# 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

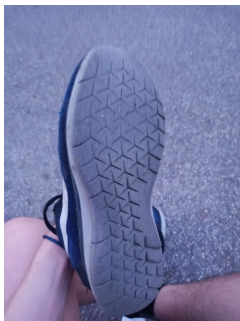# Table of Contents

## Motivation
Low Stakes

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

# Motivation
Low Stakes

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

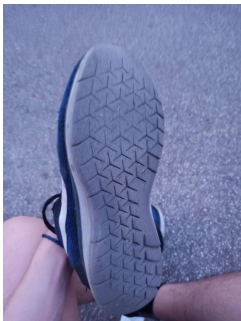I think my running shoes
are getting *worn*

## Motivation
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!

# Motivation
## 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!



My phone can *identify* my current shoes

# Motivation
## 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!



My phone can
*identify* my current shoes

*Nike Free RN Distance 2*

## Motivation
### 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!



My phone can
*identify* my current shoes

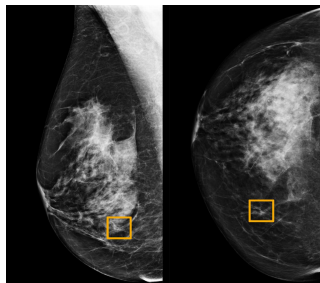*Nike Free RN Distance 2*



***How* could my phone
identify that model?**

Now let's consider riskier situations:

## Motivation
### Raising the stakes

Now let's consider riskier situations:

# Motivation
Raising the stakes

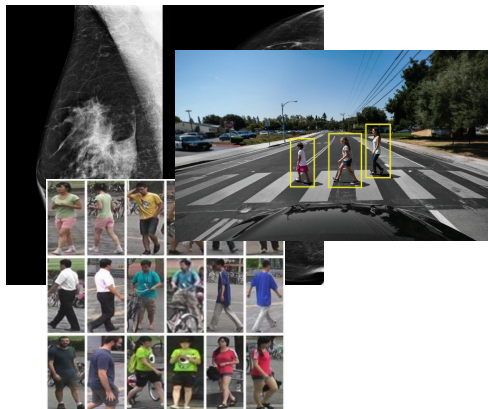Now let's consider riskier situations:

## Motivation
Raising the stakes

Now let's consider riskier situations:

## Motivation
Straight to the point



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

## Motivation
Straight to the point



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

## Motivation
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.

Computation,
Computer          $\rightarrow$     Explainable AI     $\rightarrow$          Thesis
Vision and AI                                                         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
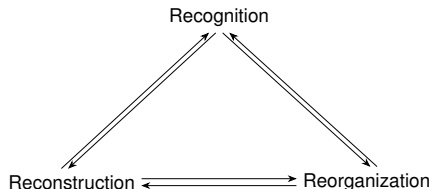Computer Vision

Replication of human vision
capabilities.

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
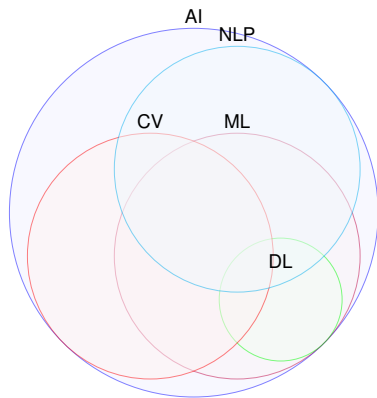tasks requiring human
intelligence [2].

# Computation, Computer Vision and AI
## Artificial Intelligence

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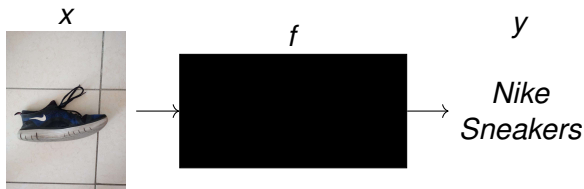tasks requiring human
intelligence [2].

Subfields:

- Machine Learning *(ML)* &
  Deep Learning (*DL*).
- Computer Vision(*CV*).
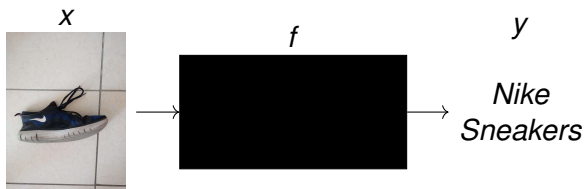- Natural Language
  Processing (*NLP*).
- Robotics.

