**C14200**: FROM ZERO TO ONE - DEEP LEARNING WITH PYTORCH
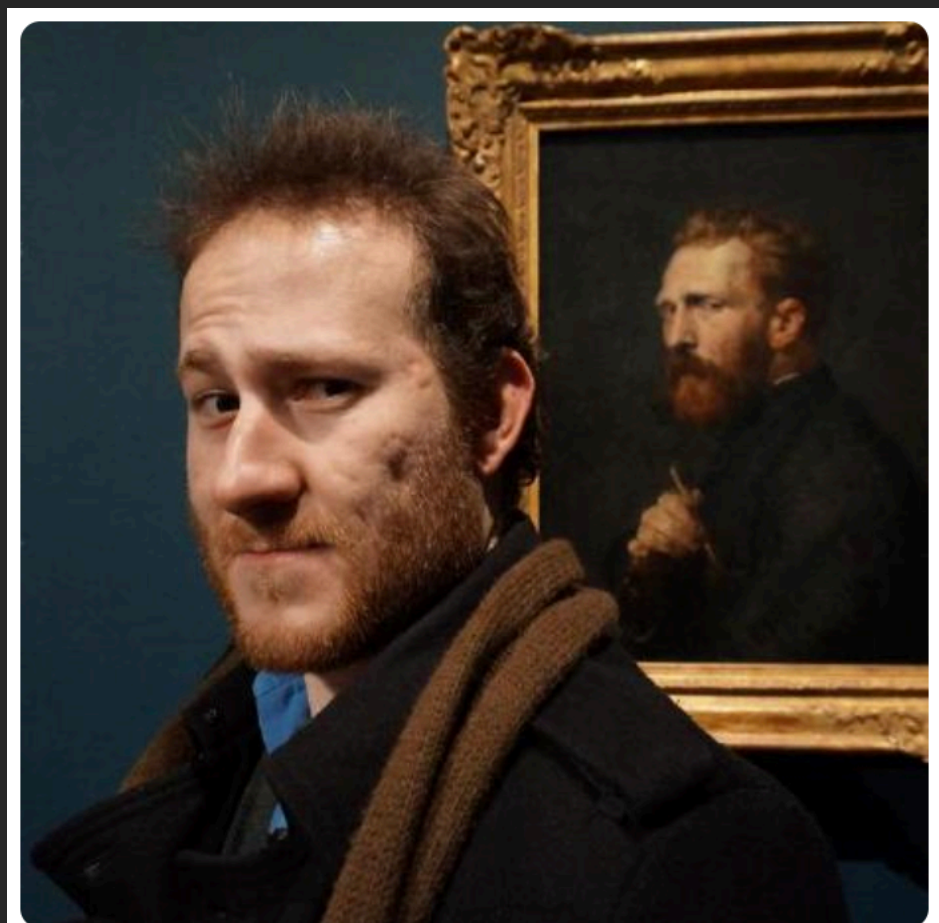
JOE SPISAK
PRODUCT MANAGER

