

20/11/2024

Convolutional Neural Networks for Object Detection

Giacomo Boracchi

giacomo.boracchi@polimi.it

Artificial Neural Networks and Deep Learning

Politecnico di Milano,

<https://boracchi.faculty.polimi.it/>

Object Detection



Object Detection, the problem

Assign to an input image $I \in \mathbb{R}^{R \times C \times 3}$:

- multiple labels $\{l_i\}$ from a fixed set of categories $\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$, each corresponding to an instance of that object
- the coordinates $\{(x, y, h, w)_i\}$ of the bounding box enclosing each object

$$I \rightarrow \{(x, y, h, w, l)_1, \dots, (x, y, h, w, l)_N\}$$

Object Detection Task

$$I: \rightarrow \{(1 \times R \times R \times R \times R)\}$$

Given a fixed set of categories and an input image which contains an unknown and varying number of instances

Draw a bounding box on each object instance

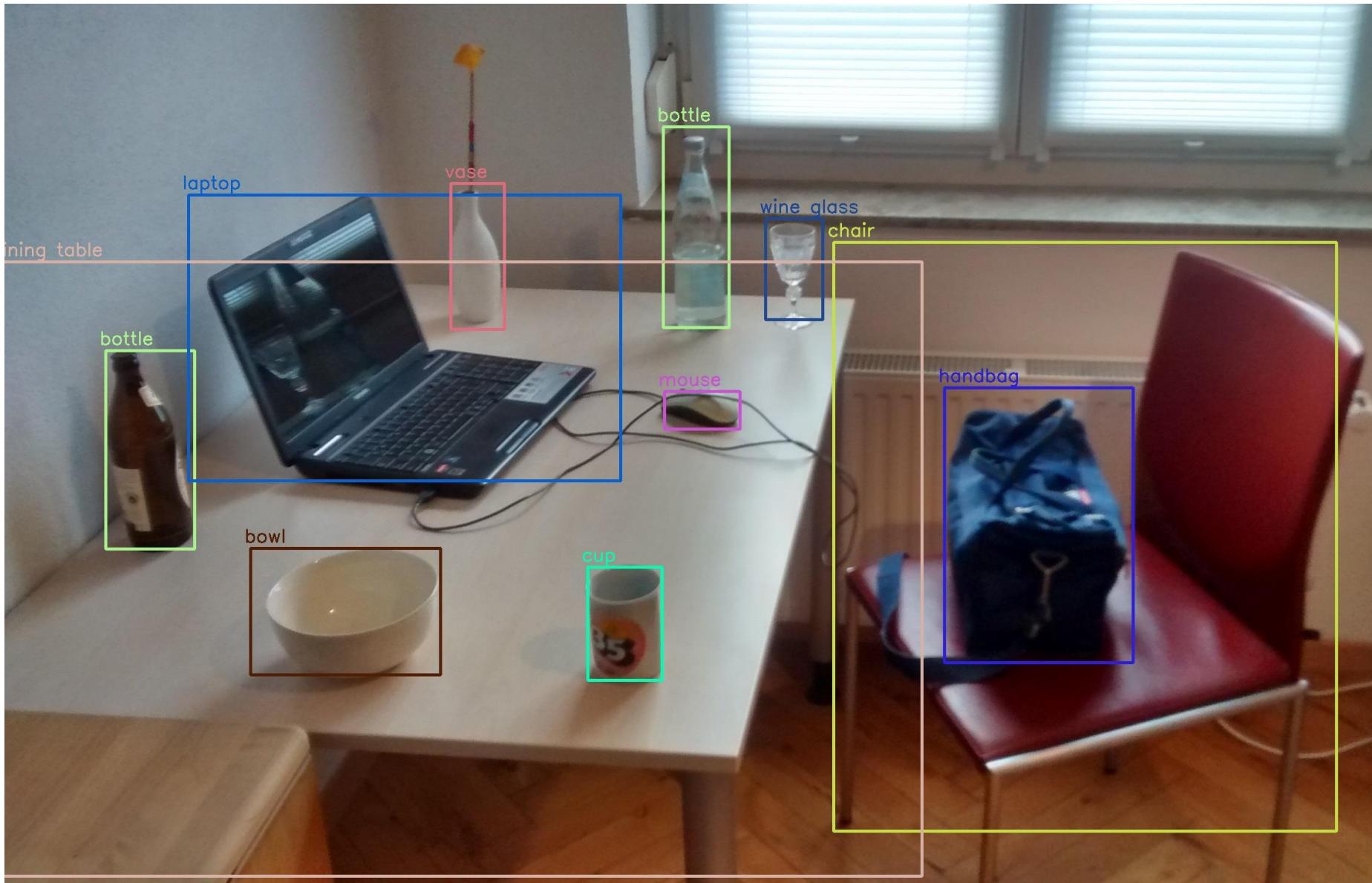
A training set of annotated images with labels and bounding boxes for each object is required

Each image requires a varying number of outputs

MAN: (x, y, h, w)
KID: (x, y, h, w)
GLOVE: (x, y, h, w)



Annotated Dataset for Object Detection

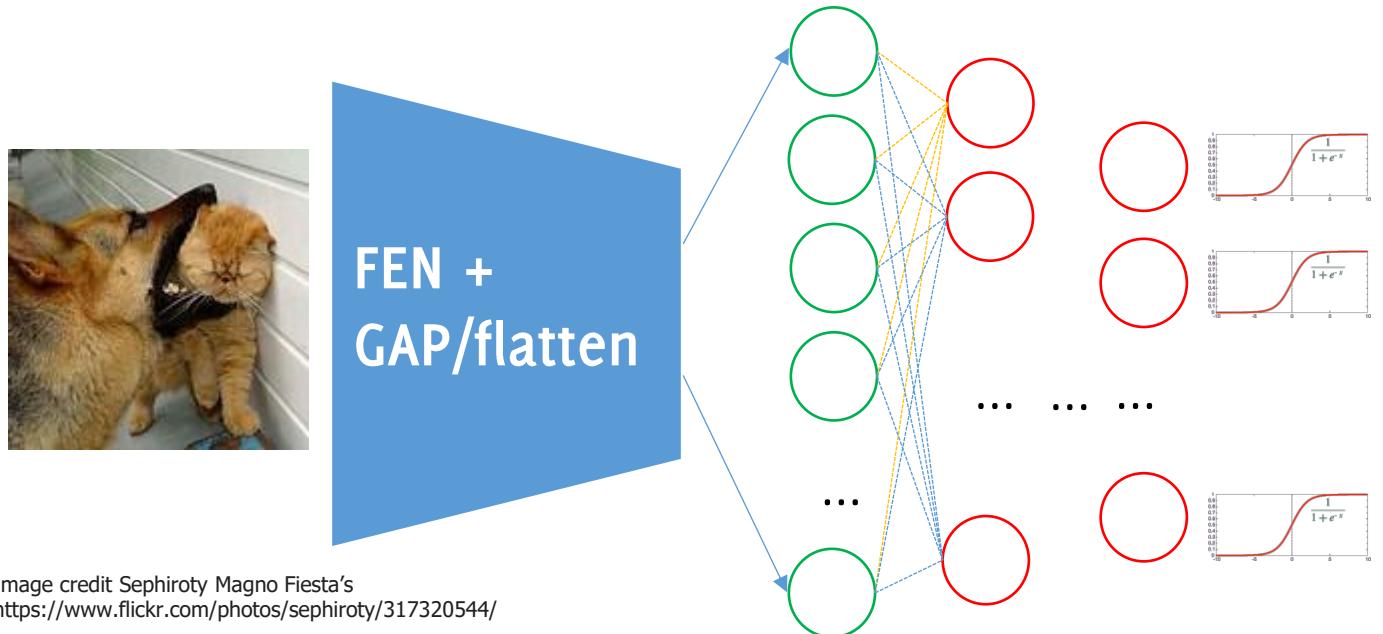


Excursus: Multi-label Classification

So far, all the models we have seen provide a fixed-sized output

CNNs (like other NN) can be modified to return **multiple classification output** by:

- Replace softmax activation in the output layer with sigmoid activation in each of the L neurons.
- The i – th neuron will predict the probability $p_i \in [0,1]$ of class membership for the input in a non-exclusive manner.
- The model must be trained with the binary cross-entropy loss function.



$$p_1 \in [0,1]$$

$$p_2 \in [0,1]$$

$$p_L \in [0,1]$$

These probabilities do not sum to 1, an input image might belong to multiple categories, depending on the content. This picture will likely be assigned to class «dog» and «cat»

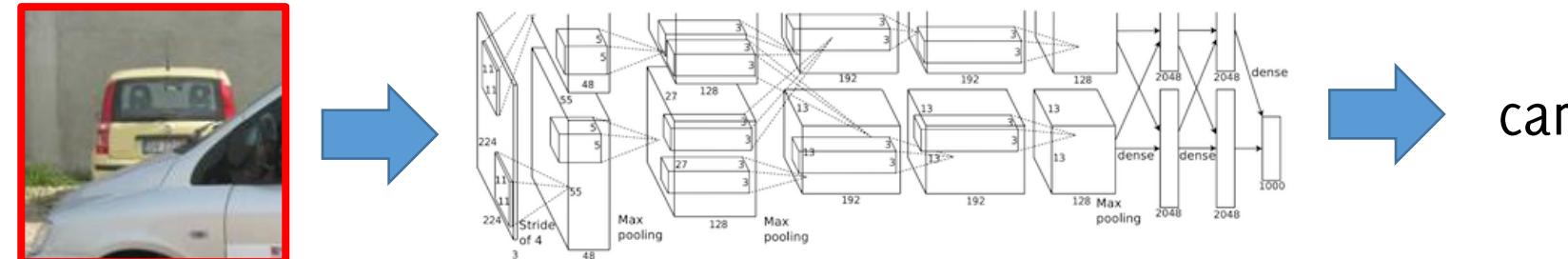
The Straightforward Solution: Sliding Window

1000 x 2000 pixels



- Similar to the sliding window for semantic segmentation
- A pretrained model is meant to process a fixed input size (e.g. 224 x 224 x 3)
- Slide on the image a window of that size and classify each region.

Adopt the whole machinery seen so far to each crop of the image



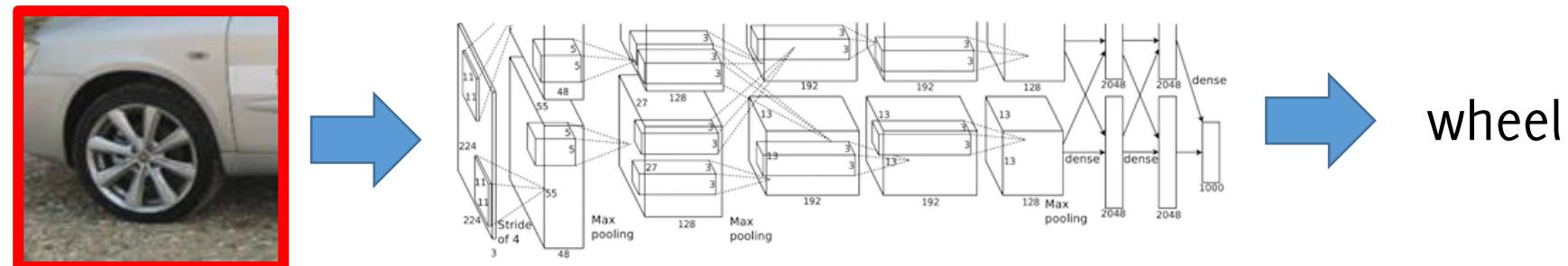
The Straightforward Solution: Sliding Window

1000 x 2000 pixels



- Similar to the sliding window for semantic segmentation
- A pretrained model is meant to process a fixed input size (e.g. 224 x 224 x 3)
- Slide on the image a window of that size and classify each region.

Adopt the whole machinery seen so far to each crop of the image



The Straightforward Solution: Sliding Window

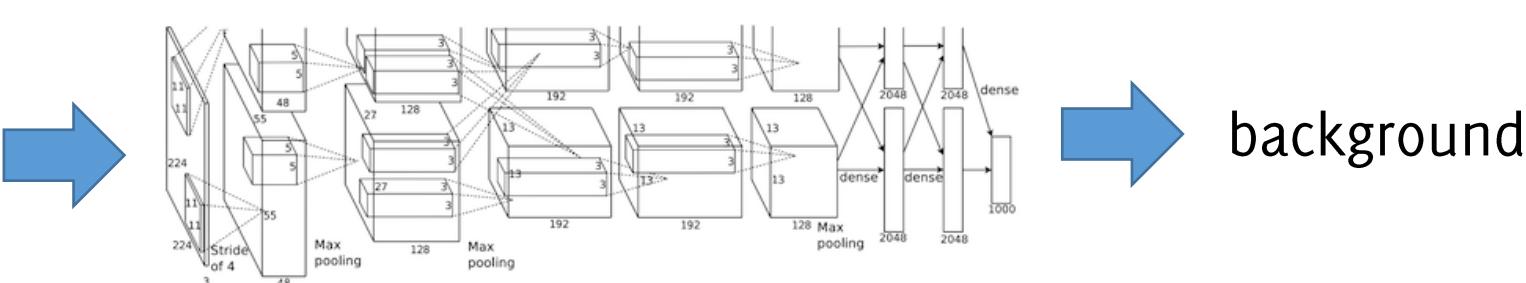
1000 x 2000 pixels



- Similar to the sliding window for semantic segmentation
- A pretrained model is meant to process a fixed input size (e.g. 224 x 224 x 3)
- Slide on the image a window of that size and classify each region.

