# Advanced Deep Learning - Introduction to Explainable Artificial Intelligence (XAI)

## XAI methods & metrics

STASSIN Sédrick

Professor: MAHMOUDI Sidi Ahmed
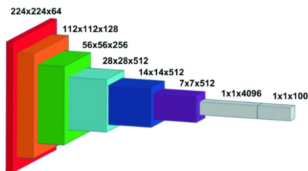
UMONS
Université de Mons

POLYTECH MONS

# INTRODUCTION

Input → **Trained** DL Model (Black Box 🙁) → Output

*Internal behavior of the code is unknown*

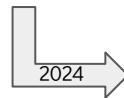**Black Box** problem of deep neural networks

European Commission

2024 → The **AI Act**

Trustworthy AI systems must be considered:

- **Lawful** – Operating within the limits of law

- **Ethical** – Fair models that do not discriminate

- **Robust** – Delivering <u>reliable</u> results in all considered situations

# XAI

**E**xplainable **A**rtificial **I**ntelligence
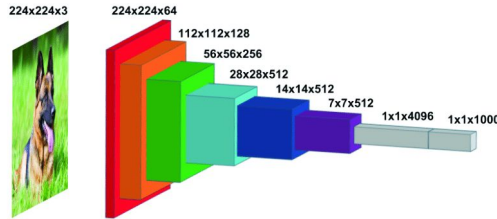
# INTRODUCTION

How does this system work ?

Can I trust this AI model ?
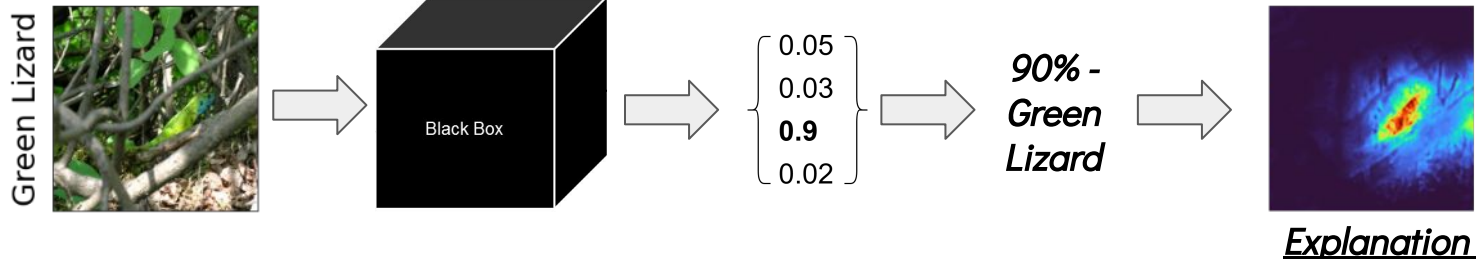
Under what circumstances ?

How to evaluate this explainability ?

