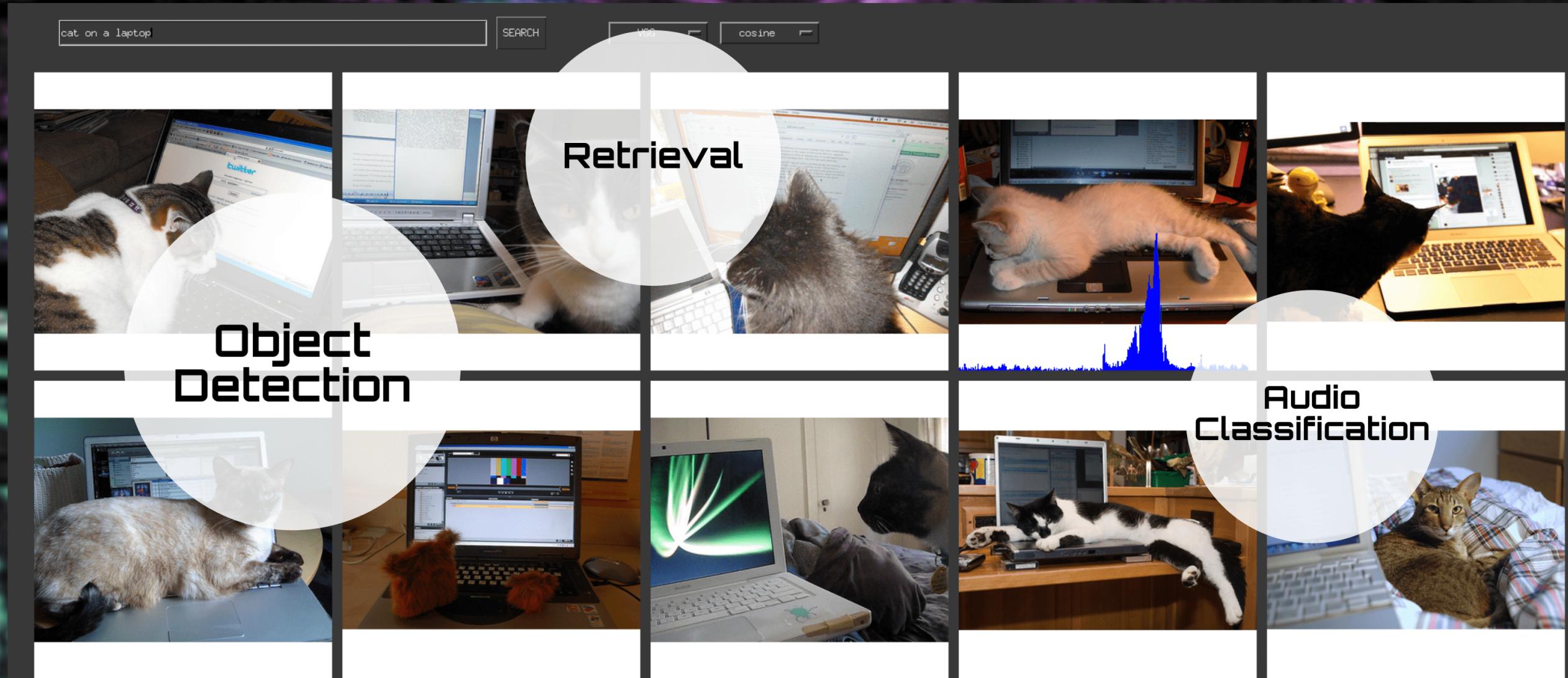


DSIM Project

Content-based IR

Cervero Riccardo 794126
Kasela Pranav 846966
Moiraghi Federico 799735





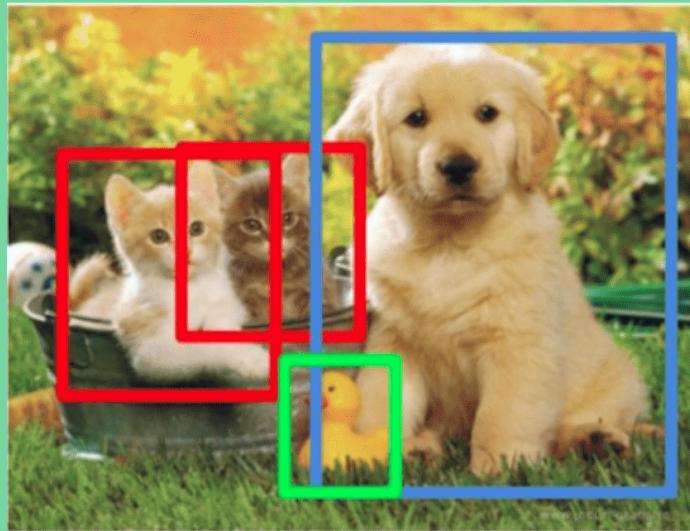
coco

Common Objects in Context



- ➡ Detection & Semantic Segmentation
- ➡ 328k images
- ➡ 80 object categories

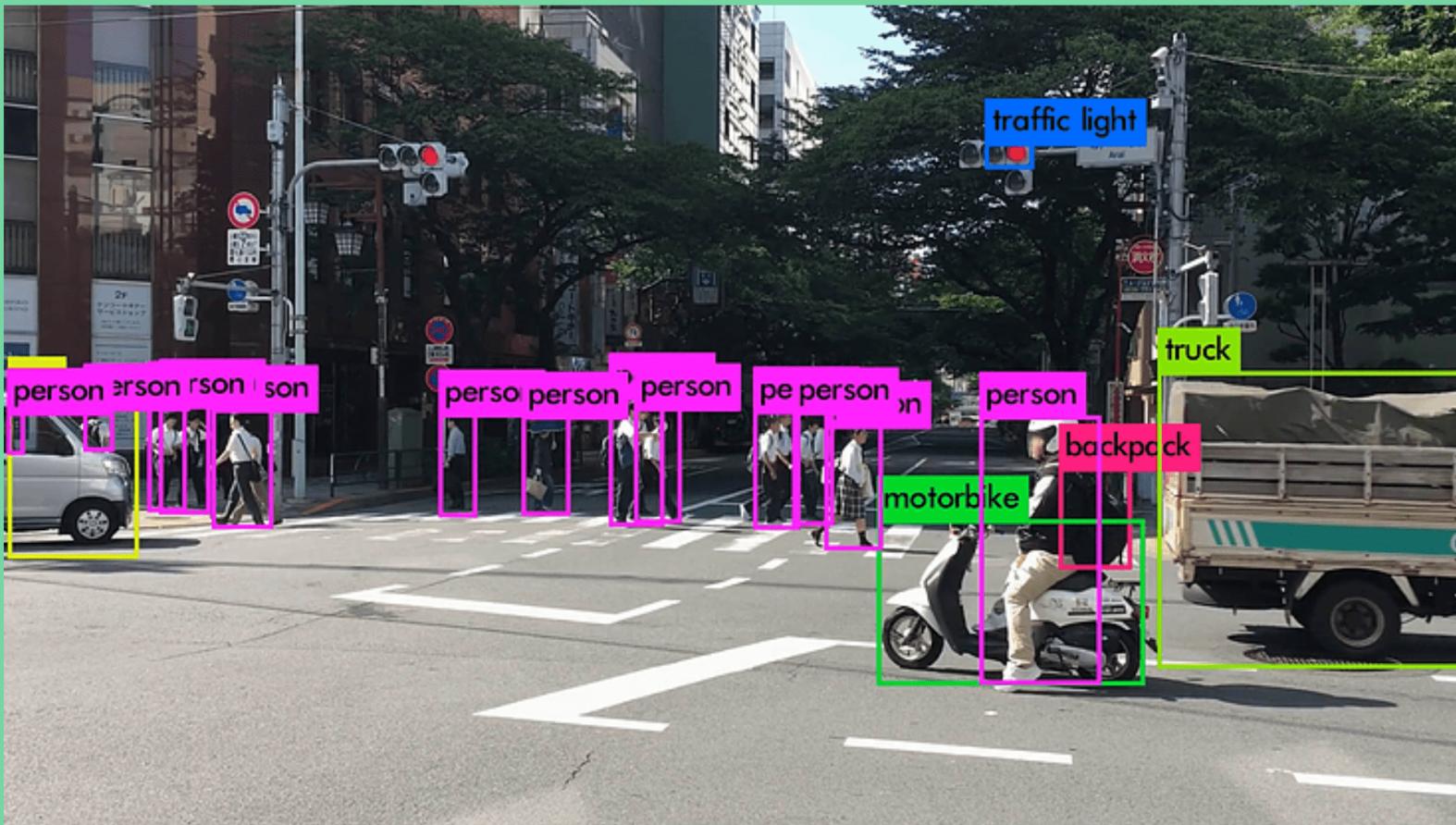
**Images are "multi-labelled"
by tagging their elements**



