# Visual Explainable AI for Convolutional Neural Networks

ECE 1512: DIGITAL IMAGE PROCESSING AND APPLICATIONS

SEMESTER: WINTER 2022                    PROJECT "A" TUTORIAL

PRESENTER: AHMAD SAJEDI

Based on NeurIPS'2020/AAAI'21 tutorial on Explainable AI: https://explainml-tutorial.github.io/

# Outline

A. Tutorial on visual explainable AI (XAI)
  ◦ Motivation
  ◦ Primer on explainability in Artificial Intelligence (AI)
  ◦ Approaches for visual explanation generation

B. Project "A" Description
  ◦ Project Goal
  ◦ Datasets and Models
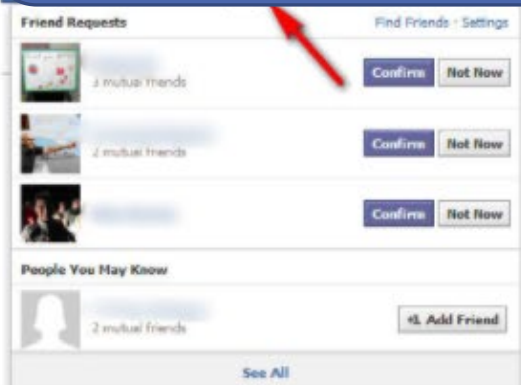  ◦ Evaluation Metrics

C. Your Questions!

# Outline

A. Tutorial on visual explainable AI (XAI)
  - Motivation
  - Primer on explainability in Artificial Intelligence (AI)
  - Approaches for visual explanation generation

B. Project "A" Description
  - Project Goal
  - Datasets and Models
  - Evaluation Metrics

C. Your Questions!

# Motivation



## Machine Learning is everywhere!

https://explainml-tutorial.github.io/
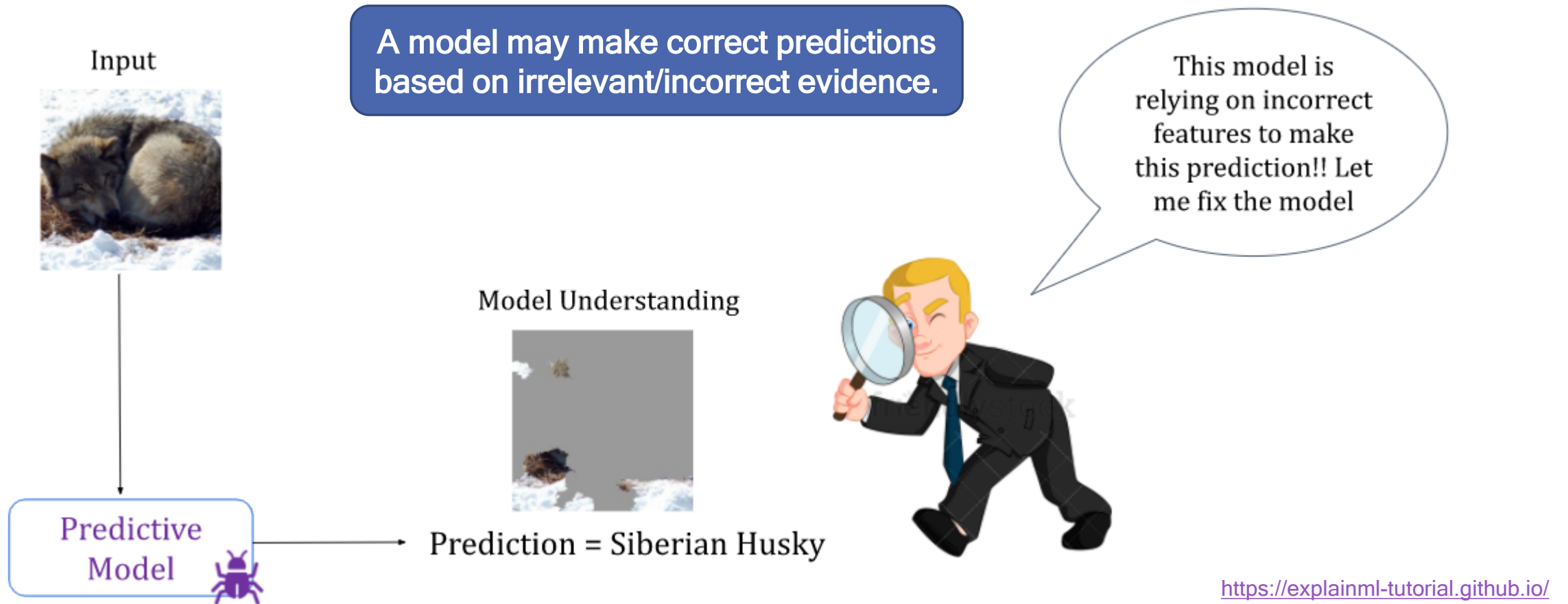
# Motivation

https://explainml-tutorial.github.io/

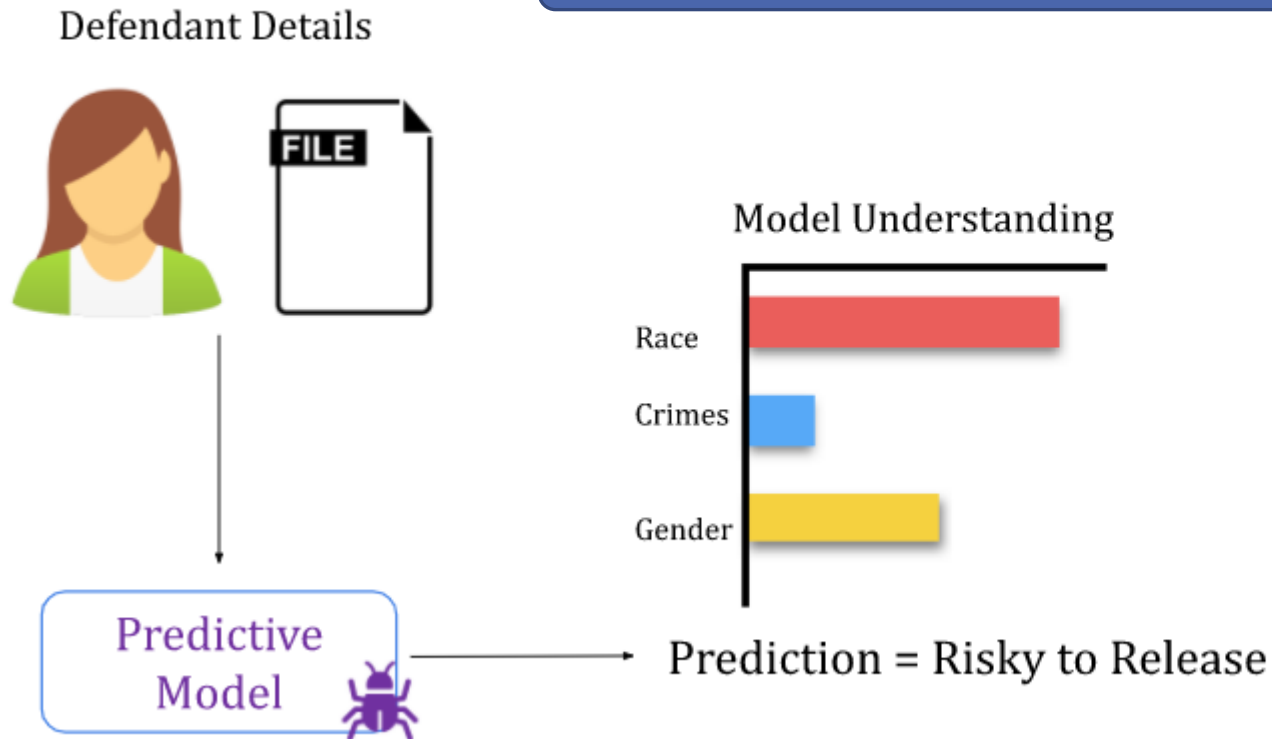Model understanding is critical in several domains, such as:

- Autonomous driving
- Healthcare
- Criminal justice

# Why model understanding?

# Why model understanding?

A model may suffer from dataset bias.

https://explainml-tutorial.github.io/

# Why model understanding?

A model may be reliable in some cases, while being unreliable in some other cases

**Patient Data**

25, Female, Cold
32, Male, No
31, Male, Cough
.
.
.
.

**Model Understanding**

If gender = female,
  if ID_num > 200, then sick

If gender = male,
  if cold = true and cough = true, then sick

**Predictions**

Healthy
Sick
Sick
.
.
Healthy
Healthy
Sick

Predictive Model

This model is using irrelevant features when predicting on female subpopulation. I should not trust its predictions for that group.

https://explainml-tutorial.github.io/

# Why model understanding?

Utility:

o Debugging

o Detecting dataset biases

o Assessing the mode's applicability in real world

o Realizing if and when to trust model predictions

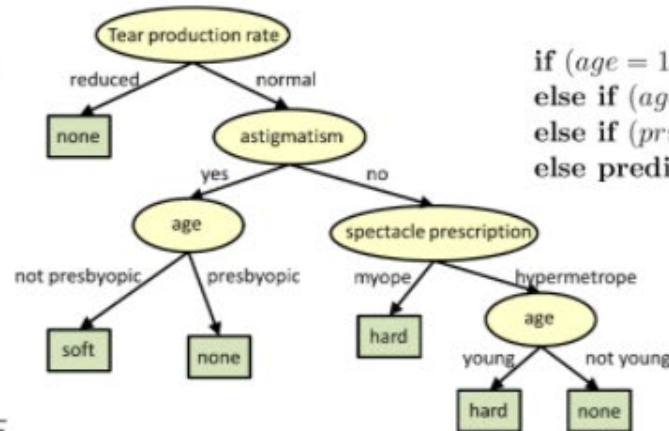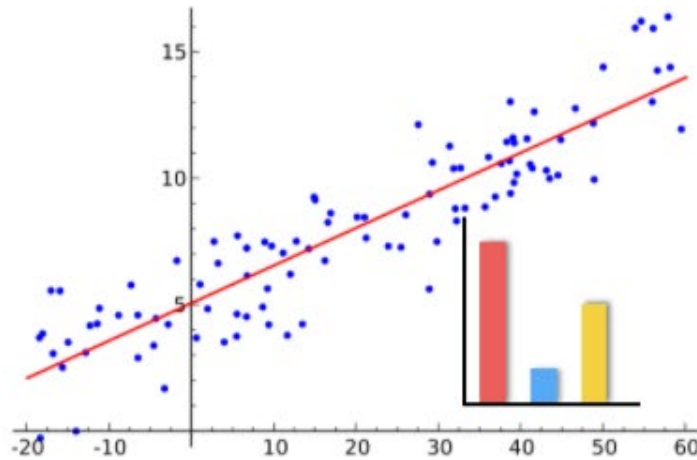# Primer on Explainability in Artificial Intelligence (AI)

# How to achieve model understanding?

# Primer on Explainability in Artificial Intelligence (AI)

How to achieve model understanding?

