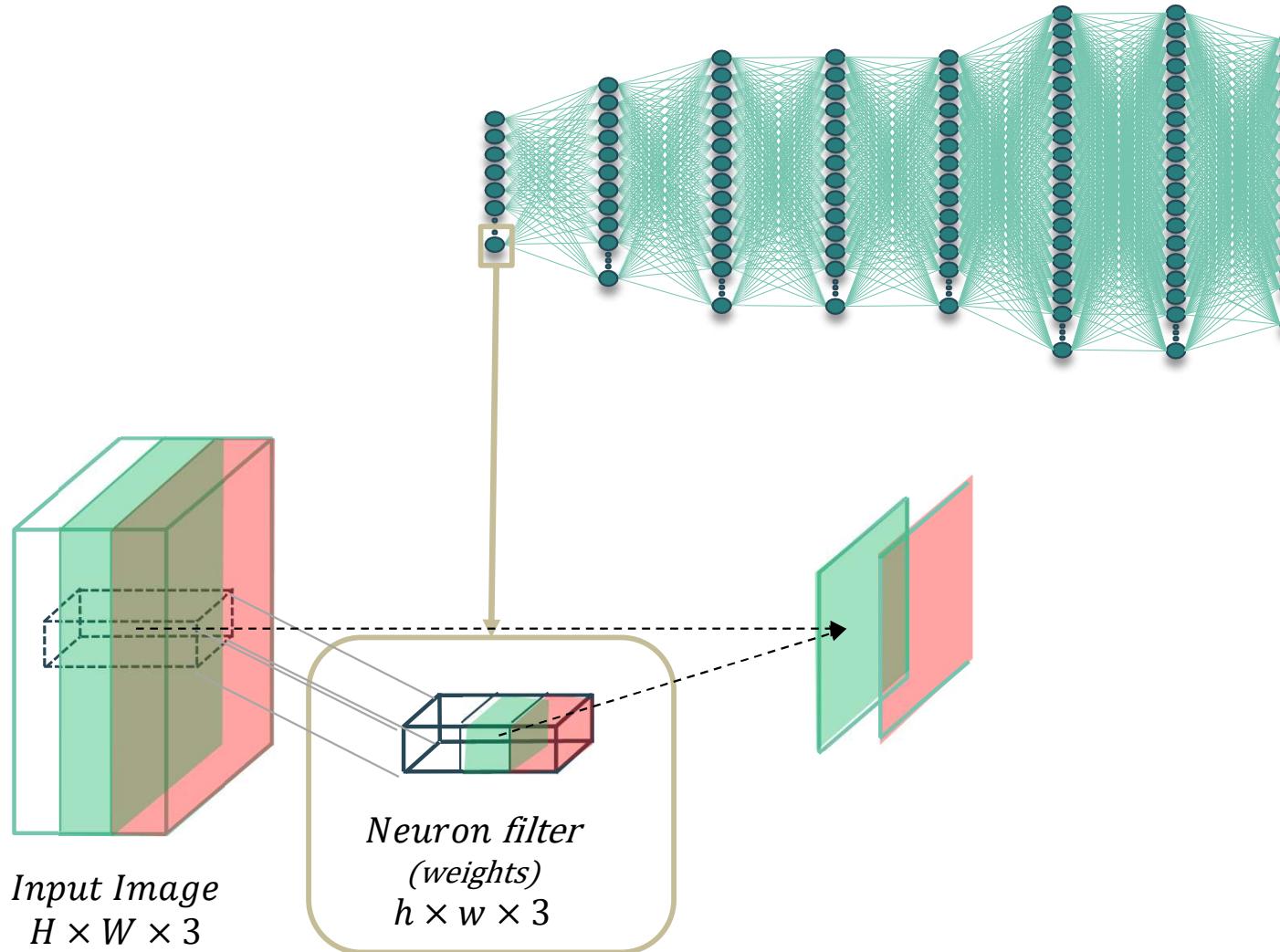
What is the Neuron Output?

each neuron has an activation for every image pixel



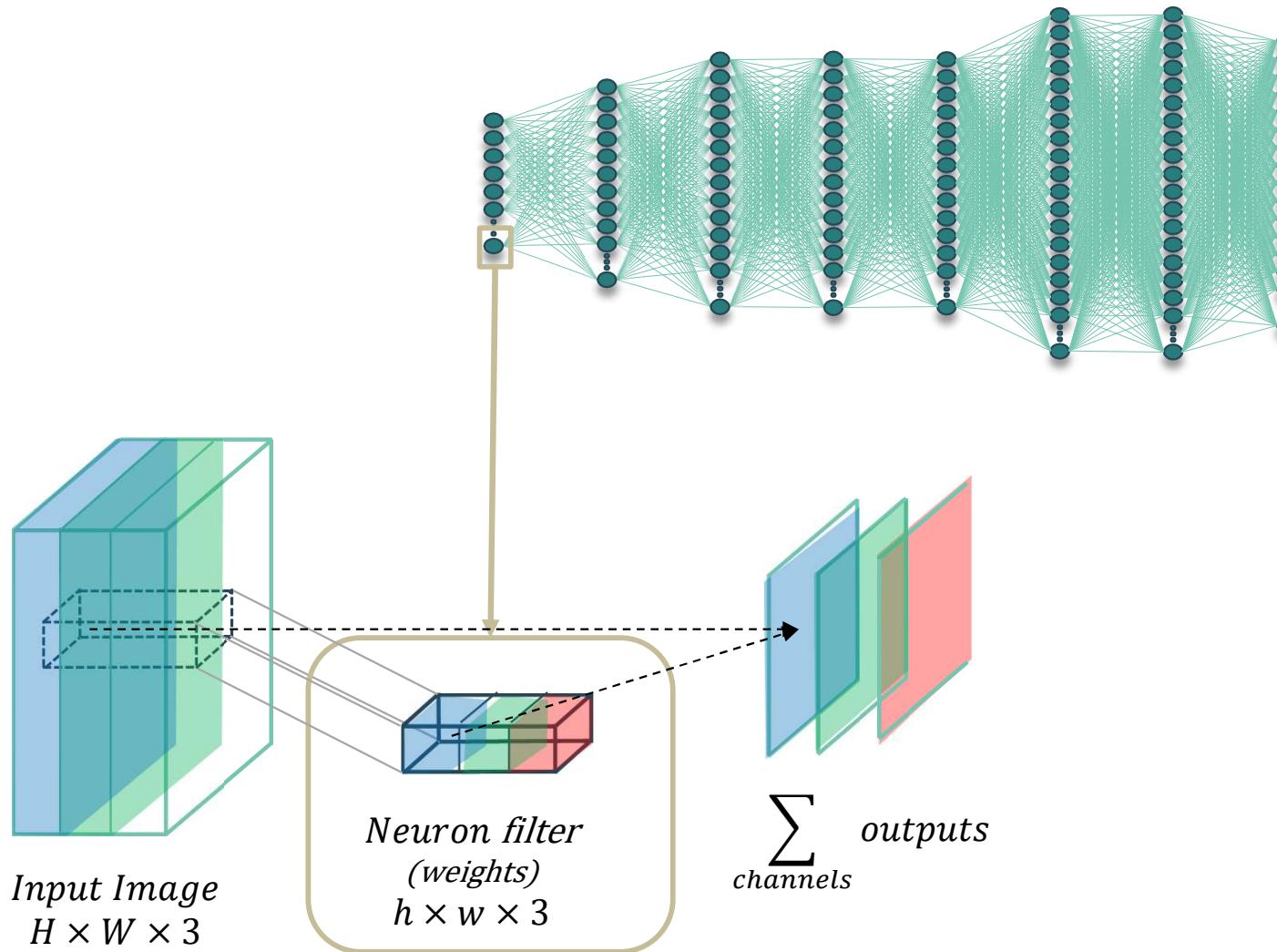
What is the Neuron Output?

each neuron has an activation for every image pixel



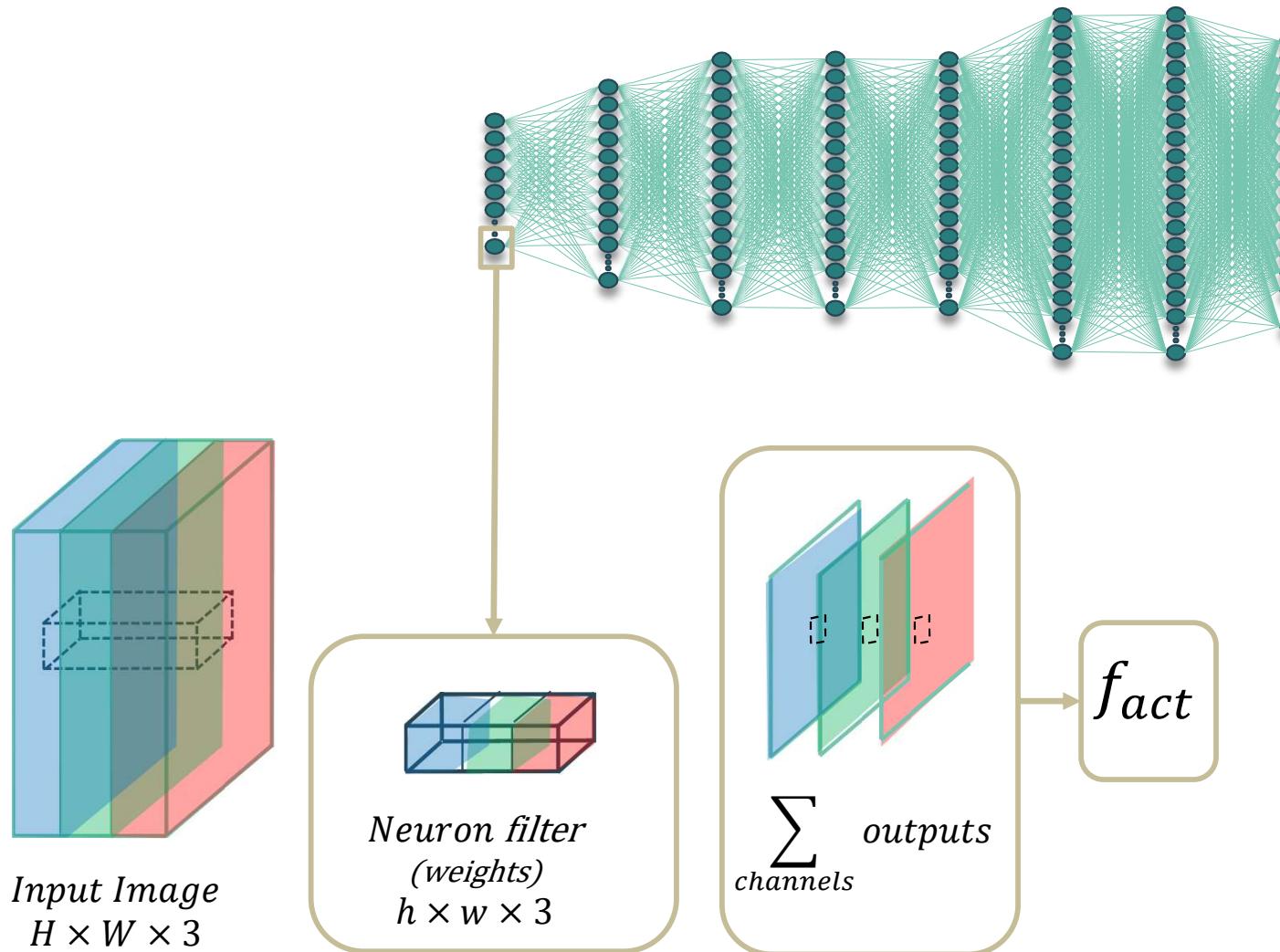
What is the Neuron Output?

each neuron has an activation for every image pixel



What is the Neuron Output?