CAT, DOG, DUCK

YOU ONLY LOOK ONCE 3v

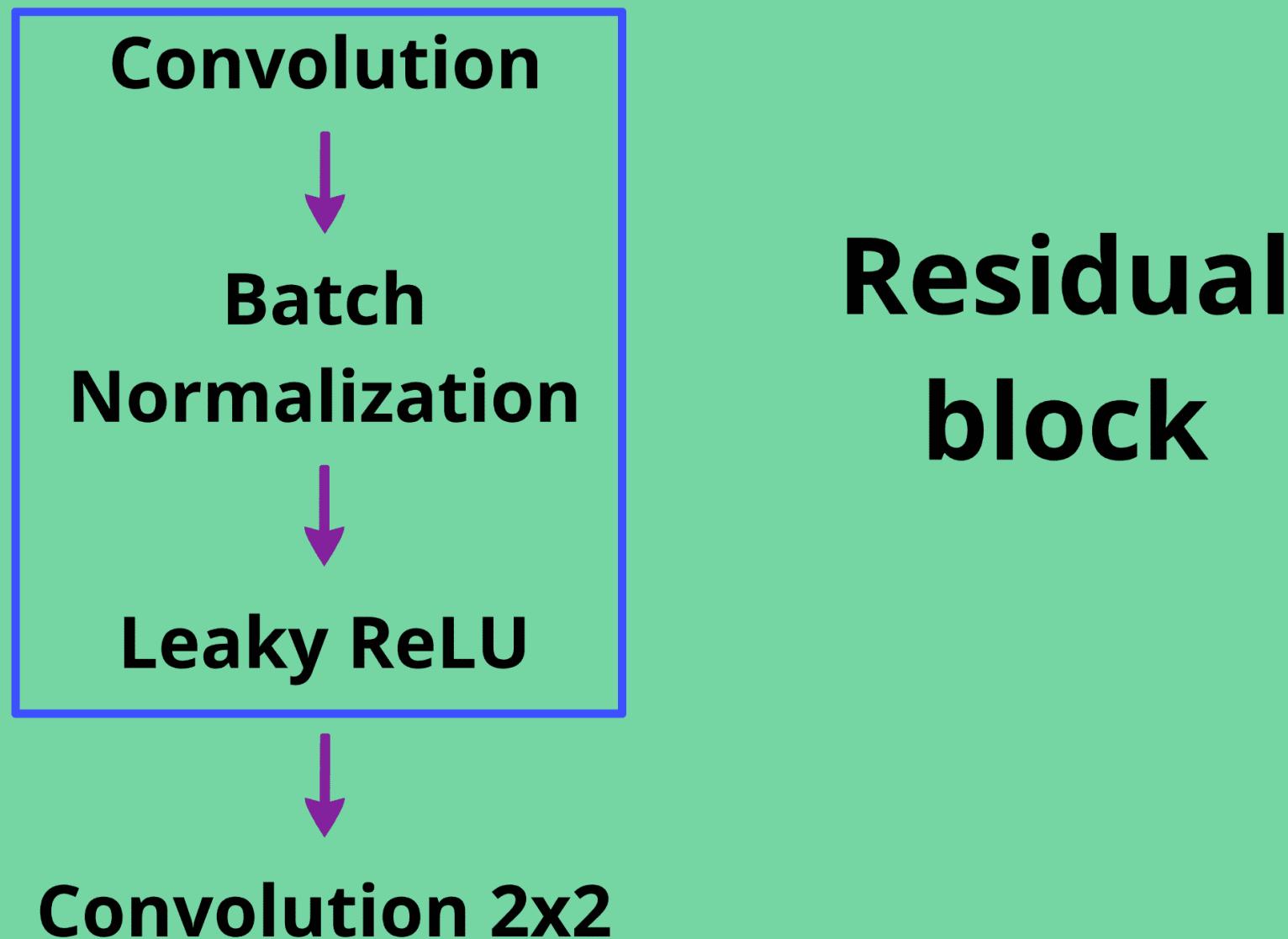
Pre-trained very fast algorithm to detect multiple object location inside images and videos



YOLOv3 Architecture

-  **Fully Convolutional Neural Network**
-  **Not a R-CNN, but "looks once"**
-  **Deeper feature extractor**
-  **Invariant to size but need to fix dimensions**
-  **Feature pyramid network**

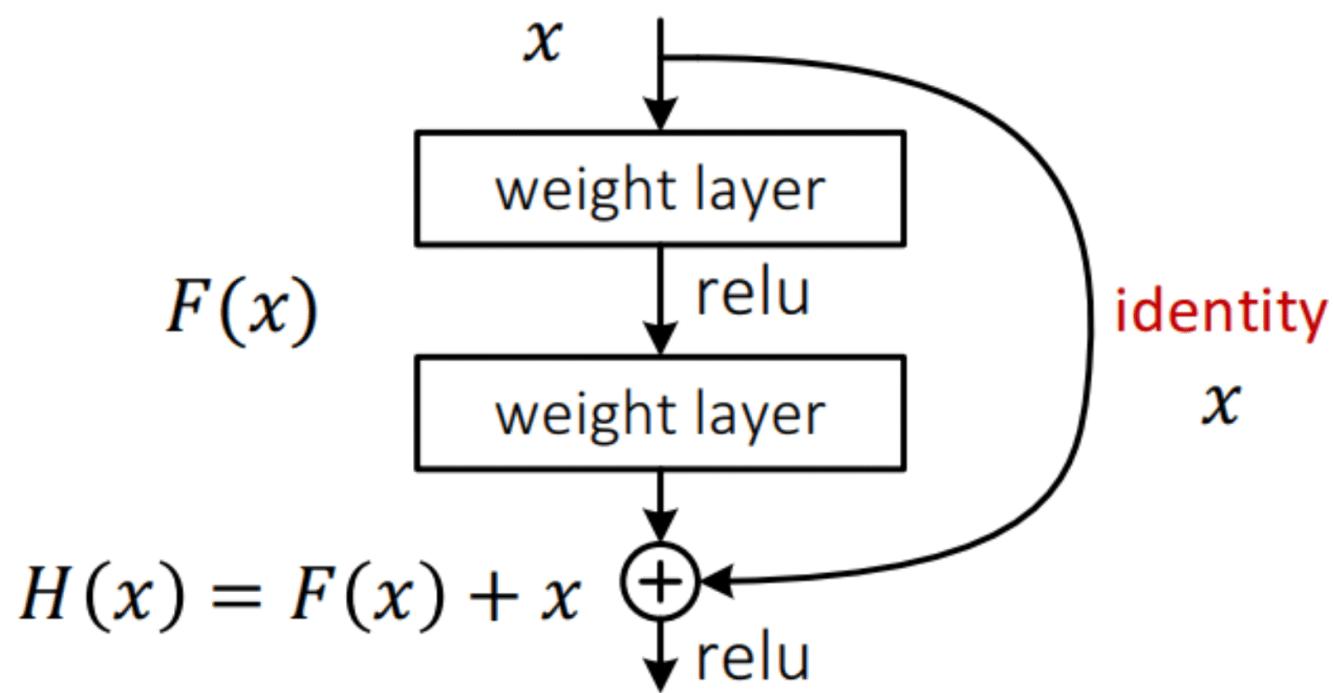
DarkNet Feature Extractor



ResNet-alike Architecture

Deep Residual Learning

- Residual net



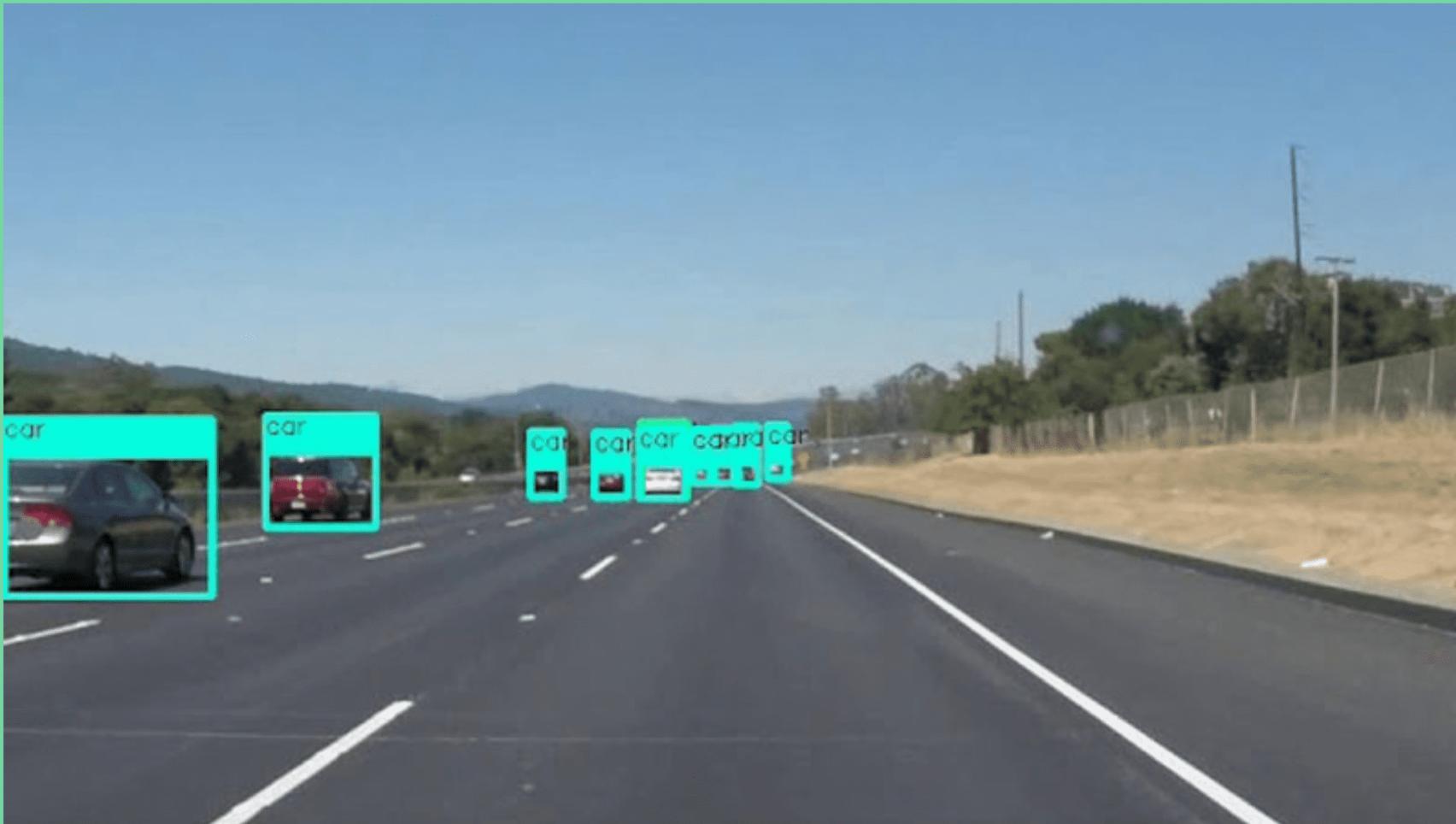
Lower learning complexity

Detecting differently-sized objects

To get objects locations and class, network classifier produces a grid called "feature map"