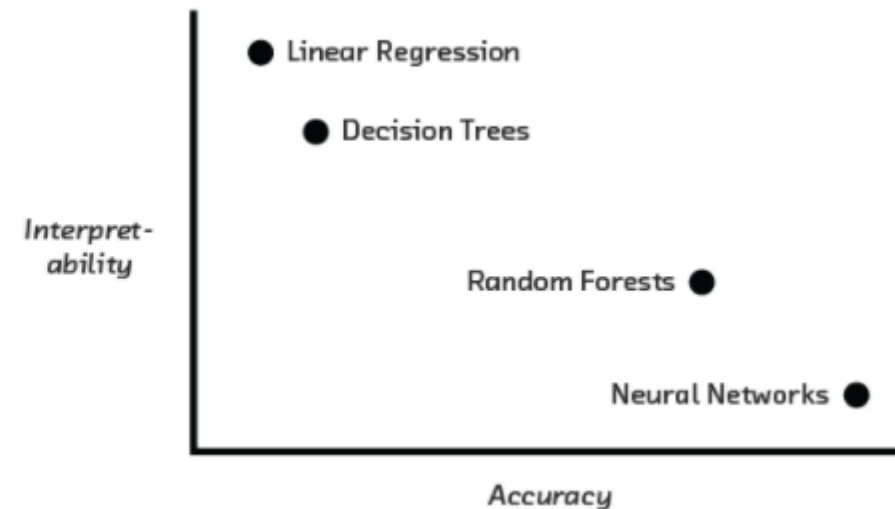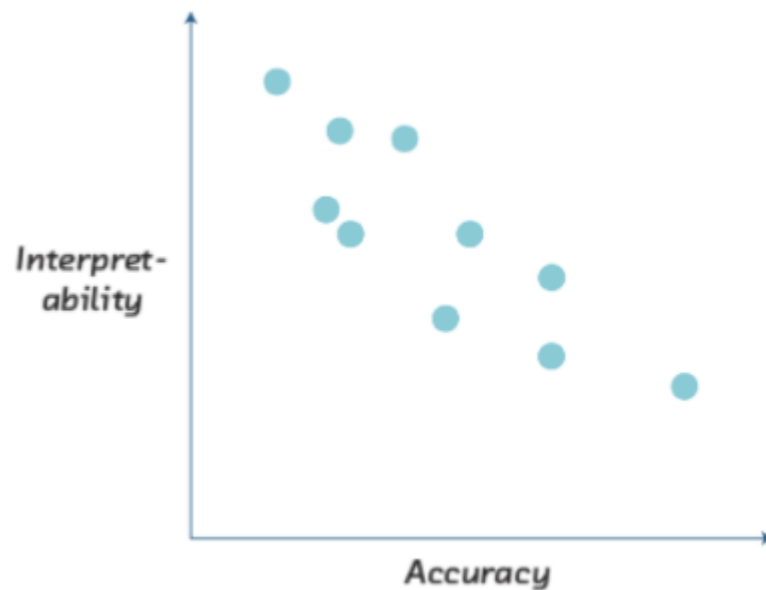Option 1: Build inherently interpretable models (e.g., decision trees, linear classifiers).



https://explainml-tutorial.github.io/

# Primer on Explainability in Artificial Intelligence (AI)

> **Drawbacks of the option 1:**
> A trade-off <u>may</u> exist between interpretability and prediction accuracy.
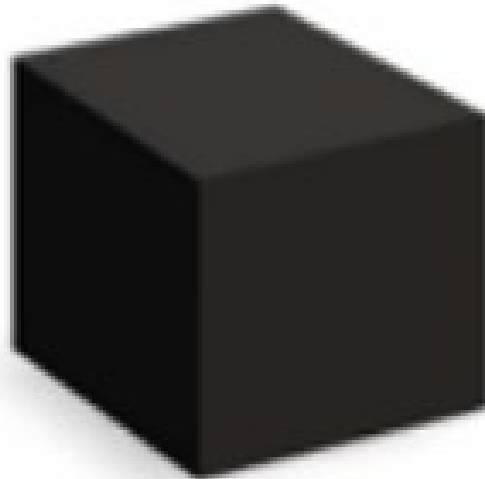


https://explainml-tutorial.github.io/

# Primer on Explainability in Artificial Intelligence (AI)

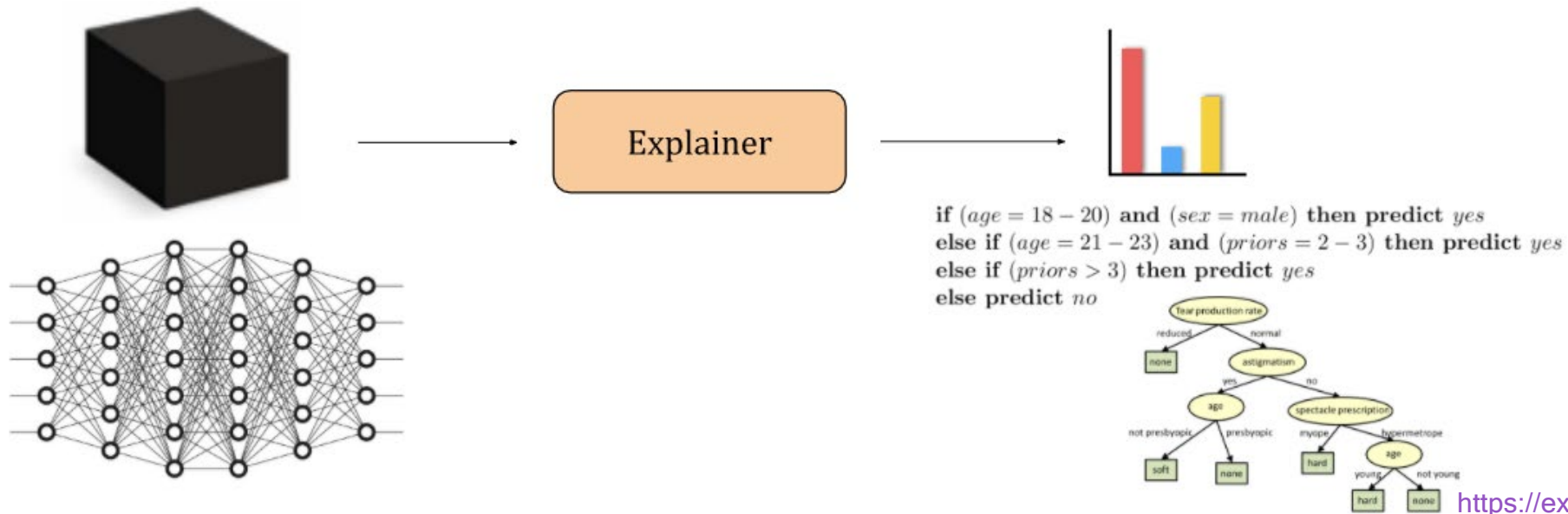If you an build interpretable models that are adequately accurate, DO IT!

Otherwise, *post-hoc explanations* come to rescue!

# Primer on Explainability in Artificial Intelligence (AI)

How to achieve model understanding?

Option 2: Explain cumbersome built models in a *"post-hoc"* manner.


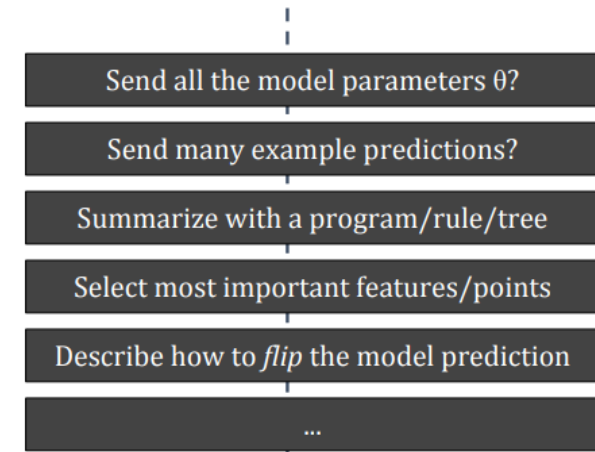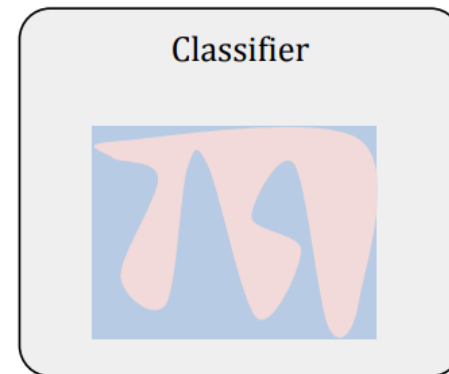
if $(age = 18 - 20)$ and $(sex = male)$ then predict $yes$
else if $(age = 21 - 23)$ and $(priors = 2 - 3)$ then predict $yes$
else if $(priors > 3)$ then predict $yes$
else predict $no$

https://explainml-tutorial.github.io/

# What is an explanation?

**Definition:**
- Interpretable description of the model's behavior

**Properties:**
- *Faithful*: from the model's end
- *Understandable*: from the user's end
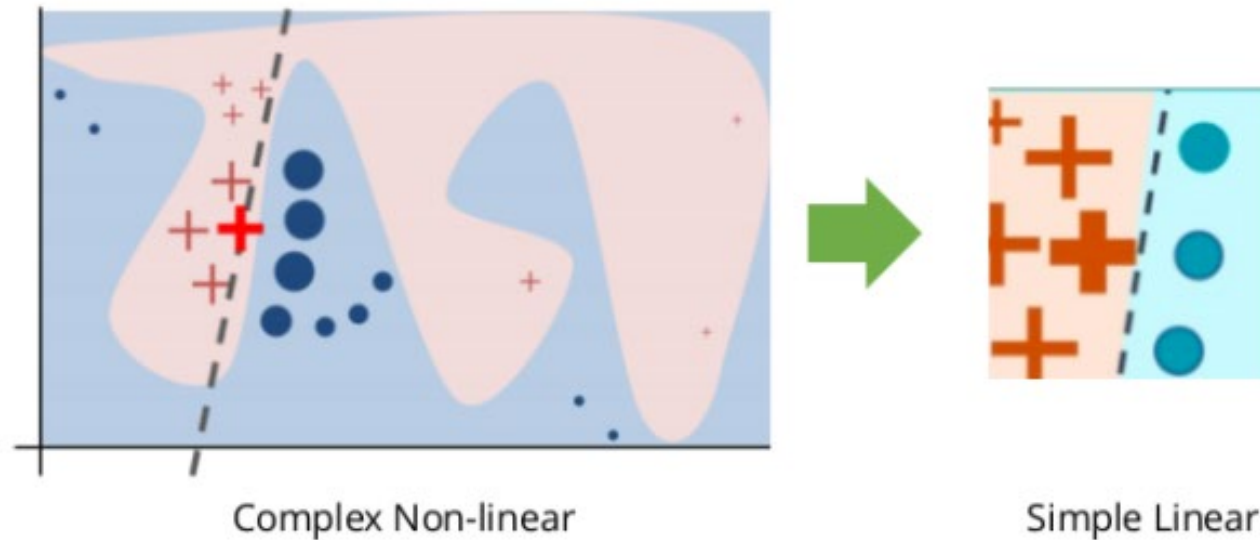


Classifier

Send all the model parameters θ?
Send many example predictions?
Summarize with a program/rule/tree
Select most important features/points
Describe how to *flip* the model prediction
...

User

https://explainml-tutorial.github.io/

# Local Explanation vs. Global Explanations



Complex Non-linear

Simple Linear

Local Explanations: Interpretable description of the model's behavior *in a target neighborhood.*

# Local Explanation vs. Global Explanations

**Local Explanations**

Explain individual predictions

Bring to light biases in the *local neighborhood* of a given instance

Assess if the correct predictions are made based on correct evidence

**Global Explanations**

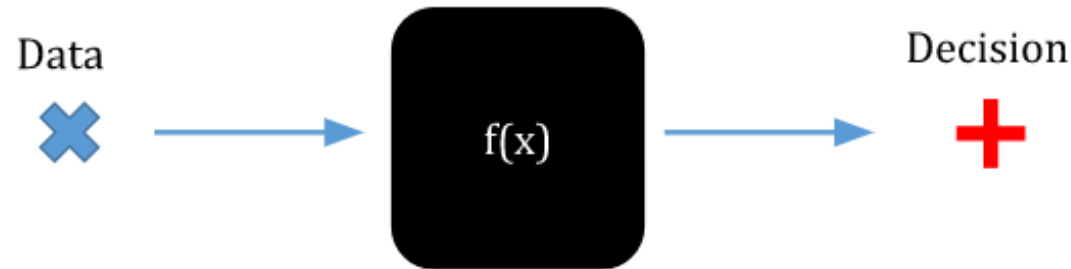Explain complete behavior of the model

Shed light of *big picture* biases affecting larger subgroups

Assess the suitability of the model for deployment in a high level

# Model-Agnostic approaches

No access to the internal structure…

Data → f(x) → Decision

Not restricted to specific models

Practically easy: not tied to PyTorch, Tflow, etc.