FRANCISCO MASSA
RESEARCH ENGINEER

## WHO AM I?

**Francisco Massa**
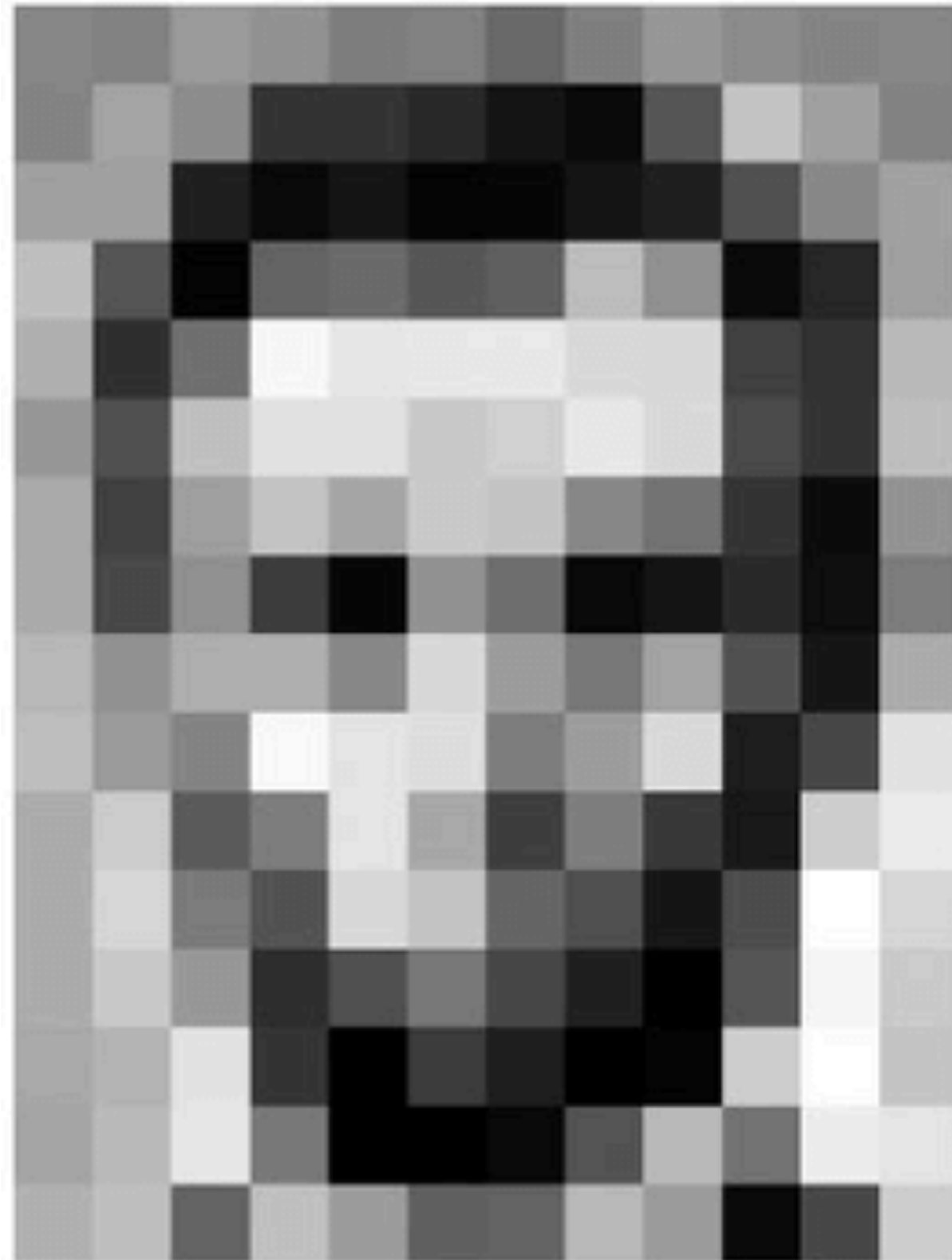fmassa

CURRENT: RESEARCH ENGINEER - PYTORCH

PREVIOUS:
- COMPUTER VISION RESEARCH ENGINEER @TWITTER

- PHD STUDENT AT ECOLE DES PONTS - FRANCE

AGENDA
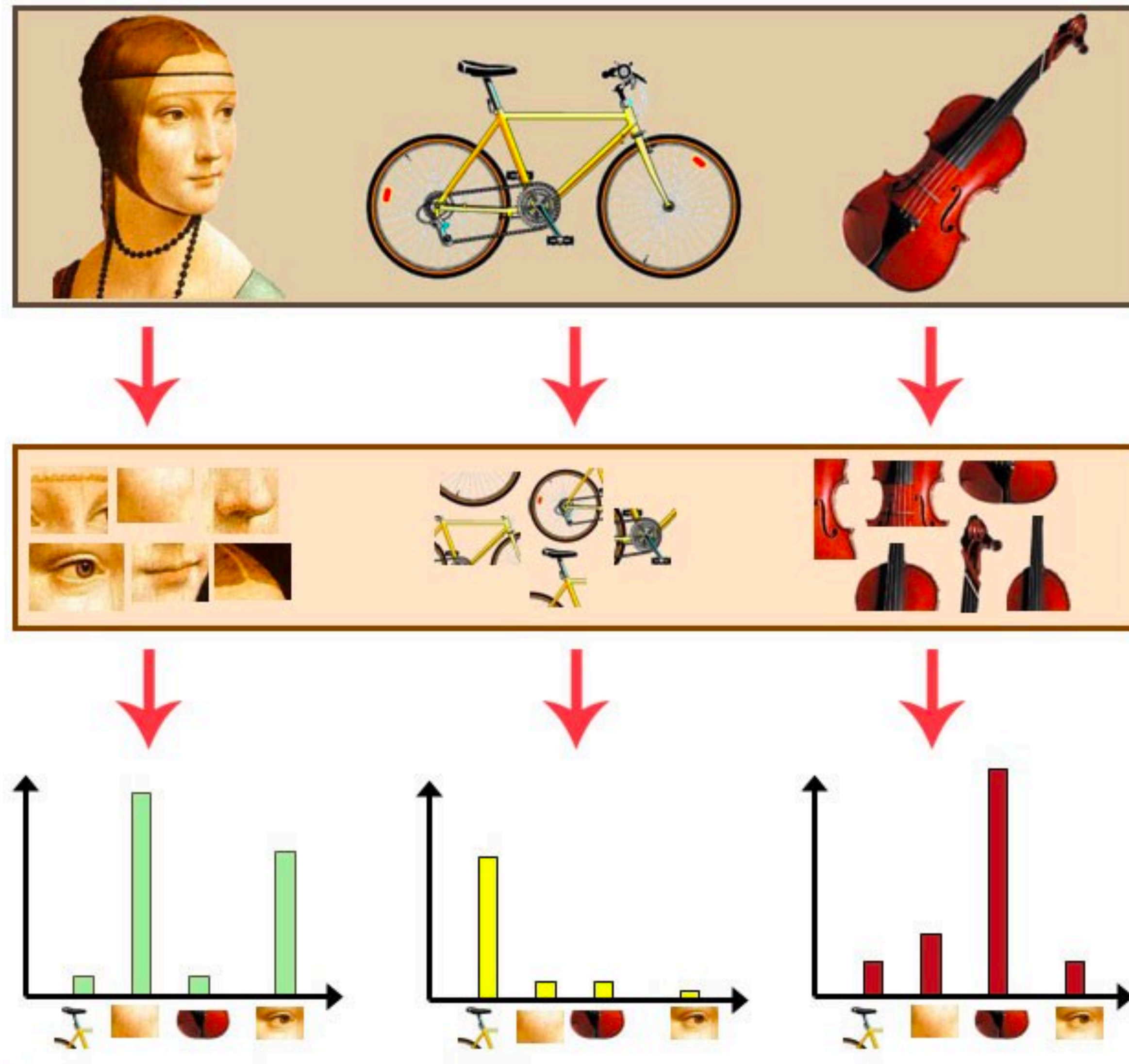
# PRIMER ON COMPUTER VISION

How can we make the computer understand what is in the image?

# OBJECTS AS A BAG OF VISUAL WORDS



Prototype patches

(visual words)

Histogram of

visual words

## ImageNet Image Classification Challenge



1M Images

1000 categories

## ImageNet Image Classification Challenge

## Ranking of the best results from each team



Key ingredients:

- Deep Convolutional Neural Networks

- Lot's of training data

- ReLU and dropout

- GPUs

## What is a convolution?



| 7 | 2 | 3 | 3 | 8 |
|---|---|---|---|---|
| 4 | 5 | 3 | 8 | 4 |
| 3 | 3 | 2 | 8 | 4 |
| 2 | 8 | 7 | 2 | 7 |
| 5 | 4 | 4 | 5 | 4 |

\*

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

| 6 | | |
|---|---|---|
| | | |
| | | |

7x1+4x1+3x1+
2x0+5x0+3x0+
3x-1+3x-1+2x-1
= 6

Image from https://medium.com/datadriveninvestor/convolutional-neural-networks-3b241a5da51e

Neural networks that uses the convolution operator



Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Image from https://pythonmachinelearning.pro/introduction-to-convolutional-neural-networks-for-vision-tasks/

# CONVOLUTIONAL NEURAL NETWORKS TODAY

## More layers work better



Image from https://medium.com/@pierre_guillou/understand-how-works-resnet-without-talking-about-residual-64698f157e0c

# GOING BEYOND IMAGE CLASSIFICATION

**Classification**

CAT

**Classification + Localization**

CAT

Single object

**Object Detection**

CAT, DOG, DUCK

**Instance Segmentation**

CAT, DOG, DUCK

Multiple objects

**Semantic Segmentation**

GRASS, CAT, TREE, SKY

No objects, just pixels