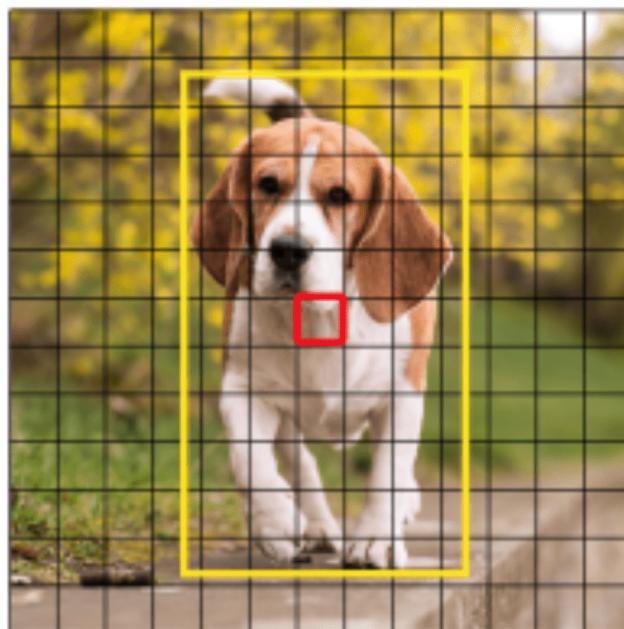
Detecting differently-sized objects

But as the network goes deeper, its feature map gets less fine-grained and could lose ability to catch small details!

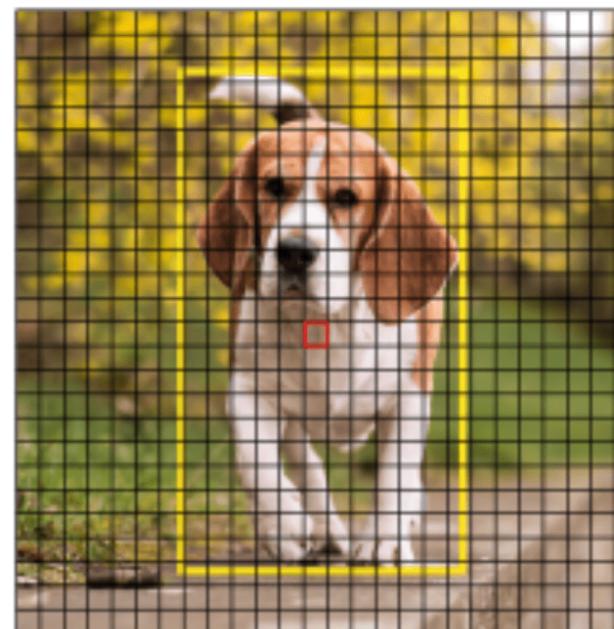


YOLO thus outputs instead of a 1D feature map...

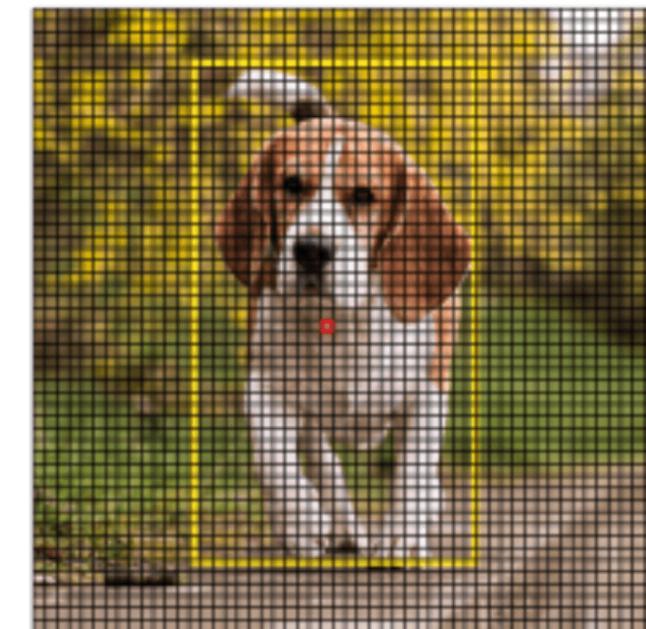
**...a 3D feature map able to find "large",
"medium" and "small" object**