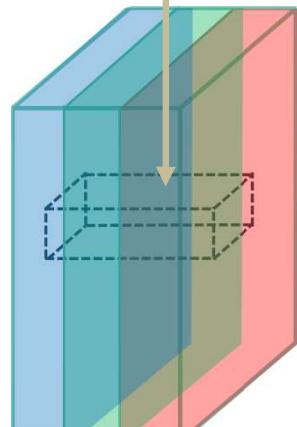
each neuron has an activation for every image pixel



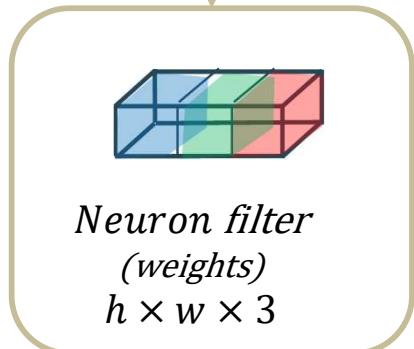
What is the Neuron Output?

each neuron has an **activation map** for every receptive field

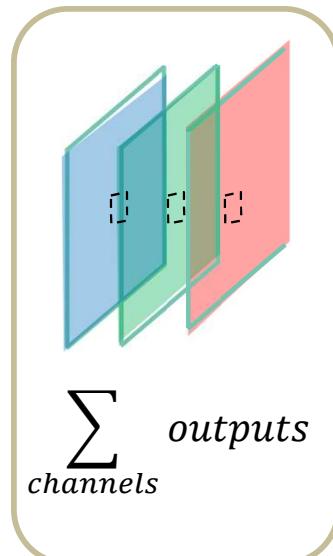
Image receptive field: is the image patch that provokes a specific activation of a neuron for a given image.



Input Image
 $H \times W \times 3$



*Neuron filter
(weights)
 $h \times w \times 3$*

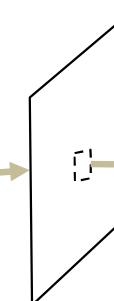


\sum_{channels} outputs

Neuron Activation map

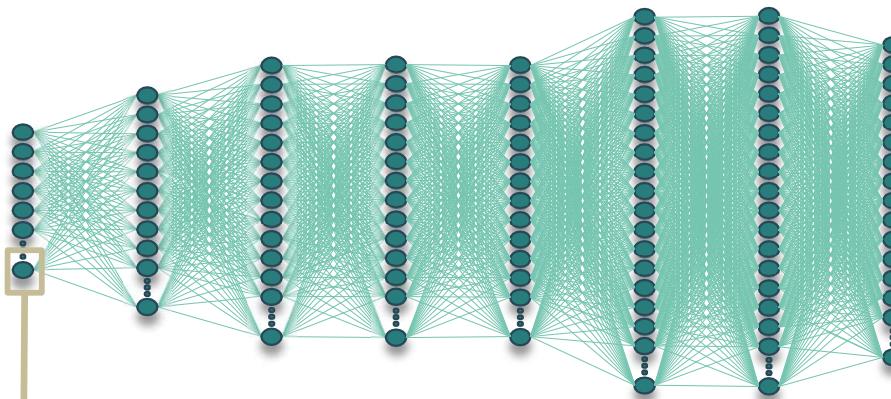
f_{act}

*Single channel output
 $(H - h + 1) \times (W - w + 1)$*



The feature encoded by the neuron has been found in the **image receptive field**

Pattern Matching interpretation in Representation Layers



High activation

Considering this Pattern Matching interpretation,

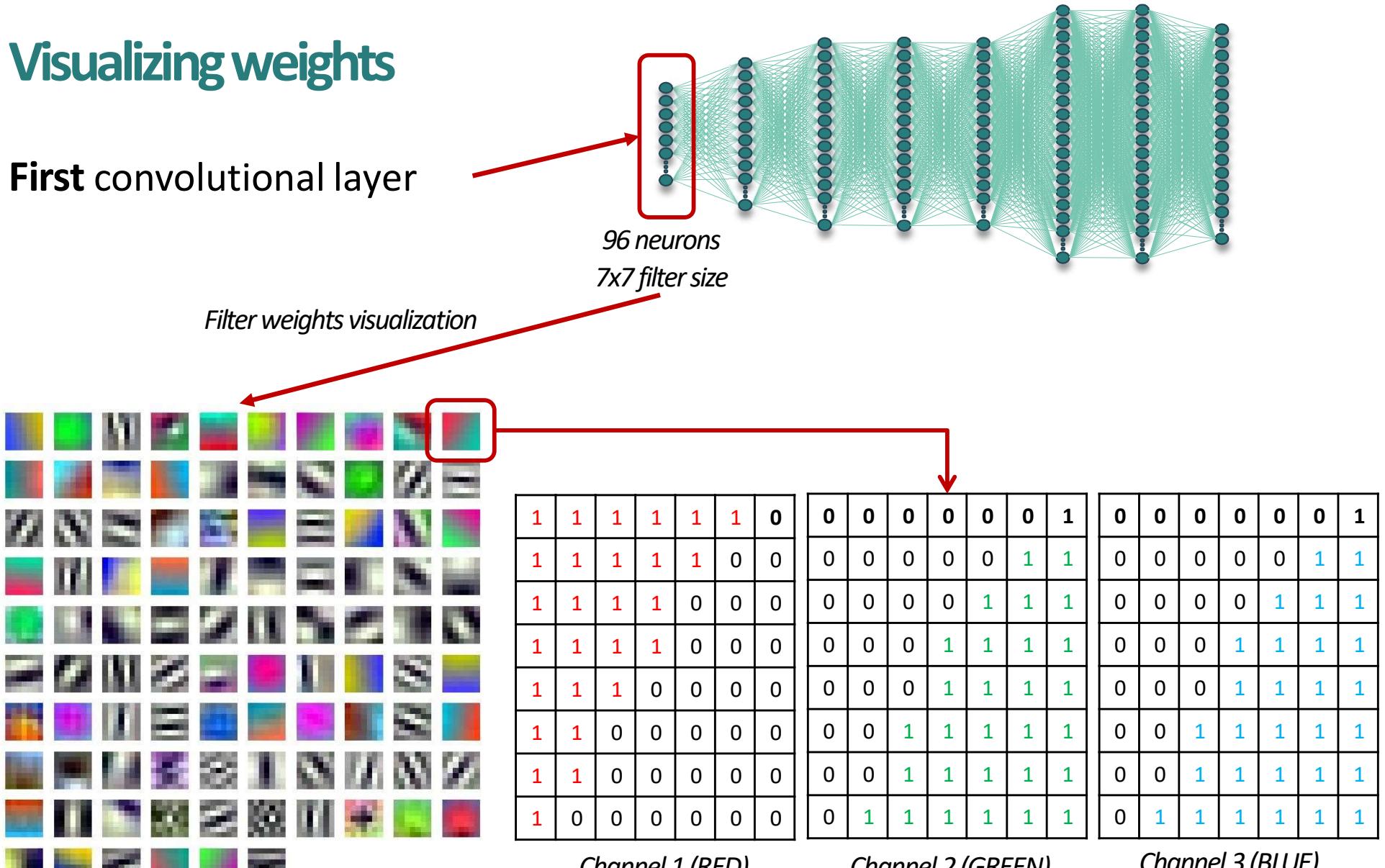
Idea: The feature encoded by the filter associated to a neuron has been found in the image

*if we visualize the weights,
we can visualize what is detected in the image when a
neuron activates*

Solution: Let's visualize the neuron weights

Visualizing weights

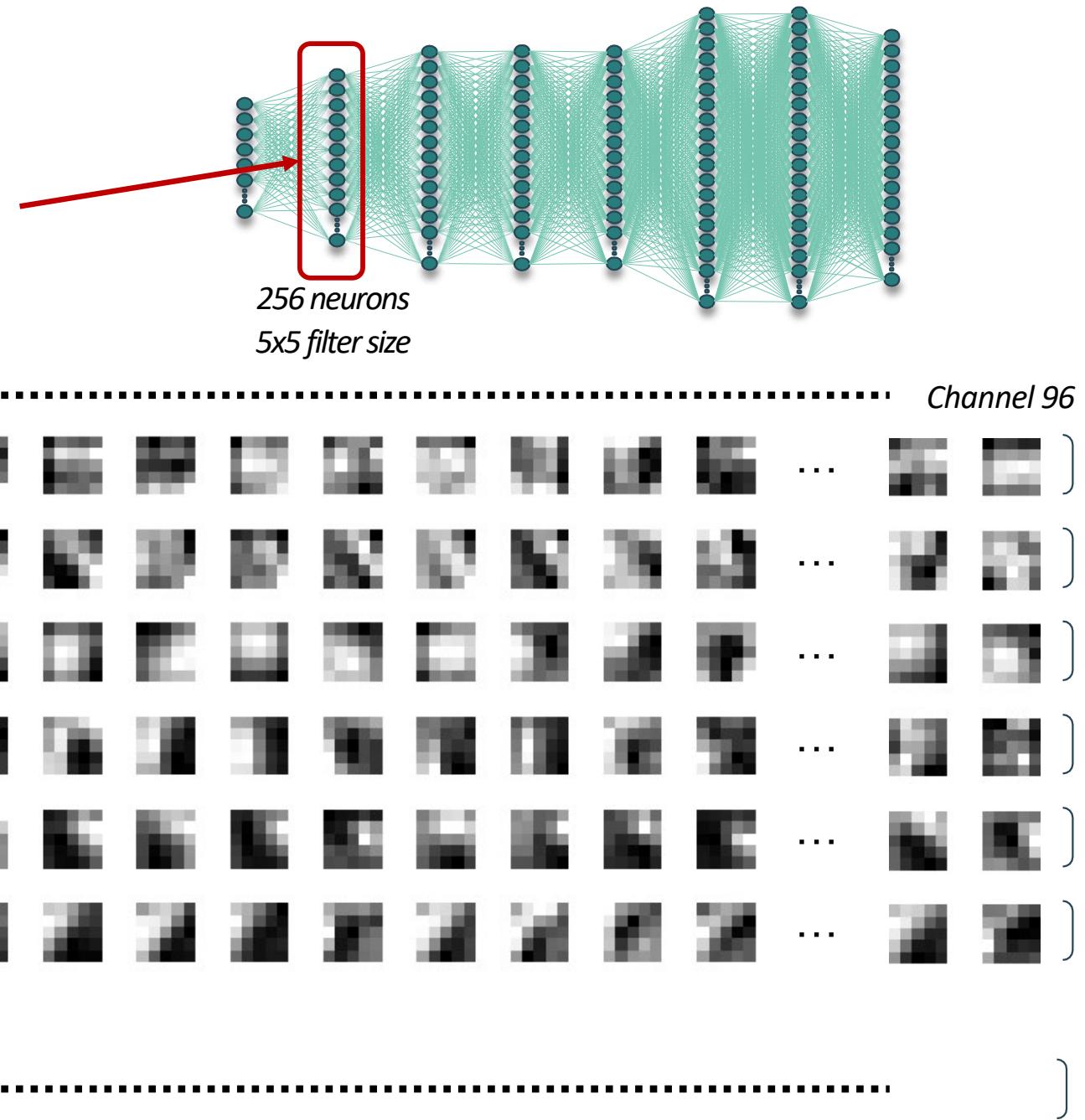
First convolutional layer



*Interpretation: Red-Cyan 45° Edge Detector
"Convolution = Pattern Matching on RGB image space"*

Visualizing weights

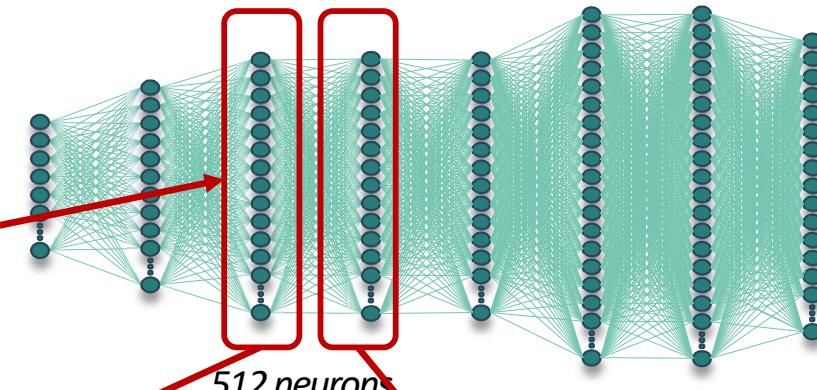
2nd convolutional layer



Interpretation? , we can not plot an image of 96 channels!!