Adopt the whole machinery seen so far to each crop of the image

The background class has to be included!



Many drawbacks...

Cons:

- **Very inefficient!** Does not re-use features that are «shared» among overlapping crops
- How to choose the crop size?
- Difficult to detect objects at different scales!
- A huge number of crops of different sizes should be considered....

Plus:

- No need of retraining the CNN



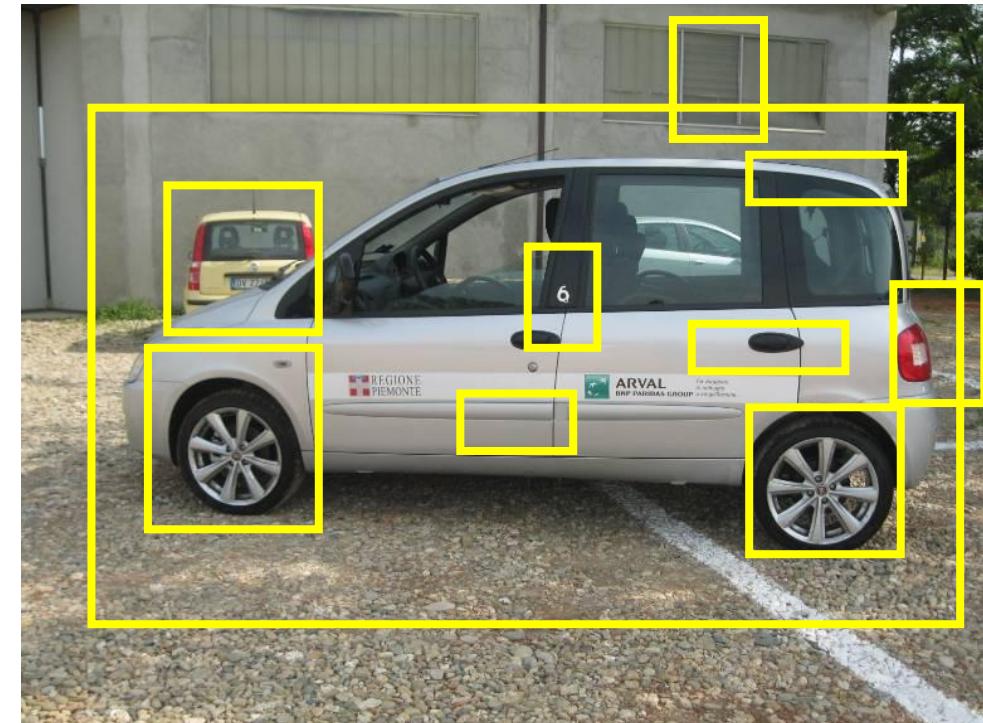
Region Proposal

Region proposal algorithms (and networks) are meant to identify bounding boxes that correspond to a candidate object in the image.

⚠️ Algorithms with very high recall (but low precision) were there before the deep learning advent

The idea is to:

- Apply a region proposal algorithm
- Classify by a CNN the image inside each proposal regions





This CVPR2014 paper is the Open Access version, provided by the Computer Vision Foundation.
The authoritative version of this paper is available in IEEE Xplore.

Rich feature hierarchies for accurate object detection and semantic segmentation

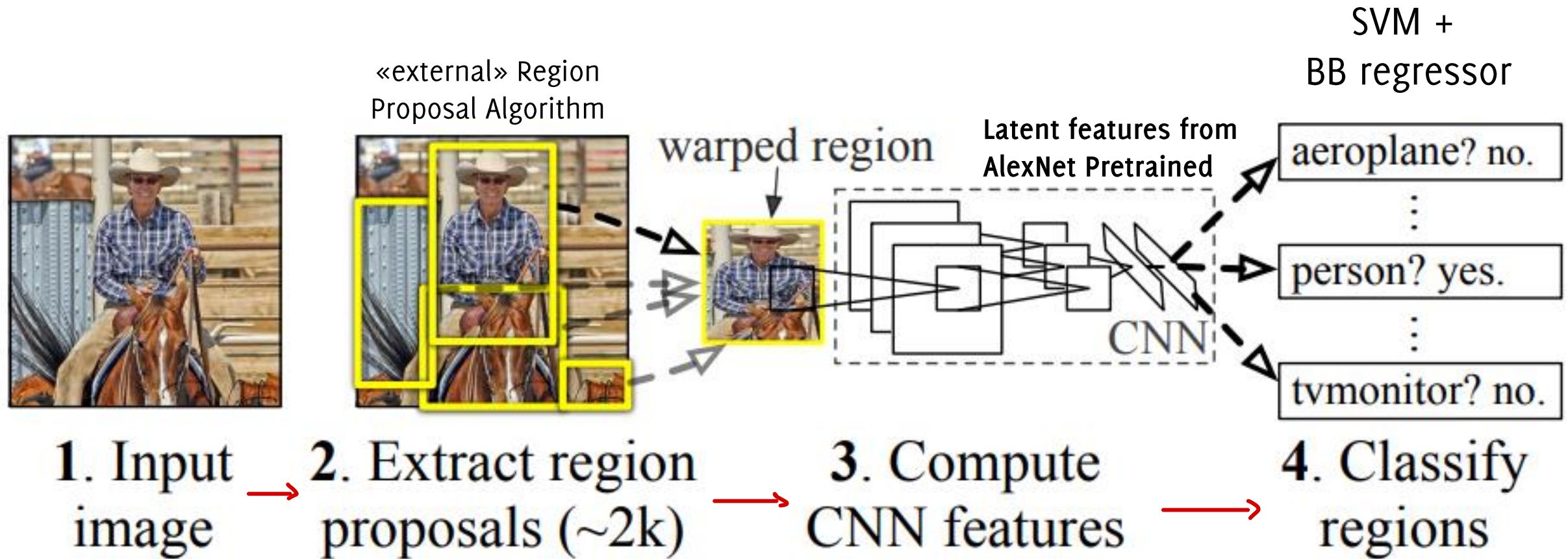
Ross Girshick¹ Jeff Donahue^{1,2} Trevor Darrell^{1,2} Jitendra Malik¹

¹UC Berkeley and ²ICSI

{rbg, jdonahue, trevor, malik}@eecs.berkeley.edu

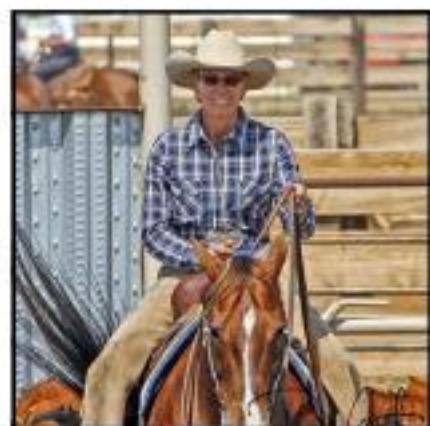
R-CNN

Object detection by means of region proposal (R stands for regions)



R-CNN

Object detection by means of region proposal



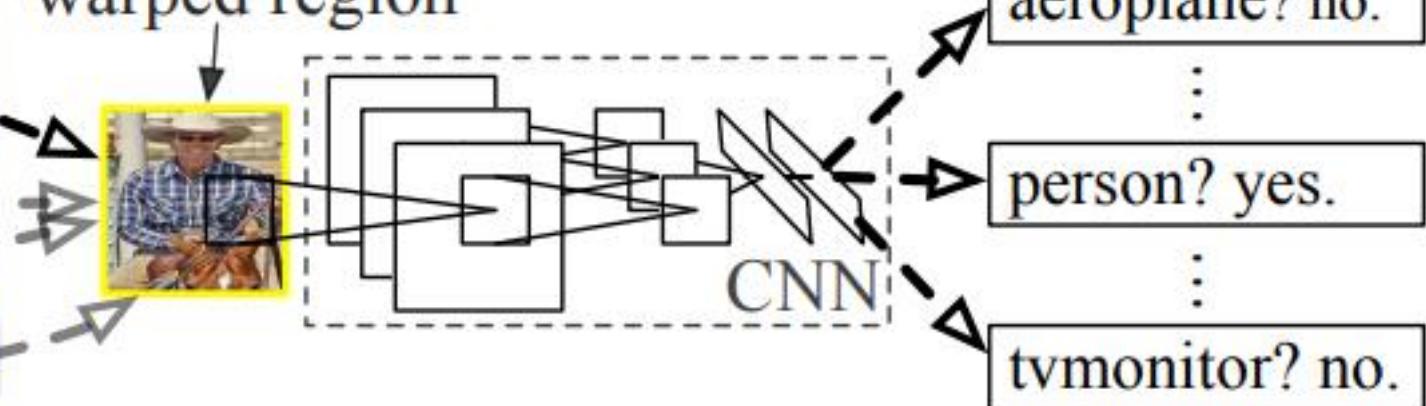
1. Input
image



2. Extract region
proposals (~2k)

Warping to a predefined
size is necessary since the
CNN has a FC layer

warped region

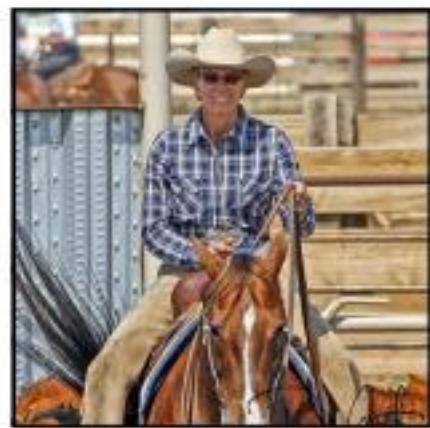


3. Compute
CNN features

4. Classify
regions

There is no learning in the region proposal
algorithm, very high recall (e.g. Selective Search)

R-CNN

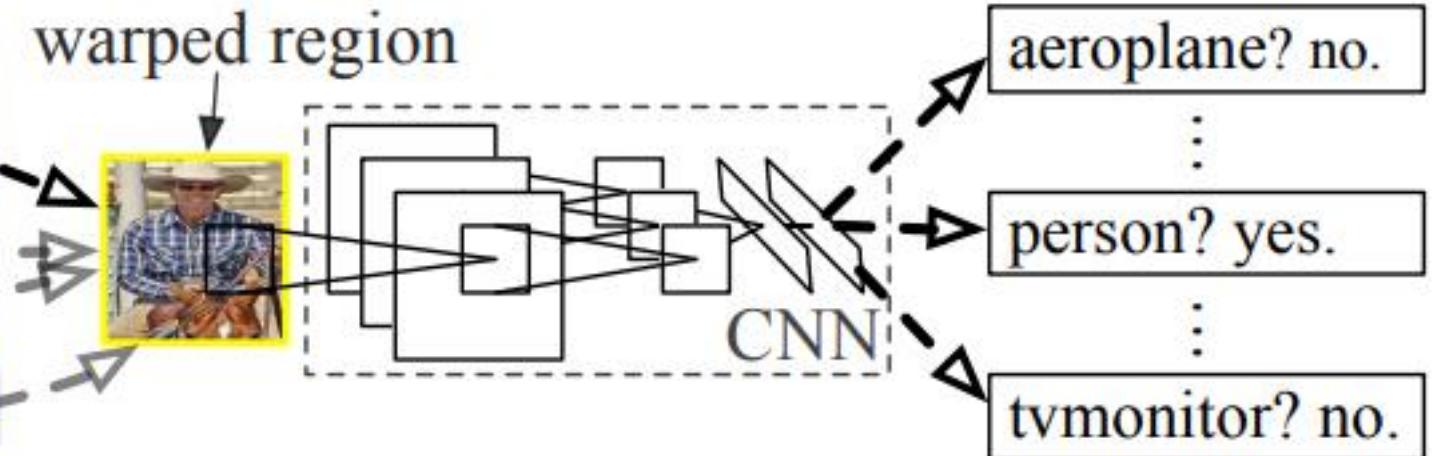


1. Input image



2. Extract region proposals (~2k)

warped region



3. Compute CNN features

SVM +
BB regressor

SVM trained to minimize the classification error over the extracted ROI

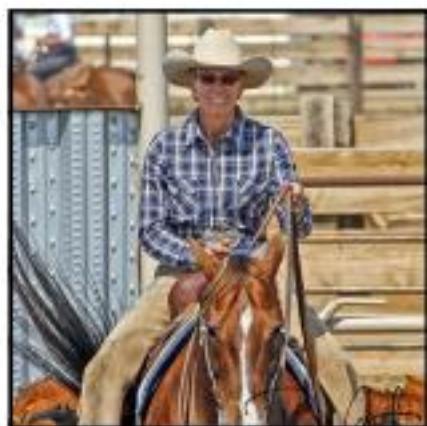
aeroplane? no.