13 x 13

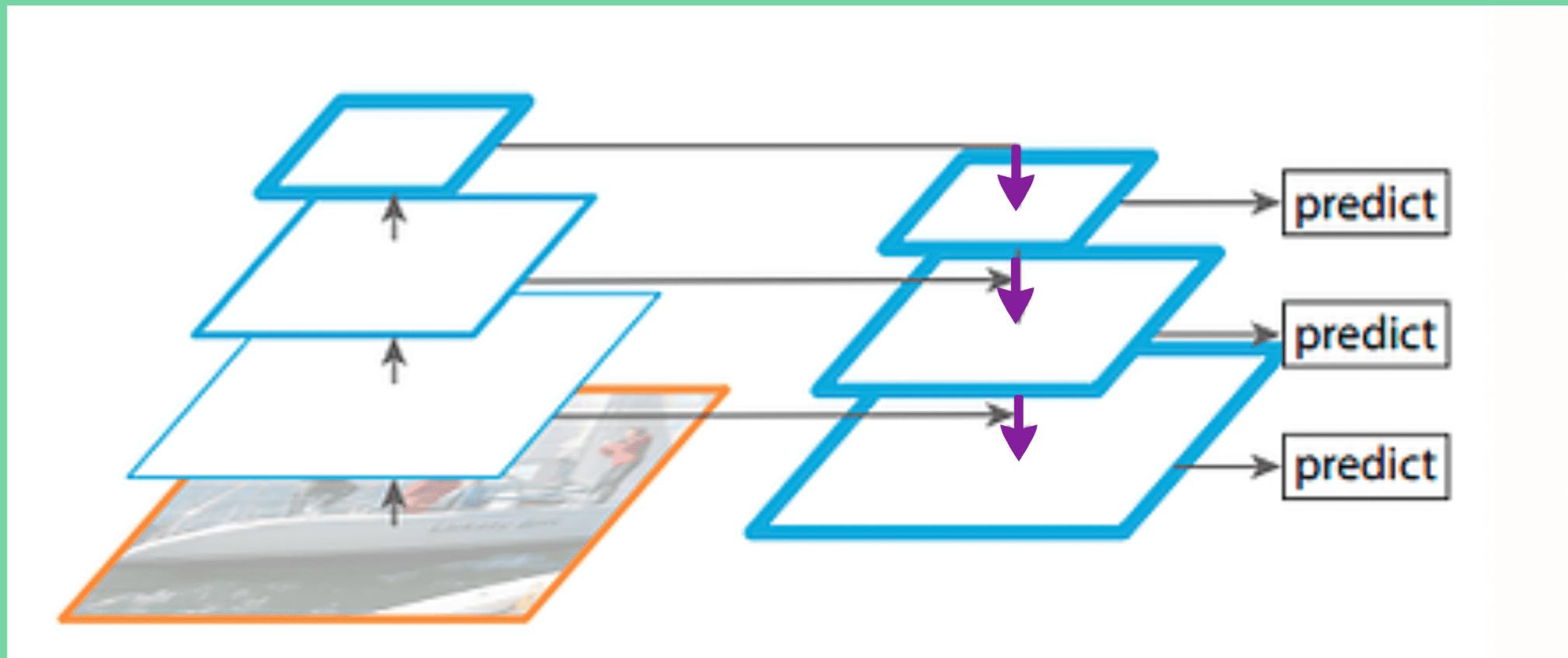


26 x 26



52 x 52

Feature pyramid network



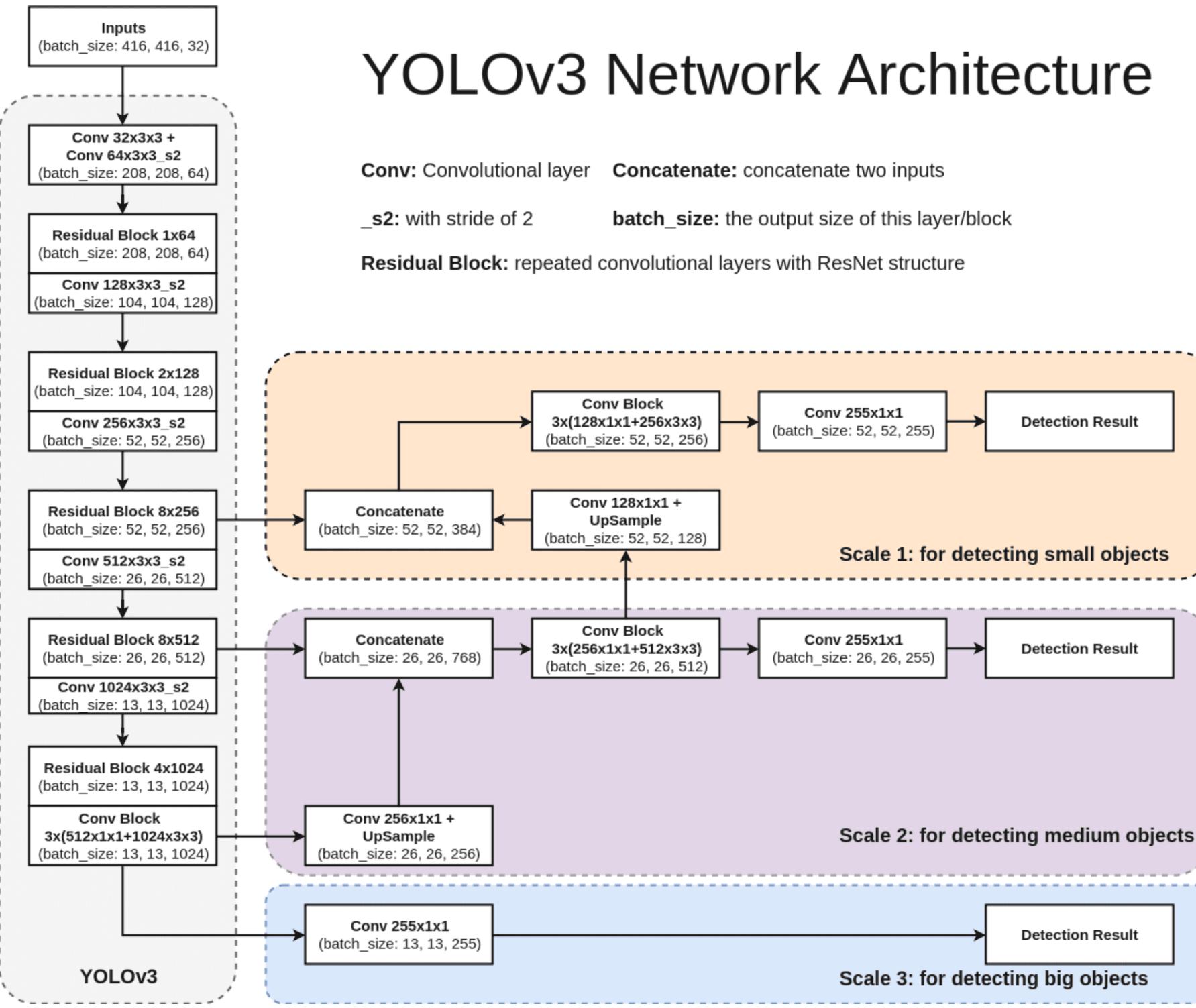
Concatenation



Upsampling ↓



YOLOv3 Network Architecture



Conv: Convolutional layer **Concatenate:** concatenate two inputs

_s2: with stride of 2 **batch_size:** the output size of this layer/block

Residual Block: repeated convolutional layers with ResNet structure

Inference

**Given an input image, for each object detected,
YOLOv3 outputs:**



One bounding box to identify location

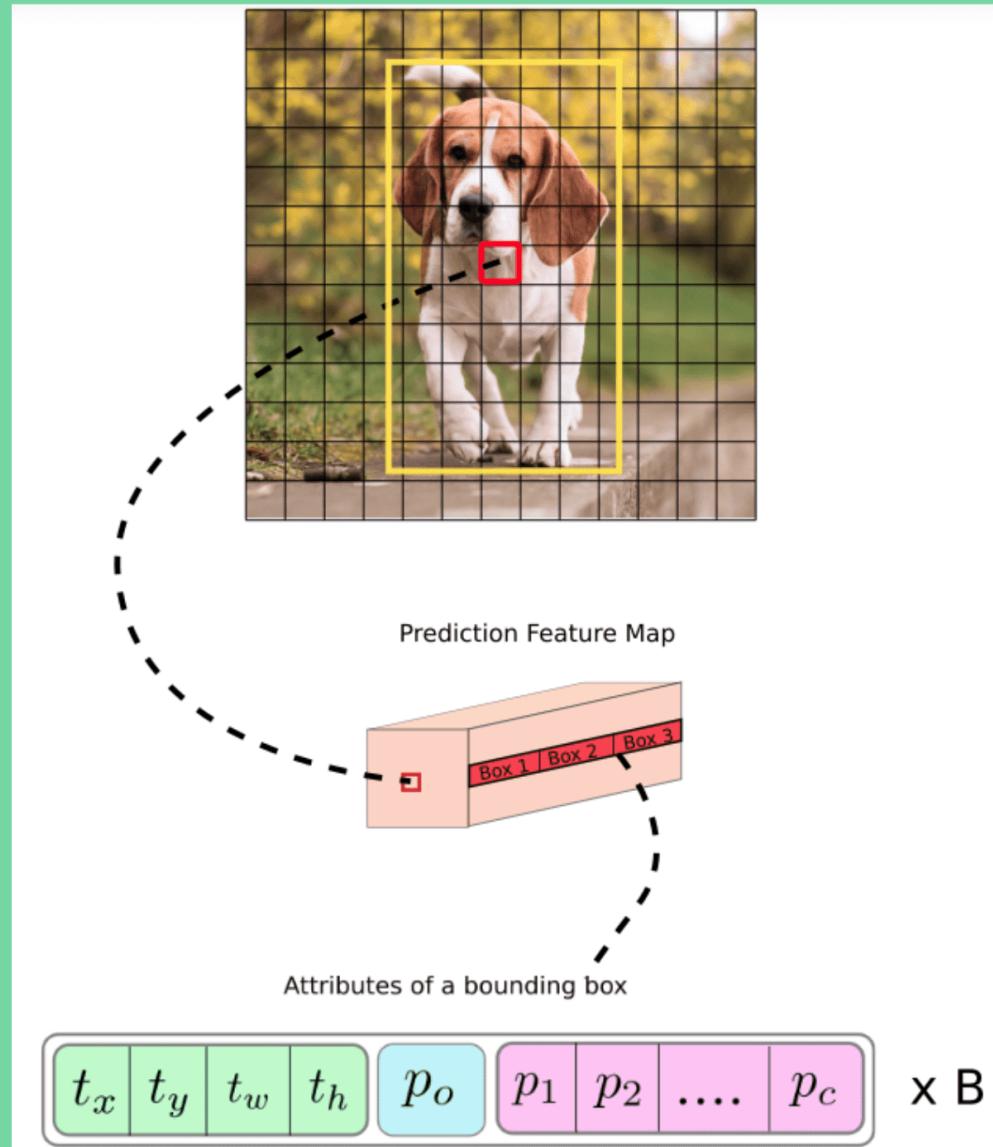


Label of the object



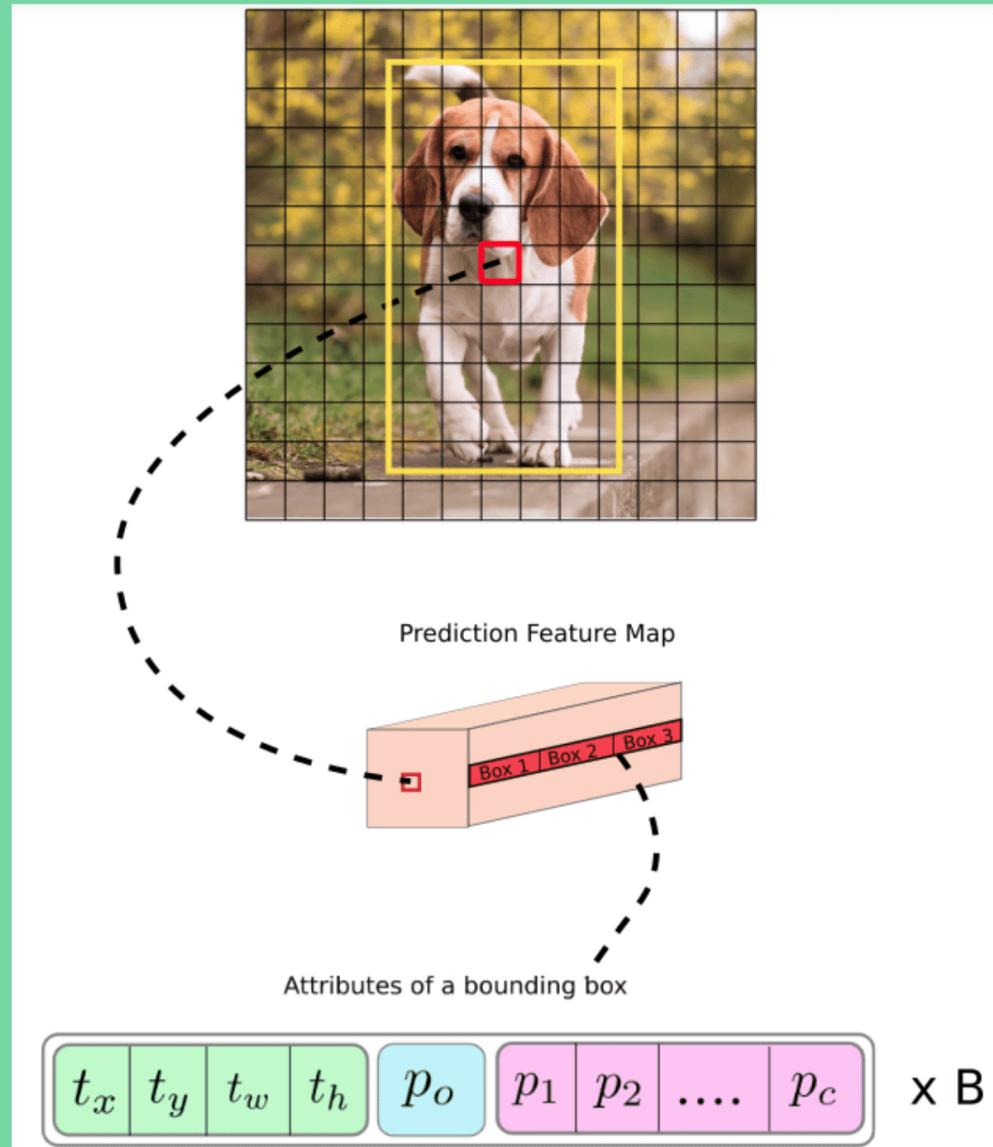
Class probability

Generation of the input



Each entry of each feature map predicts 3 bounding boxes, with $5+C$ attributes : center coordinates, dimensions, "objectness score" (by sigmoid) and C , which is the set of class probabilities (by sigmoid).