# EXplainable Artificial Intelligence (XAI)

**XAI in computer vision – Images**

**applied to CNN**



- Most used - attribution-based methods = *saliency methods*



*Explanation*

# CONTENTS

**02**

# BACKGROUND
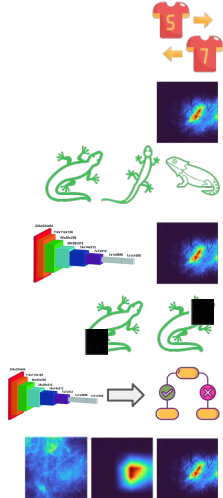
XAI Taxonomy

XAI Model Taxonomy

**UMONS**
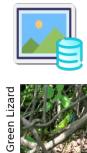Université de Mons

**Taxonomy:**
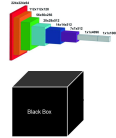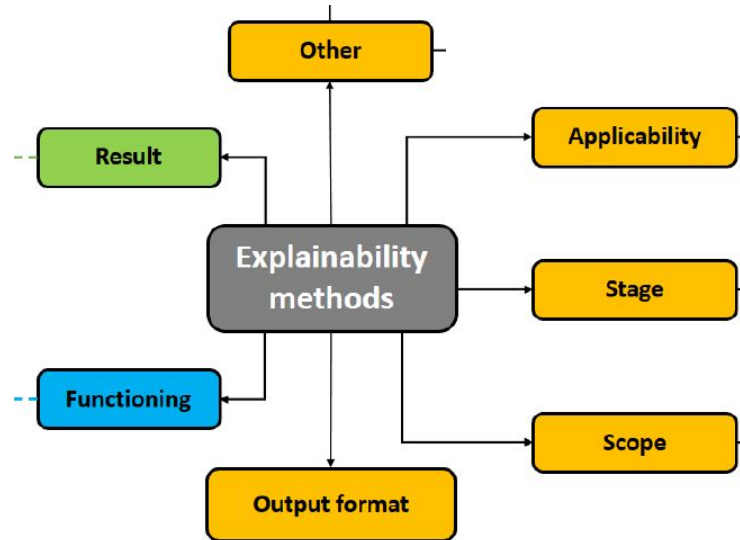
- **Conceptual**
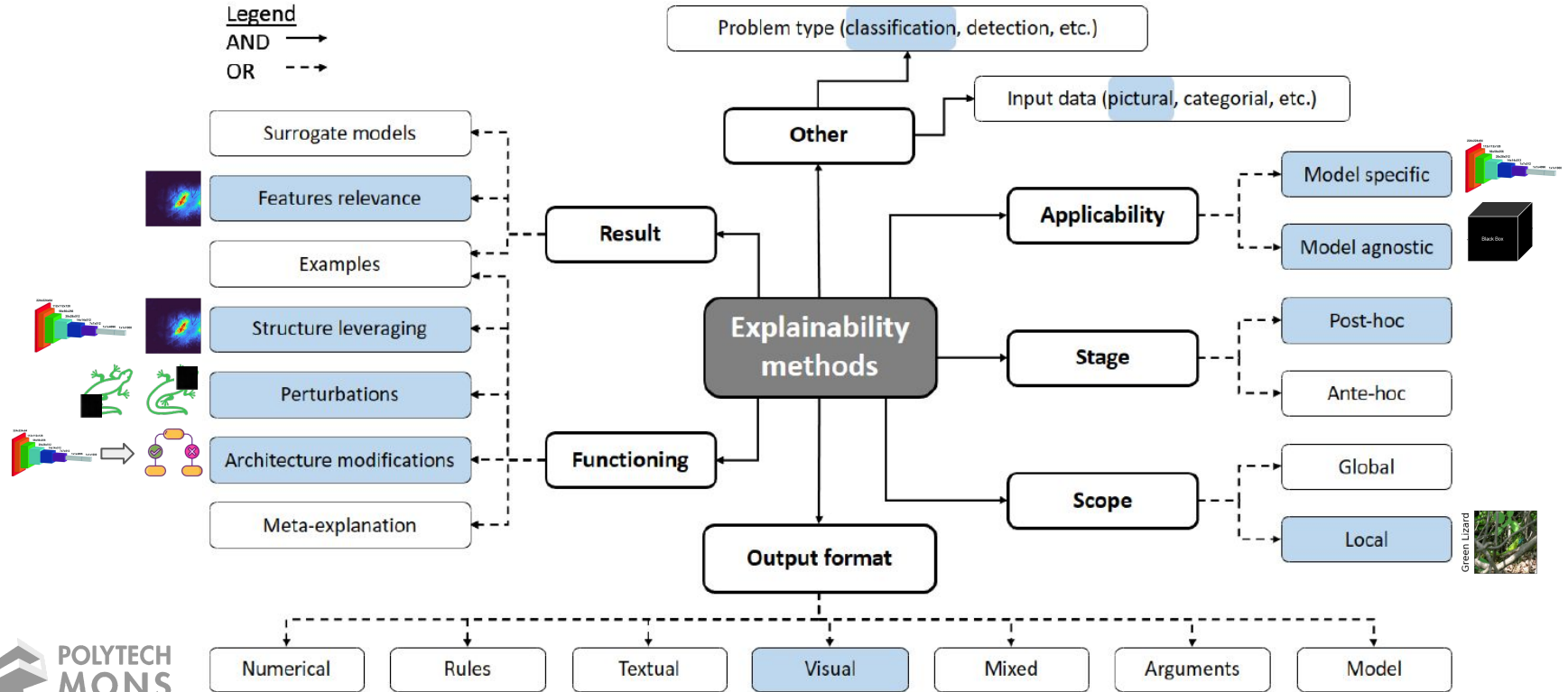- **Result-based**
- **Function-based**