⋮
person? yes.

⋮
tvmonitor? no.

4. Classify regions

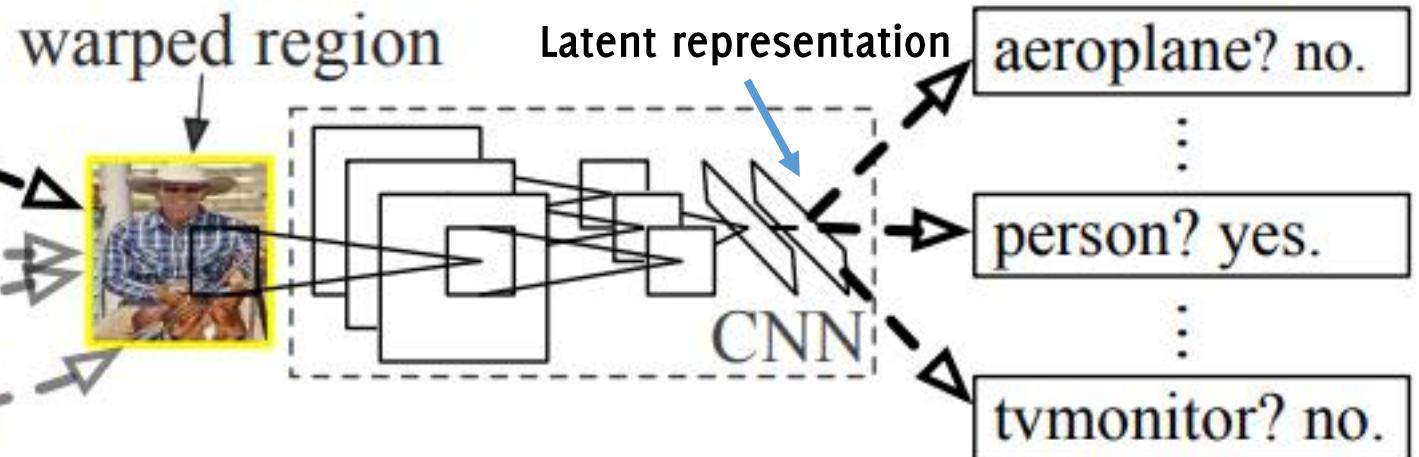
R-CNN



1. Input image



2. Extract region proposals (~2k)



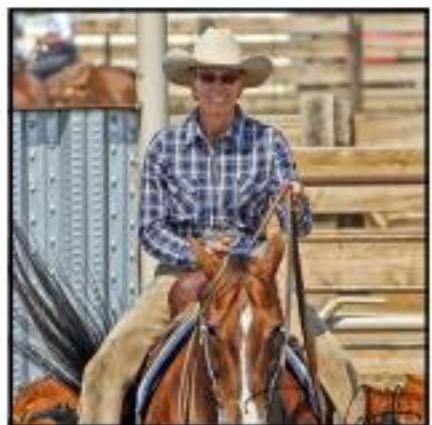
3. Compute CNN features

4. Classify regions

R-CNN

The pretrained CNN is fine-tuned over the classes to be detected (21 vs 1000 of Alexnet) by placing a FC layer after feature extraction.
No end-to-end training of the SVM

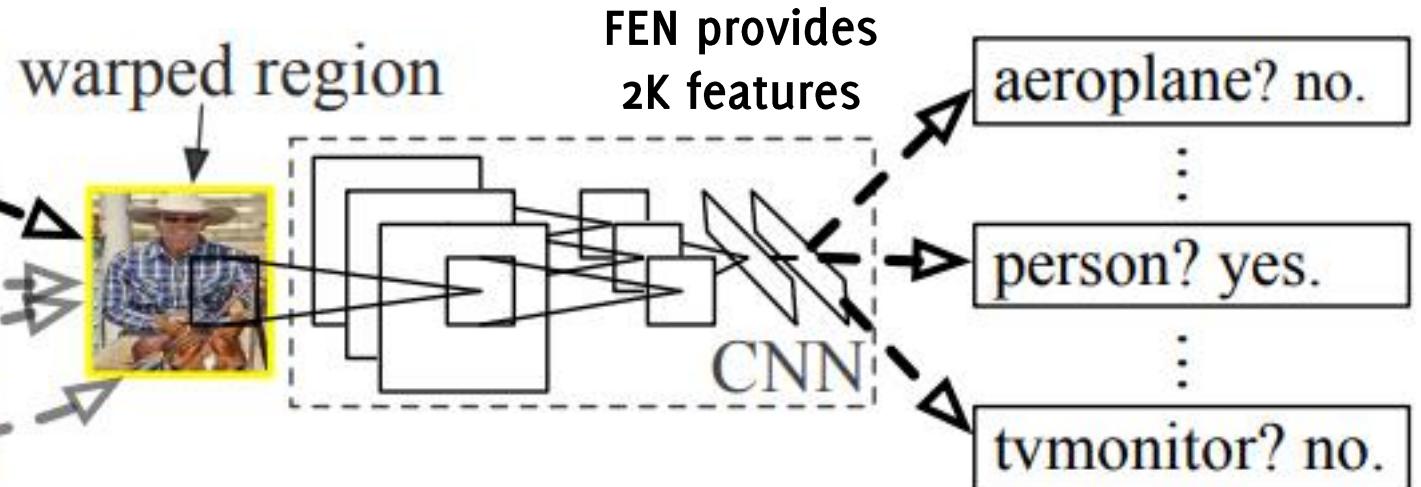
SVM +
BB regressor



1. Input image



2. Extract region proposals (~2k)



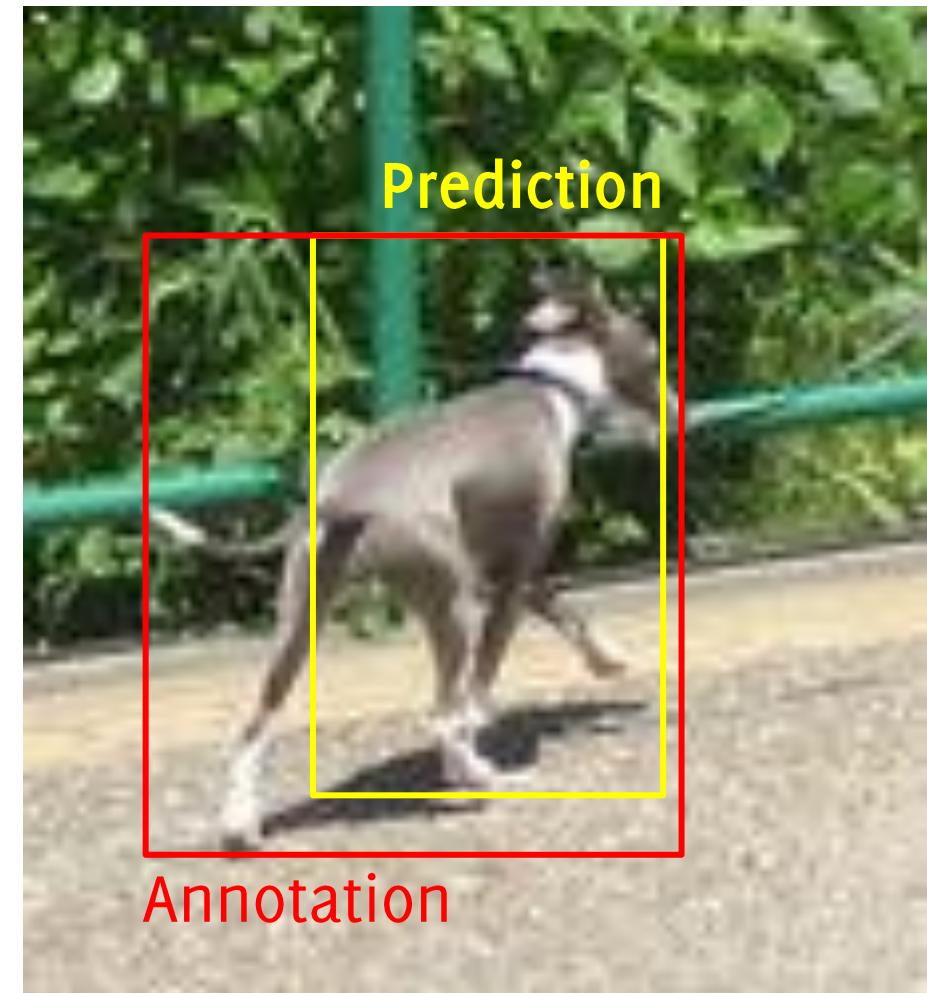
3. Compute CNN features

4. Classify regions

Include a background class to get rid of those regions not corresponding to an object

The Loss Function (over bounding boxes)

We need to quantitatively assess the network performance over each and every image in the test set



The Loss Function (over bounding boxes)

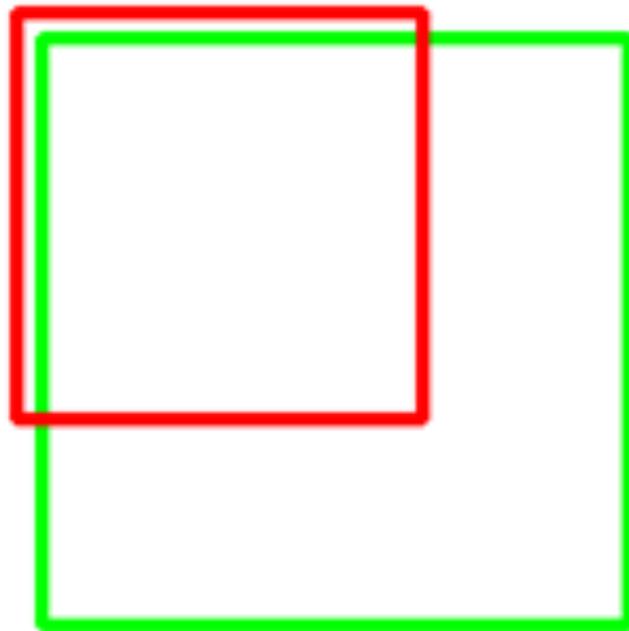
We need to quantitatively assess the network performance over each and every image in the test set

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



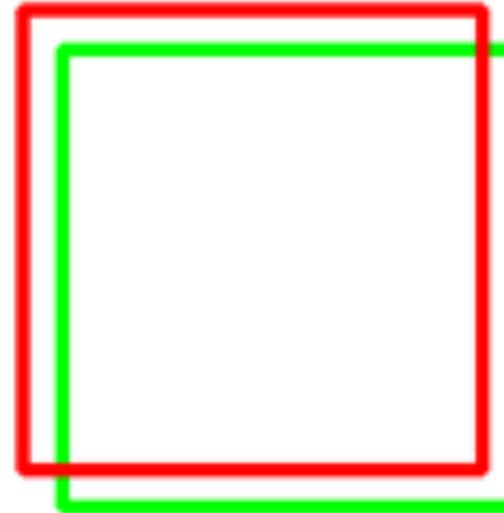

A number quantitatively assessing detection performance

IoU: 0.4034



Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

The Loss Function

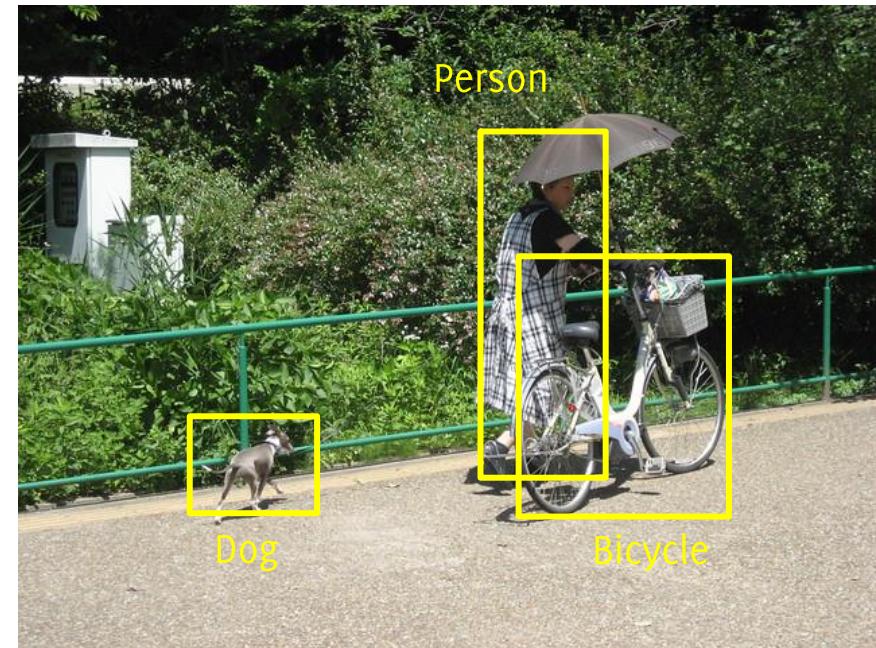
The loss function compares the detection results for an image x : $NN(x)$ against the annotations y and gives you a number

$$\mathcal{L}(y, \hat{y})$$

That is indicating «*how happy we are with the predictions*», the lower the better. This considers:

- how many missed items → classified as background but it was an object.
- how many false positives → classified as object but it was background.
- how close correct detections are to the annotation (IoU)

about bounding box.