Generation of the input



The 3 bounding boxes used for each scale are default type, called anchors

Generation of the input (algorithm)

For each anchor of each grid cell, compute probability of contain every class

Select the maximum among all classes and across 3 anchor boxes

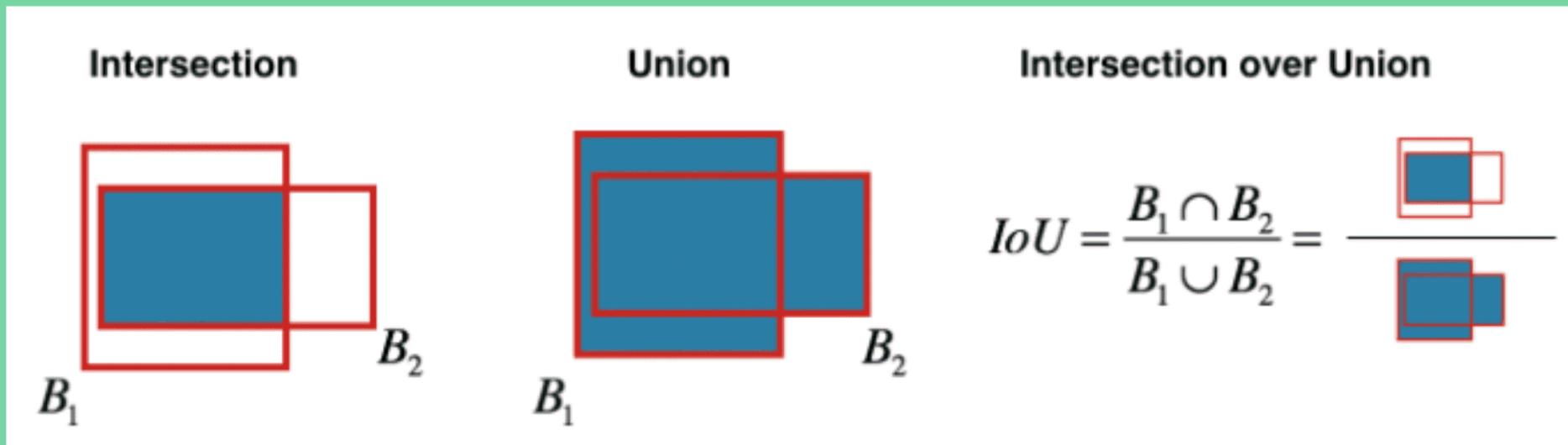
Non-Maximum Suppression

Remove boxes with an objectness score under 0.5

Return optimal boxes to detect all the objects

Non-maximal suppression with IoU

**Suppresse overlapping boxes using
Intersection over Unit metric**



Remove boxes with an IoU > 0.45

Crop features

**Tagged objects are then
cut out from the image
into crops**

**For IR purposes, these
crops need to be
described**

Crop Descriptors



Original image ID



Label



Importance score



Feature extracted with Transfer Learning

Feature Extraction

VGG16

(cut before classifier)

YOLOv3

(cut before detection layer)

SIFT

Audio Model

Data Format:

- 2 seconds
- 44100 Hz
- Normalized

981 samples - half
with background
noises

- 1 multi-classification
 - 3 binary
- Multi-classification for 'sì',
'no' and 'forse' and the 3
binary for the speaking user

1st Model

2nd Model

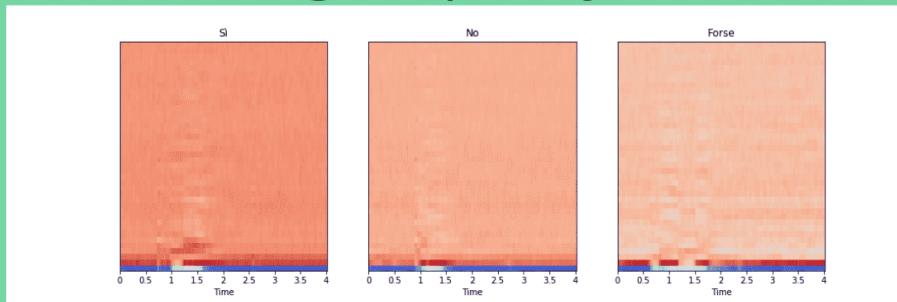
Combine

1st Model

Features Extracted mean of :

- MFCCS

Derives from a inverse
fourier transformation of
the log frequency



Model Specs

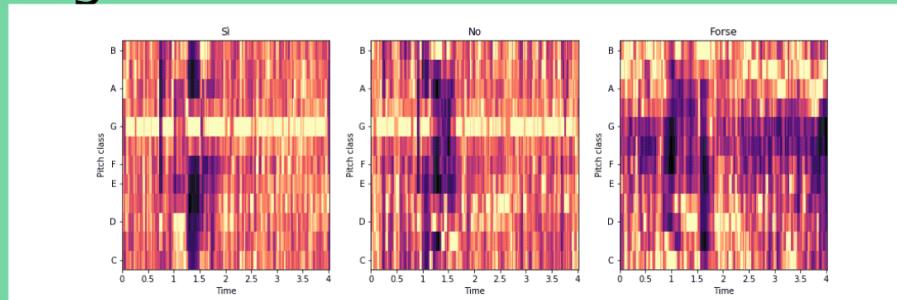
1st Model

Features Extracted mean of :

- MFCCS
- Chroma (stft)

It indicates the pitch class of the audio (musical notes)

Might be useful in user identification



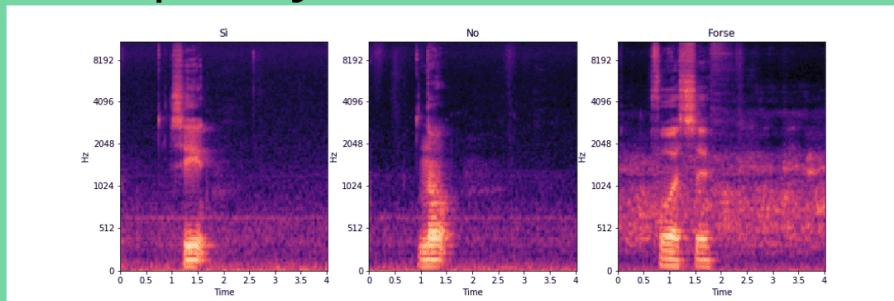
Model Specs

1st Model

Features Extracted mean of :

- MFCCS
- Chroma (stft)
- Melspectrogram

A spectrogram created from the Mel Frequency (same as MFCCS).



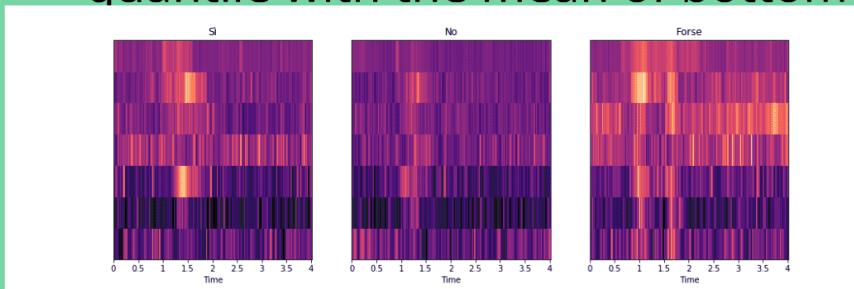
Model Specs

1st Model

Features Extracted mean of :

- MFCCS
- Chroma (stft)
- Melspectrogram
- Contrast (stft)

Energy contrast for each band of the spectrogram (ratio between mean of top quantile with the mean of bottom one)



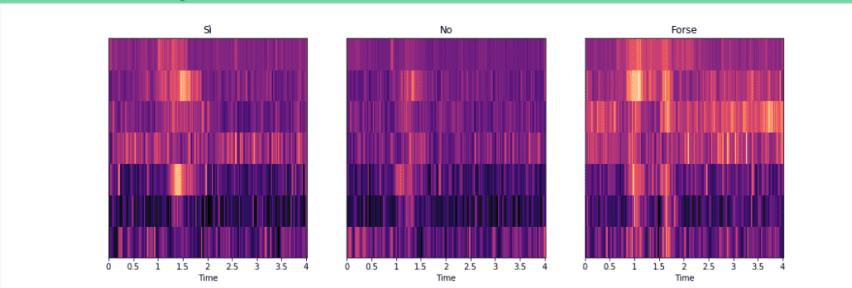
Model Specs

1st Model

Features Extracted mean of :