Legend
AND →
OR --→

# XAI Taxonomy

## Selected Method Taxonomy

**03**

# EXPLAINABILITY AND BIAS DETECTION

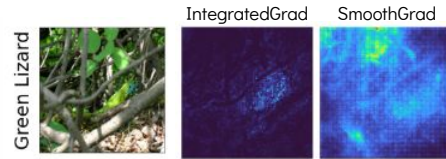**03a** XAI Methods

**03b** Covid-19 Use Case

# **03a** Saliency **methods**

## Gradient-based methods :
[Smilkov et al., 2017] [Sundarajan et al., 2017] ...
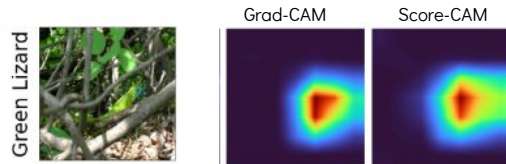
- Use the gradient (e.g. back-propagated) to picture the derivative of the model output w.r.t the input image.



## CAM-based methods :
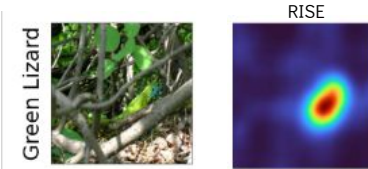[Selvaraju et al., 2017] [Wang et al., 2020] ...

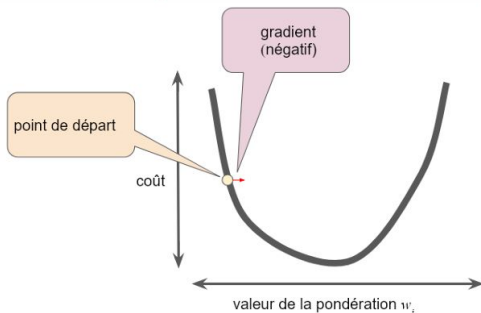- Produce a weighted sum of the activations from a convolutional layer



## Perturbation-based methods :
[Petsiuk et al., 2018] ...

- Study the output model response to small changes in the input.

## Gradients

[Simonyan et al., 2013]

### Neural Network Backpropagation



- **Gradient :** vecteur ayant deux caractéristiques : direction et magnitude
- Il indique la direction de la croissance maximale de la fonction de perte
- L'algorithme de descente de gradient fait **un pas dans le sens inverse** afin de **réduire la perte** aussi rapidement que possible.

