

# Learning Discriminative Representations to Interpret Image Recognition Models

Thèse de Doctorat

Felipe Torres Figueroa

École Centrale de Marseille

Laboratoire d'Informatique et de Systèmes (LIS)

Marseille, September 2024



# Table of Contents

- Introduction
- 1 Background
- 2 Opti-CAM: Optimizing saliency maps for interpretability
- 3 CA-Stream: Attention-based pooling for interpretable image recognition
- 4 A learning paradigm for interpretable gradients
- References

# Motivation

## Low Stakes

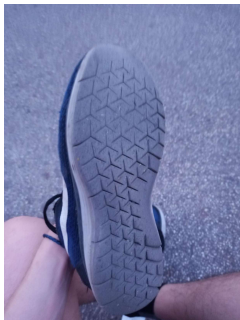
My go to exercise is running, **but...**

# Motivation

## Low Stakes

My go to exercise is running, **but...**

I think my running shoes  
are getting *worn*



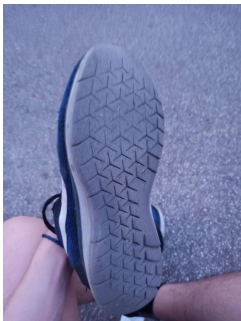




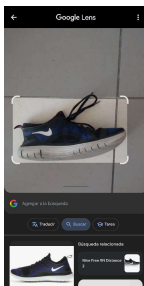
## Low Stakes

My go to exercise is running, **but...**

I think my running shoes  
are getting *worn*



I want a replacement,  
*but* I know about  
machines, not shoes!



## Nike Free RN Distance 2



My phone can  
*identify* my current shoes





# Motivation

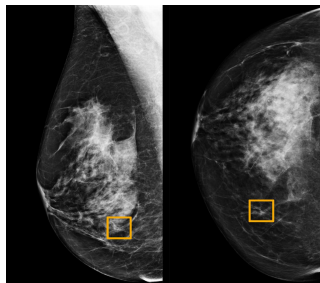
## Raising the stakes

Now let's consider riskier situations:

# Motivation

## Raising the stakes

Now let's consider riskier situations:



# Motivation

## Raising the stakes

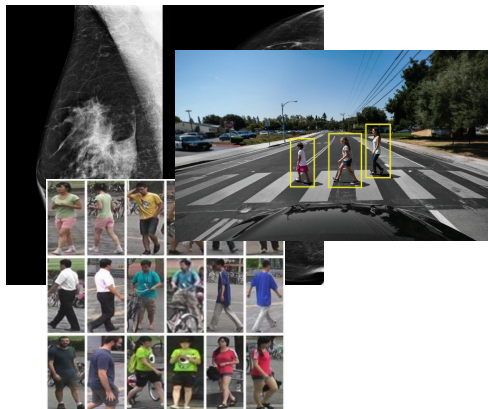
Now let's consider riskier situations:



# Motivation

## Raising the stakes

Now let's consider riskier situations:



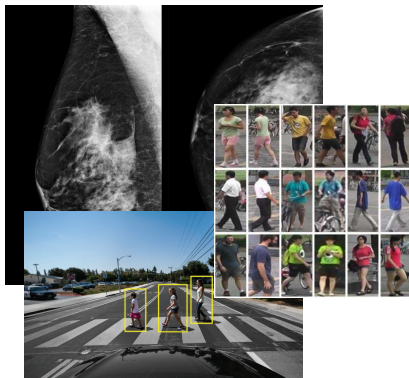
# Motivation

Straight to the point

- How do we **know how** a system works?



## Straight to the point



- How do we **know how** a system works?
- How do we **know how** safe a system is?

## Straight to the point



- How do we **know how** a system works?
- How do we **know how** safe a system is?
- If a system fails, **who** is accountable?

We must **understand** the behaviour of these models.



## Step by step

# Computation, Computer Vision and AI



## Explainable AI



## Thesis objectives

# Computation, Computer Vision and AI

## Computation

Better known as *Computer Science*.



*Alan Turing* forefather of current computer science.

# Computation, Computer Vision and AI

## Computation



*Alan Turing* forefather of current computer science.

Better known as *Computer Science*.

Study of:

- Algorithms.
- Data structures.
- Design of hardware and software.

# Computation, Computer Vision and AI

## Computer Vision

Replication of human vision capabilities.

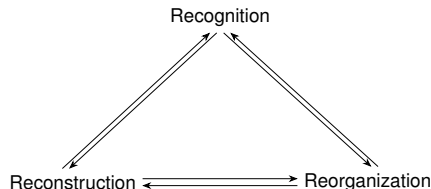
# Computation, Computer Vision and AI

## Computer Vision

Replication of human vision capabilities.

Three fundamental tasks[1]:

- Recognition.
- Reconstruction.
- Reorganization.