- MFCCS
- Chroma (stft)
- Melspectrogram
- Contrast (stft)
- Tonnetz

Computes some tonal features



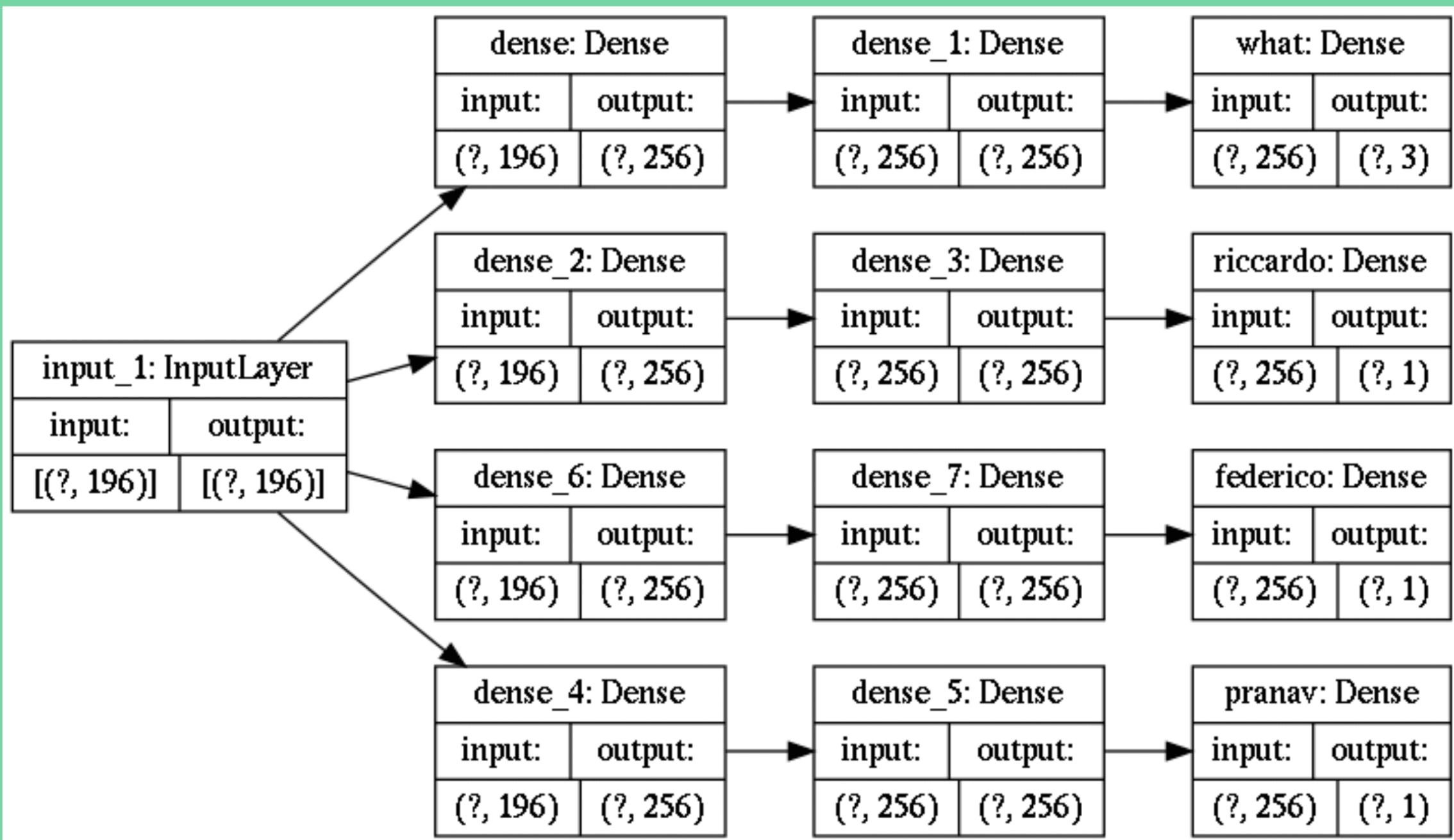
Model Specs

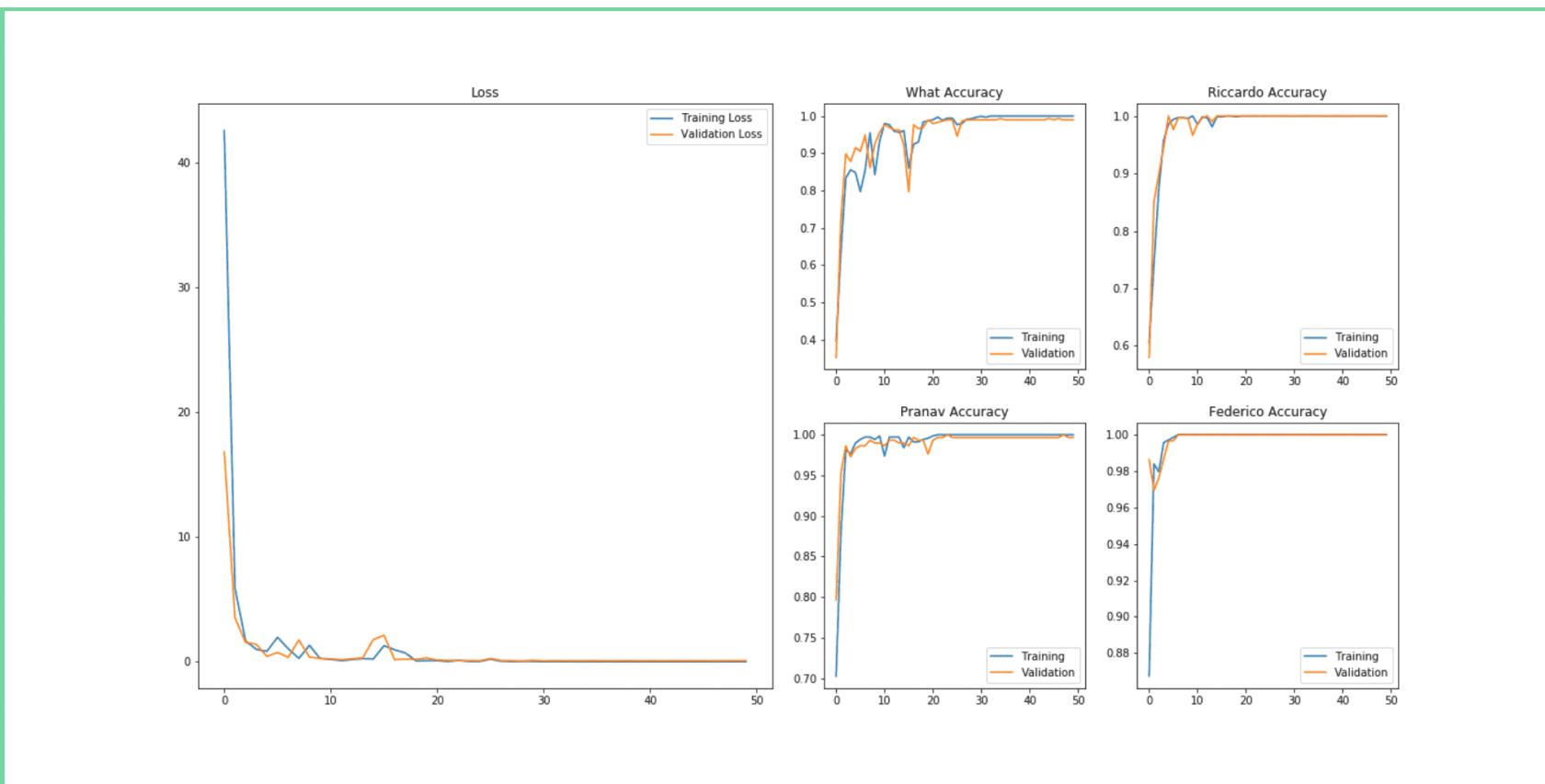
1st Model

Features Extracted mean of :

- MFCCS
- Chroma (stft)
- Melspectrogram
- Contrast (stft)
- Tonnetz
- energy
- standard deviation
- average