Université de Mons — Sidi Ahmed Mahmoudi — Cours IA. Chapitre 5 — 40

Loss function

$$w_{t+1} = w_t - \alpha \frac{\partial \mathcal{L}}{\partial w_t}$$

learning rate — weight

### Formula

Loss function

$$R^c = \frac{\partial \mathcal{L}_c(x)}{\partial x}$$

Relevance of class c

Input

**SmoothGrad:** Add *Gaussian noise* to *n* samples of the image - Compute gradients - Average

**Integrated Gradients:** *Create n* images - ranging linearly from a *baseline image* to the *input image* - Average the gradients 13
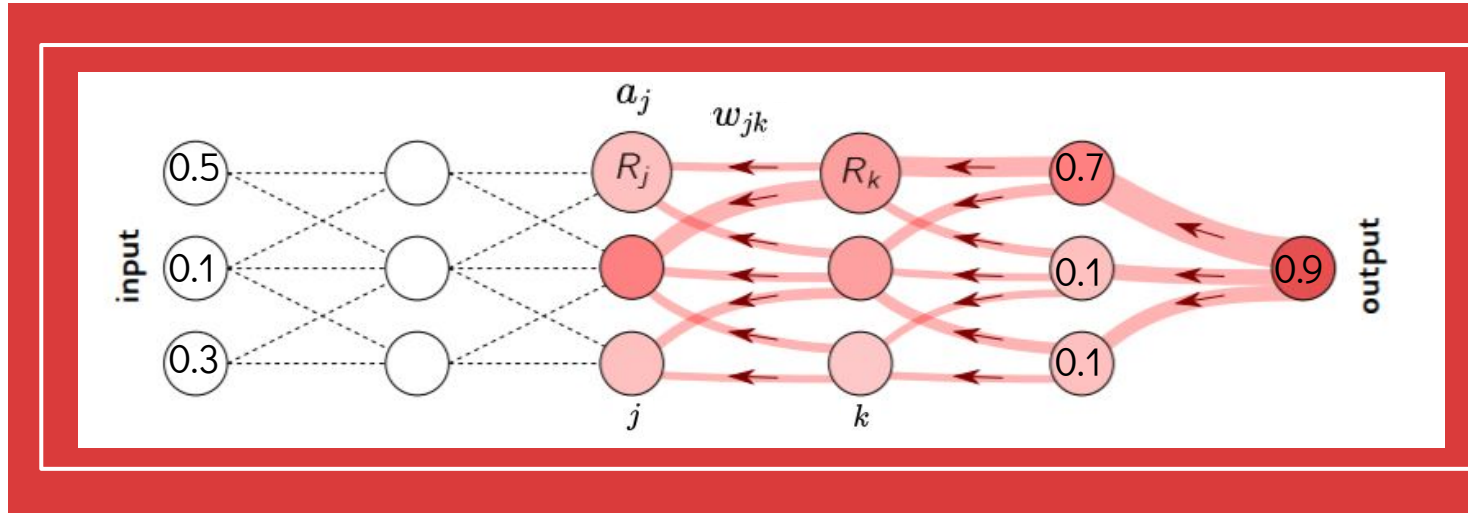
**Conservation property**

What has been received by a neuron must be redistributed to the lower layer in equal amount

$j$  $k$  Consecutive layers

$R_j$  $R_k$  Neuron relevances

$a_j$  Neuron $j$ activation

$w_{jk}$  Weight between neuron $j$ and $k$
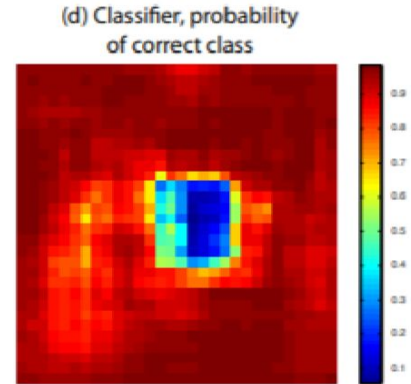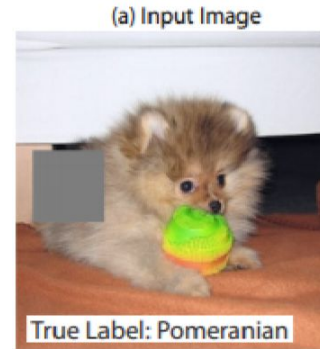


$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$$

14

### Main steps

- Choose an image, a size of square

- For all possible positions of the square in the image

  ○ Occlude

  ○ Compute score

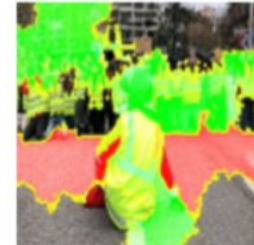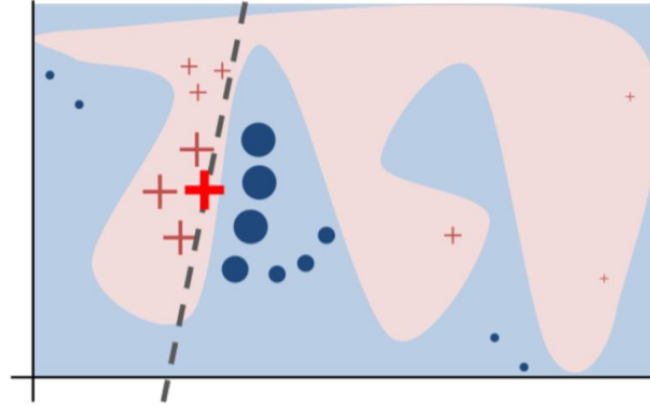⟹ Pixels are important if the class score drops significantly