Visualizing weights

3rd and 4th convolutional layers



Channel 1 Channel 256

Filter 1 [...]

Filter 2 [...]

Filter 3 [...]

Filter 512 [.....]

Channel 1 Channel 512

Filter 1 [...]

Filter 2 [...]

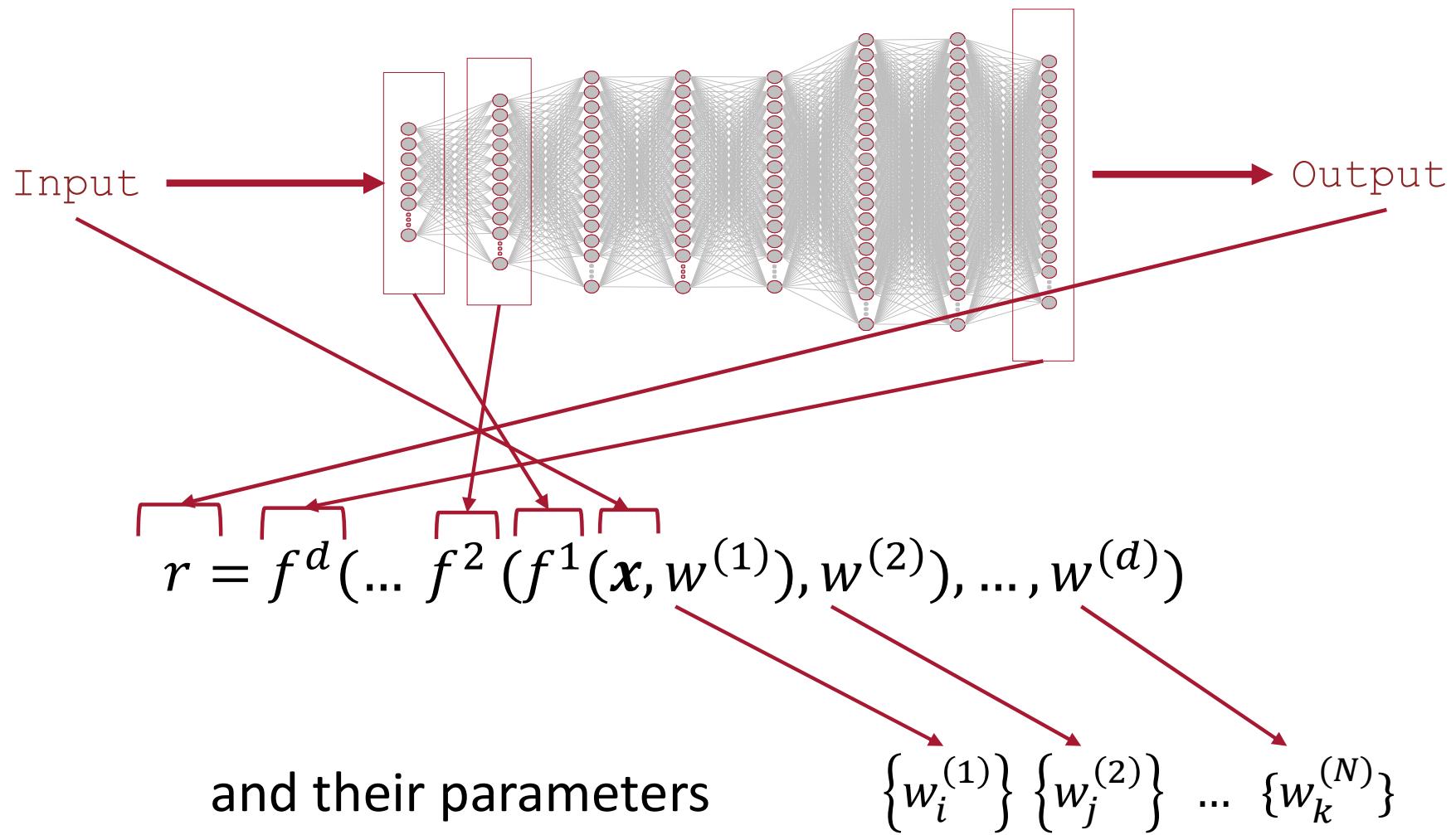
Filter 3 [...]

Filter 512 [.....]