**Model
Specs**



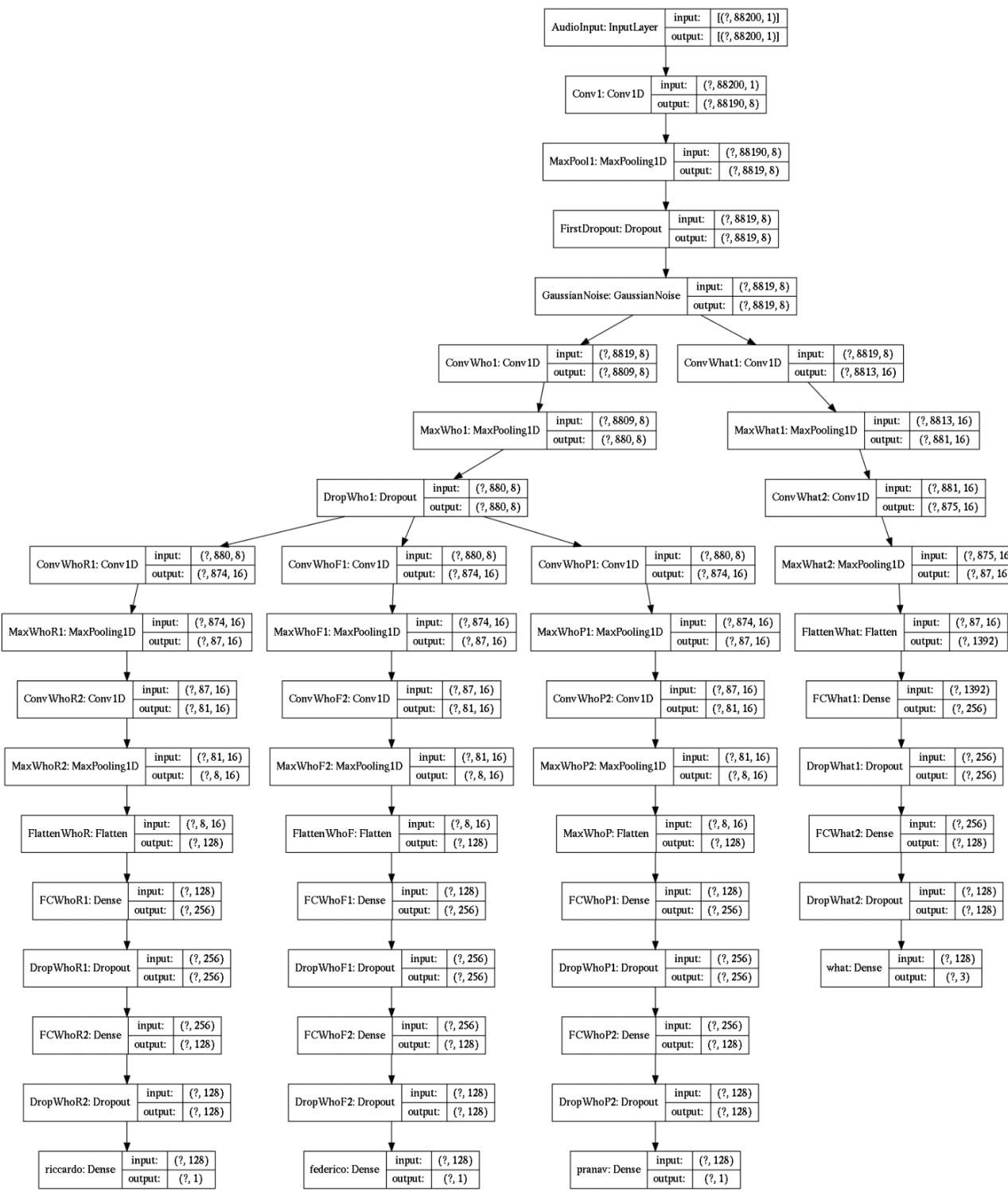


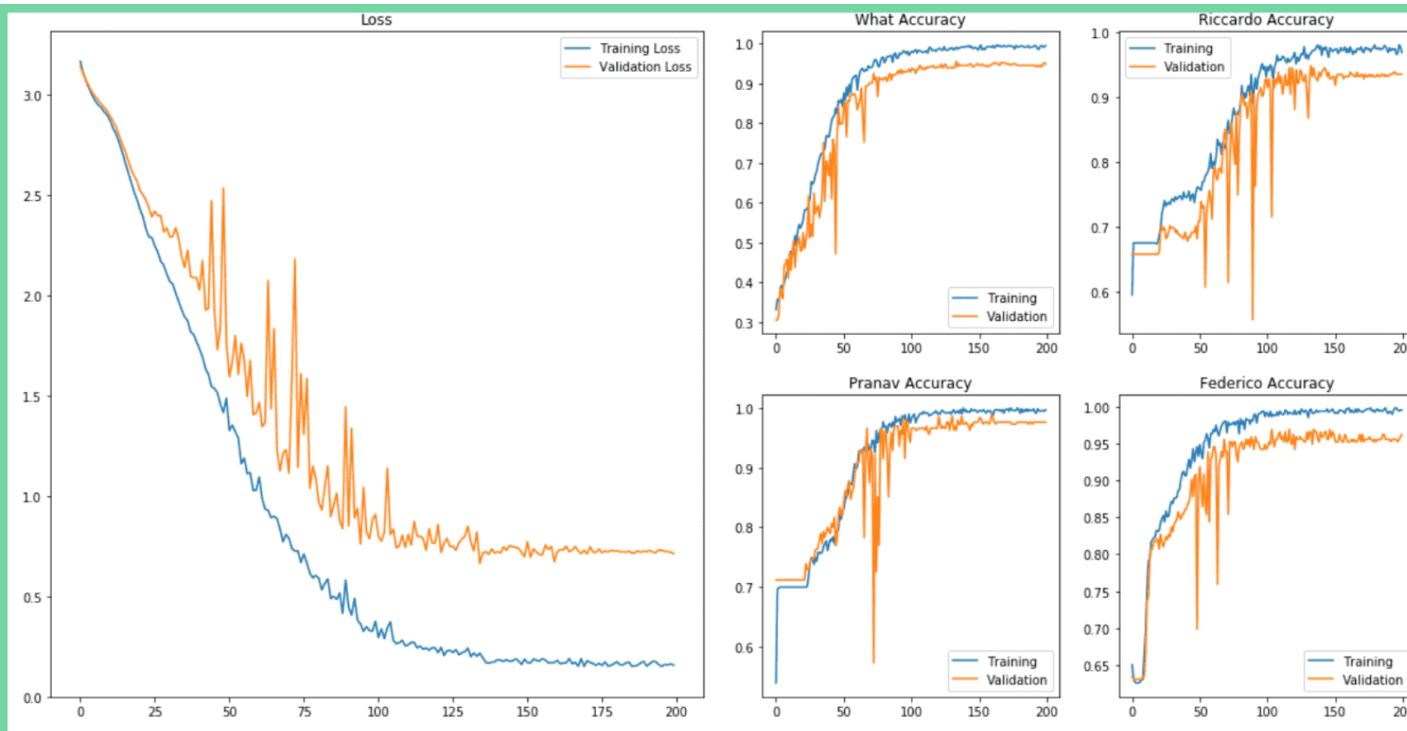
**Very accurate when the
speaker is one of us but
doesn't generalize well to a
new speaker!!!**

2nd Model

- Way more complex as a model, uses the convolutions to extract features
- Input consists of 88200 neurons
- Gaussian noise to reduce overfitting
- Combines the losses adding them
- Uses two callbacks:
 - ReduceLROnPlateau (patience 5)
 - EarlyStopping (patience 100)
- A lot of shared features between models

Model Specs





Way more noisy training, and seem
to over fit the training data, but
achieves good performance.

No AutoML :(
So many parameters to optimize.

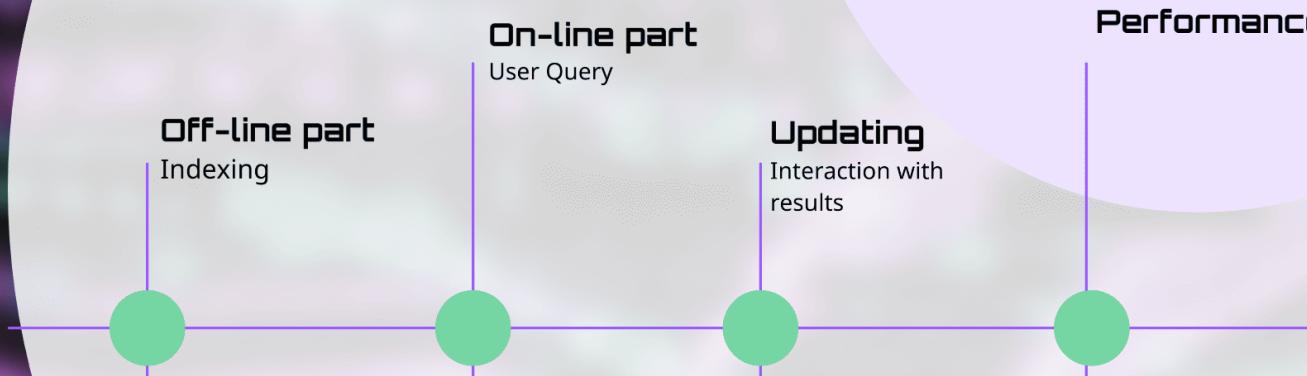
Combine the results

- First model is very accurate
- Second model is a little better at generalizing

Best choice seems to use both of them, so we take the average of both predictions and if the score is bigger than a fixed threshold it is accepted and rejected and classified as external class otherwise.

$$\frac{2 * \text{CNN model} + 1 * \text{No CNN model}}{3}$$

Information Retrieval



Off-line part: indexing

Crop Extraction
(with label)



Feature Extraction
from the Crop



Storing vectors and
Crop Size on HD

Off-line part: indexing

crop size

file

features

```
data
└── SIFT
└── VGG
└── YOLO
    ├── aeroplane.csv
    ├── apple.csv
    ├── backpack.csv
    ├── banana.csv
    ├── baseball bat.csv
    ├── baseball glove.csv
    ├── bear.csv
    ├── bed.csv
    ├── bench.csv
    ├── bicycle.csv
    ├── bird.csv
    ├── boat.csv
    ├── book.csv
    ├── bottle.csv
    ├── bowl.csv
    ├── broccoli.csv
    ├── bus.csv
    ├── cake.csv
    ├── car.csv
    ├── carrot.csv
    ├── cat.csv
    ├── cell phone.csv
    ├── chair.csv
    ├── clock.csv
    ├── cow.csv
    ├── cup.csv
    └── diningtable.csv

./val2017/00000002149.jpg 0.332186768198829 "[4.4942074,0.09536005,0.7241454,2.3763955,1.0973221,0$  
./val2017/00000002149.jpg 0.384890954325527 "[0.53094274,0.2661042,0.9301314,4.7975087,0.30232587,$  
./val2017/00000002149.jpg 0.09596750585480093 "[0.0,0.0,0.0,0.0,0.0,0.0,0.53095,0.0,0.0,0.16182597,$  
./val2017/000000482719.jpg 0.071547916666667 "[5.7086277,3.946248,1.1103738,1.2445958,0.0,0.0,0.529$,  
./val2017/000000303566.jpg 0.53461699530516 "[0.5908617,5.0579515,2.6,2.3,0.014211853,2.0619097,0$,  
./val2017/000000125936.jpg 0.016788990825688074 "[0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,$  
./val2017/000000125936.jpg 0.02163914373088685 "[0.14218193,0.0,0.0,0.050794527,0.19629297,0.0,0.0,0.$  
./val2017/000000125936.jpg 0.007486238532110091 "[0.0825237,0.0,0.0,0.11099987,0.22813958,0.0,0.0,0.0,$  
./val2017/000000125936.jpg 0.008807339449541285 "[0.018619329,0.0,0.0,0.039218664,0.1920746,0.0,0.0,0.$  
./val2017/000000409424.jpg 0.03486328125 "[0.3605157,3.4485998,1.0421764,3.7517483,0.0,0.635046$,  
./val2017/000000216277.jpg 0.2702180989583333 "[3.5326602,1.1985853,3.5543804,0.8625804,0.14950421,0$,  
./val2017/000000217753.jpg 0.00984375 "[9.633409,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,$  
./val2017/000000537506.jpg 0.0999 "[1.1631529,5.1990805,6.5065155,26.134287,0.0,0.0,0.0,0.0,$  
./val2017/000000188592.jpg 0.018125 "[0.0,0.0,0.6,8044505,0.0,0.0,0.0,1.1987189,0.0,0.0,0.0,0.0,$  
./val2017/000000189078.jpg 0.08773652694610778 "[3.4091957,4.169929,1.3720208,0.0,0.0,0.0,0.0,0.0,0.0,$  
./val2017/000000562059.jpg 0.10235655737704918 "[2.650178,8.126871,18.168177,0.0,1.3693635,0.0,1.0580$,  
./val2017/000000562059.jpg 0.12172131147540984 "[0.0012675449,0.038269904,0.0,0.0007994231,0.01963814$,  
./val2017/000000119452.jpg 0.019166666666666665 "[0.0,0.0,0.0,0.0,0.0,0.0,0.0,3.6931825,0.0,2.9052842,$  
./val2017/000000119452.jpg 0.013151041666666667 "[0.0,0.0,0.0,0.0,3.5074928,0.0,0.0,5.957371,0.0,19.35$,  
./val2017/000000388903.jpg 0.2027710843373494 "[0.0,0.0,0.9,134814,0.0,0.07538362,0.0,0.46320564,2.347$,  
./val2017/000000546626.jpg 0.08265 "[0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.02764116,0.0,0.0,0.0,0.0,$  
./val2017/000000185599.jpg 0.33 "[1.599073,1.3286817,0.00995423,0.4212316,0.22010246,0$,  
./val2017/000000193926.jpg 0.02458071278825996 "[1.9731374,0.0,0.0,0.0,0.0,0.0,0.6,4.478349,0.0,0.0,1$,  
./val2017/000000486438.jpg 0.1002195550351288 "[6.423902,0.0,0.4,0.23174,20.958826,0.07910835,0.0,5.461$,  
./val2017/00000044260.jpg 0.015779411764705882 "[0.119758904,0.0,0.1090616,0.072967544,0.22672345,0.0$,  
./val2017/00000044260.jpg 0.013191176470588236 "[41.481777,0.0,0.0,0.0,3.390477,0.0,0.0,14.671722,0.0$,  
./val2017/000000203095.jpg 0.03017578125 "[0.0,0.0,0.0,0.0,0.5880337,0.0,0.0,0.0,0.0,0.0,0.0,0.0,$  
./val2017/000000113589.jpg 0.064453125 "[0.63864505,0.0,0.0,0.16,6.67279,0.0,0.0,0.0,1.1338391,0$,  
./val2017/000000113589.jpg 0.1639388020833333 "[11.348594,0.36317176,1.8471438,3.5356371,2.1124234,0$,  
./val2017/00000060855.jpg 0.1454035087719298 "[4.278553,0.34758127,1.1238421,0.26897964,0.0,0.0,0.0$,  
./val2017/00000060855.jpg 0.20057142857142857 "[4.524483,4.2390027,2.3613894,1.3759605,2.7867727,5.4$,  
./val2017/000000509131.jpg 0.026356617647058822 "[16.463696,0.0,0.0,0.8880062,0.0,0.0,0.10.29312,0.0$,  
./val2017/000000421721.jpg 0.026356617647058822 "[14.020125,0.0,0.0,0.702722,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,$  
Treemacs
~D/P/DSIM-Project/data/VGG/apple.csv 1:0 Top
LF UTF-8 CSV
<Scroll_Lock> is undefined
```