(a) Input Image

True Label: Pomeranian

(d) Classifier, probability of correct class

## Principles

- Fit a linear model to *n* perturbed samples of an image

  - Predicted by the complex model

  - With distances as weights

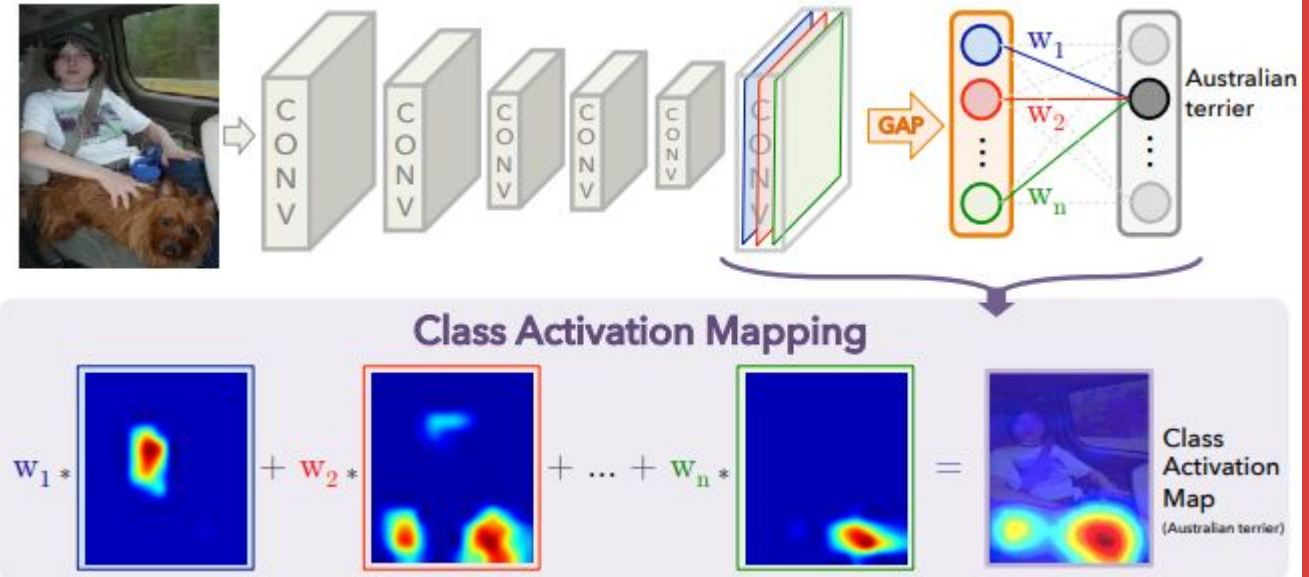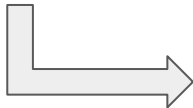- The linear model using the *m* best features provides a locally faithful explanation for the image

# CAM-based Methods

## Class Activation Mapping (CAM)

[Zhou et al. 2016]

Class Activation map for class **c**

$$R^c = \sum_{n=1}^{N} w_{n,c} A_n$$



**Class Activation Mapping**

$w_1 *$ $+$ $w_2 *$ $+ \ldots + w_n *$ $=$ Class Activation Map (Australian terrier)
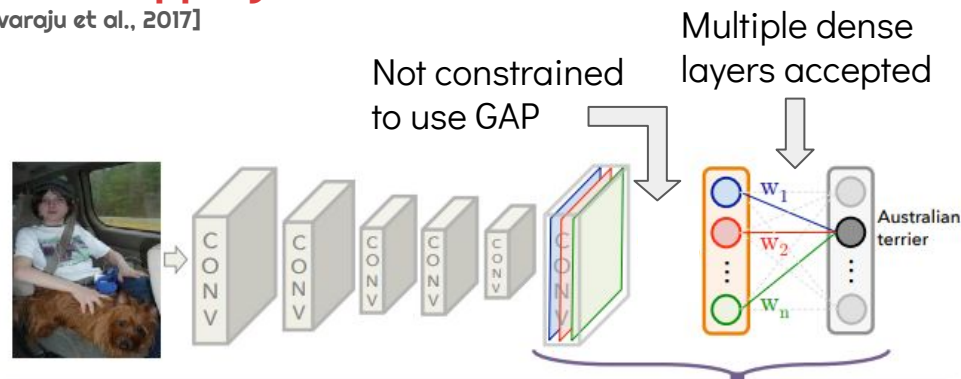
# CAM-based Methods
## Grad-Class Activation Mapping (Grad-CAM)
[Selvaraju et al., 2017]