*Let's try to see the meaning
of these weights*

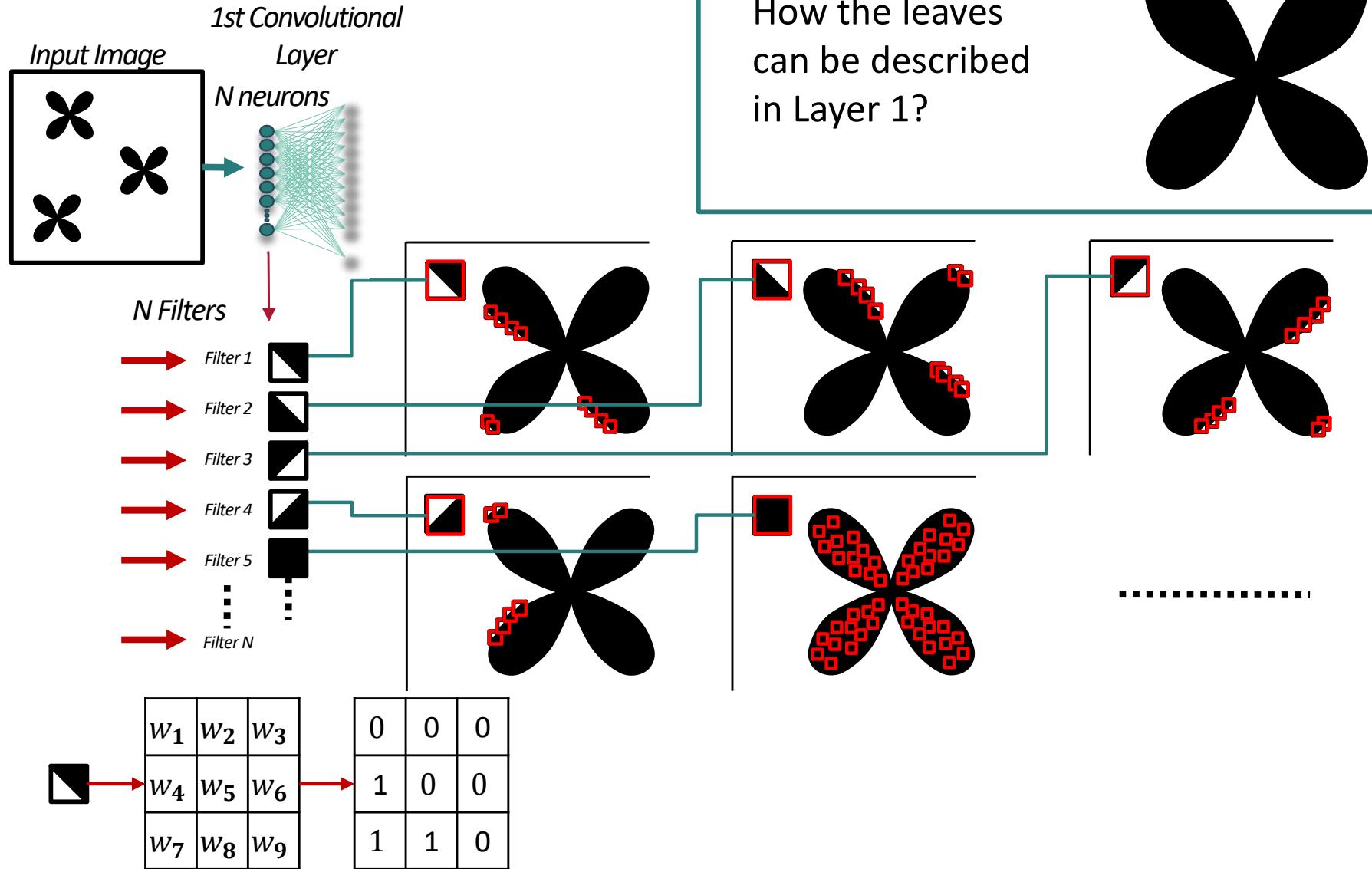
Understanding the weights

means to understand the composition of multiple layers



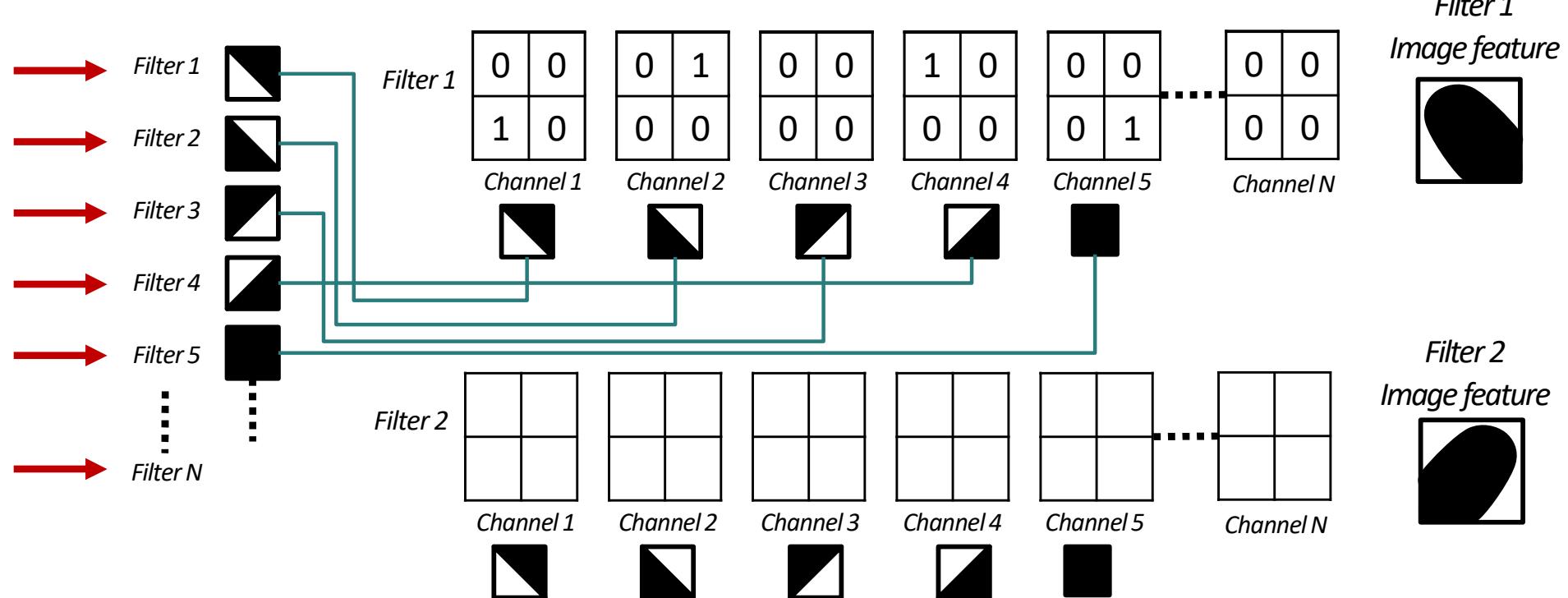
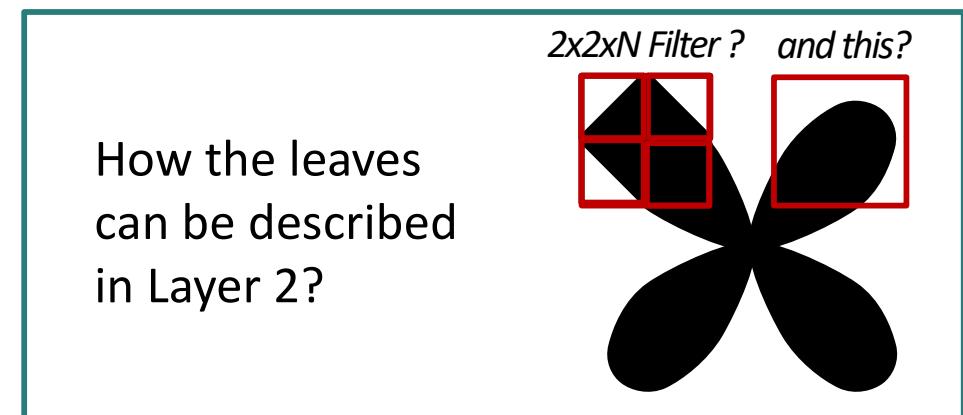
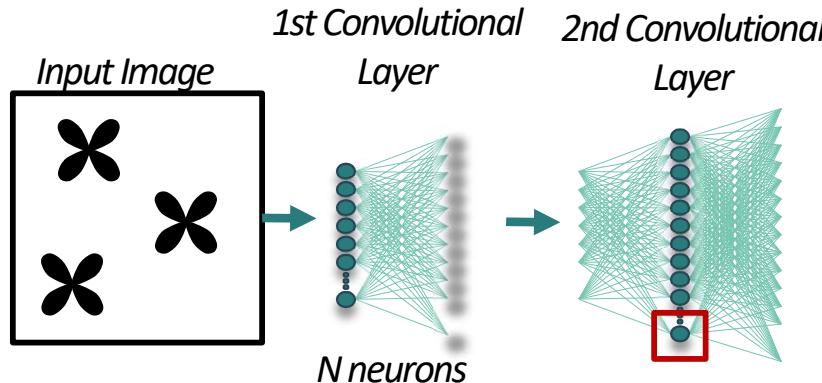
Understanding Composition of Convolutions

Example:



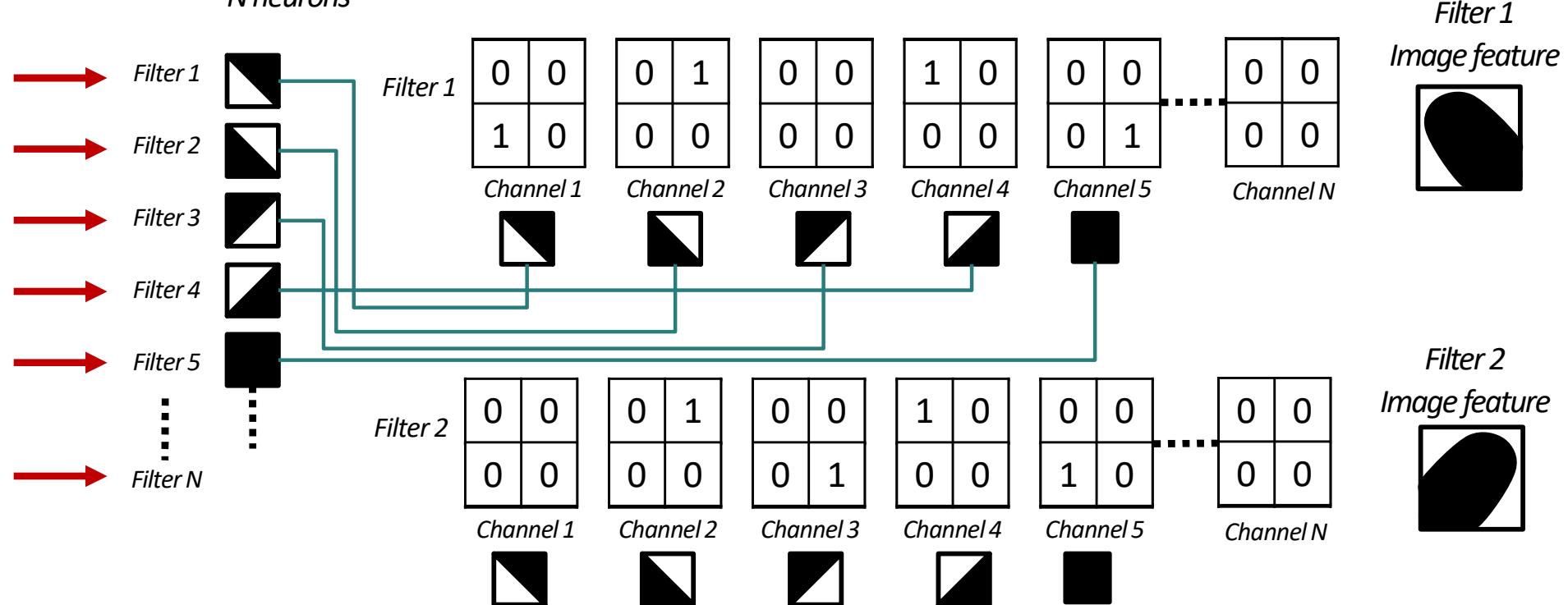
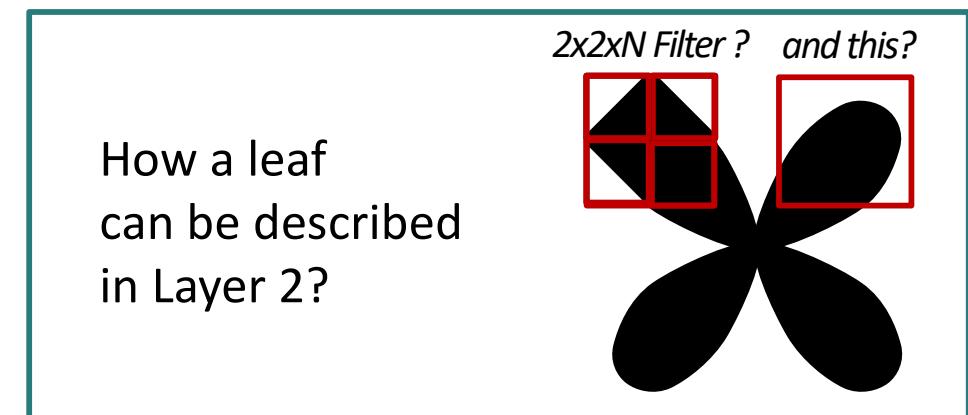
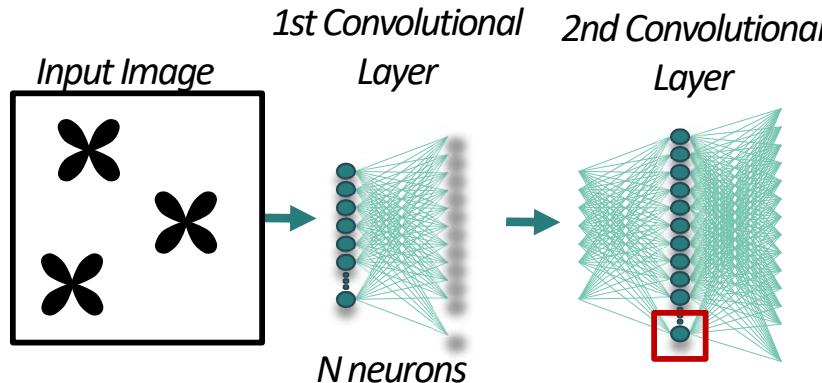
Understanding Composition of Convolutions

Example:

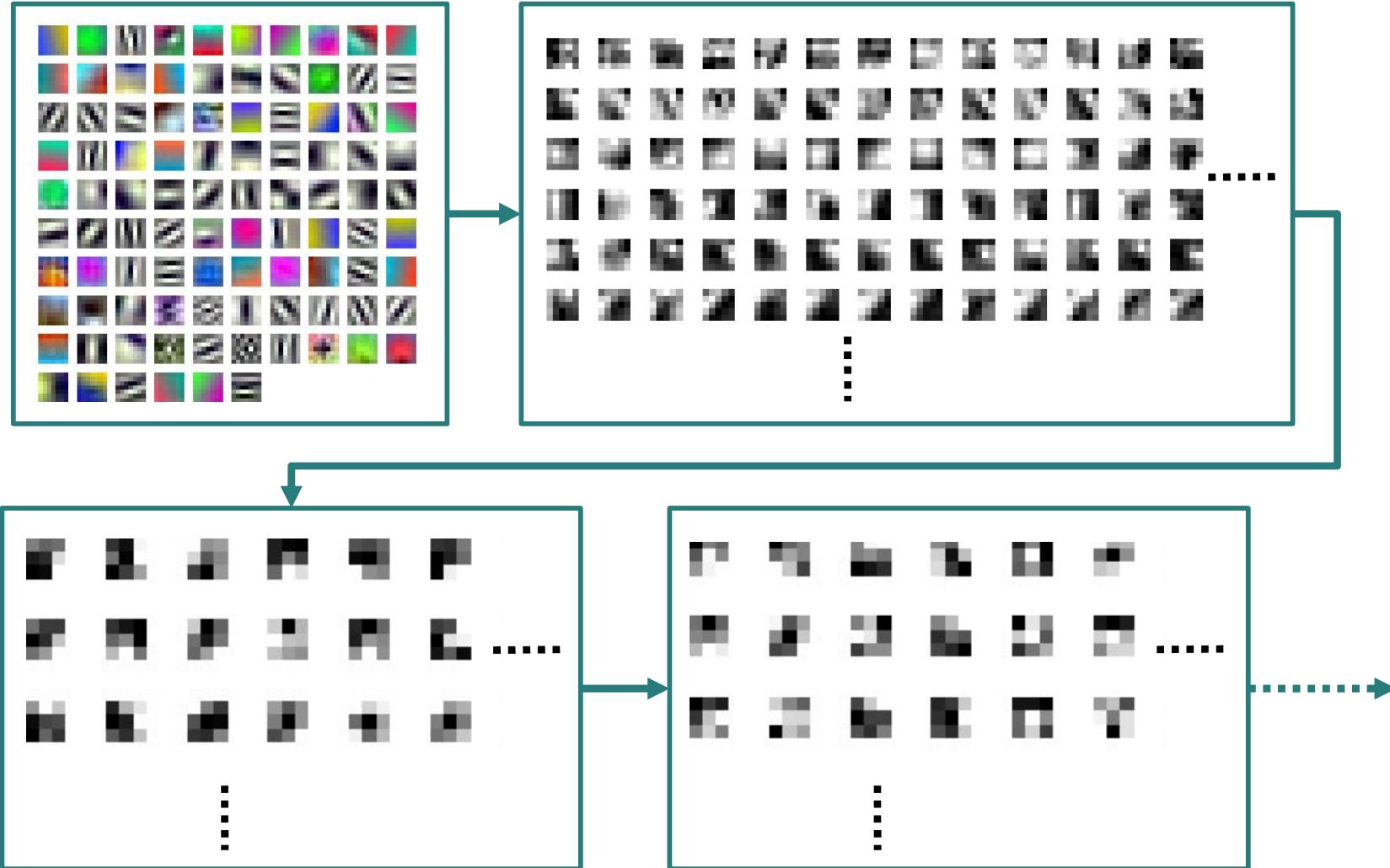


Understanding Composition of Convolutions

Example:



Understanding Composition of Convolutions through layers



It is a powerful representation, but deeply non-understandable

Conclusion: Composition of Convolutions is a powerful representation of complex spatial features where the hierarchy of layers is maximizing the representation power

But, explaining based on visualizing weights ...

Pros:

- Easy and Understandable **for the 1st layer**

Contras:

- No interpretation for the rest of layers

We need to use other tools to understand ...

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Index of this Lecture:

Post-hoc analysis

What: Methodologies and tools to analyse pretrained models

How: Firstly you train a model regularly, then use the tools to understand the model's behaviour

Why: It can provide usefull information about the model (biases, fails,...)

Ad-hoc modelling

What: Models specially created to be understandable

How: you create a new model from scratch tuning specific parameters to make it easy to understand.

Why: This models are easier to understand for us, thus easier to correct

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

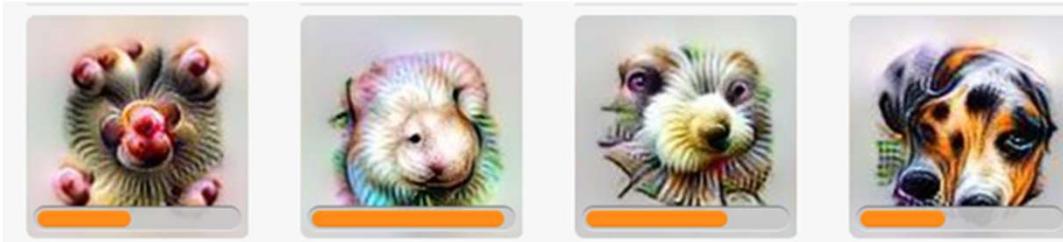
How color is represented in a CNN? and parallelisms with HVS

Neuron Analysis

IDEA: Understand a network by visualizing the individual neurons

There are two main methodologies used to visualize the neuron preference:

- **Inverting-based methods:** Generate the image that produces a specific activation



Olah, et al., "The Building Blocks of Interpretability", Distill, 2018

- **Activation maximization methods:** Find the images that maximally activate a neuron

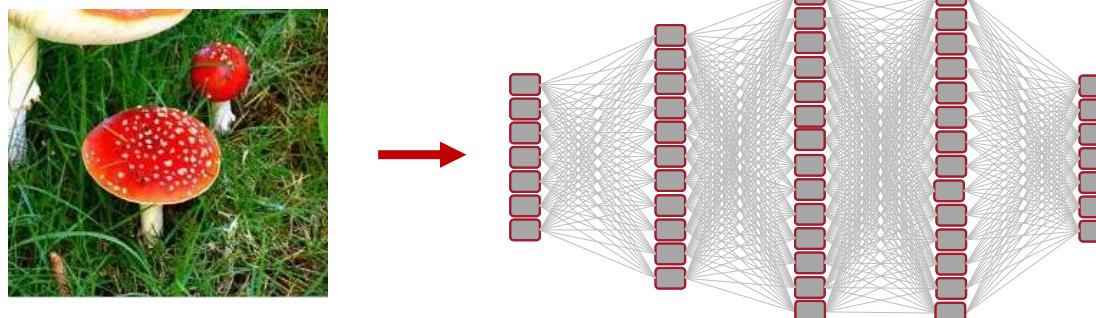


Rafegas, I.,et.al (2020).
Understanding trained CNNs by indexing neuron selectivity. *Pattern Recognition Letters*, 136, 318-325.

Neuron Analysis

Example

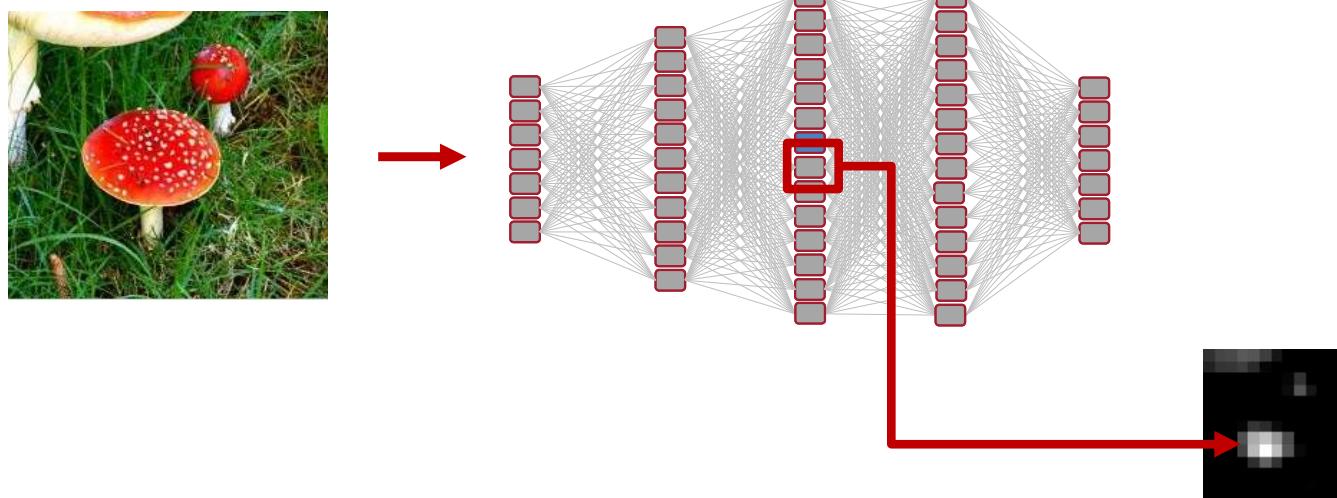
- 1) Run an image through a neural network.



Neuron Analysis

Example

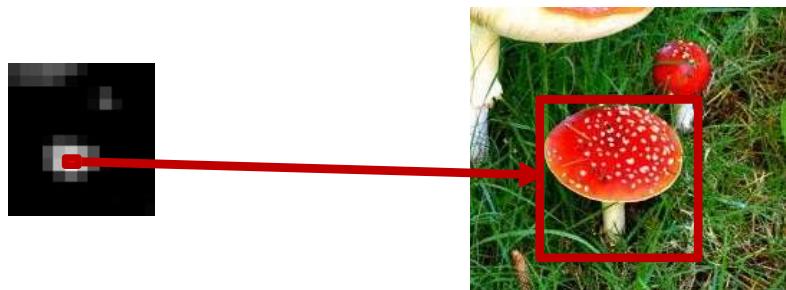
- 1) Run an image through a neural network.
- 2) Obtain the output of a neuron (Activation Map)



Neuron Analysis

Example

- 1) Run an image through a neural network.
- 2) Obtain the output of a neuron (Activation Map)
- 3) Find the part of the image (with size of the receptive field) that triggered the **maximum** activation in the feature map.

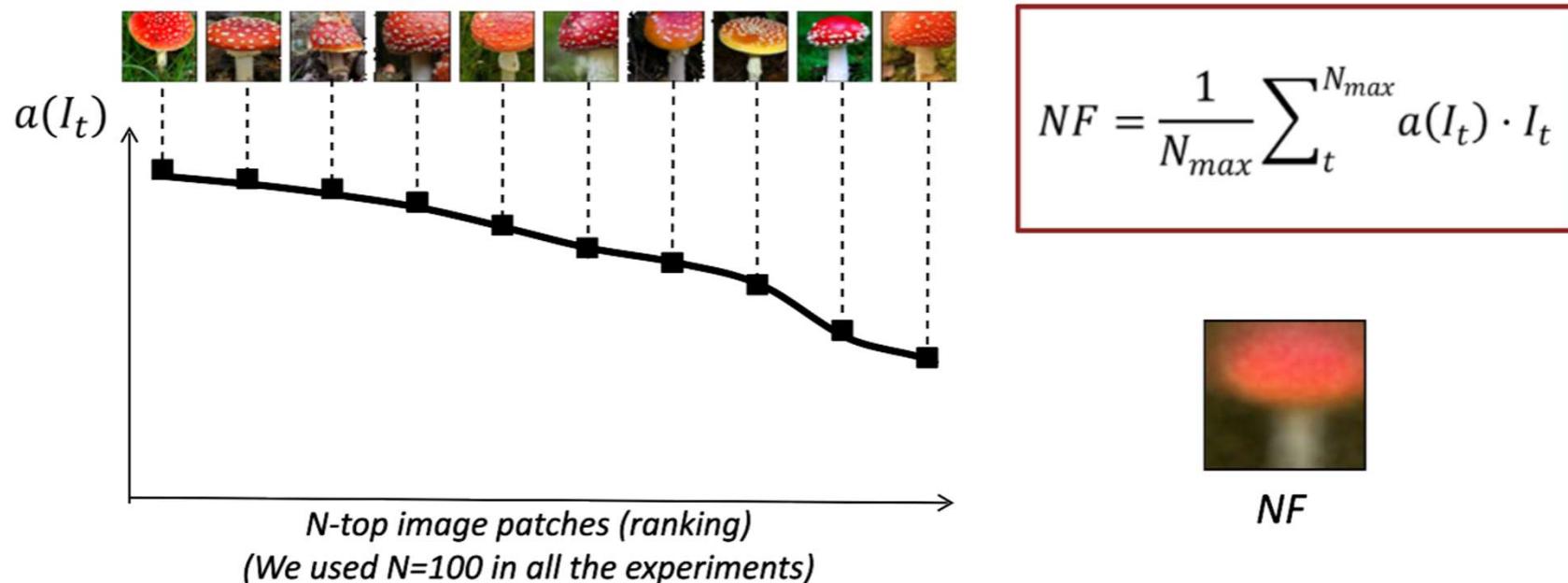


Receptive field
(image patch)
of a neuron
for a given image

Neuron Analysis

Example

- 1) Run an image through a neural network.
- 2) Obtain the output of a neuron (Activation Map)
- 3) Find the part of the image (with size of the receptive field) that triggered the **maximum** activation in the feature map.
- 4) Find the top activating images and perform a weighted average (Neuron Feature).