R-CNN Limitations

- Ad-hoc training objectives and not an end-to-end training
 - Fine-tuning network with softmax classifier (log loss) before training SVM
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressions (least squares)
- Region proposals are from a different algorithm and that part has not been optimized for the detection by CNN
- Training is slow (84h), takes a lot of disk space to store features
- Inference (detection) is slow since the CNN has to be executed on each region proposal (no feature re-use)
 - 47s / image with VGG16

Rmk: efficiency in object detection network is key! Otherwise you might want to train a segmentation network instead!



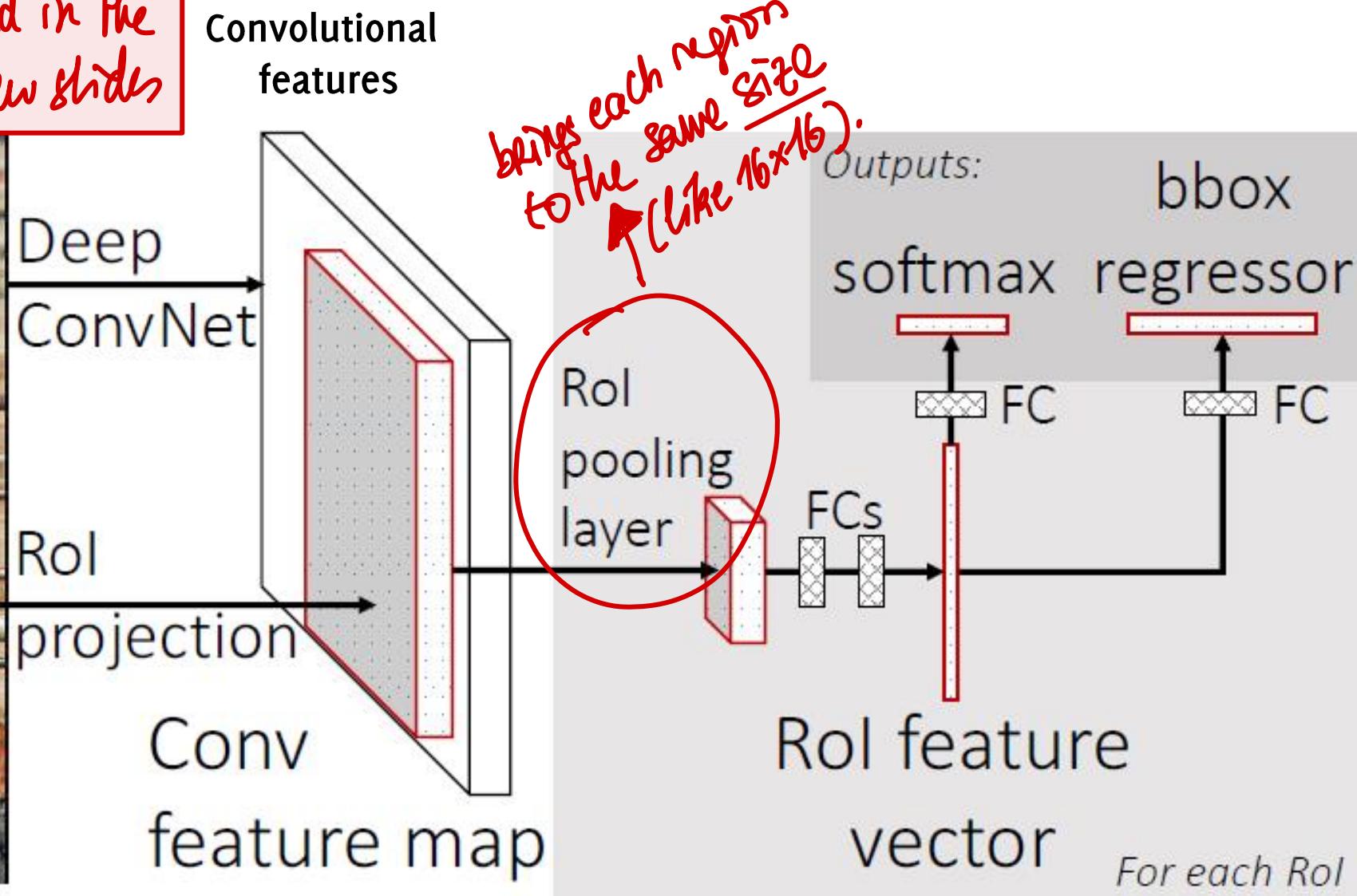
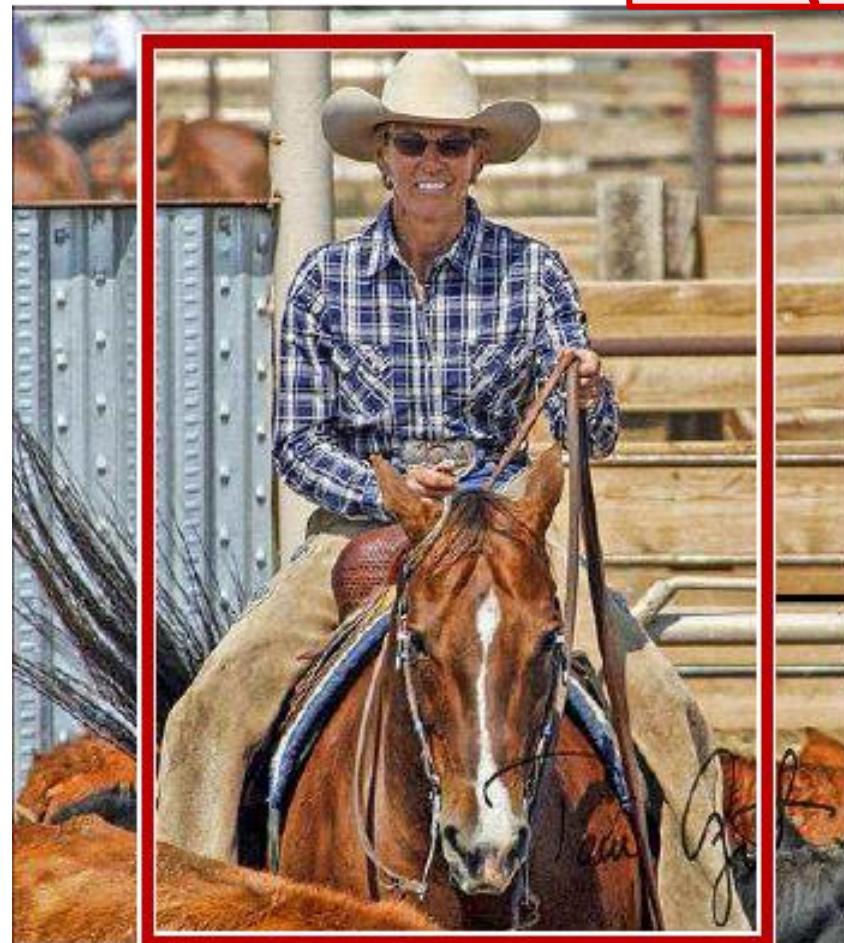
This ICCV paper is the Open Access version, provided by the Computer Vision Foundation.
Except for this watermark, it is identical to the version available on IEEE Xplore.

Fast R-CNN

Ross Girshick
Microsoft Research
rbg@microsoft.com

Fast R-CNN

Explained in the
next few slides

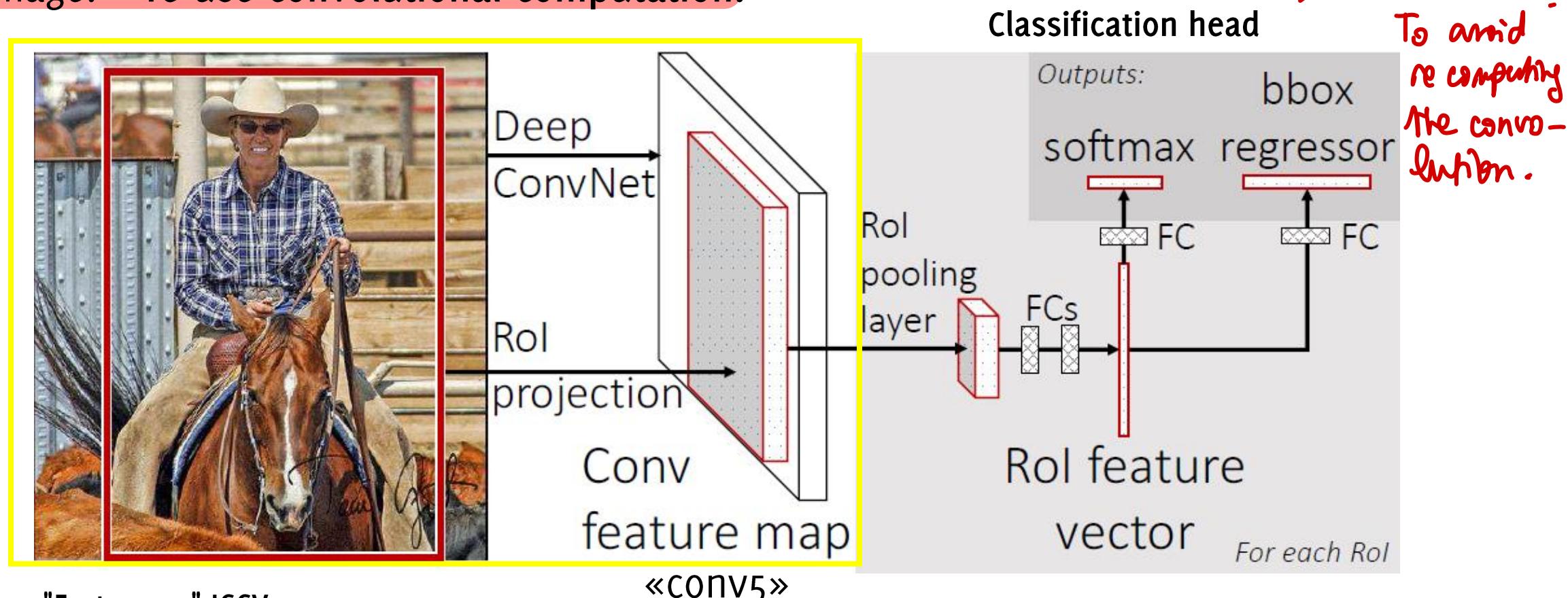


Region from external RPA

⚠ We get rid of the SVM.
→ replace it by NN!