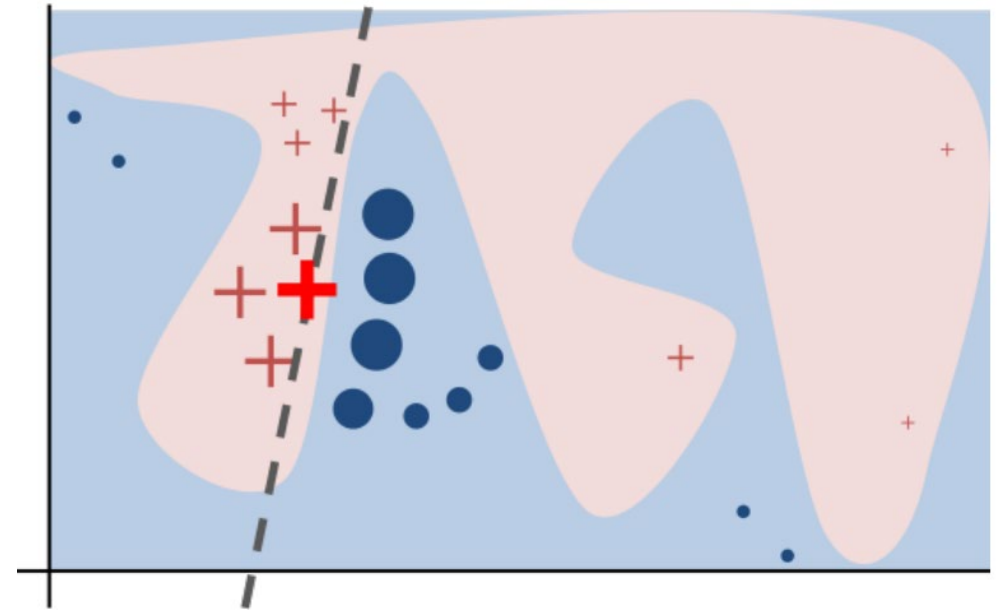Study models that you don't have access to!

# LIME

**Definition:**

- Local Interpretable Model-Agnostic Explanations

**Idea:**

- Identifying the important dimensions and quantifying their relative importance.

**Utility:**

- Initially proposed for tabular data, but also applicable to image/text data.

Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "" Why should i trust you?" Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016.

# LIME



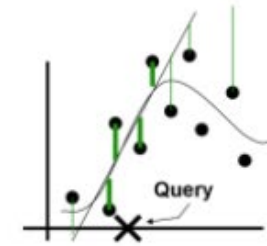Original Image

P(labrador) = 0.21

LIME is quite customizable:
- How to perturb?
- Distance/similarity?
- How *local* you want it to be?
- How to express explanation

| Perturbed Instances | P(Labrador) |
|---|---|
| | 0.92 |
| | 0.001 |
| | 0.34 |

Locally weighted regression

Query

Explanation

https://explainml-tutorial.github.io/

# SHAP

**Definition:**

- Shapely Additive Explanations

**Idea:**

- Computing *marginal contribution* of each input feature towards the prediction, averaged over all possible permutations.

**Utility:**

- Initially proposed for tabular data, but also applicable to image/text data.



$$P(y) = 0.9$$

$$P(y) = 0.8$$

$$M(x_i, O) = 0.1$$

Lundberg, Scott M., and Su-In Lee. "A Unified Approach to Interpreting Model Predictions." *Advances in Neural Information Processing Systems* 30 (2017): 4765-4774.
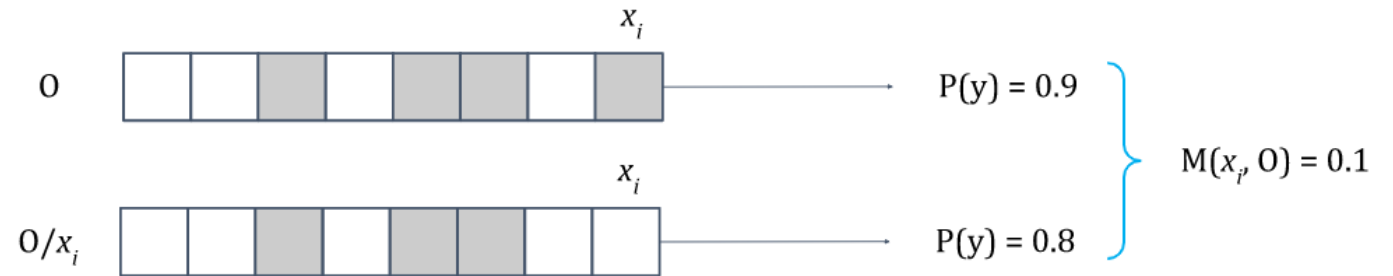
# SHAP

**Definition:**
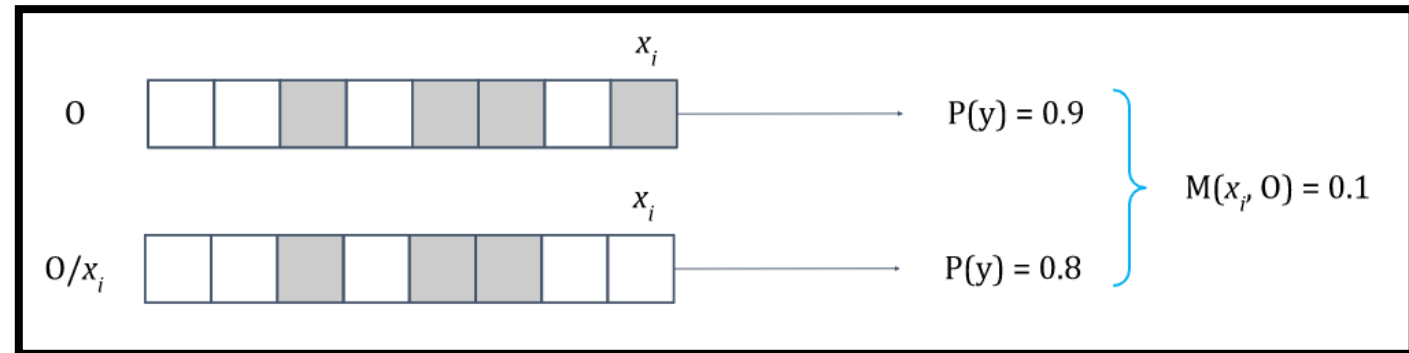- Shapely Additive Explanations

**Idea:**
- Computing *marginal contribution* of each input feature towards the prediction, averaged over all possible permutations.

**Utility:**
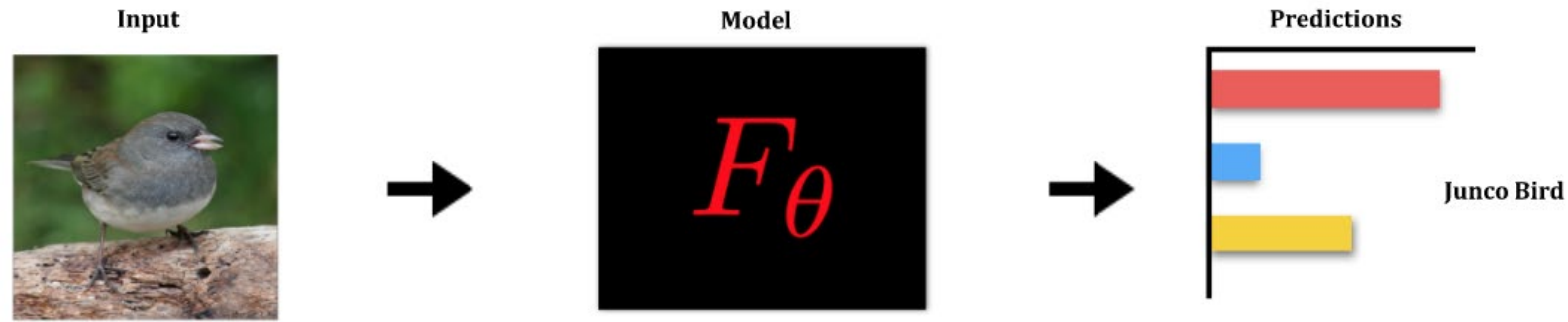- Initially proposed for tabular data, but also applicable to image/text data.

$$g(x_i) = \mathbb{E}_O \{M(x_i, O)\}$$



$g(x_i)$ :: Importance of the feature $x_i$

$M(x_i, O)$ :: Marginal contribution of $x_i$ over the permutation set $O$

# Backpropagation-based approaches

Assumption: we are explaining a *differentiable* model.



$$F : \mathbb{R}^d \to \mathbb{R}^c \qquad \text{Model}$$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - -

$$F_i : \mathbb{R}^d \to \mathbb{R} \qquad \text{class specific logit}$$

# Backpropagation-based approaches

Simplest approach: Vanilla Gradient

Input   Model   Predictions

Shortcomings:
o   Generating visually noisy and hard-to-understand explanation maps
o   Gradient saturation

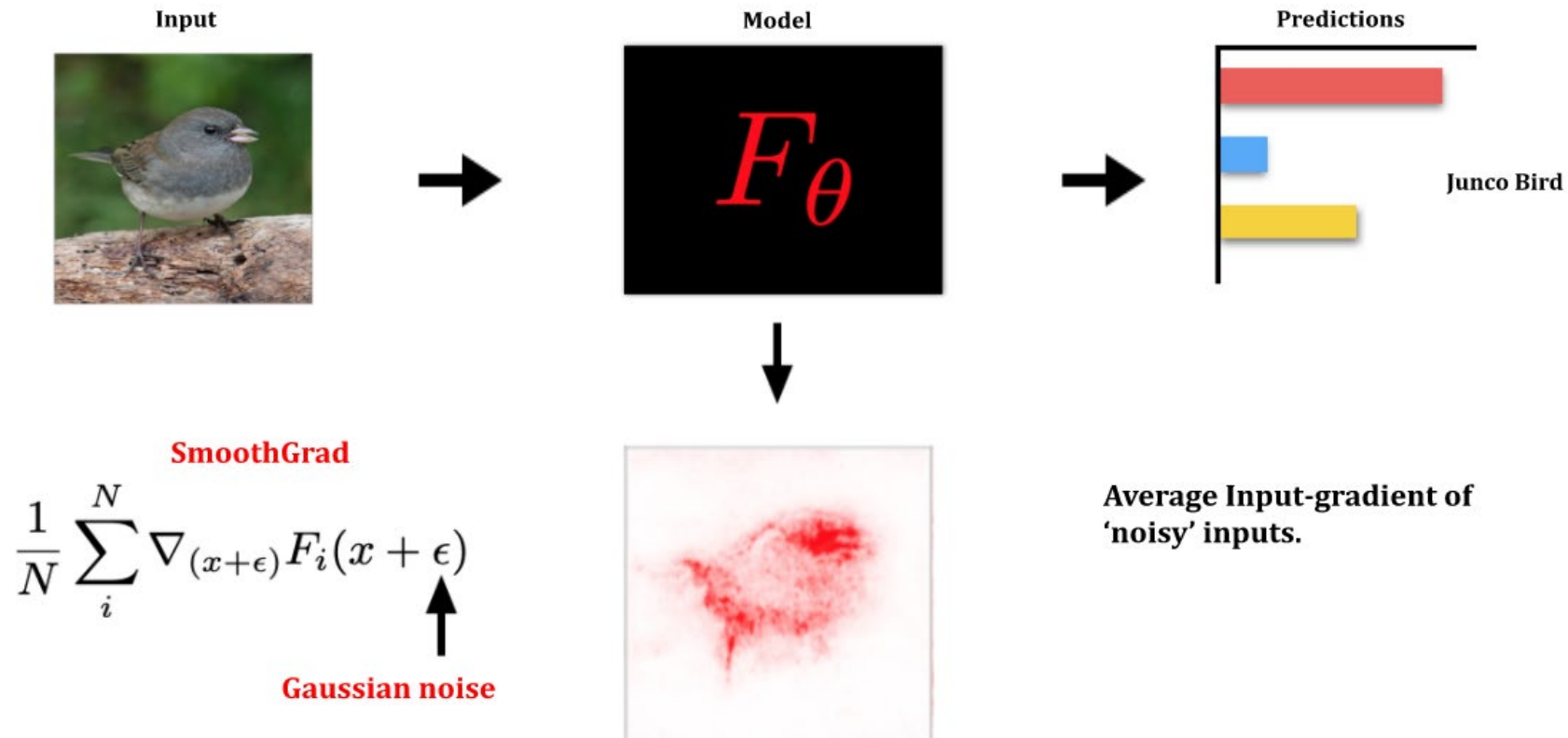Further backpropagation-based methods are proposed to address the above shortcomings.