Explainable AI

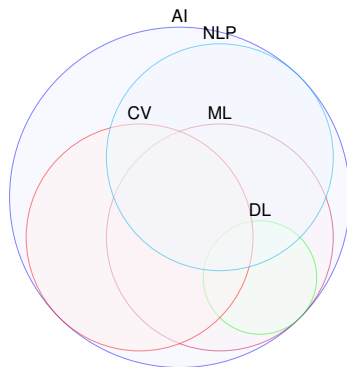We are interested in understanding models,
behaving like a black box model:



$x$    $f$    $y$

*Nike*
*Sneakers*

Explainable AI

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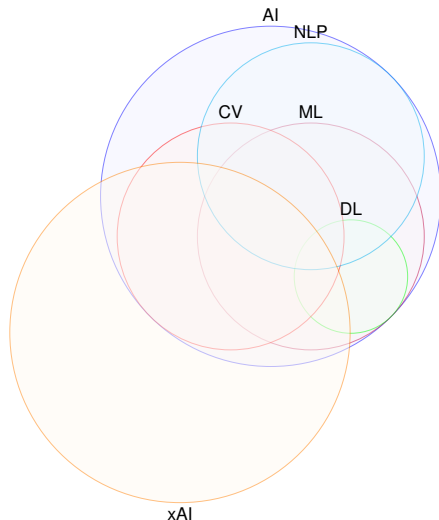behaving like a black box model:



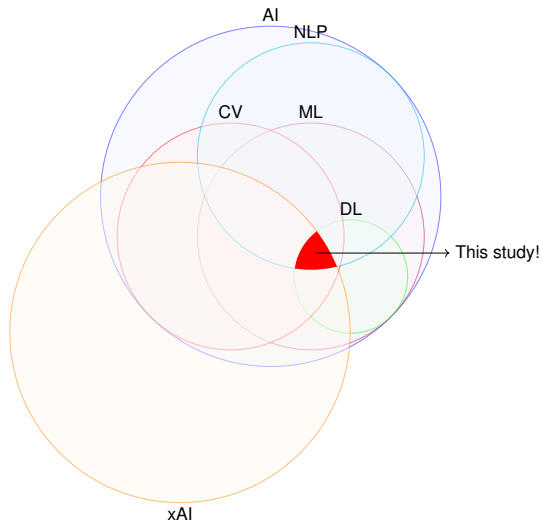**We want to *know why* $f(x) \to y$**

## Fitting it all together

# Fitting it all together

# Fitting it all together

## Thesis Objectives

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

## Thesis Objectives

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

In particular:

- Development of low cost/complexity explainability approaches.

## Thesis Objectives

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

In particular:

- Development of low cost/complexity explainability approaches.
- Establishment of a fixed evaluation protocol.

## Thesis Objectives

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

In particular:

- Development of low cost/complexity explainability approaches.
- Establishment of a fixed evaluation protocol.
- Differenciation 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.
- Desiredata of Interpretability Study.

**Image Recognition Models**

- Traditional Models.
- Convolutional Neural Networks (CNN).
- Hybrid Architectures.

**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.
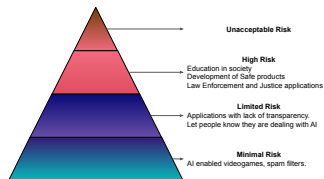  - Three fundamental tasks. [1]

## Preliminaries
David Marr's approach



Addressing vision on three
levels:

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

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

# Preliminaries
CV Currently

# Preliminaries
## Desiredata of Interpretability Study

Introduction
Background
Opti-CAM
CA-Stream
Gradient
References

# Image Recognition Models
Classic Models

# Image Recognition Models
Convolutional Neural Networks

# Image Recognition Models
Self Attention Architectures

# Image Recognition Models
Hybrid Architectures

Introduction
Background
Opti-CAM
CA-Stream
Gradient
References

# Interpretability
Transparency

Introduction
0000000000000

Background
0000000000000●00

Opti-CAM
0

CA-Stream
00

Gradient
0

References
00

# Interpretability
Post-Hoc Interpretability

# Interpretability
Class Activation Methods

Introduction
Background
Opti-CAM
CA-Stream
Gradient
References

# Interpretability
Evaluating Interpretability

## Table of Contents

## Table of Contents

## Table of Contents

## Table of Contents

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.