## Main differences

1. Generalization to multiple dense layers

   a. **Gradient** w.r.t feature maps of the last convolutional layer

2. Not constrained to models with GAP

   a. **GAP** is used to obtain neuron importance weights

3. **ReLU** activation to only keep positive contributions

Not constrained to use GAP

Multiple dense layers accepted

Australian terrier

global average pooling

$$w_{n,c} = \frac{\partial y_c}{\partial A_{n,i,j}}$$

gradients via backprop

$$R^c = \left( \sum_{n=1}^{N} w_{n,c} A_n \right)$$

linear combination

18

## Problem

Covid-19  Image Classication

## Dataset

2621 training (90%) — 284 testing (10%)
3 classes

## Models

VGG16 – VGG19 –  ResNet50  – DenseNet121 – DenseNet201

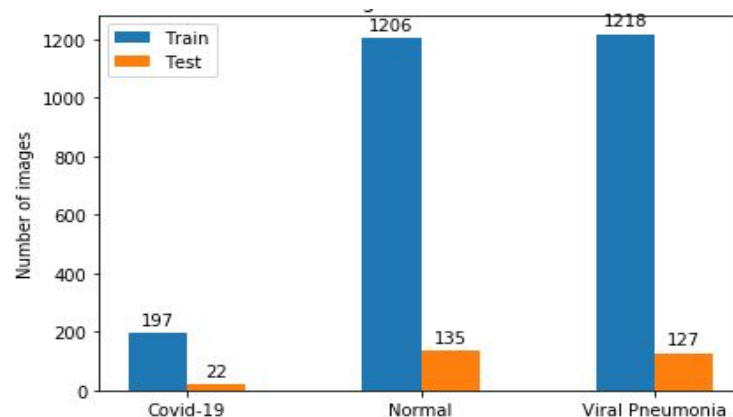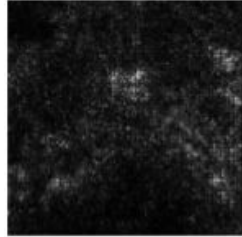## Analysis

1- XAI on predicted class
2-  Model comparison
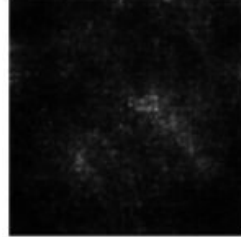


Image distribution among existing classes

19

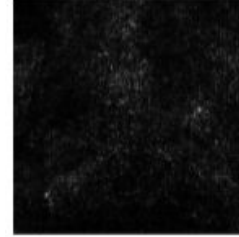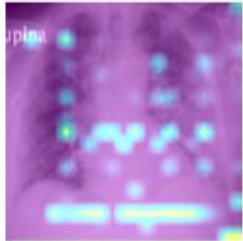## XAI : Method comparison - Predicted class



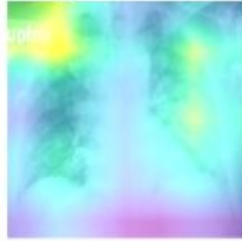(a) Input  (f) Gradient  (h) SmoothGrad  (i) Integrated

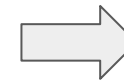(e) Occlusion  (j) GradCAM  (m) PresetAFlat  (p) LIME

DenseNet121 explanation for a Covid-19 x-ray image

- <u>Gradients + Occlusion:</u> noisy results

- <u>Best visual results:</u> LRP (Preset) + LIME + GradCAM
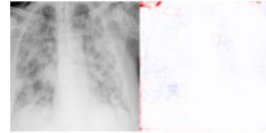
  → Letter Detection ?