Fast R-CNN

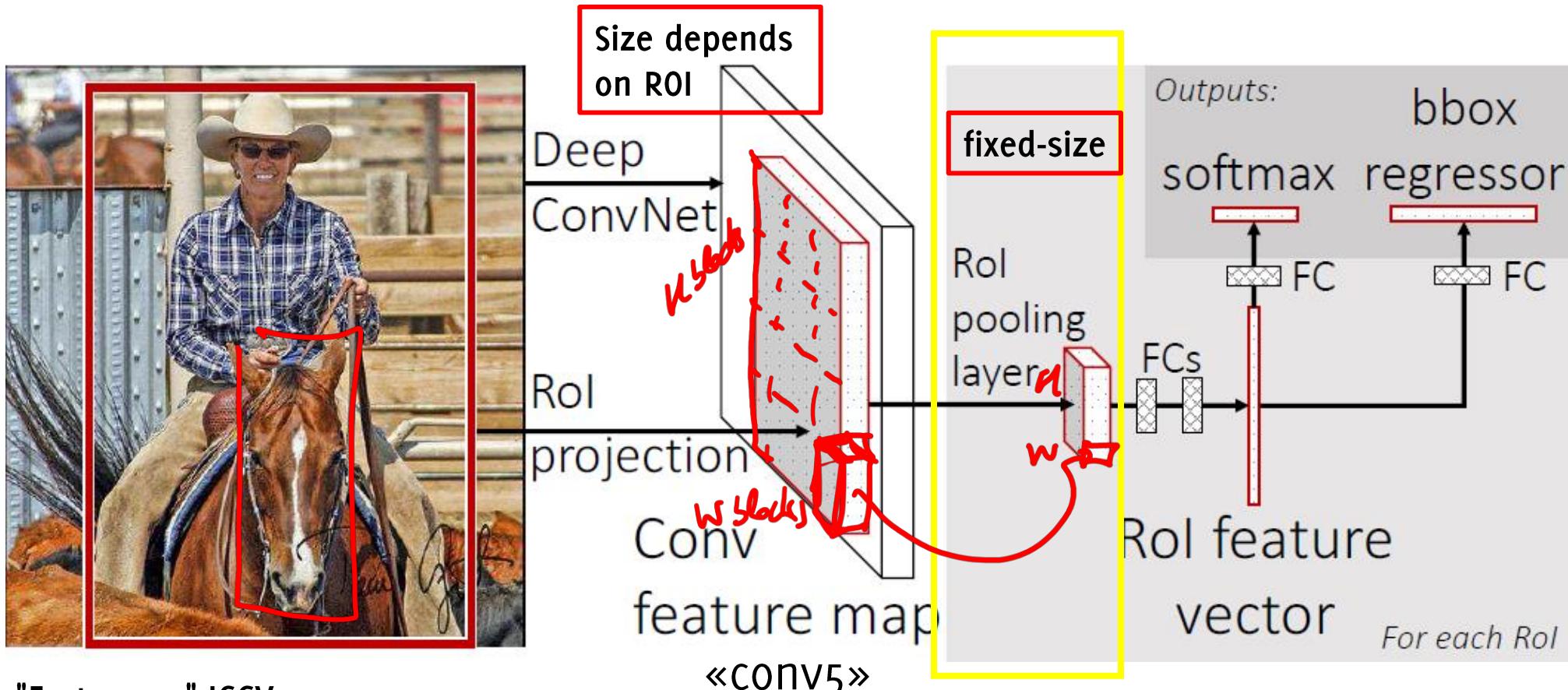
1. The whole image is fed to a CNN that extracts feature maps.
2. Region proposals are identified from the image and projected into the feature maps. Regions are directly cropped from the feature maps, instead from the image: →re-use convolutional computation.



Fast R-CNN

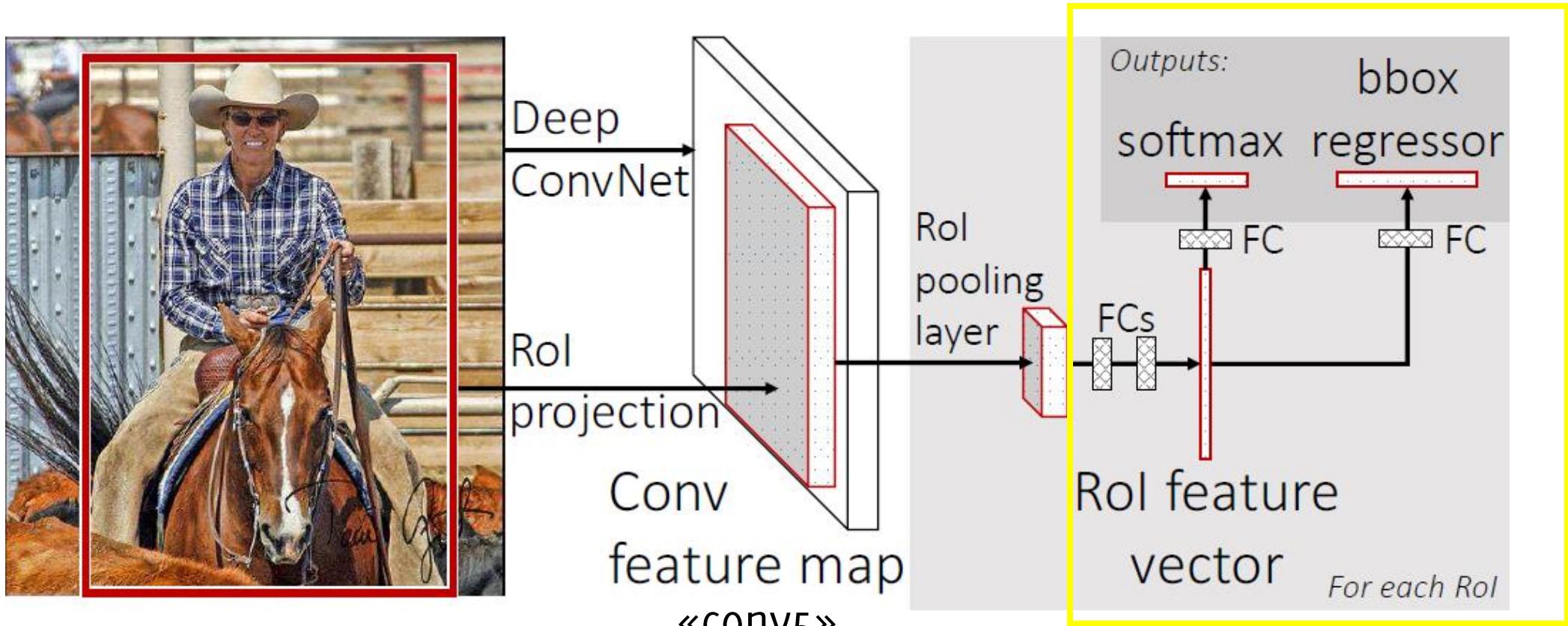
3. Fixed size is still required to feed data to a fully connected layer.

ROI pooling layers: extract a **fixed size $H \times W$** activation from each region proposal. Each ROI in the feature maps is divided in a $H \times W$ grid and then maxpooling over the grid provides a fixed size input (vectorized) to the next step.



Fast R-CNN

4. The FC layers estimate both classes and BB location (bb regressor)
A convex combination of the two is used as a **multitask loss** to be optimized (as in R-CNN, but no SVM here).
↑ cf slide 21.
5. Training performed in an **end-to-end** manner, convolutional part executed only once



Fast R-CNN

In this new architecture it is possible to **back-propagate through the whole network**, thus train the whole network in an end-to-end manner.

It becomes **incredibly faster than R-CNN during testing**.

Now that convolutions are not repeated on overlapping areas, **the vast majority of test time is spent on ROI extraction** (e.g. selective search).

Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks

Shaoqing Ren* **Kaiming He** **Ross Girshick** **Jian Sun**

Microsoft Research

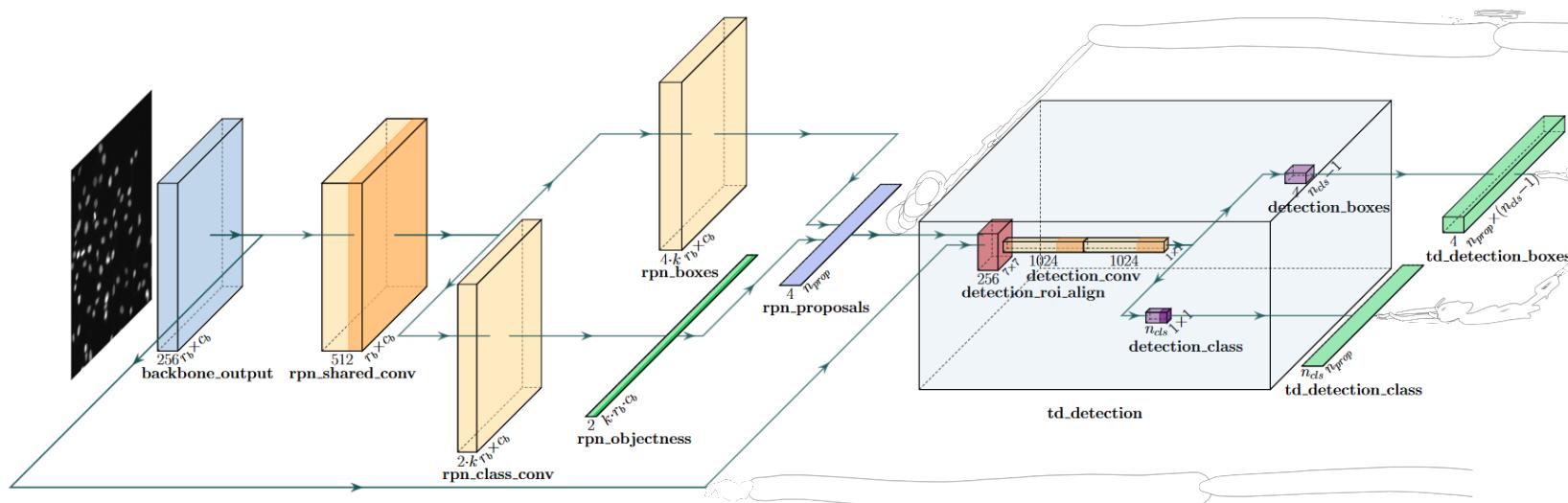
{v-shren, kahe, rbg, jiansun}@microsoft.com

Faster R-CNN

Conv backbone
Classification

What is new : you manage to make learnable the ROI extraction .

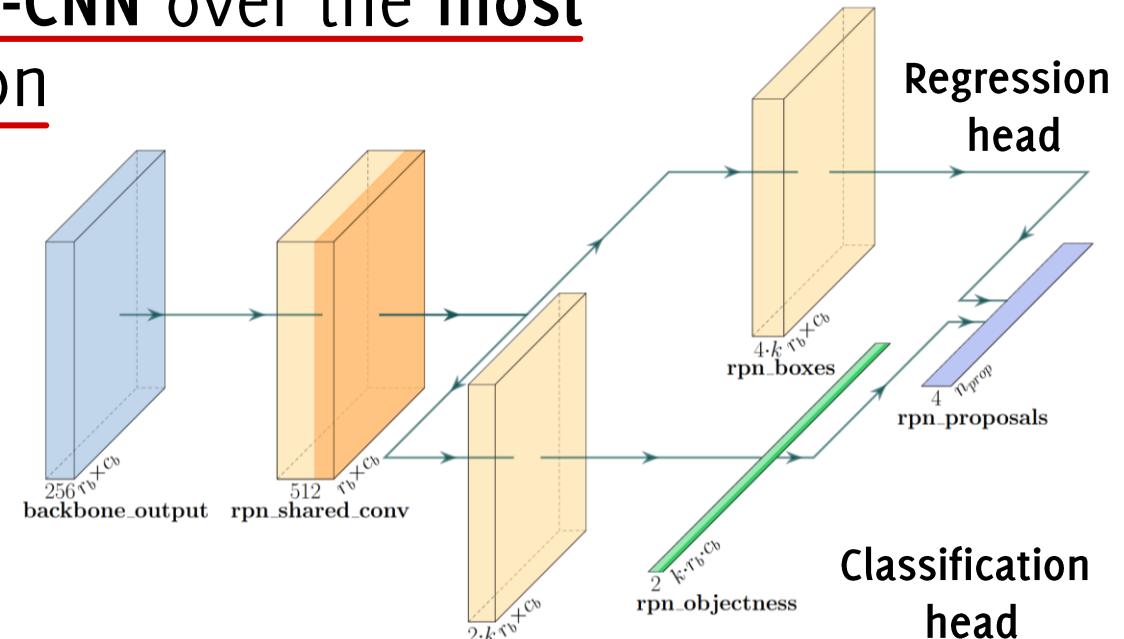
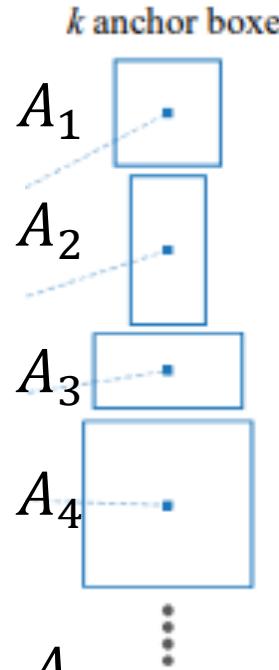
Detection head