Logit

Input

Simonyan, Karen, Andrea Vedaldi, and Andrew Zisserman. "Deep inside convolutional networks: Visualising image classification models and saliency maps

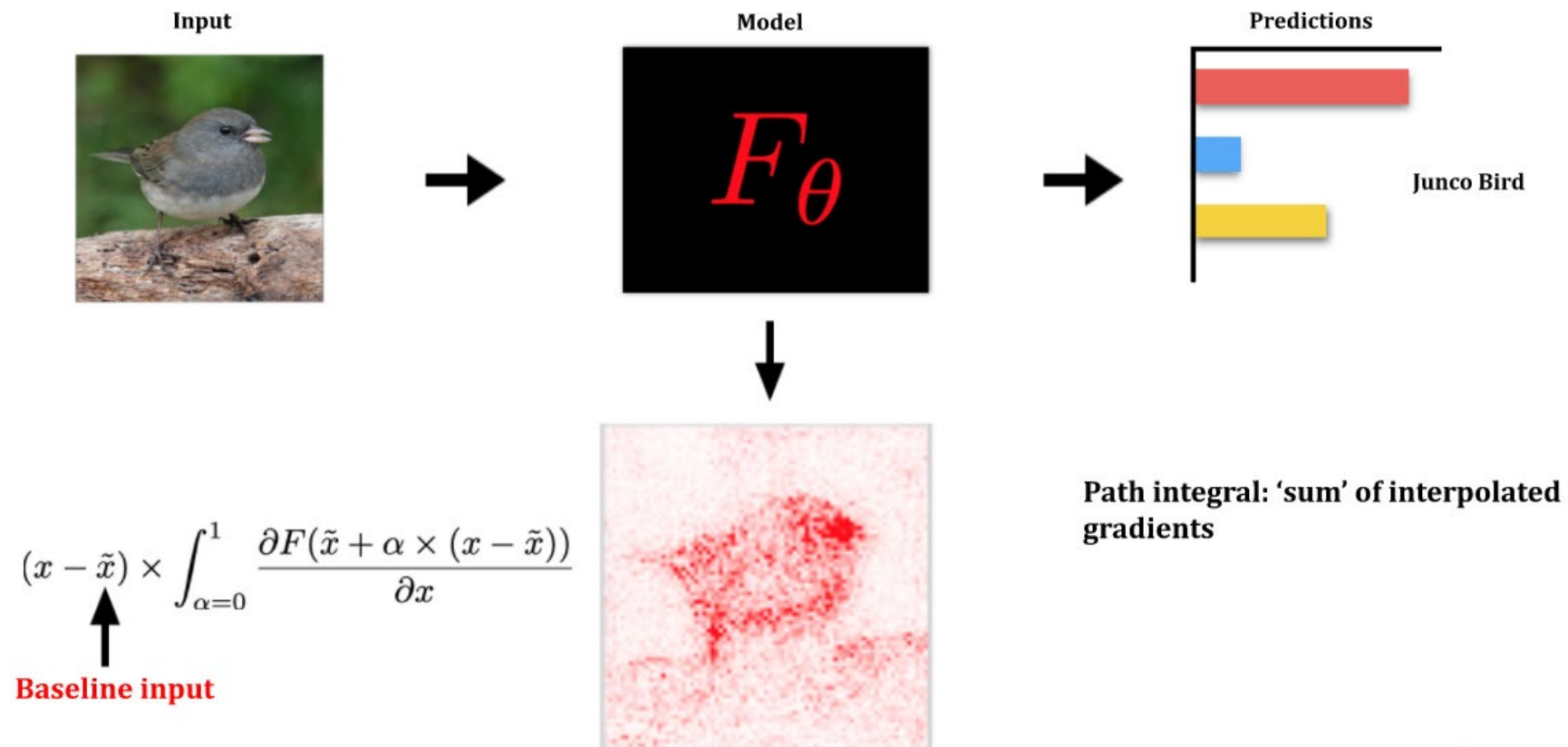# Backpropagation-based approaches

SmoothGrad



Smilkov, Daniel, et al. "Smoothgrad: removing noise by adding noise." *arXiv preprint arXiv:1706.03825* (2017).

# Backpropagation-based approaches

Integrated Gradients



$$(x - \tilde{x}) \times \int_{\alpha=0}^{1} \frac{\partial F(\tilde{x} + \alpha \times (x - \tilde{x}))}{\partial x}$$

**Baseline input**

Path integral: 'sum' of interpolated gradients

Sundararajan, M., Taly, A., & Yan, Q. (2017, July). Axiomatic attribution for deep networks. In *International Conference on Machine Learning* (pp. 3319-3328). PMLR.
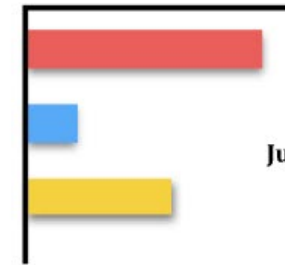
# Backpropagation-based approaches

Gradient ⊙ Input



Shrikumar, Avanti, Peyton Greenside, and Anshul Kundaje. "Learning important features through propagating activation differences." *International Conference on Machine Learning*. PMLR, 2017.

# : Class Activation Mapping

**What CAM stands for?**

o Class Activation Mapping

**Intuition:**
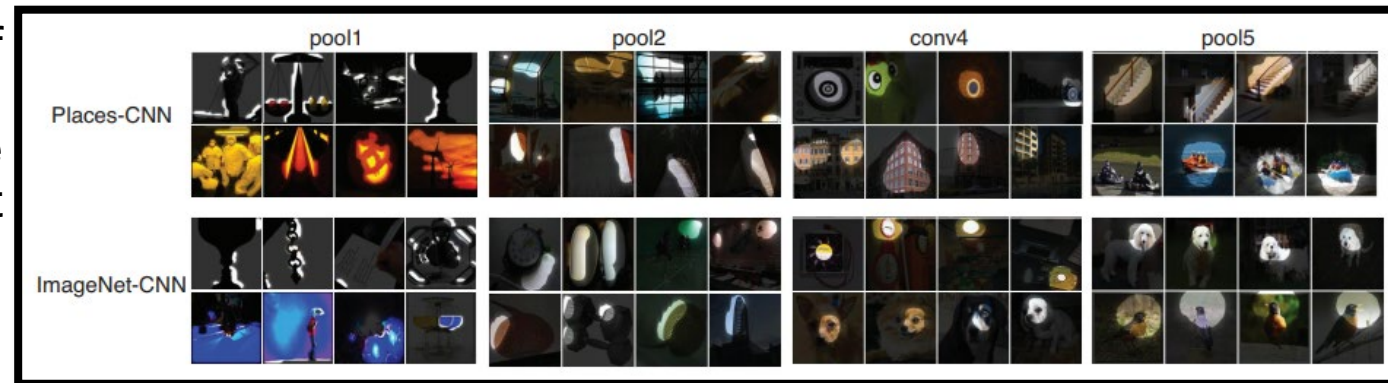
o The convolutional units of various layers of convolutional neural networks (CNNs) actually behave as object detectors despite no supervision on the location of the object was provided (Zhou et al. 2014)

**Utility:**

o Model-specific explanation method

o Specialized for convolutional neural networks (CNNs).



Zhou, Bolei, et al. "Object Detectors Emerge in Deep Scene CNNs." *ICLR*. 2015.
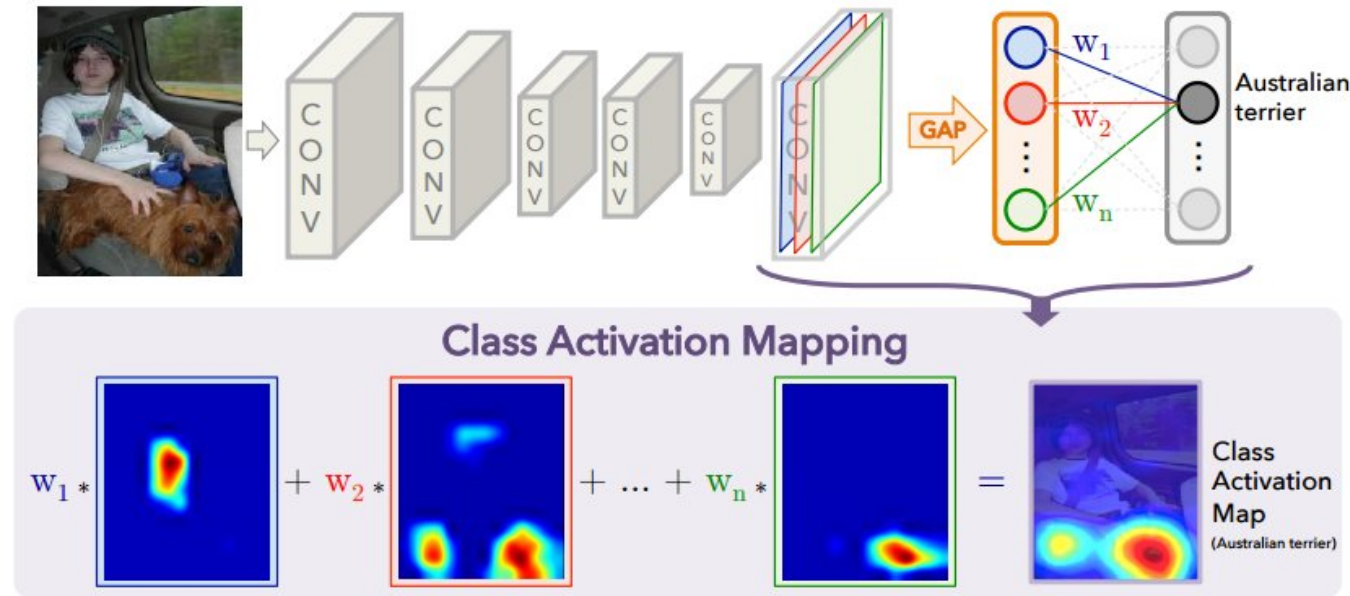
# CAM-based methods

**Aim**

- Visualizing the features extracted by the CNN.
- NOT explanation.