# Computation, Computer Vision and AI

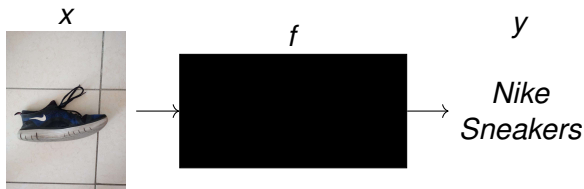
## Artificial Intelligence

Systems capable of performing tasks requiring human intelligence [2].



# Explainable AI

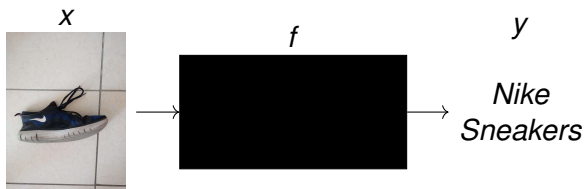
We are interested in understanding models, behaving like a black box model:





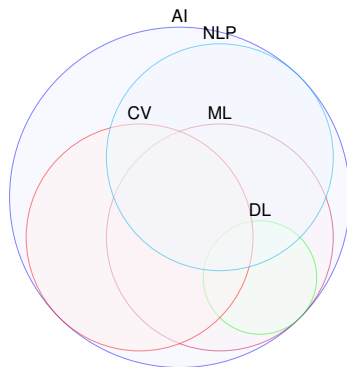
# Explainable AI

We are interested in understanding models, behaving like a black box model:

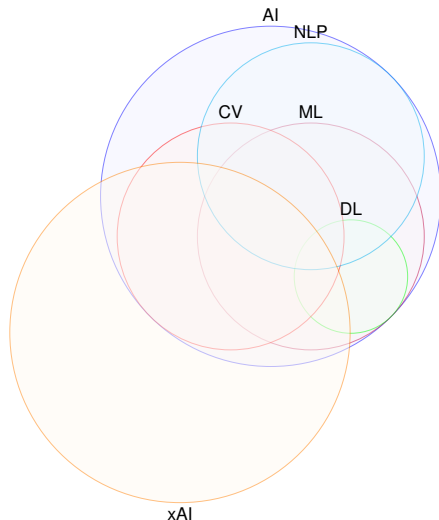


We want to *know why*  $f(x) \rightarrow y$

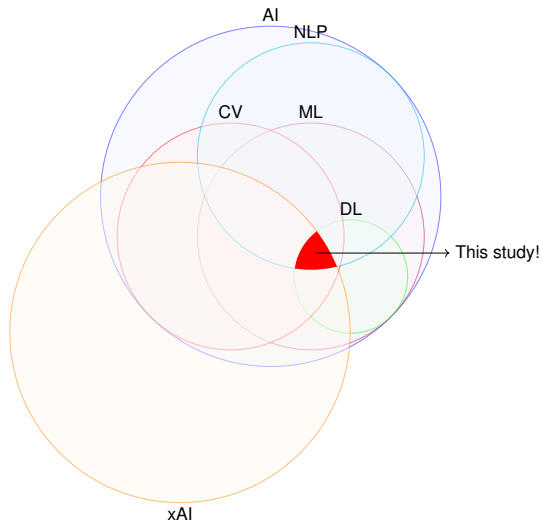
# Fitting it all together



# Fitting it all together



# Fitting it all together



# Thesis Objectives

**Improvement of recognition and interpretable properties of model predictions.**

# Thesis Objectives

**Improvement of recognition and interpretable properties of model predictions.**

In particular:

- Development of low cost/complexity explainability approaches.
- Establishment of a fixed evaluation protocol.
- Differentiation of human based and machine explanations.

# Table of Contents

- Introduction
- 1 Background**
- 2 Opti-CAM: Optimizing saliency maps for interpretability
- 3 CA-Stream: Attention-based pooling for interpretable image recognition
- 4 A learning paradigm for interpretable gradients
- References

# Background

To familiarize with this work, we split it into three points:



# Background

To familiarize with this work, we split it into three points:

## Preliminaries

- Approaching Vision.
- David Marr's approach.
- CV currently.
- Desired data of Interpretability Study.

## Image Recognition Models

- Traditional Models.
- Convolutional Neural Networks (CNN).
- The Current Landscape.

## Interpretability

- Transparency.
- Post-Hoc Interpretability.
  - Class Activation Methods.
- Evaluating Interpretability.

# Preliminaries

## Approaching Vision

# Preliminaries

David Marr's approach



Addressing vision on three levels:

- Algorithmic.
- Implementation.
- Computational.
  - Three fundamental tasks. [1]

## David Marr's approach



- Algorithmic.
- Implementation.
- Computational.
  - Three fundamental tasks. [1]

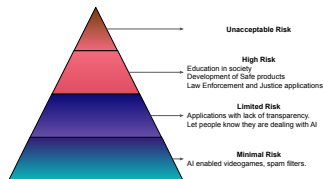
**Computer Vision focuses on the last level.**

# Preliminaries

CV Currently

# Preliminaries

## Desired data of Interpretability Study



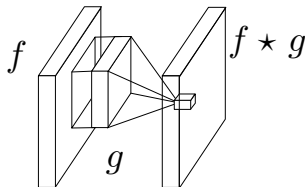
## Classic Models

# Image Recognition Models

Based on the **convolution operation**.