RPN: The Region Proposal Network

Explained in the
next few slides

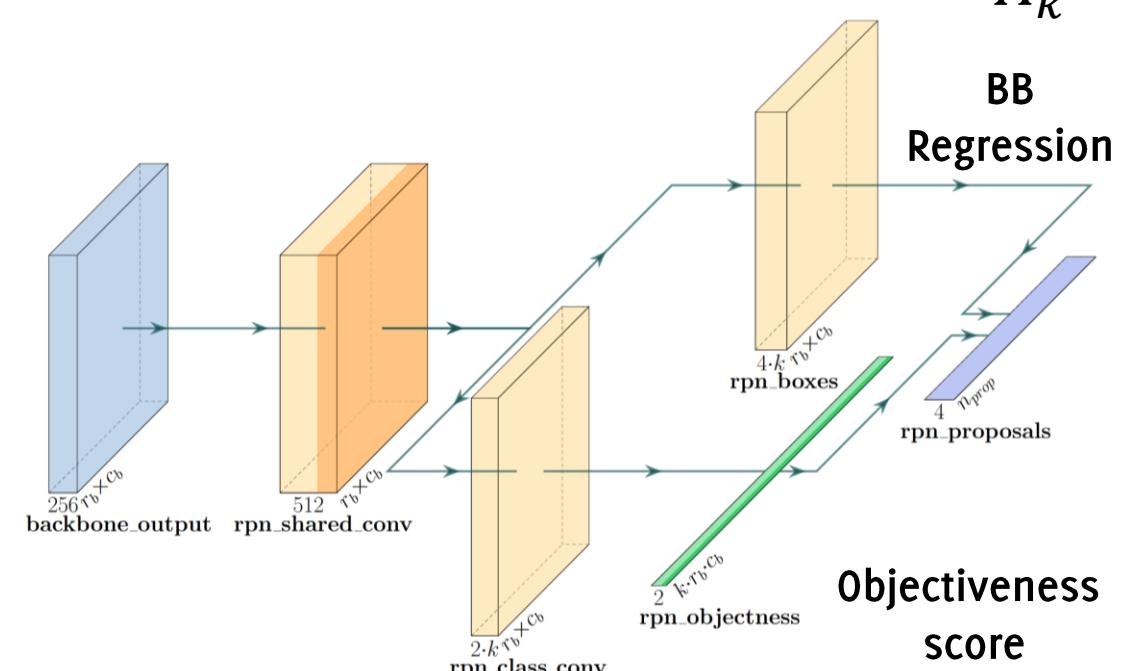
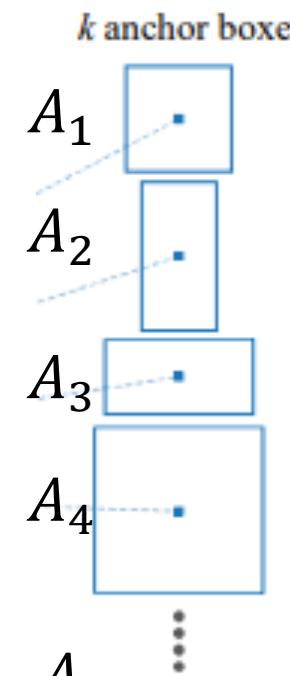
- Instead of the ROI extraction algorithm, **train a region proposal network (RPN)**, which is a F-CNN (3×3 filter size)
- **RPN operates on the same feature maps used for classification**, thus at the last conv layers
- RPN can be seen as an **additional (learnable) module** that **improves efficiency and focus Fast R-CNN over the most promising regions for object detection**



RPN: Region Proposal Network

Goal: Associate to each spatial location k anchor boxes, i.e. ROI having different scales and ratios (e.g. $k = 3 \times 3$, 3 sizes of the anchor side, 3 height/width ratios). Assume the feature maps are $r_b \times c_b$.

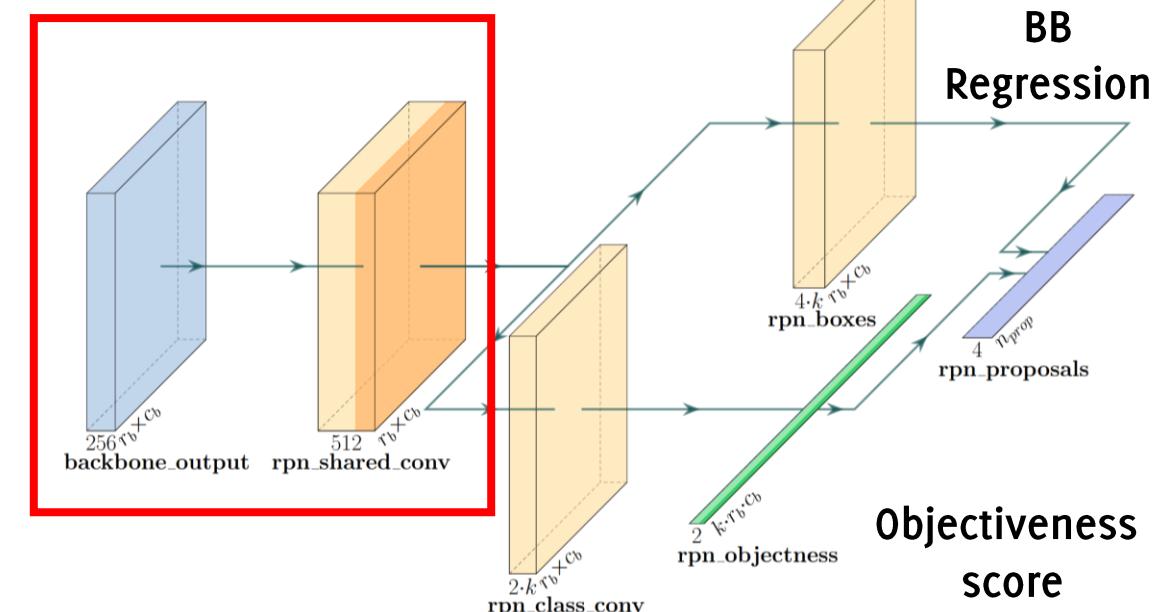
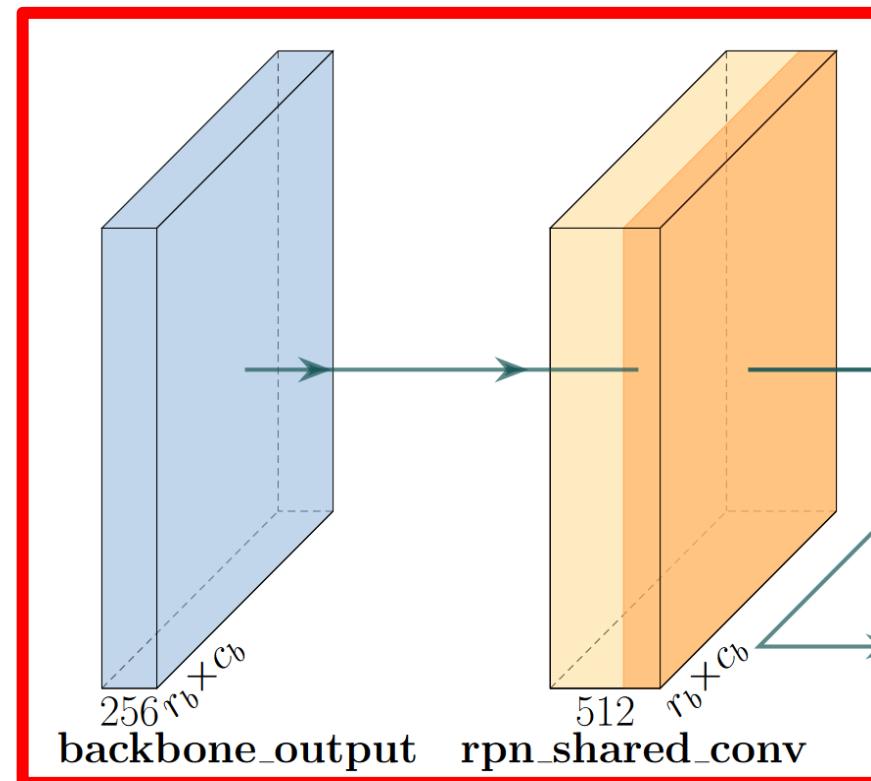
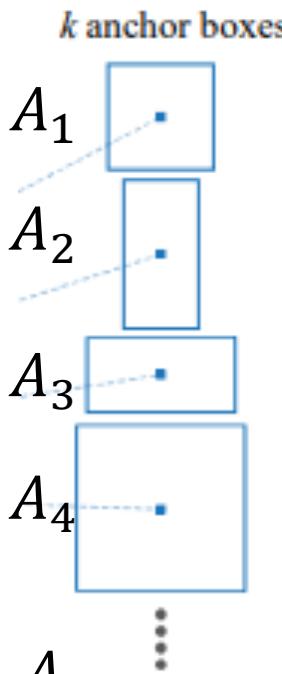
The network outputs $r_b \times c_b \times k$ candidates anchor and estimate objectiveness scores for each anchor



RPN: Intermediate Layer

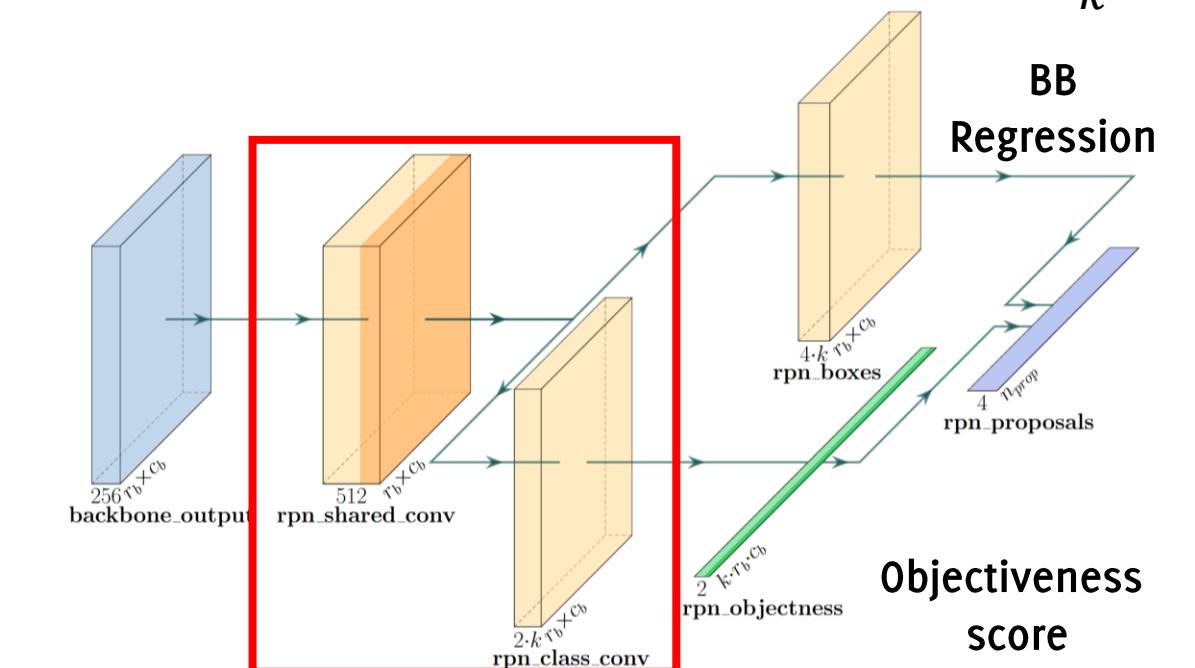
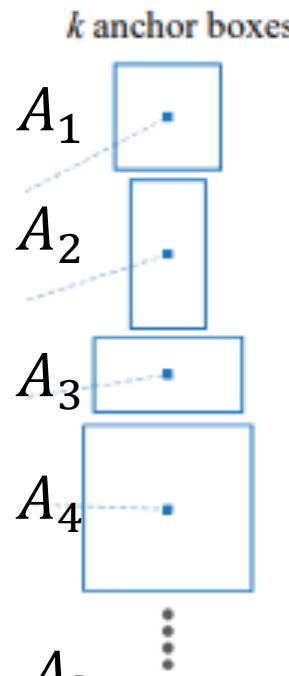
Intermediate layer: conv2D that takes as input the last layer of the feature-extraction network and uses 256 filters of size 3×3 .