Neuron Analysis

USEFUL TOOLS

- **Distill.pub** --> Web-journal dedicated to the understanding CNN at a neuron level
<https://distill.pub/2018/building-blocks/>
- **Network Dissection** --> Tool to get activation information of each neuron
<http://netdissect.csail.mit.edu/>

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Data Inspection

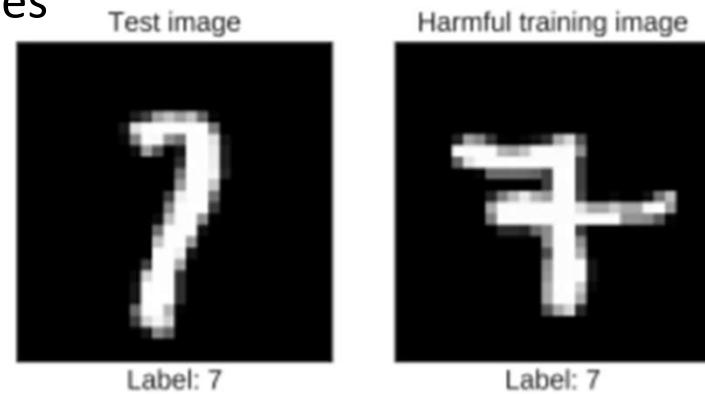
IDEA: Study the dataset to understand possible training bias

Study harmful Images on predictions

1. Find the most similar images in the dataset
2. Positive influence if they share label / Negative otherwise

Why is it usefull: help identify mis-annotated labels and outliers existing in the data

- Incorrectly labeled images
- Similar images belonging to different classes
- Context over object

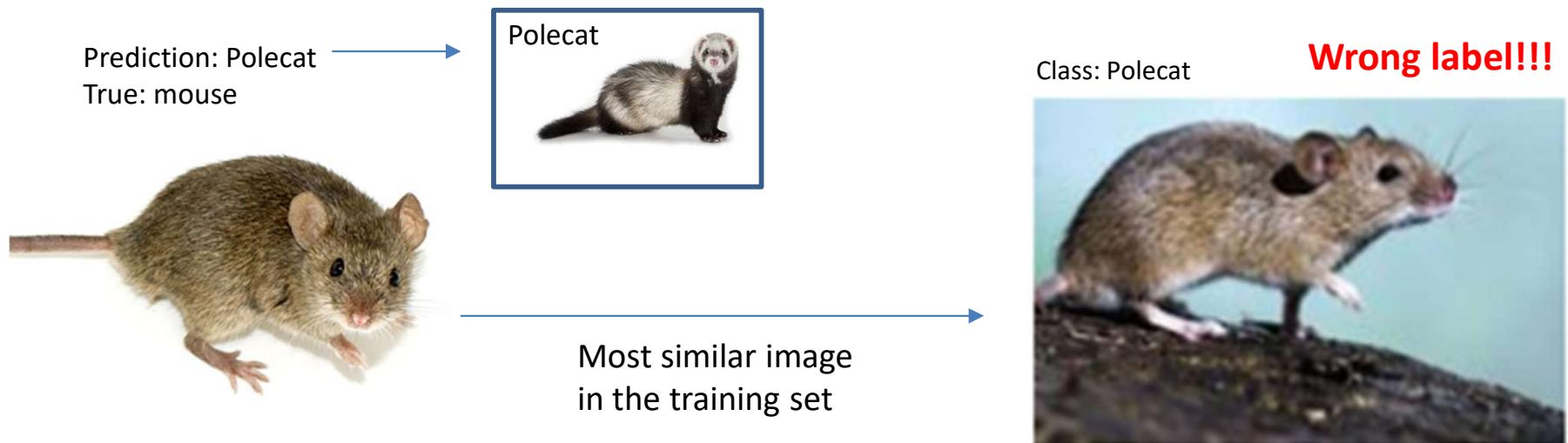


Koh, P. W., & Liang, P. (2017, July). Understanding black-box predictions via influence functions. In *International Conference on Machine Learning* (pp. 1885-1894). PMLR.

Data Inspection

Example

- Find incorrectly classified test data
- Use algorithm to find similar images in the train set (handcrafted descriptor KNN)
- Check if the error comes from wrongly labeled train data / similar class



Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

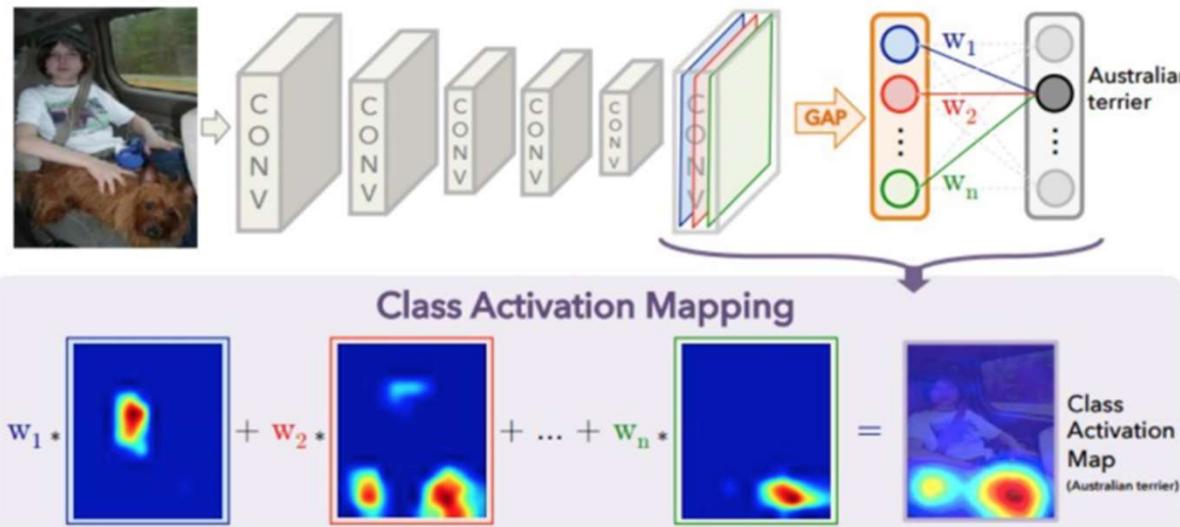
How color is represented in a CNN? and parallelisms with HVS

Saliency based

IDEA: Visualize the attributes of input data being more relevant to a prediction.
Relevant attributes are highlighted as attention maps

Saliency maps can be used for understanding and for improving CNN:

- **Class Activation Maps (CAM):** For each possible class prediction, highlights the regions of the image that contributed the most
- **Use saliency maps to improve performance:** Add saliency information as an additional input of the CNN to ensure that the CNN is taking into account only the selected areas of interest

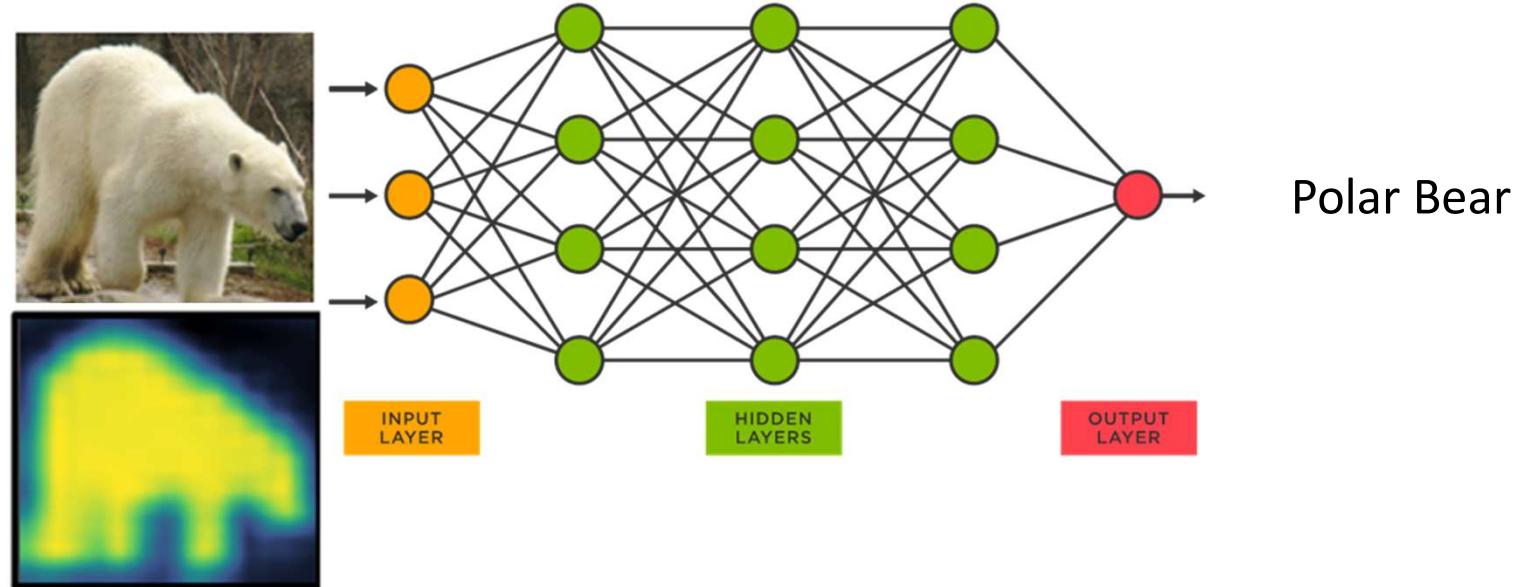


Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In Proceedings of the IEEE international conference on computer vision (pp. 618-626).