A representation  $f \star g$  is computed for a feature map  $f$  and a kernel  $g$ .

### First approach with *Neocognitron* [3]



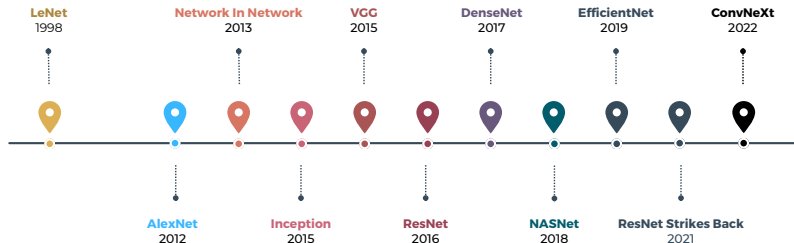
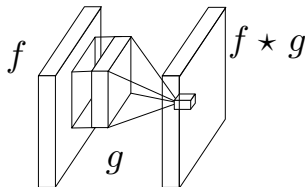


# Image Recognition Models

## Convolutional Neural Networks

Based on the **convolution operation**.  
A representation  $f \star g$  is computed for a feature map  $f$  and a kernel  $g$ .

First approach with *Neocognitron* [3]

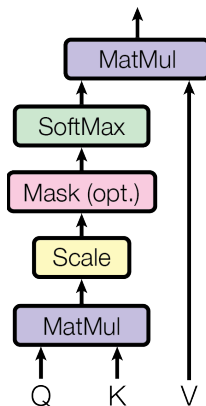


# Image Recognition Models

## Self Attention Architectures

Updates a representation, using the relevance of each element relative to others in an embedding.

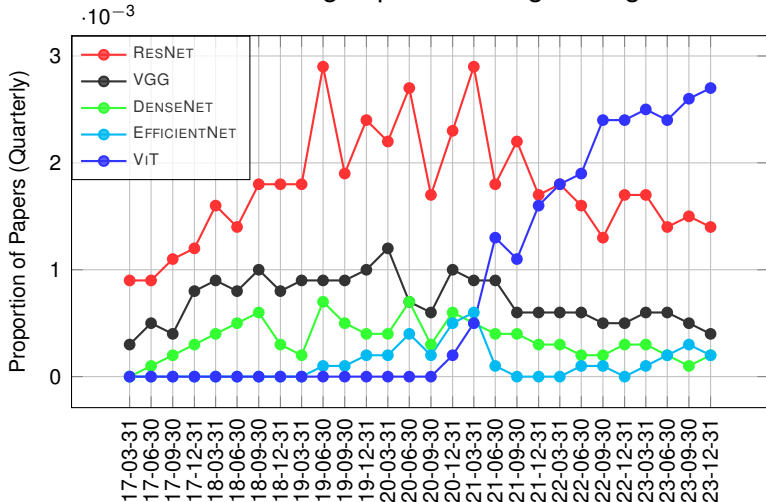
An input is projected to three spaces ( $QKV$ ), the weights control the relevance of each element.



# Image Recognition Models

## The Current Landscape

Transformers had a strong impact on image recognition.



# Interpretability

## Transparency

# Interpretability

## Post-Hoc Interpretability

# Interpretability

# Interpretability

## Evaluating Interpretability

# Table of Contents

- Introduction
- 1 Background
- 2 **Opti-CAM: Optimizing saliency maps for interpretability**
- 3 CA-Stream: Attention-based pooling for interpretable image recognition
- 4 A learning paradigm for interpretable gradients
- References



# Table of Contents

- Introduction
- 1 Background
- 2 Opti-CAM: Optimizing saliency maps for interpretability
- **3 CA-Stream: Attention-based pooling for interpretable image recognition**
- 4 A learning paradigm for interpretable gradients
- References






# Table of Contents

- Introduction
- 1 Background
- 2 Opti-CAM: Optimizing saliency maps for interpretability
- 3 CA-Stream: Attention-based pooling for interpretable image recognition
- 4 A learning paradigm for interpretable gradients**
- References



# References I

-  J. Malik, P. Arbeláez, J. Carreira, K. Fragkiadaki, R. Girshick, G. Gkioxari, S. Gupta, B. Hariharan, A. Kar, and S. Tulsiani, “The three r’s of computer vision: Recognition, reconstruction and reorganization,” *Pattern Recognition Letters*, vol. 72, pp. 4–14, 2016.
-  J. McCarthy *et al.*, “What is artificial intelligence,” 2007.
-  K. Fukushima, “Cognitron: A self-organizing multilayered neural network,” *Biological cybernetics*, vol. 20, no. 3-4, pp. 121–136, 1975.