Modifies the dimension of feature map, each region is embedded in a lower dimensional vector of size (output size $r_b \times c_b \times 512$)



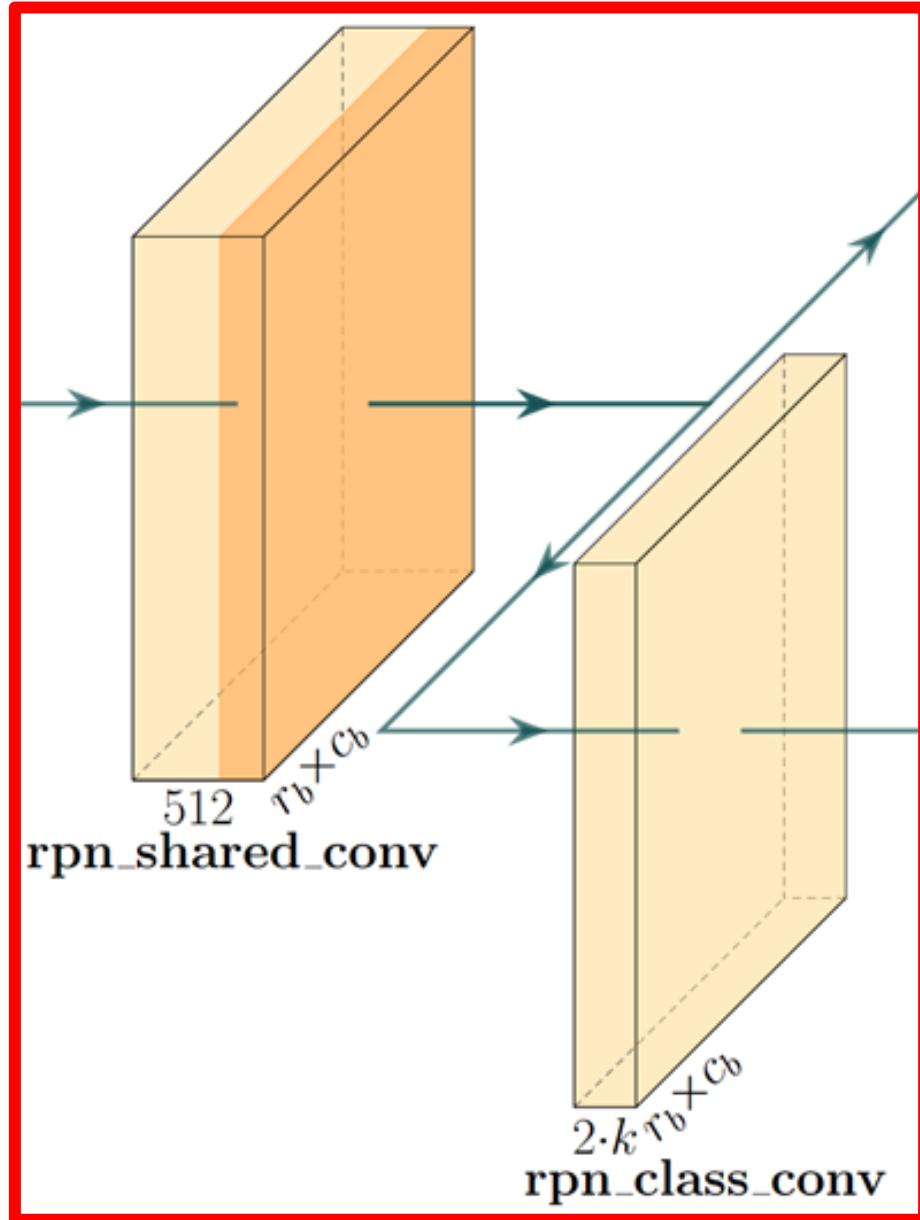
RPN Estimating k Anchors

- The **cls** (classification) network is trained to predict the **object probability**, i.e. the probability for an anchor to contain an object [*contains* / *does not contain*] → $2k$ probability estimates per each spatial location
- Made of a stack of 1×1 convolutional layers
- Each of these **k probability pairs** $[p, 1 - p]$ corresponds to a specific anchor (having a specific dimension) and expresses the probability for that anchor in that spatial location to contain *any object*

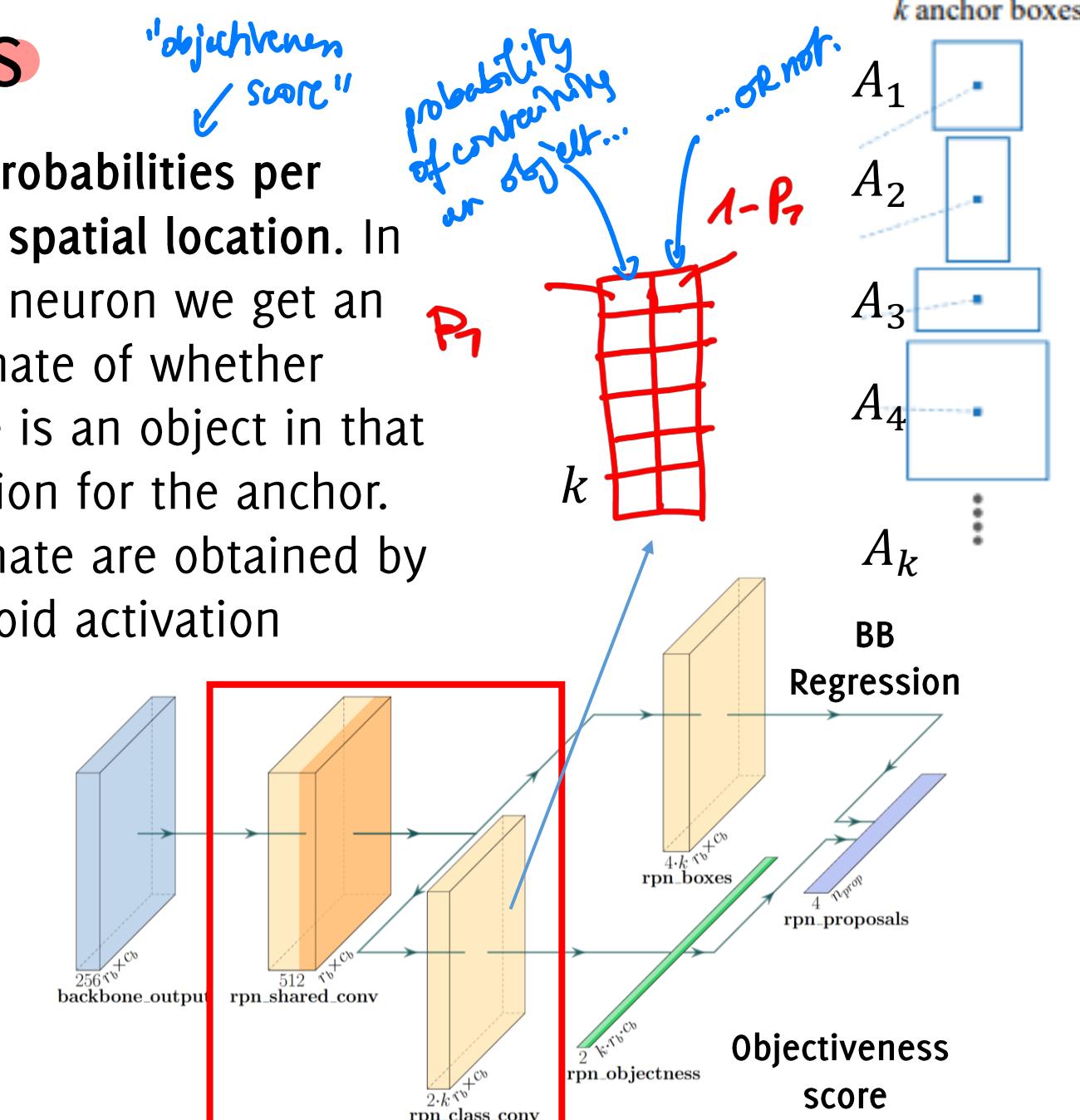


**Objectiveness
score**

RPN Estimating k Anchors

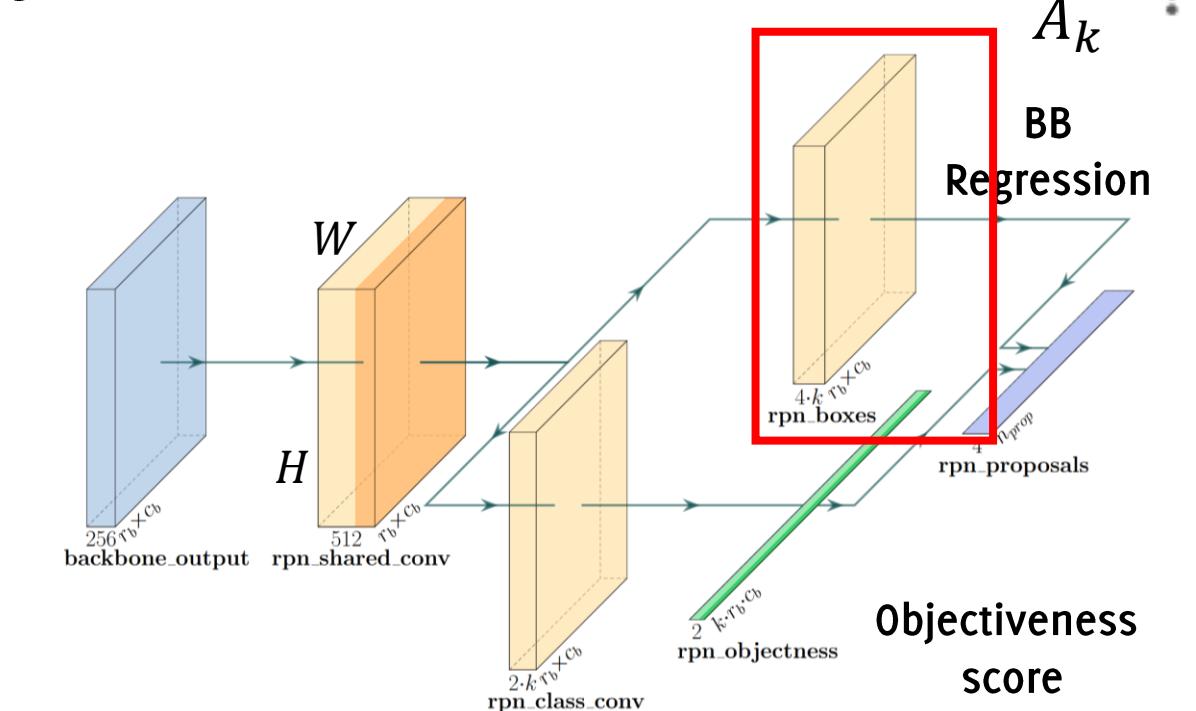
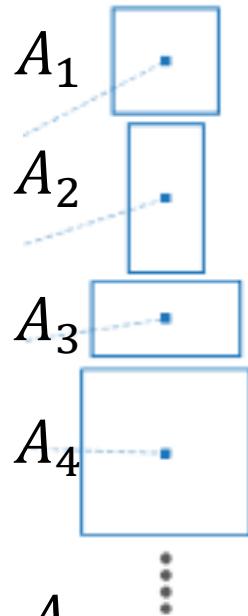


2 k probabilities per each spatial location. In each neuron we get an estimate of whether there is an object in that location for the anchor. Estimate are obtained by sigmoid activation