**Methodology**

- Replacing the classifier units of the target model with a Global Average Pooling (GAP) layer, followed by an output layer.
- Training the model.
- Linear combination of the feature maps in the last convolutional layer, using the weights in the output layer.
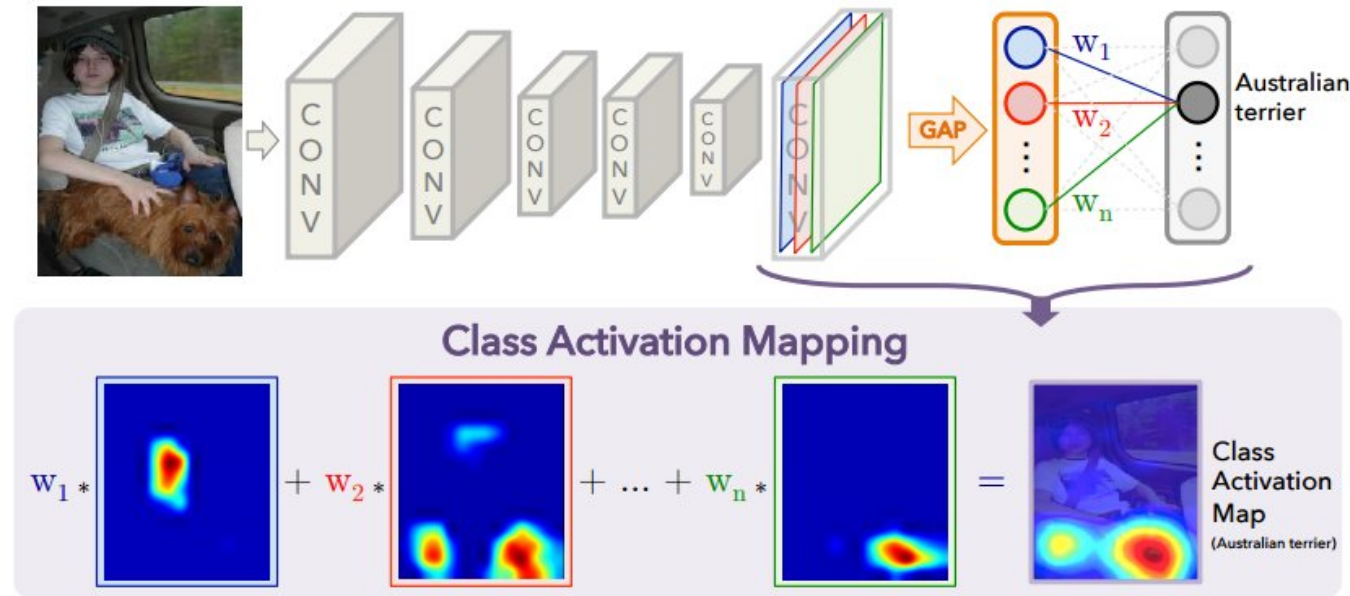


Zhou, Bolei, et al. "Learning deep features for discriminative localization." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2016.
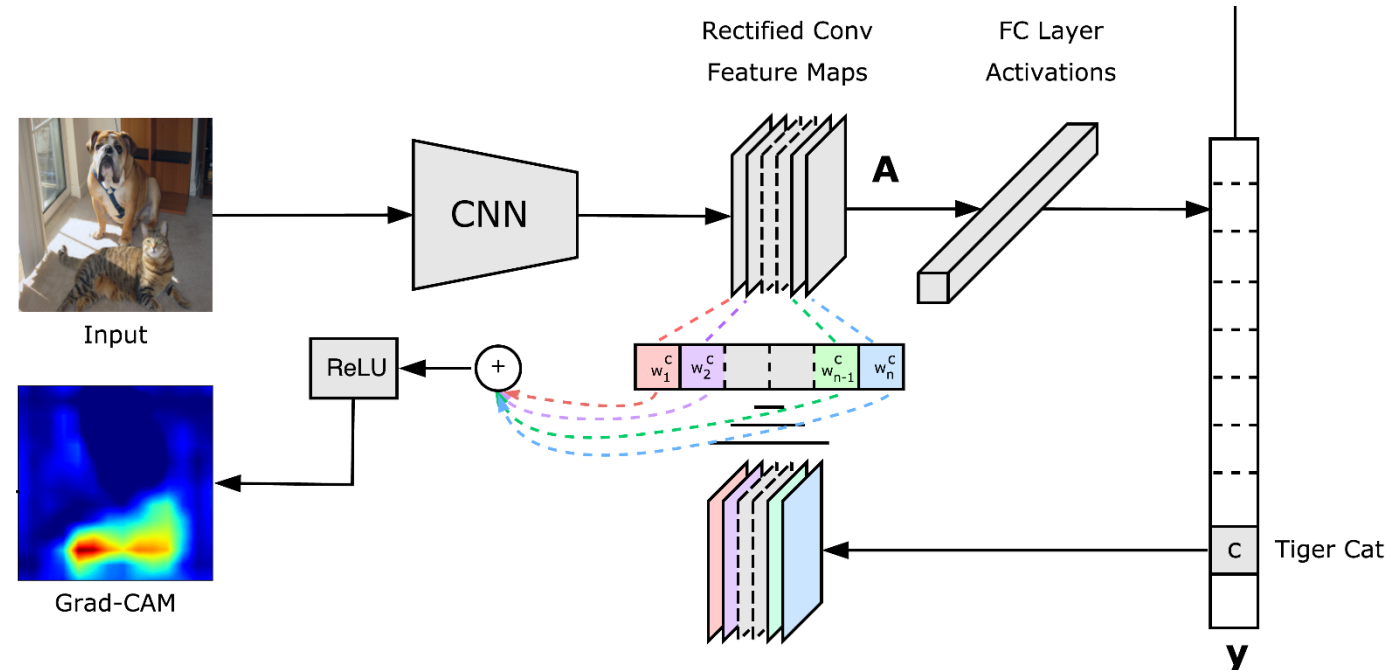
# CAM-based methods

**Limitations**

- Fails to visualize CNNs with complicated classifiers.
- Inapplicable for CNN explanation.
- Performance degradation because of the alteration in the CNN's classifier units.
- Inapplicable to the models trained for more complex image recognition tasks.

# CAM-based methods

## Grad-CAM

o Addresses the mentioned limitations in CAM.

o Uses "average gradient" values to score the feature maps in the last convolutional layer.

o Backpropagates the signal from the output, to calculate the average gradient values
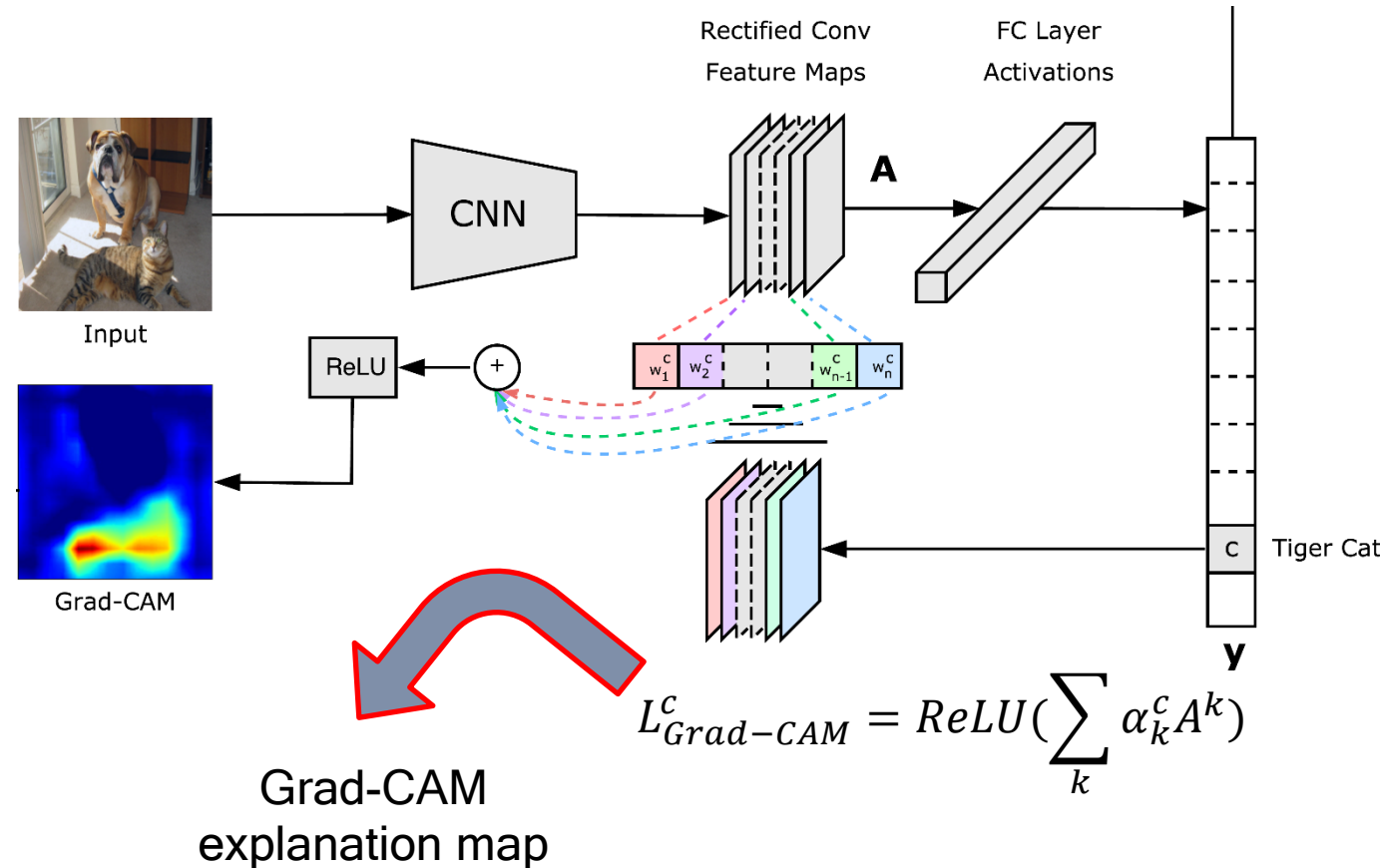
o Runs in a forward pass and a backward pass.



Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." *Proceedings of the IEEE international conference on computer vision*. 2017.

# Grad-CAM

## Grad-CAM

- Addresses the mentioned limitations in CAM.
- Uses "average gradient" values to score the feature maps in the last convolutional layer.
- Backpropagates the signal from the output, to calculate the average gradient values
- Runs in a forward pass and a backward pass.

Model's prediction for class $c$

Feature map values

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y_c}{\partial A_{ij}^k}$$

Grad-CAM explanation map

$$L_{Grad-CAM}^c = ReLU\left(\sum_k \alpha_k^c A^k\right)$$

Rectified Conv Feature Maps

FC Layer Activations

CNN

Input

ReLU

$+$

$w_1^c$ $w_2^c$ $w_{n-1}^c$ $w_n^c$

Grad-CAM

A

c    Tiger Cat

y

# Perturbation-based methods: RISE
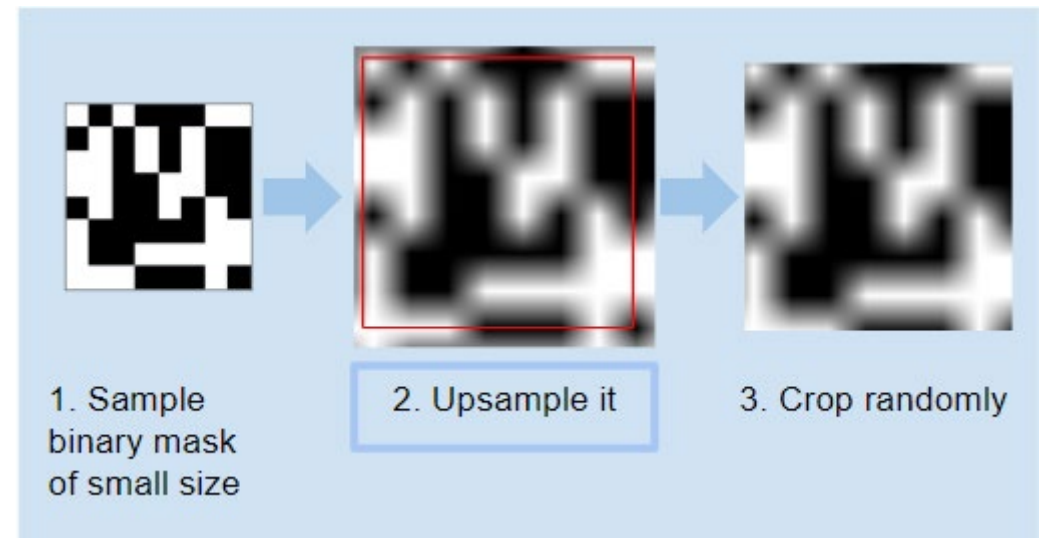
**Definition:**
- Randomized Input Sampling for Explanation

**Idea:**
- Feeding the target model with the perturbed copies of the input image.
- Perturbation is performed using a set of smooth *random masks*. (why smooth?)