Saliency based

Example:

- Check the saliency map of an incorrectly classified image
- Use saliency to correct what should be important



Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

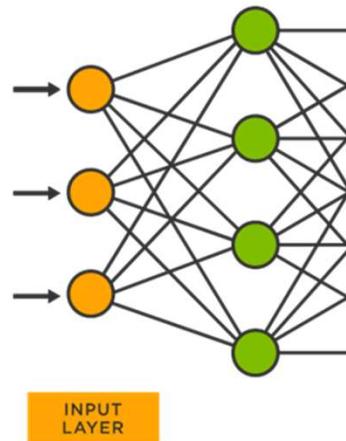
How color is represented in a CNN? and parallelisms with HVS

Proxy model

IDEA: Create an alternative model that performs similarly to the DNN

There are three main methodologies:

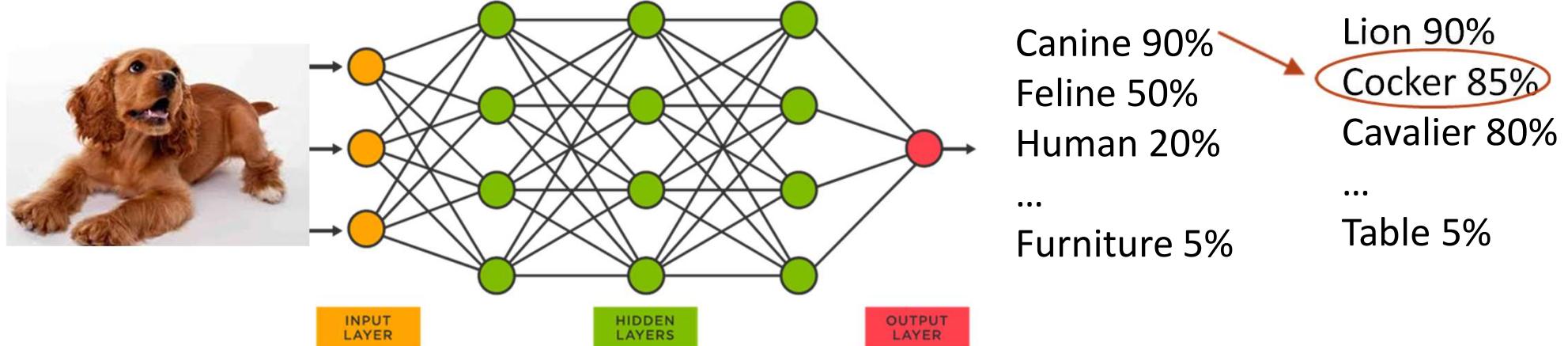
- **Simplistic models:** Create a decision tree or a set of rules that mimics the behaviour of the CNN
- **Knowledge distillation:** Use soft labels to better understand the classification process (I.e: Breed, specie, animal). Can be used to improve training.
- **Local Interpretable Modelagnostic Explanation:** Cut the NN at a certain point to study intermediate representations



Proxy model

Example (Knowledge distillation)

- Use soft labels to locate the error in classification
- Use soft label information in the loss function to improve training



Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Modifications

IDEA: Neural networks outputs additional information providing some insight of why a decision was taken

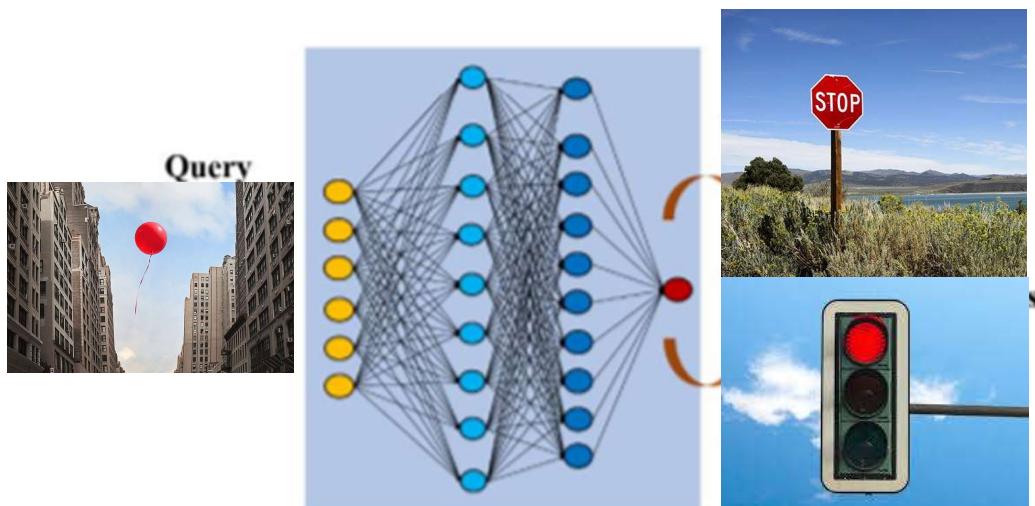
There are two main methodologies:

- Explaining by case
- Explaining by text

Explaining-by-Case

IDEA: Give similar examples of training instances that have a similar label

- Given a query the CNN outputs a label as well as nearest neighbors of that activation values in the training set
- Visual information provided by this function may help the user understand why a certain decision was made

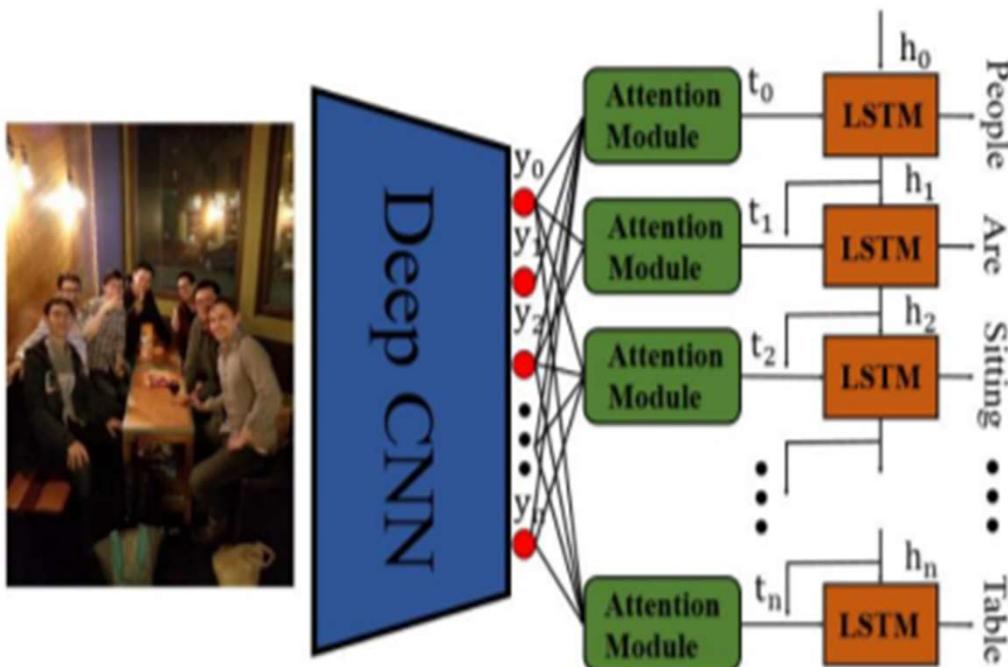


S. Wachter, B. Mittelstadt and C. Russell,
“Counterfactual Explanations without
Opening the Black Box: Automated
Decisions and the GDPR,” Harv. JL &
Tech., vol. 31, no. 841, 2017

Explaining-by-Text

IDEA: Give a text explanation of the output instead of just a label

- Combination of Attention modules and LSTM:
 - Attention: Detects the most important objects in the image
 - LSTM: Produces a sentence based on ordered attention



A. Karpathy, L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," In CVPR, pp. 3128-3137, 2015

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Theoretical Analysis

IDEA: Study of the DNN architecture from a mathematical point

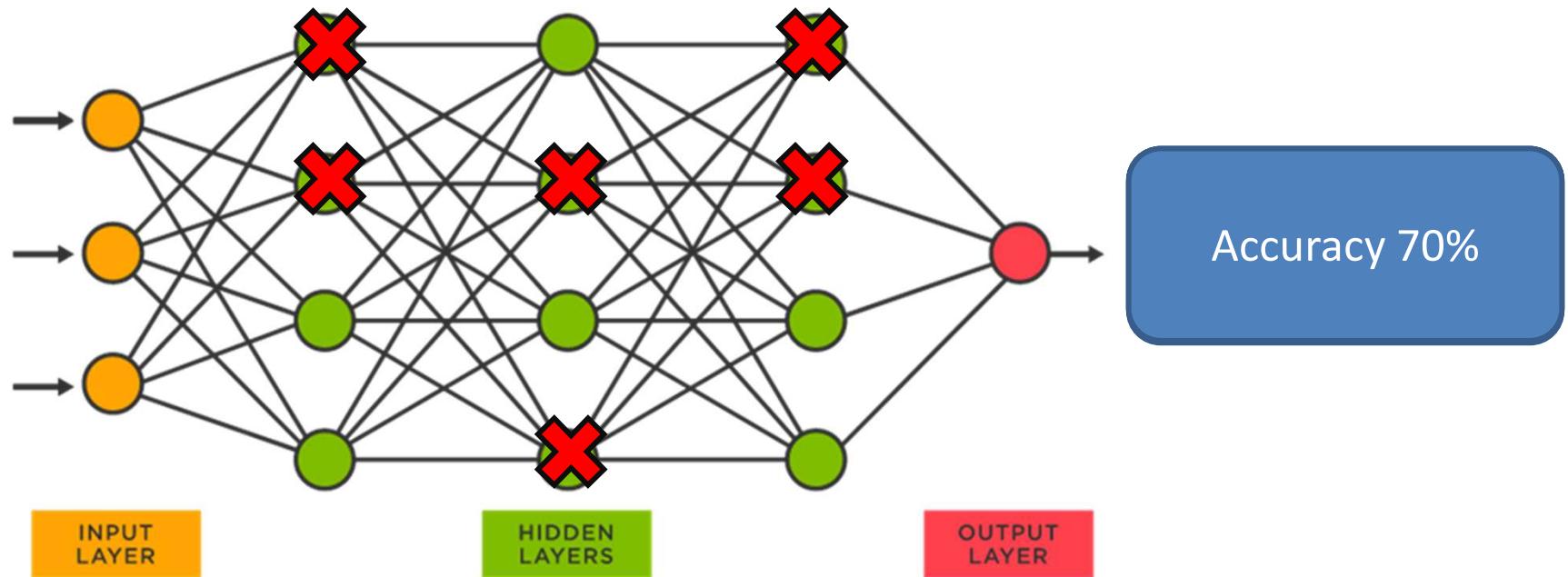
There are three main areas of study

- **Representation:** Deeper networks are more expressive than shallow ones
- **Optimization:** The number of parameters in a network exceeds the number of data instances (try to find the minimum number of parameter with best performance)
- **Generalization:** Explain why a deep network can generalize well despite the number of parameters is greater than the number of data samples

Theoretical Analysis

Example: Ablation

- Use ablation to find the optimal number of parameters
- In Pytorch: **import** torch.nn.utils.prune **as** prune
 - https://pytorch.org/tutorials/intermediate/pruning_tutorial.html



Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Interpretable representation