# CNN OVER REGIONS (R-CNN)

# COMPUTER VISION WITH TORCHVISION

A library built to facilitate research and experimentation in the field of Computer Vision

**DATASETS**

COMMON DATASETS

**MODELS**

PRE-TRAINED MODELS

**OPS**

EFFICIENT OPERATORS

**TRANSFORMS**

DATA TRANSFORMATION

**IO**

EFFICIENT VIDEO READER

**REFERENCES**

TRAINING SCRIPTS

DATASETS

TRANSFORMS

MODELS

OPS

IO



IMAGENET 64X64

# DATASETS

```
import                    for

T.Com
    T                 224
    T              p(
    T
    T                 86
                  , 0      ])
)
```

An Analysis of Deep Neural Network Models for Practical Applications
Alfredo Canziani, Adam Paszke, Eugenio Culurciello

DATASETS

TRANSFORMS

MODELS

OPS

IO

```python
import torchvision.ops

torchvision.ops.box_iou(...)
torchvision.ops.roi_align(...)
torchvision.ops.nms(...)
torchvision.ops.roi_pool(...)
```

DATASETS

TRANSFORMS

MODELS

OPS

IO

```python
import torchvision.io

torchvision.io.read_video(filename,
                          start_pts=0,
                          end_pts=None)


torchvision.io.read_video_timestamps(filename)


torchvision.io.write_video(filename,
                           video_array,
                           fps,
                           video_codec='libx264',
                           options=None)
```

```python
import torchvision

model = torchvision.models.detection.maskrcnn_resnet50_fpn(pretrained=True)
# set it to evaluation mode, as the model behaves differently
# during training and evaluation
model.eval()

image = PIL.Image.open('/path/to/an/image.jpg')
image_tensor = torchvision.transforms.functional.to_tensor(image)

# pass a list of (potentially different sized) tensors
# to the model, in 0-1 range. The model will take care of
# batching them together and normalizing
output = model([image_tensor])
# output is a list of dict, containing the post processed predictions
```

HANDS ON WITH TORCHVISION