**Utility:**
- Initially proposed for image data, but also applicable to tabular /text data.
- Model-agnostic explanation method.



1. Sample binary mask of small size
2. Upsample it
3. Crop randomly

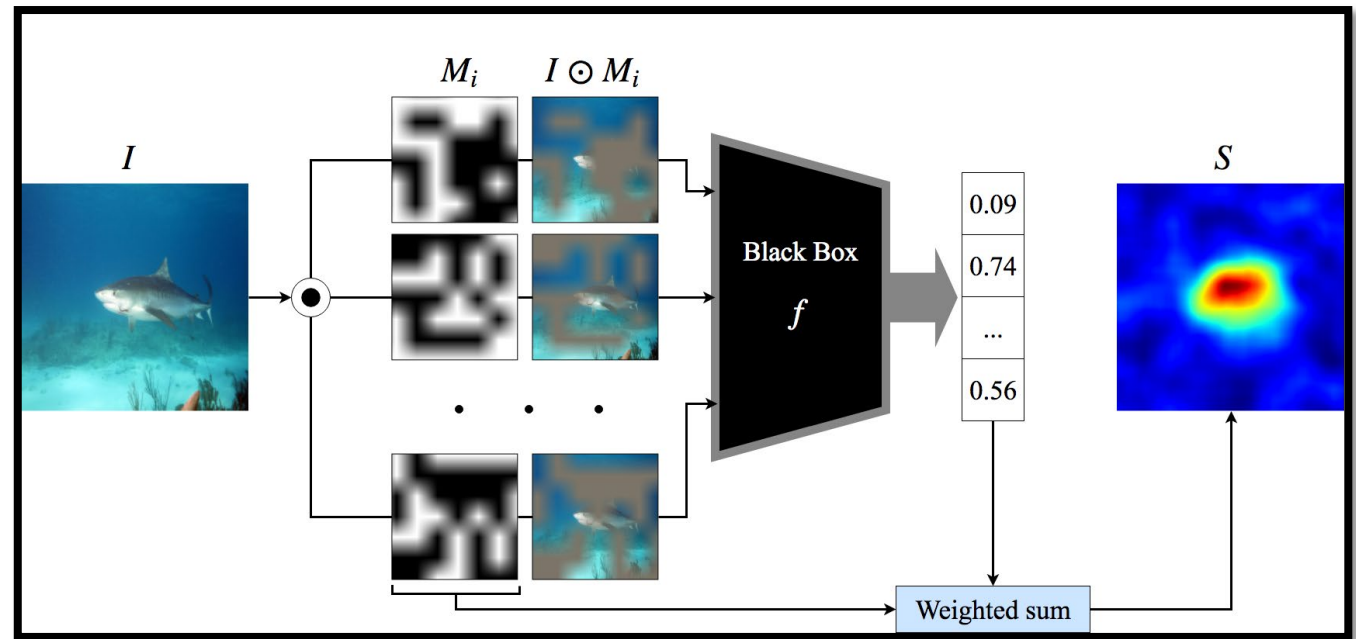https://cs-people.bu.edu/vpetsiuk/rise/#

# Perturbation-based methods: RISE

**Methodology:**
- Generating a set of random masks
- Masking the input with the random masks
- Feeding the model with the masked images
- Weighted linear combination of the masks.



Petsiuk, Vitali, Abir Das, and Kate Saenko. "Rise: Randomized input sampling for explanation of black-box models." *arXiv preprint arXiv:1806.07421* (2018).
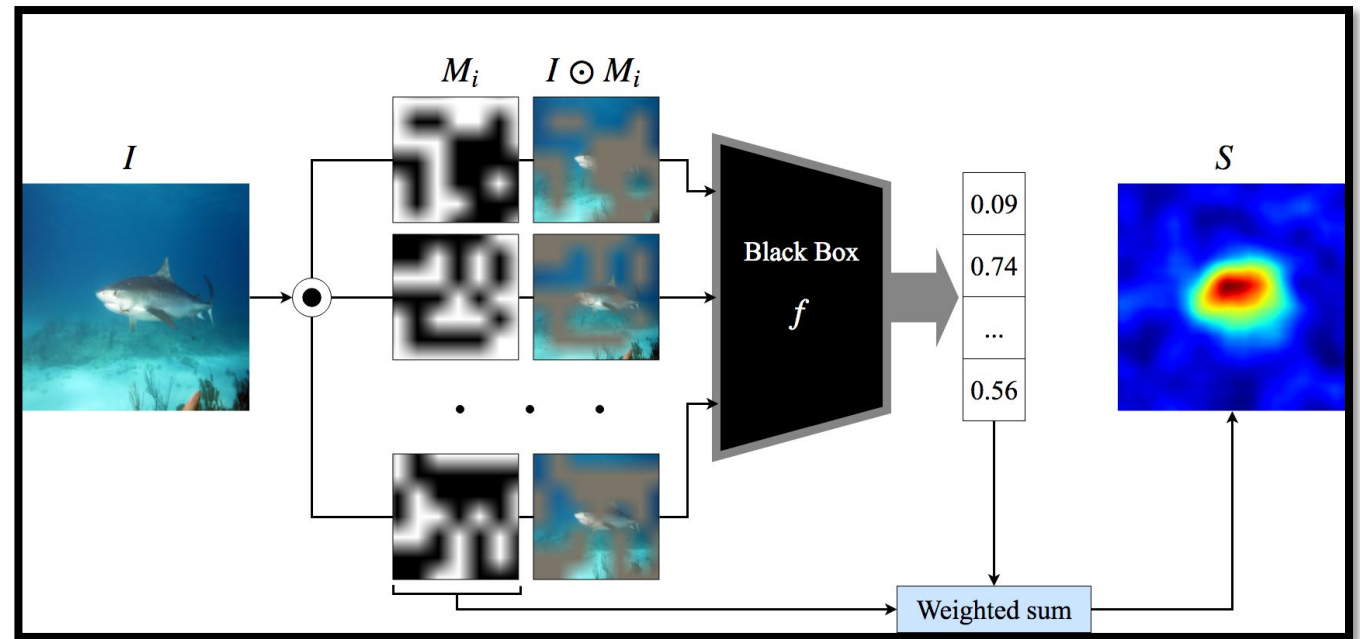
# Perturbation-based methods: RISE

**Pros:**

o Applicability of the method to the AI models beyond the family of CNNs.

o Shows the superior preciseness of perturbation rather than backpropagation, in forming explanation map.

**Cons:**

o Low visual quality of RISE explanation maps.

o Increase of failure chance, while dealing with small object instances

o Slow runtime, as it passes numerous (4000-8000) masked images through a model.
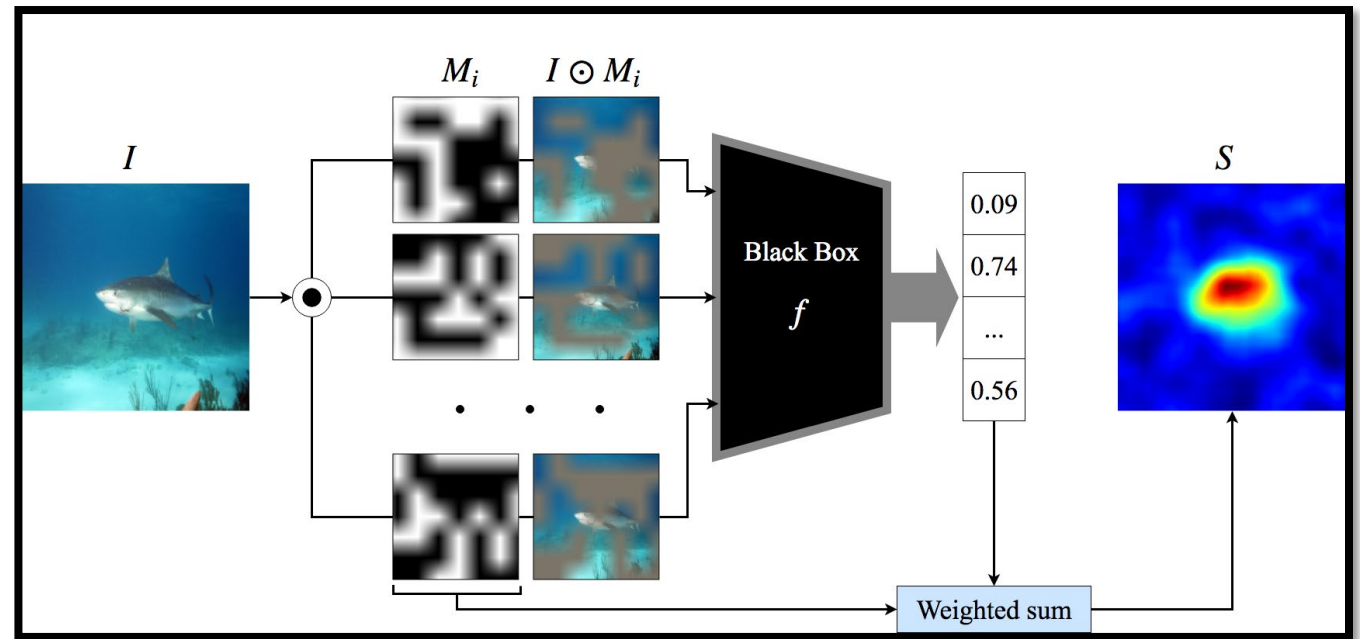
# Perturbation-based methods: RISE

**Pros:**

o Applicability of the method to the AI models beyond the family of CNNs.

o Shows the superior preciseness of perturbation rather than backpropagation, in forming explanation map.

**Cons:**

o Low visual quality of RISE explanation maps.

o Increase of failure chance, while dealing with small object instances

o Slow runtime, as it passes numerous (4000-8000) masked images through a model.
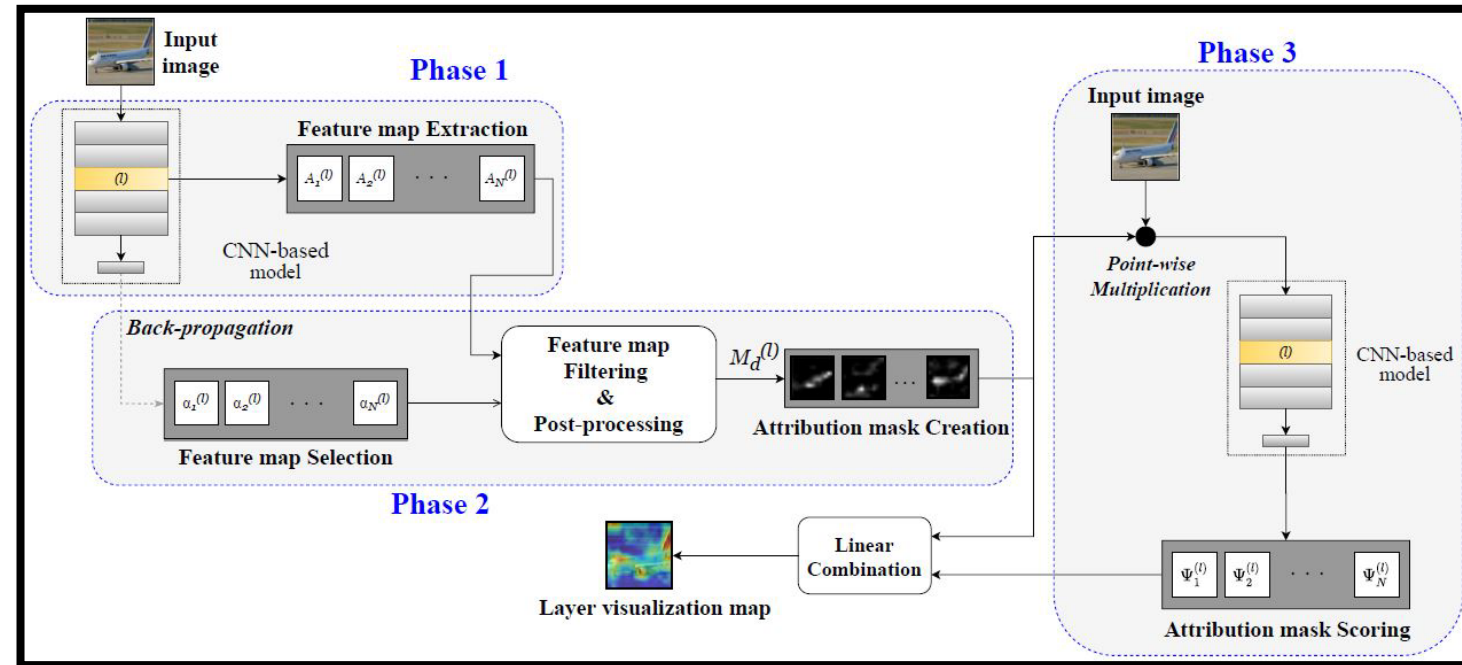
# Perturbation-based methods: SISE