➡ **XAI helps to detect biases in models**

**Model bias ?**

**Models**

- DenseNet201
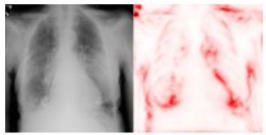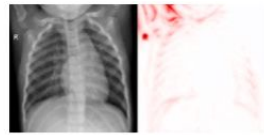
- ResNet50

- VGG16

- VGG19



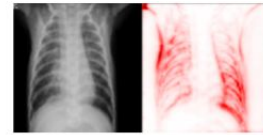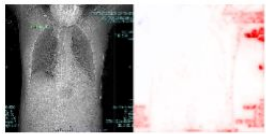(a) Covid-19 Image
(b) Normal Image
(c) Viral Pneumonia Image

(a) Covid-19 Image
(b) Normal Image
(c) Viral Pneumonia Image

(a) Covid-19 Image
(b) Normal Image
(c) Viral Pneumonia Image

(a) Covid-19 Image
(b) Normal Image
(c) Viral Pneumonia Image

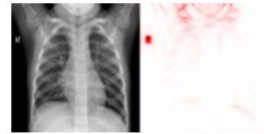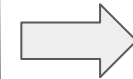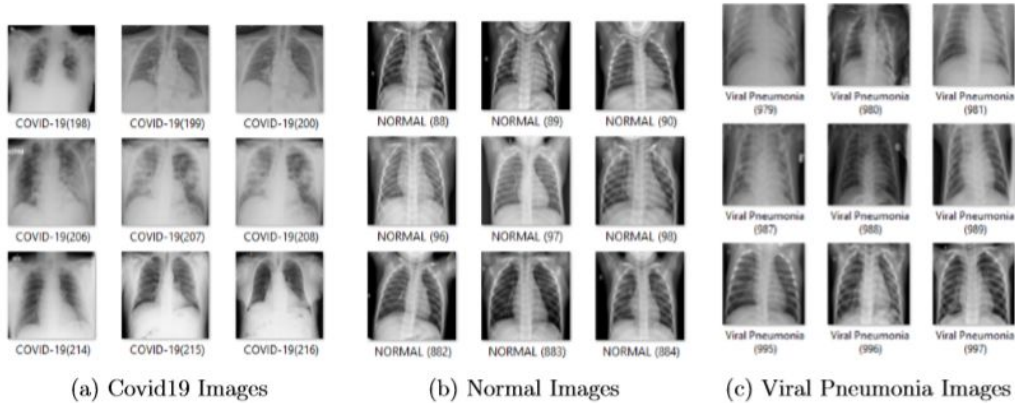- Every model uses zones out of interest

- VGG-16 seems less reliant

➡ **XAI can guide model selection**

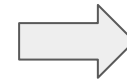- Normal Images: Detection towards the head ?

Visual comparison between 9 images by class

Biases due to:

- Camera positioning

- Arm positioning and patient characteristic

**XAI helps to understand the sources of bias**

# 036 COVID-19 CLASSIFICATION

**CT-Scan Images classification**

**XAI : Method comparison -** Predicted class

## Dataset - Model

- 349 Covid-19 CT
- 297 Normal CT

VGG-16
**Classification**

89% test accuracy

Lung Segmentation preprocessing
- 233 Covid-19 CT
- 293 Normal CT

VGG-16
**Classification**

70% Test Accuracy

## Biases

Focus outside lungs (integrated)

Top right Corner Bias

Focus inside lungs (integrated)

No apparent bias



(a) Unsegmented Covid-19 image

(b) Integrated Gradients

(c) LRP PresetAFlat

(d) LIME Proxy Model

(a) Segmented Covid image

(b) Integrated Gradients

(c) LRP PresetAFlat

(d) LIME Proxy Model

Unsegmented and Segmented Covid-19 CT Scan explained with two VGG-16 model

XAI helps to verify the model learning

23

# INTRODUCTION – EVALUATION OF EXPLAINABILITY METHODS

Input image → Classification Model (Black Box) → Probabilities $\begin{bmatrix} 0.05 \\ 0.03 \\ \textbf{0.9} \\ 0.02 \end{bmatrix}$ → *90% - Green Lizard* → **Saliency** map

**How to evaluate XAI methods ?**

**Faithfulness metrics :**
[Bhatt et al., 2020] [Petsiuk et al., 2018] ...


Deletion → does the accuracy drop?