OFF-line part: indexing

feature extractor

tag

On-line Part: Query

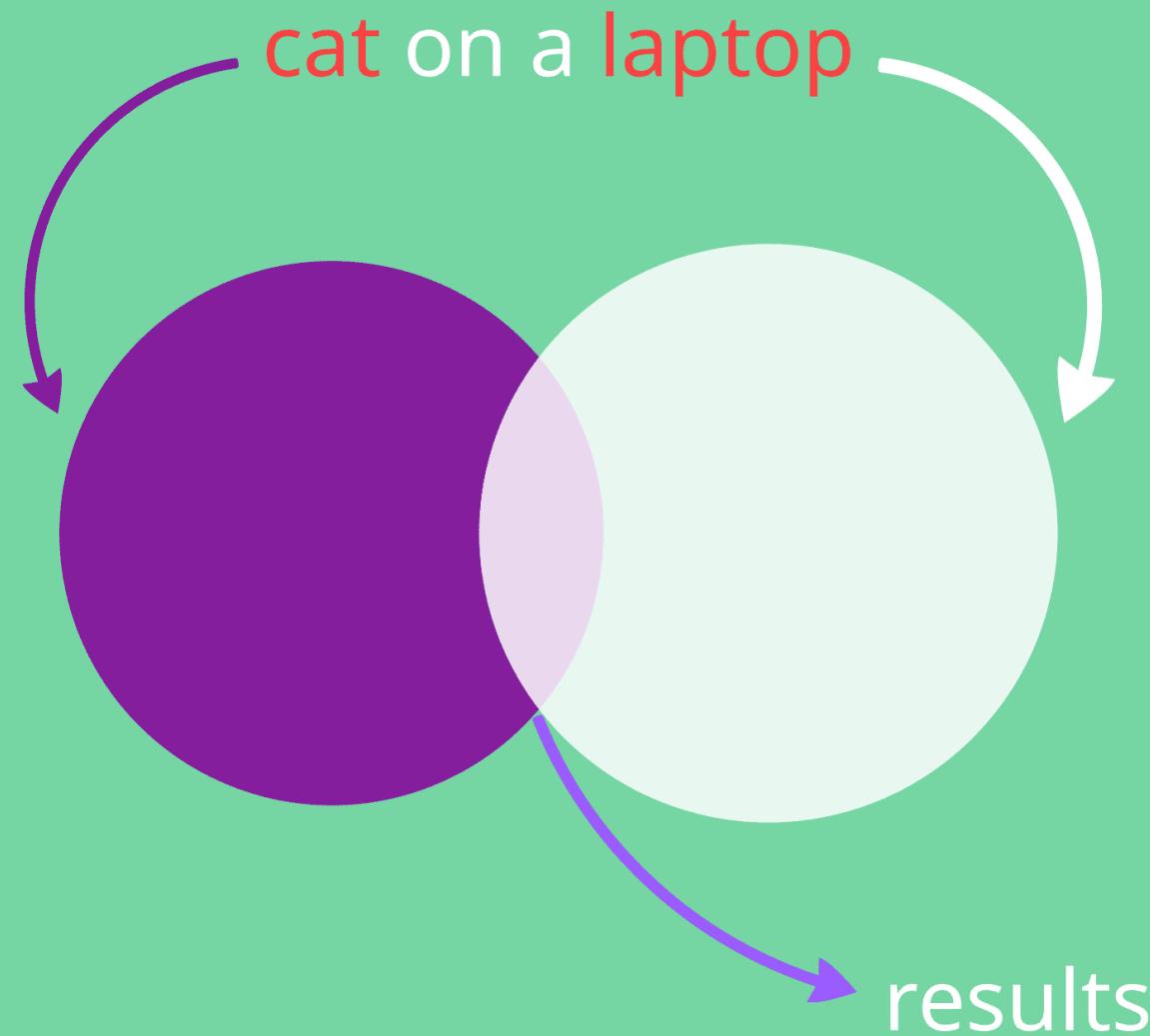
Boolean Matching
of the Text

Loading vectors

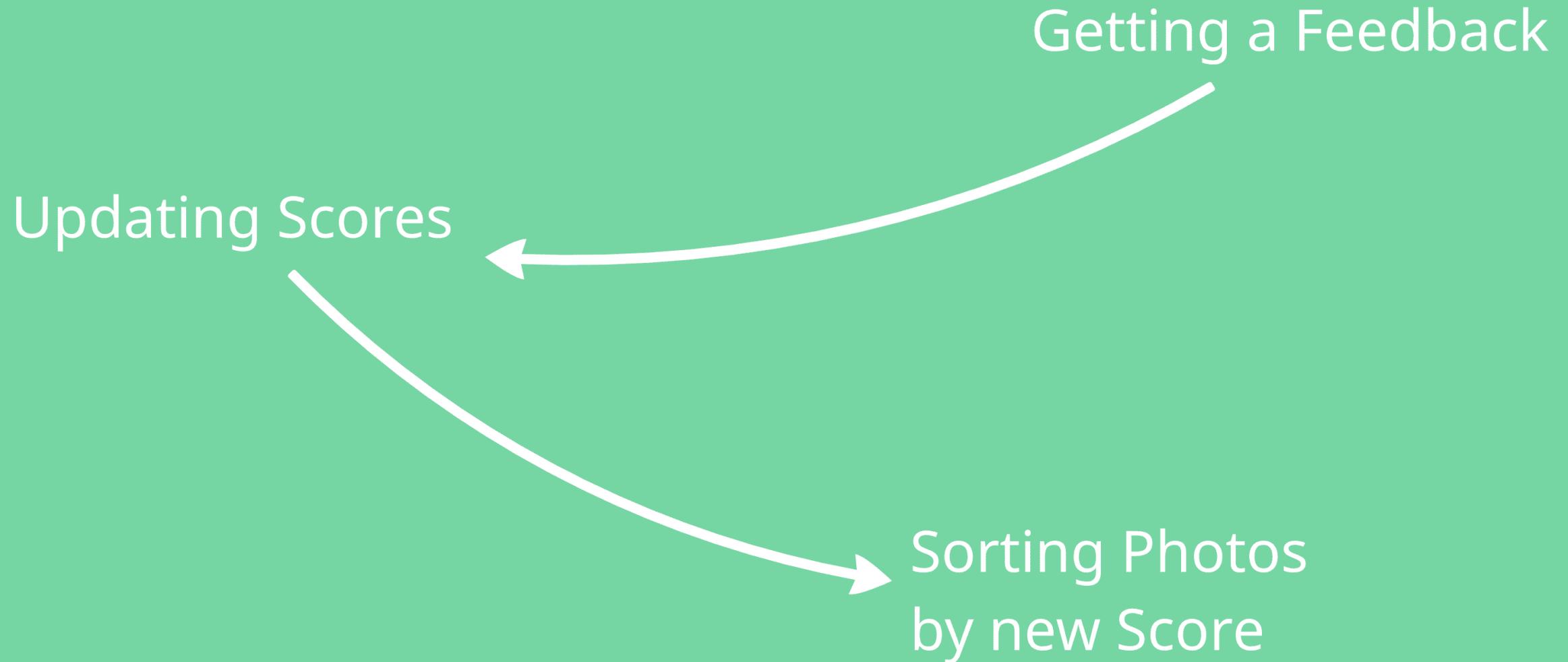


Sorting Photos
by Crop Size

On-line Part: Query



Interaction: Updating Scores



Interaction: Updating Scores

Fist matches between crops are found
(if I select a cat, I want similar cats)

Then the crop score
is updated

Total score for the photo is
the mean of its crop scores

```
selected ← selected photo
w ← weight of the update
∀ photo
    ∀ crop ∈ photo
         $dist_{crop, selected} \leftarrow \min_{crop_{selected}} \{dist(crop, crop_{selected})\}$ 
         $score_{crop} \leftarrow (1 - w) \cdot score_{crop} + w \cdot dist_{crop, selected}$ 
     $score_{photo} \leftarrow mean(score_{crop} \forall crop \text{ in photo})$ 
```

ALL DISTANCES ARE NORMALIZED

Performance

Efficiency: fast and scalable

$O(n)$

Effectiveness: most results ($> 75\%$)
are found in < 3 iterations

Results found using our Search Engine