**Definition:**
- Semantic Input Sampling for Explanation

**Idea:**
- Feeding the target model with the perturbed copies of the input image.
- Random masks in RISE are replaced with *attribution masks.*

**Utility:**
- Model-specific explanation method
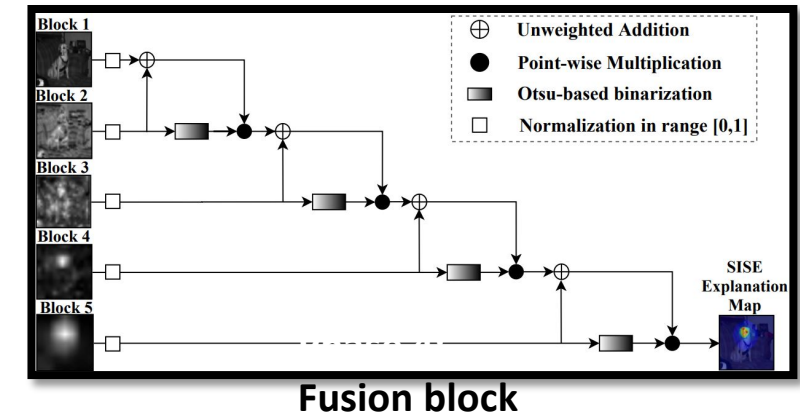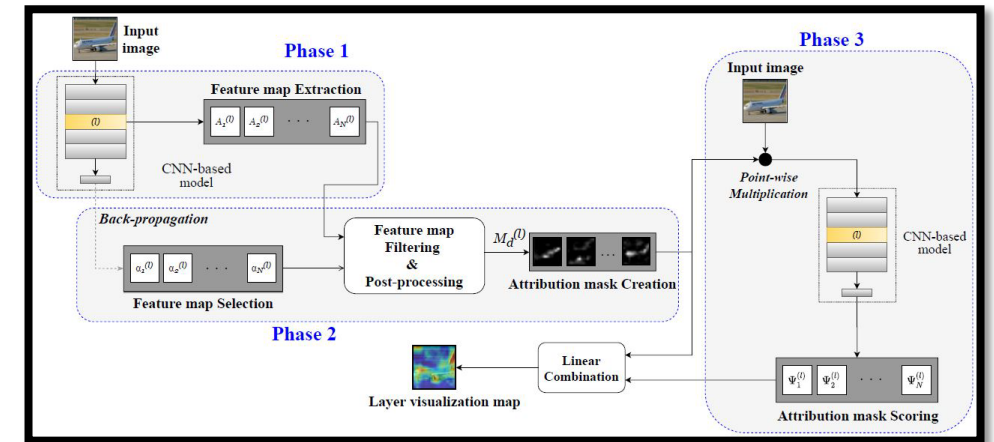- Specialized for convolutional neural networks (CNNs).



Sattarzadeh, Sam, et al. "Explaining Convolutional Neural Networks through Attribution-Based Input Sampling and Block-Wise Feature Aggregation." *arXiv preprint arXiv:2010.00672* (2020).

# Perturbation-based methods: SISE

**Methodology:** Consists four consecutive phases
- Feature map extraction
- Feature map selection
- Attribution mask scoring
- Feature aggregation



- The first phases are applied on the last layers in all convolutional blocks (why?). Corresponding to each layer, the third phase outputs a 2-dimensional map called *visualization map*.

- The visualization maps are **aggregated** in the last phase to form the desires explanation map.
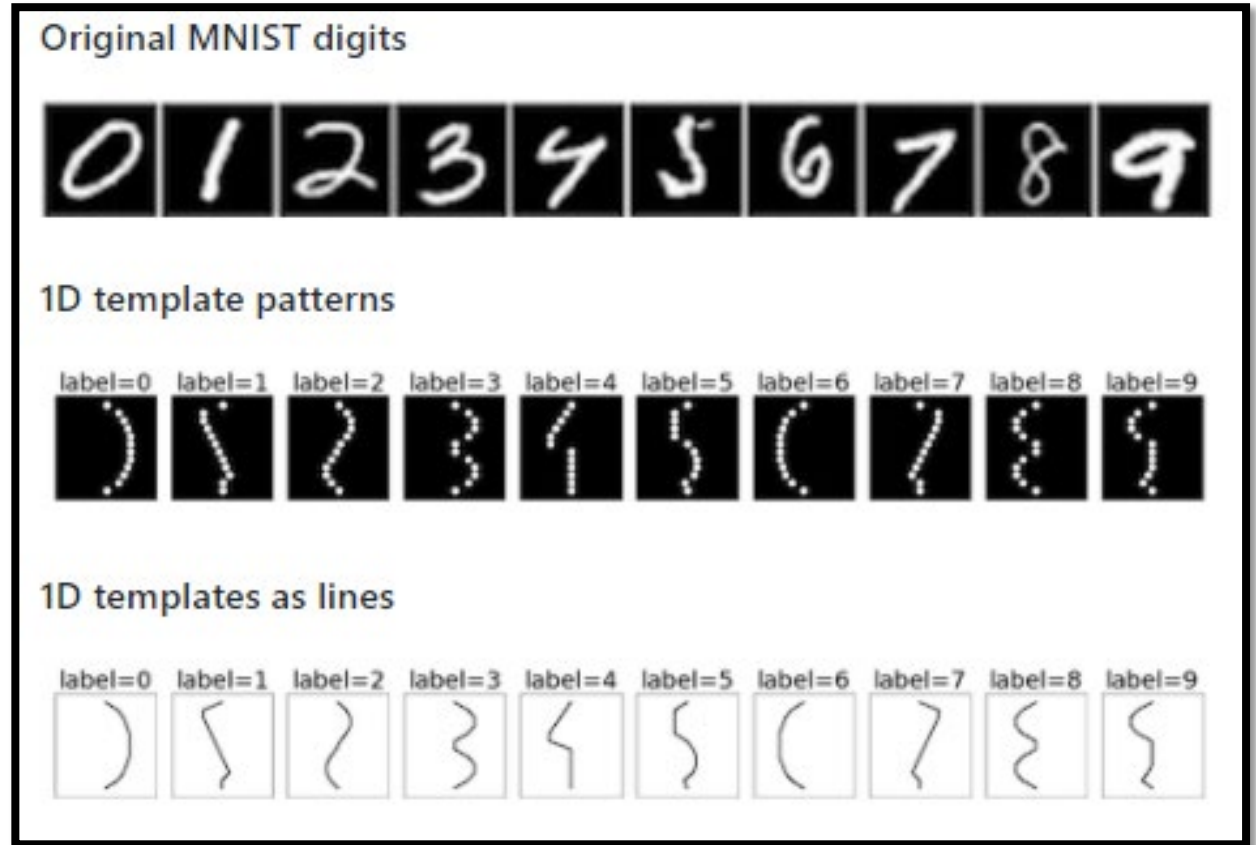


**Fusion block**

# Outline

A. Tutorial on visual explainable AI (XAI)
- Motivation
- Primer on explainability in Artificial Intelligence (AI)
- Approaches for visual explanation generation

B. Project "A" Description
- Project Goal
- Datasets and Models
- Evaluation Metrics

C. Your Questions!

# Datasets: MNIST-1D

o A 1-dimensional analogue for the famous **MNIST** digit classification dataset.

o Small-scale, low-memory, synthetic

o 4000 train data + 1000 test data (partitioned by the dataset promoters).

o Each data is generating by applying random transformations to one of the 1-D templates.

o Transformations include:

  o Random padding

  o Random translation

  o Random 1-D shearing

  o Gaussian noise addition.

https://greydanus.github.io/2020/12/01/scaling-down/

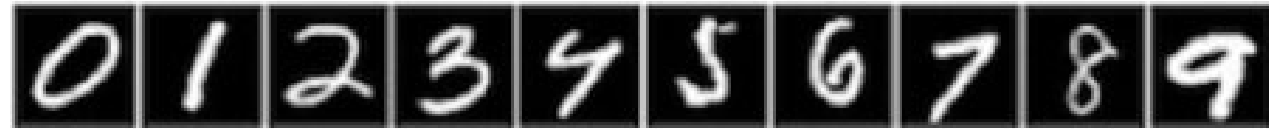Greydanus, Sam. "Scaling* down* Deep Learning." *arXiv preprint arXiv:2011.14439* (2020).

# Project Goal

- Inspect the ability of one/two explanation methods(s) in providing correct explanations for CNNs in two different scenarios.
- **Specifications:**
  - Datasets:
    - MNIST-1D: Generic digit pattern classification
    - HMT: Histopathological tissue classification
  - Models:
    - A shallow CNN trained on the MNIST-1D dataset.
    - A VGG-7 network trained on the HMT dataset.
  - Evaluation metrics:
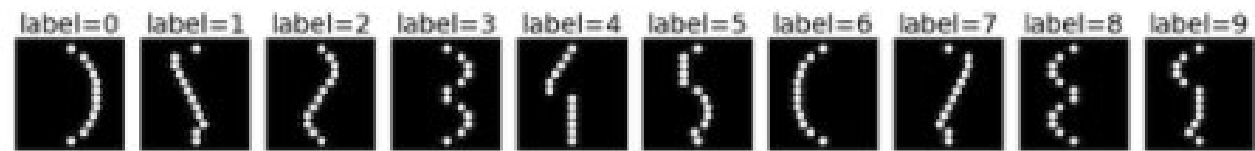    - Drop%
    - Increase%

# Datasets: MNIST-1D

- A 1-dimensional analogue for the famous **MNIST** digit classification dataset.
- Small-scale, low-memory, synthetic
- 4000 train data + 1000 test data (partitioned by the dataset promoters).
- Each data is generating by applying random transformations to one of the 1-D templates.
- Transformations include:
  - Random padding
  - Random translation
  - Random 1-D shearing
  - Gaussian noise addition.