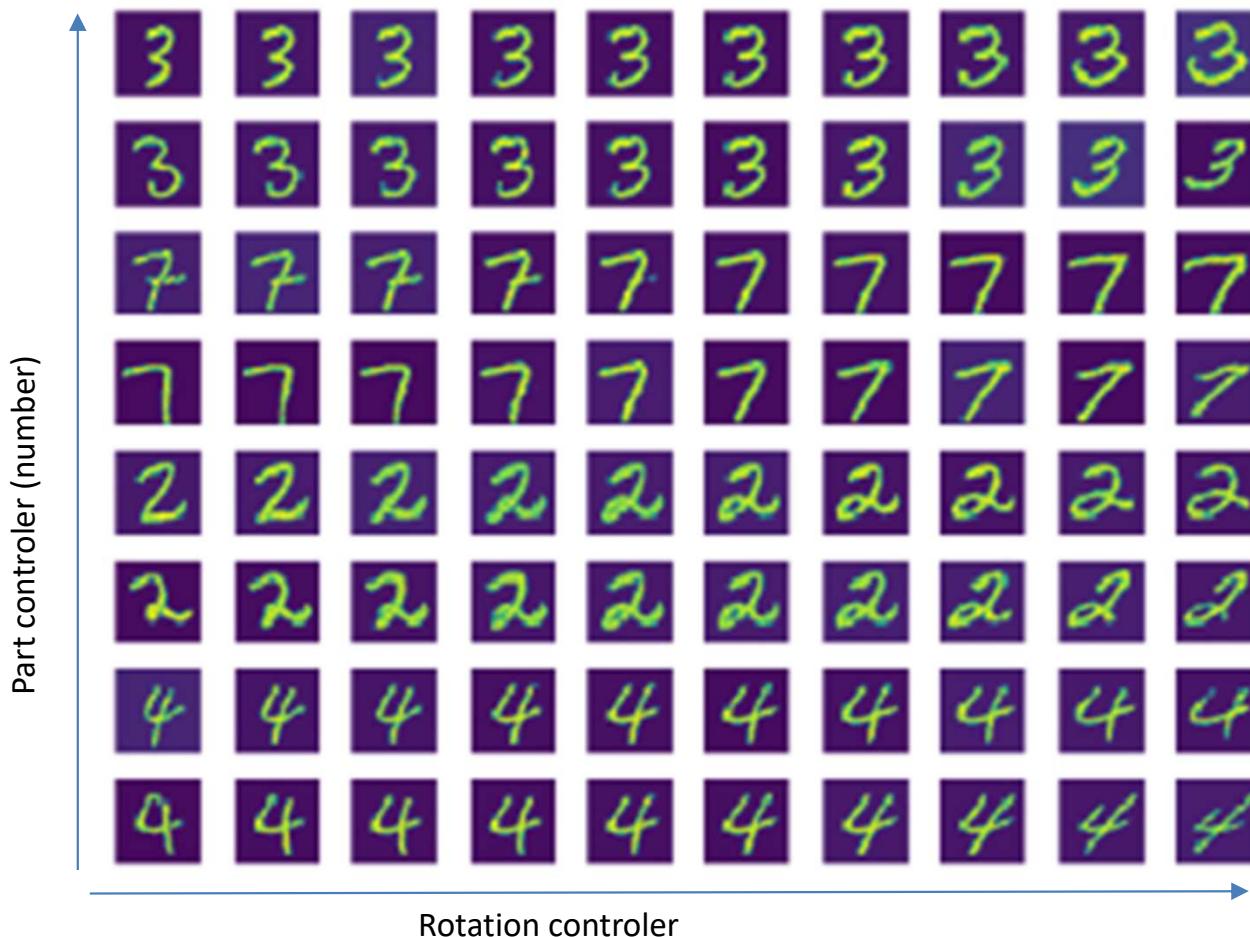
IDEA: Training the DNN ensuring that the neurons are sensitive to known or controlled stimulus by using regularization techniques

Regularization techniques steer the optimization towards more interpretable representations with the following properties:

- Decomposability: Neurons represent simple concepts
- Mathematical constraints: monotonicity or Non-negativity
- Sparsity: Maximizes the difference of internal representations
- Human-in-the-loop prior: Use handcrafted features

Interpretable representation

Example:



In an InfoGAN, two latent codes control the localized parts and rotation parts respectively

Internal features aligned with the PCA of the training data

X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel, “Infogan: Interpretable representation learning by information maximizing generative adversarial nets,” In NeurIPS, 2016

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- Model Renovation

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS

Model renovation

IDEA: Seek interpretability by means of designing and deploying more interpretable machineries into a network

There are three main methodologies

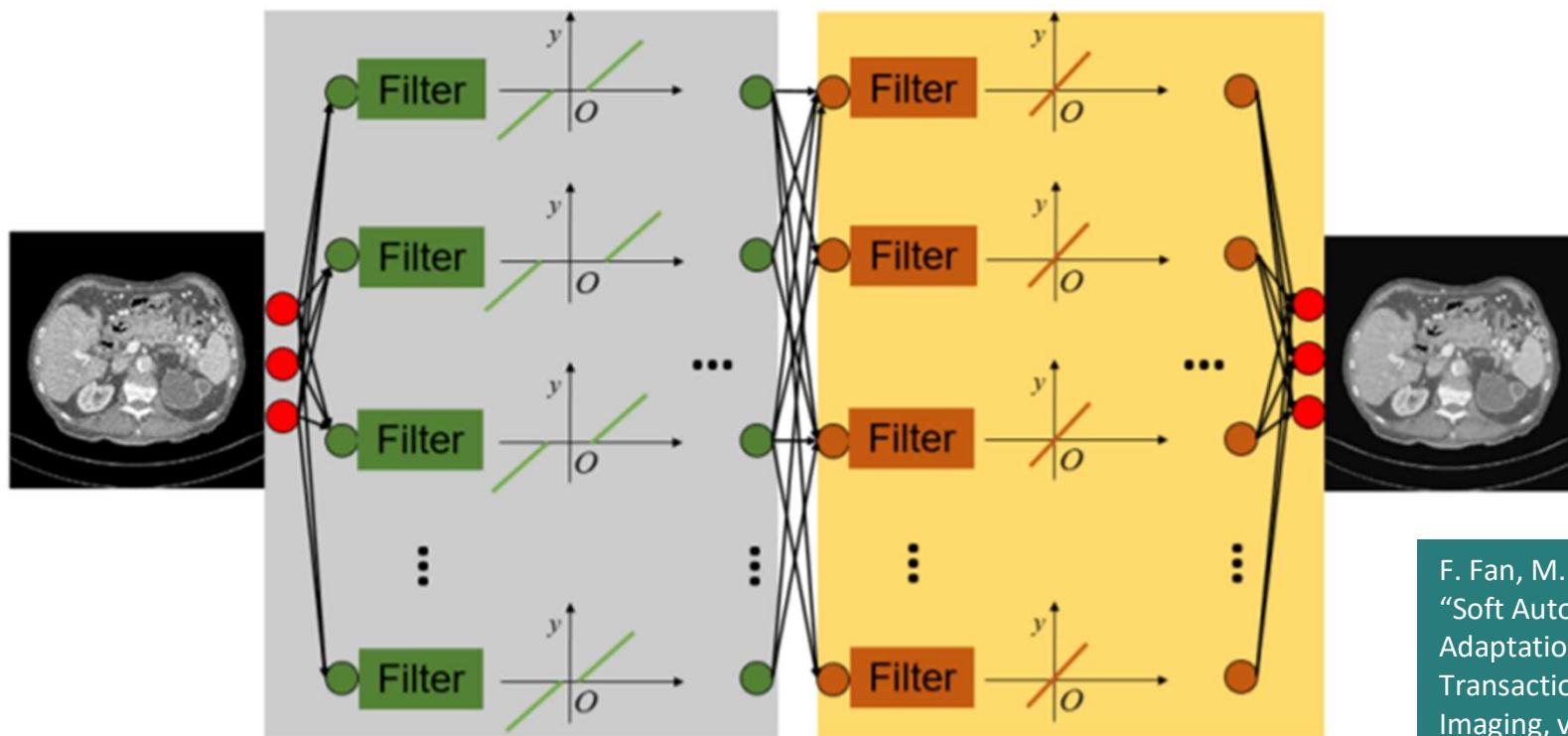
- **Neuron with purposely designed activation function**
Some parameters of the neurons can be modified externally
- **Inserted layer with a special functionality**
Intermediate layers output or interlayer connections
- **Modularized architecture**
Neural Network is a sum of simple modules

Model renovation

Example:

Example of Neurons with a soft Autoencoder

The threshold values of the RELU function can be modified at any point and act as a wavelet filter



F. Fan, M. Li, Y. Teng and G. Wang,
"Soft Autoencoder and Its Wavelet
Adaptation Interpretation," IEEE
Transactions on Computational
Imaging, vol. 6, pp. 1245-1257, 2020

Interpretation Model Properties

In order to create a good interpretation model we need to take into account this properties

Exactness: how accurate an interpretation method is.

Consistency: there is no contradiction in an explanation.

Completeness: show effectiveness in support of the maximal number of data instances and data types.

Universality: universal interpreter that deciphers many models.

Reward: What are gains from the improved understanding.

Why Is Interpretability Difficult?

- **Human Limitation:** As humans we are limited to understanding only basic features
- **Algorithmic Complexity:** The complexity and combination of algorithms makes it very difficult to follow the dataflow.
- **Commercial Barrier:** There is an effort to make algorithms hard to understand

Summary

- Neuron Analysis: Visualize individual neurons.
- Data Inspection: Inspect the data to understand the training
- Saliency based: Visualize attention maps to highlight relevant regions
- Proxy models: Use simplistic models to summarize how it works
- Modifications: Build the CNN to give extra information about the decision made.
- Theoretical Analysis: Use theory to understand how a CNN works.
- Interpretable representation: Regularize the CNN training to make neurons easier to understand
- Model Renovation: Build CNNs with components easier to understand

Example

Interpreting CNNs via Decision Trees

Quanshi Zhang[†], Yu Yang[‡], Haotian Ma[§], and Ying Nian Wu[‡]

[†]Shanghai Jiao Tong University, [‡]University of California, Los Angeles,
[§]South China University of Technology

Abstract

This paper¹ aims to quantitatively explain the rationales of each prediction that is made by a pre-trained convolutional neural network (CNN). We propose to learn a decision tree, which clarifies the specific reason for each prediction made by the CNN at the semantic level. I.e. the decision tree decomposes feature representations in high conv-layers of the CNN into elementary concepts of object parts. In this way, the decision tree tells people which object parts activate which filters for the prediction and how much each object part contributes to the prediction score. Such semantic and quantitative explanations for CNN predictions have specific values beyond the traditional pixel-level analysis of CNNs. More specifically, our method mines all potential decision modes of the CNN, where each mode represents a typical case of how the CNN uses object parts for prediction. The decision tree organizes all potential decision modes in a coarse-to-fine manner to explain CNN predictions at different fine-grained levels. Experiments have demonstrated the effectiveness of the proposed method.

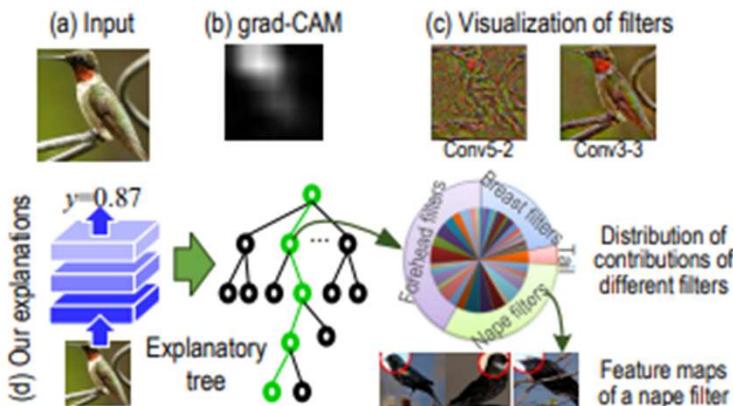


Figure 1. Different types of explanations for CNNs. We compare (d) our task of quantitatively and semantically explaining CNN predictions with previous studies of interpreting CNNs, such as (b) the grad-CAM [26] and (c) CNN visualization [23]. Given an input image (a), we infer a parse tree (green lines) within the decision tree to project neural activations onto clear concepts of object parts. Our method quantitatively explains which filters/parts (in the small/big round) are used for the prediction and how much they contribute to the prediction. We visualize numerical contributions from randomly selected 10% filters for clarity.

Index of this Lecture:

Preliminary considerations

Post-hoc analysis

- Neuron Analysis
- Data Inspection
- Saliency based
- Proxy models
- Modifications
- Theoretical Analysis

Ad-hoc modelling

- Interpretable representation
- *Model Renovation*

A case study on a single feature (*post-hoc analysis*)

How color is represented in a CNN? and parallelisms with HVS