- Measure how explanations follow the **predictive** behavior of the model

**Randomization metrics :**
[Sixt et al., 2020] ...


Randomization → Is the saliency map modified ?

- Measure change in saliency map as a function of **parameter randomization**

**Robustness metrics :**
[Yeh et al., 2019] ...


Small perturbation → Is the saliency map stable ?

- Measure the **stability** of explanations with respect to small input **perturbations**

**Complexity metrics :**
[Chalasani et al., 2020] ...


Is the saliency map concise ?

- Measure explanation **conciseness**

**Localization metrics :**
[Zhang et al., 2018] ...


Does it fit into the bounding box ?

- Measure whether explanations fit into the delimitation of a region of interest (e.g. bounding-boxes)

26

# ADL - Intro to XAI

XAI is a **valuable tool** to:

- Detect biases
- Guide model selection
- Understand the bias source
- Verify the model learning

XAI methods can be evaluated:

- Visually (subjective !)
- Through metrics (properties!)

# Gradient-based Methods

## Gradients

[Simonyan et al., 2013]

### Interpretation

*"Which pixels need to be changed the least to affect the class score the most"*

# Gradient-based Methods

## SmoothGrad - Integrated Gradients
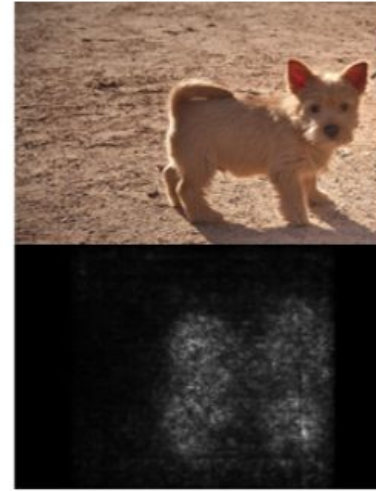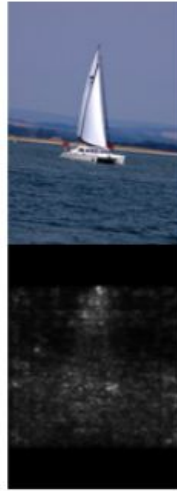**[Smilkov et al., 2017]**          **[Sundarajan et al., 2017]**

### SmoothGrad

*" Removing noise by adding noise "*

- Image $x$ - class $c$ - $n$ samples - noise parameters

- Add *Gaussian noise* to $n$ samples of the image

- Calculate the gradient of the class for each noisy image

- Average the gradients

$$\hat{M}_c(x) = \frac{1}{n} \sum_{1}^{n} M_c\big(x + \mathcal{N}\big(0, \sigma^2\big)\big)$$
$$\hookrightarrow M_c(x) = \partial S_c(x)/\partial x$$

### Integrated Gradients

*"Path from baseline to input"*

- Image $x$ - baseline $x'$ - number of steps $n$

- *Create $n$ images - ranging linearly from the **baseline** to the input image*

- *Calculate the gradient* for each image

- *Average* all calculated gradients, and multiply by (input-baseline)

$$\text{IntegratedGrads}_i(x) ::= \big(x_i - x_i'\big) \times \int_{\alpha=0}^{1} \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha$$

29

## Detected biases – Summary

| Dataset | Samples | Classes | Biases |
|---|---|---|---|
| Covid-19: X-ray [30] | 2,621 | 3 | Letters; Artefacts; Camera Positioning; Arm Positioning; Patient Characteristic (children) |
| Covid-19: CT Scan [143] | 646 / 556 | 2 | Out of lung detection |
| ChestXRay2017 [64] | 5,856 | 2 | Letters; Artefacts; Camera Positioning; Arm Positioning; Patient Characteristic (children) |
| ChestX-ray14 [134] | 112,120 | 14 | Letters; Artefacts; Horizontal Lines; Vertical Lines; Rotation; |

Table 4.2: Summary of biases detected in each dataset

- Models alone, without XAI **cannot** be trusted based on predictions