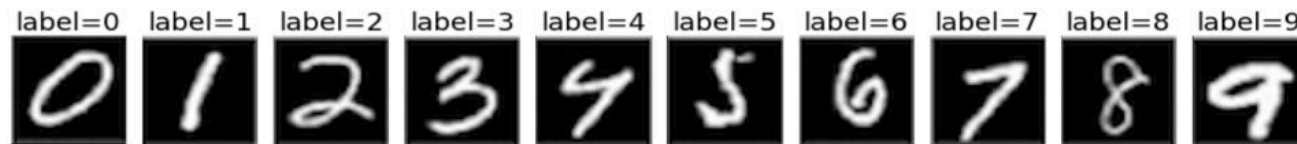Original MNIST digits

1D template patterns
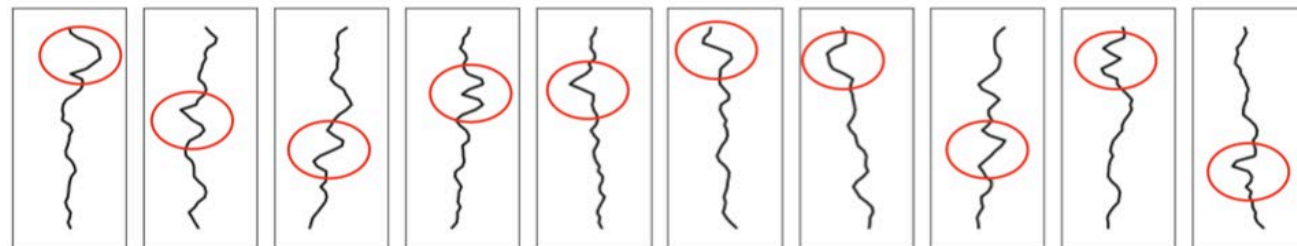
1D templates as lines

# Datasets: MNIST-1D
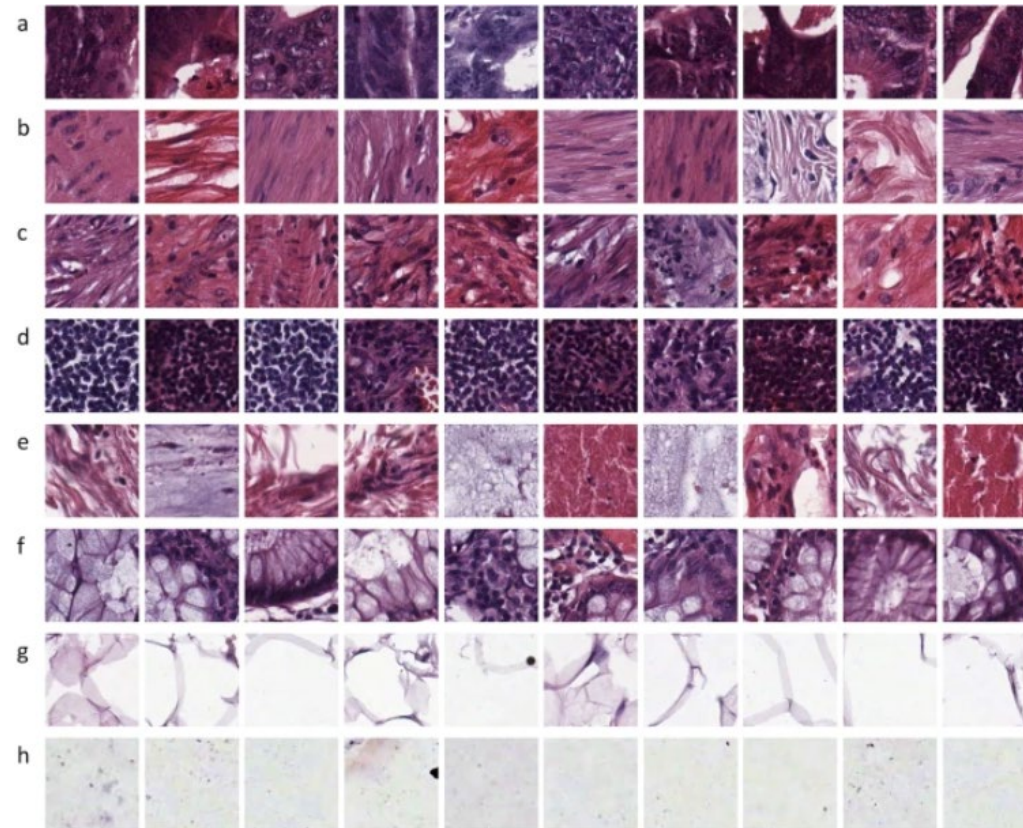
# Datasets: HMT

o Multi-class texture analysis in colorectal cancer histology.

o An equally-balanced dataset.

o classes:
- o (a) tumor epithelium,
- o (b) simple stroma,
- o (c) complex stroma,
- o (d) immune cell conglomerates,
- o (e) debris and mucus,
- o (f) mucosal glands,
- o (g) adipose tissue,
- o (h) background.

o 4504 train images + 496 test images.



https://www.nature.com/articles/srep27988

Greydanus, Sam. "Scaling* down* Deep Learning." *arXiv preprint arXiv:2011.14439* (2020).

# Evaluation Metrics

Properties of a "good" explanation method:

○ Faithfulness: correlataion of the generated explanation maps with the behavior of the model (the model's end)

○ Understandability: Clarity of the generated explanation maps (the user's end)

Model truth-based metrics:

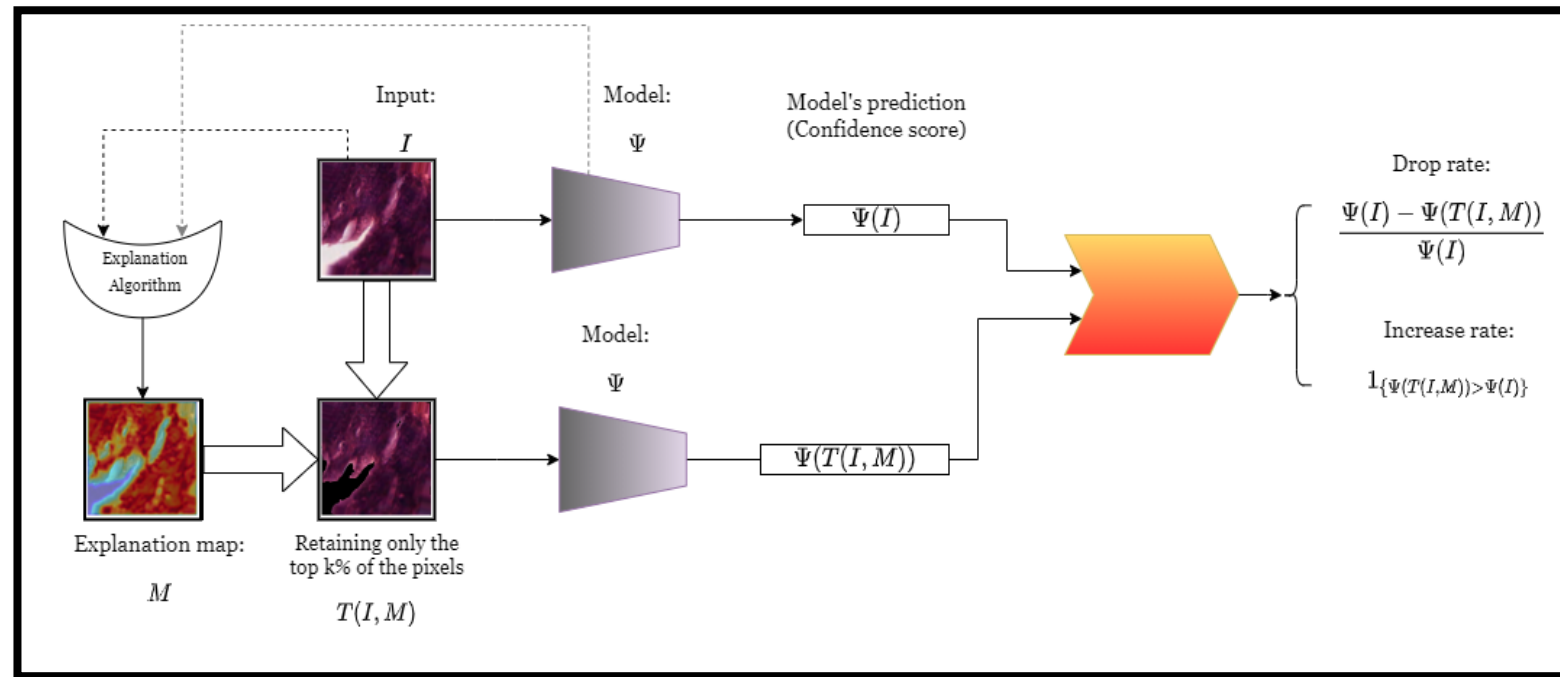○ Aim to evaluate the faithfulness of the explanation methods.

The focus of Project A

# Evaluation Metrics

**Drop%**

o When the most important features are retained and the other features are removed, the model's output should not "drop" significantly.
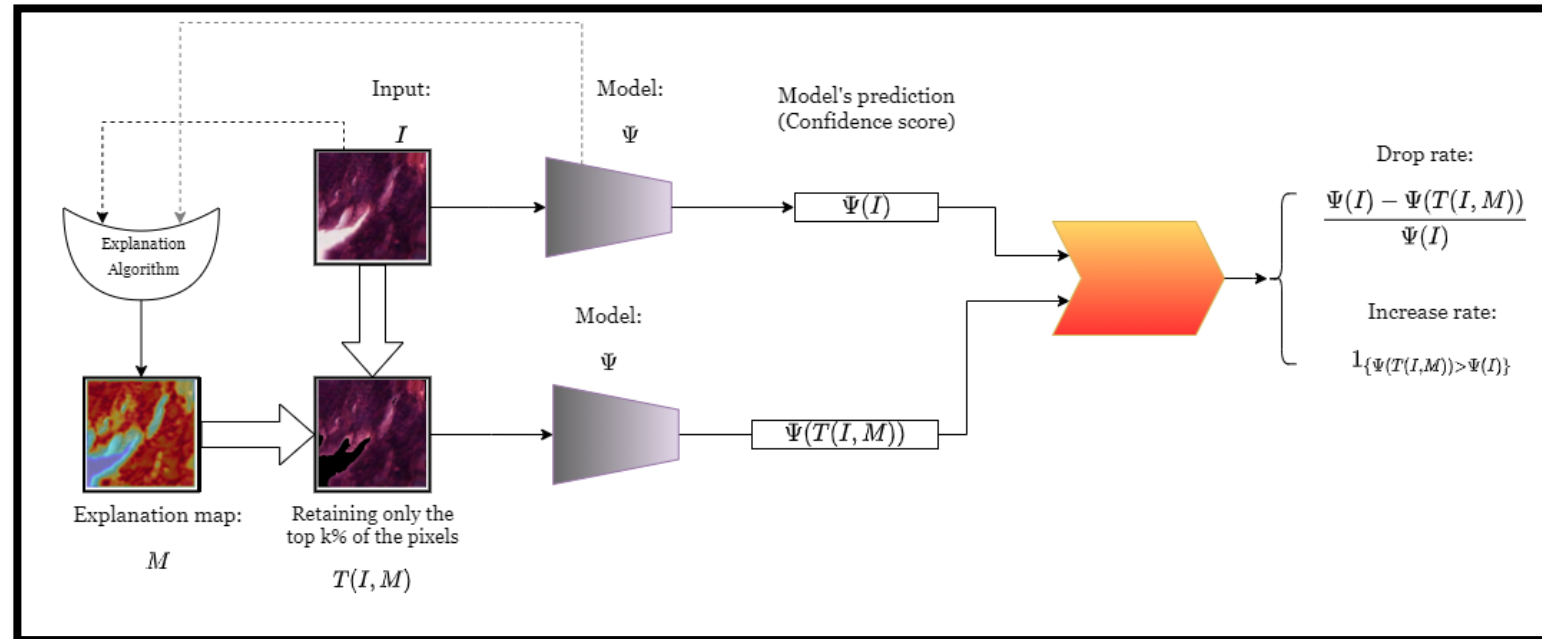
**Increase**

o When the most important features are retained and the other features are removed, the model's output may increase in some cases

# Evaluation Metrics

How Drop/Increase% evaluate your explanation algorithm?

o Retain the top-most $k\%$ pixels of the image, and set the other pixels to black.

o Pass the main image and the masked image to the model.

o Compare the model's output in the two cases.



$$k= \begin{cases} \text{The MNIST-1D dataset: } 30\% \\ \\ \text{The HMT dataset: } 90\% \end{cases}$$

# Outline

A. Tutorial on visual explainable AI (XAI)
   ◦ Motivation
   ◦ Primer on explainability in Artificial Intelligence (AI)
   ◦ Approaches for visual explanation generation

B. Project "A" Description
   ◦ Project Goal
   ◦ Datasets and Models
   ◦ Evaluation Metrics

C. Your Questions!

# THANK YOU
# Questions?