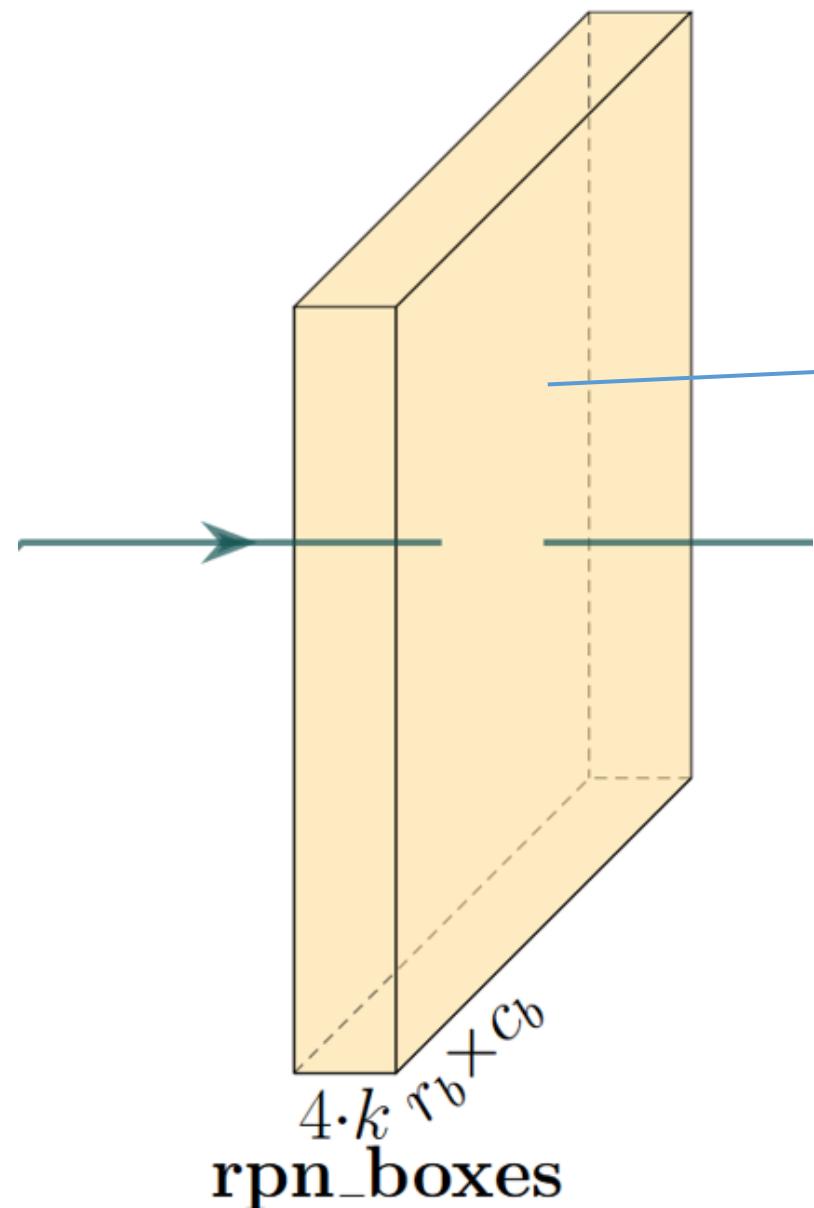


RPN Estimating k Anchors

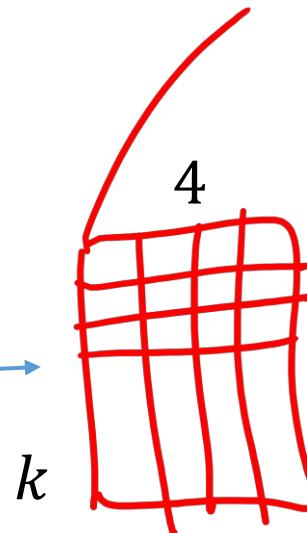
- The **reg** (regression) **network** is trained to adjust each of the k anchors to better match object ground truth blob
→ $4k$ estimates for the «delta» over the 4 bounding box coordinates (per each spatial location)
- Each of these k 4 –tuples **expresses the delta / refinements** for a specific anchor and are computed in a log measure to ease network convergence



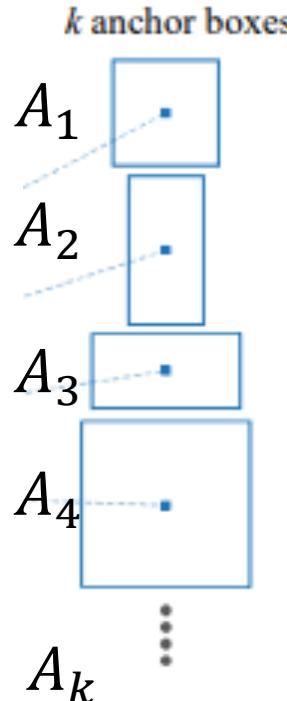
RPN Estimating k Anchors



BB
Regression



x, y, h, w



$4k$ numbers predicted per each spatial location. In each neuron we get an estimate of the 4 bounding box correction factors for each anchor.

RPN output: the proposals

The RPN returns k proposals for each location in the latent space

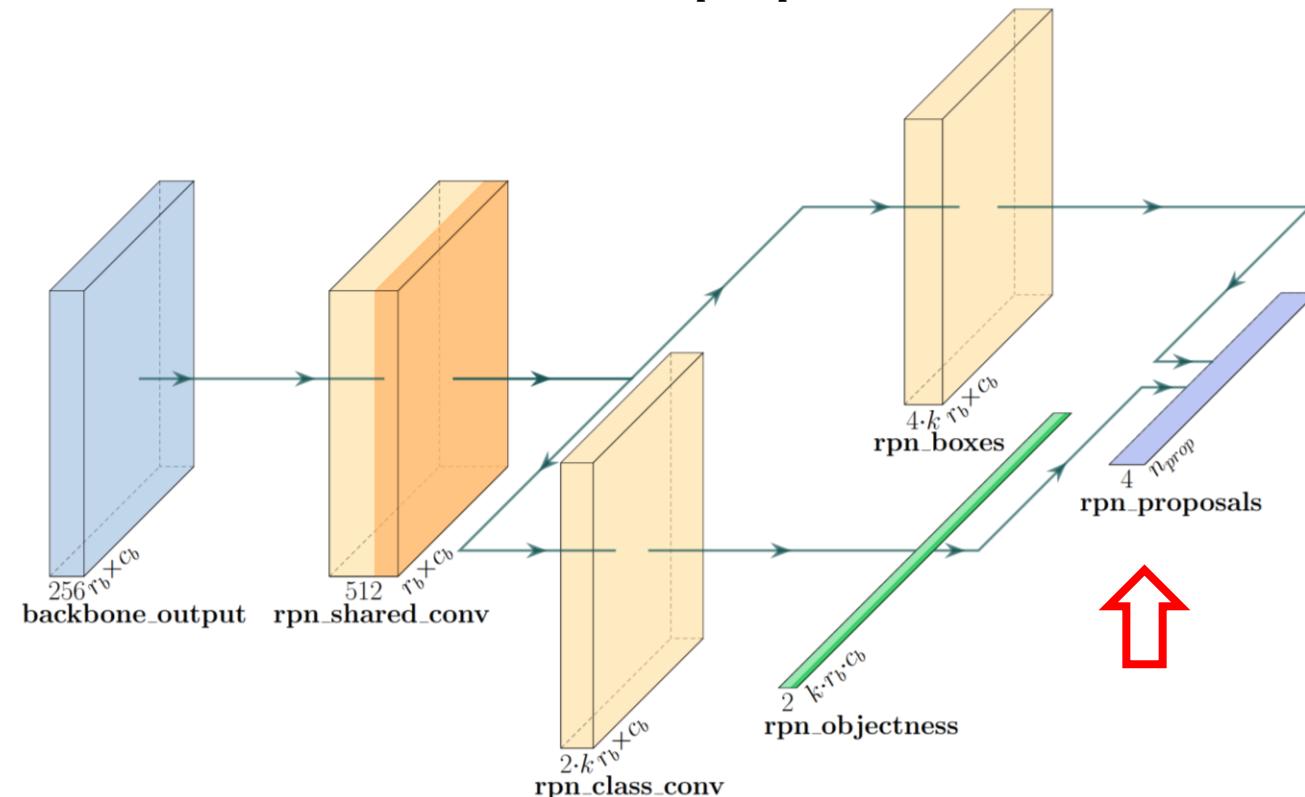
- $2k$ objectiveness scores
- $4k$ correction bounding boxes

Nonmaximum suppression: on the proposals to return only the first n_{prop} proposals

Proposals are selected analyzing:

- Non-maximum suppression based on the objectiveness score and the IoU considering corrected boxes
- Fixed threshold on the objectiveness score

Only the top n_{prop} proposals are propagated to the next layer

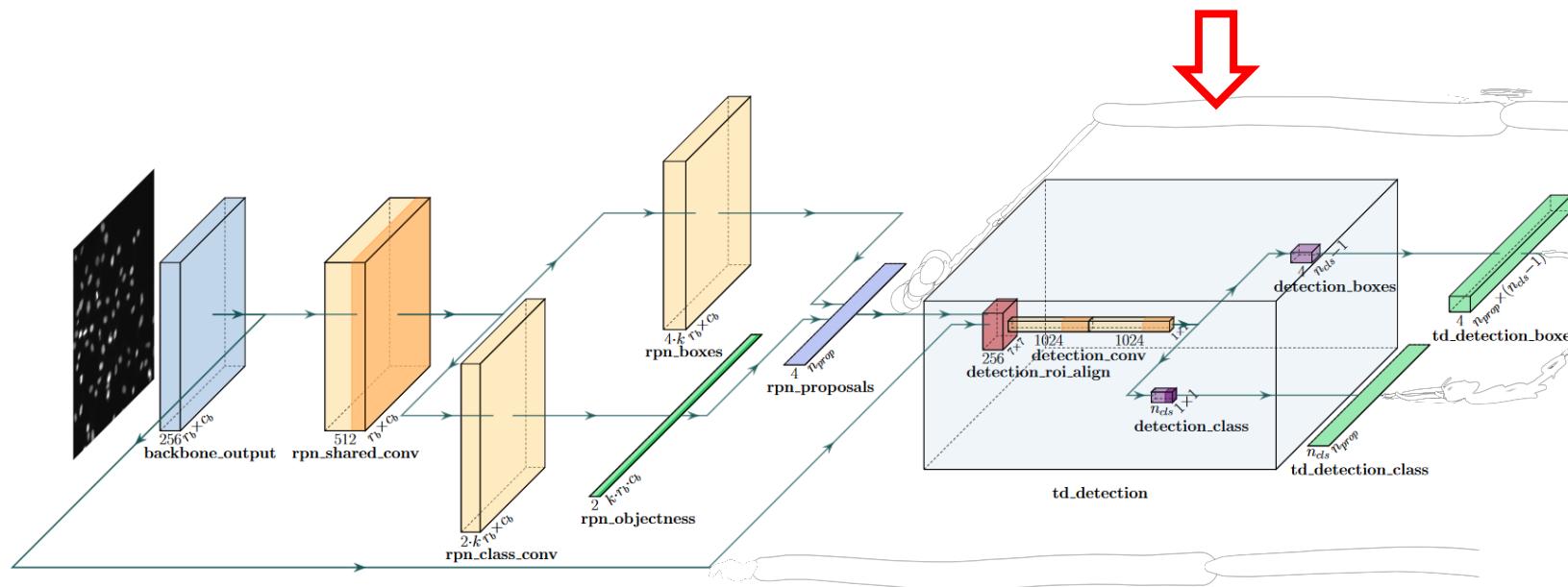


Faster R-CNN

Conv backbone Classification

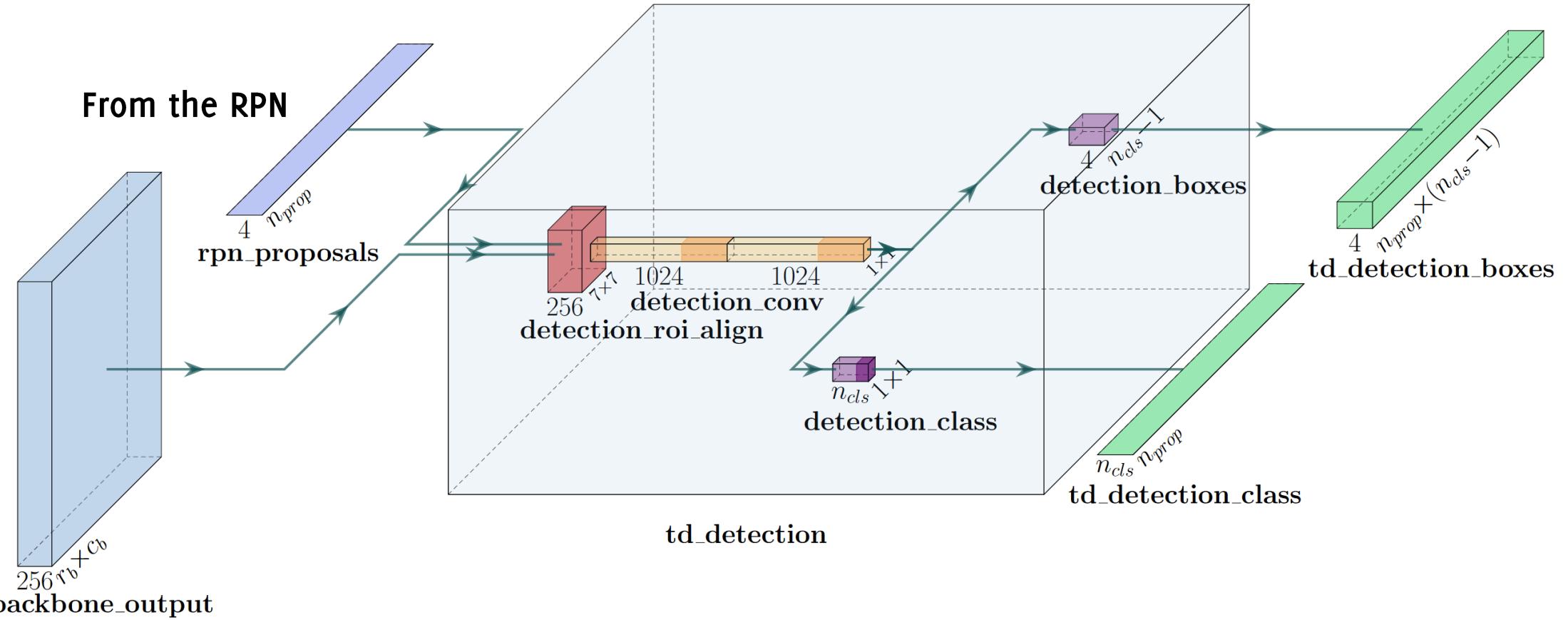
RPN

Detection head



The detection head

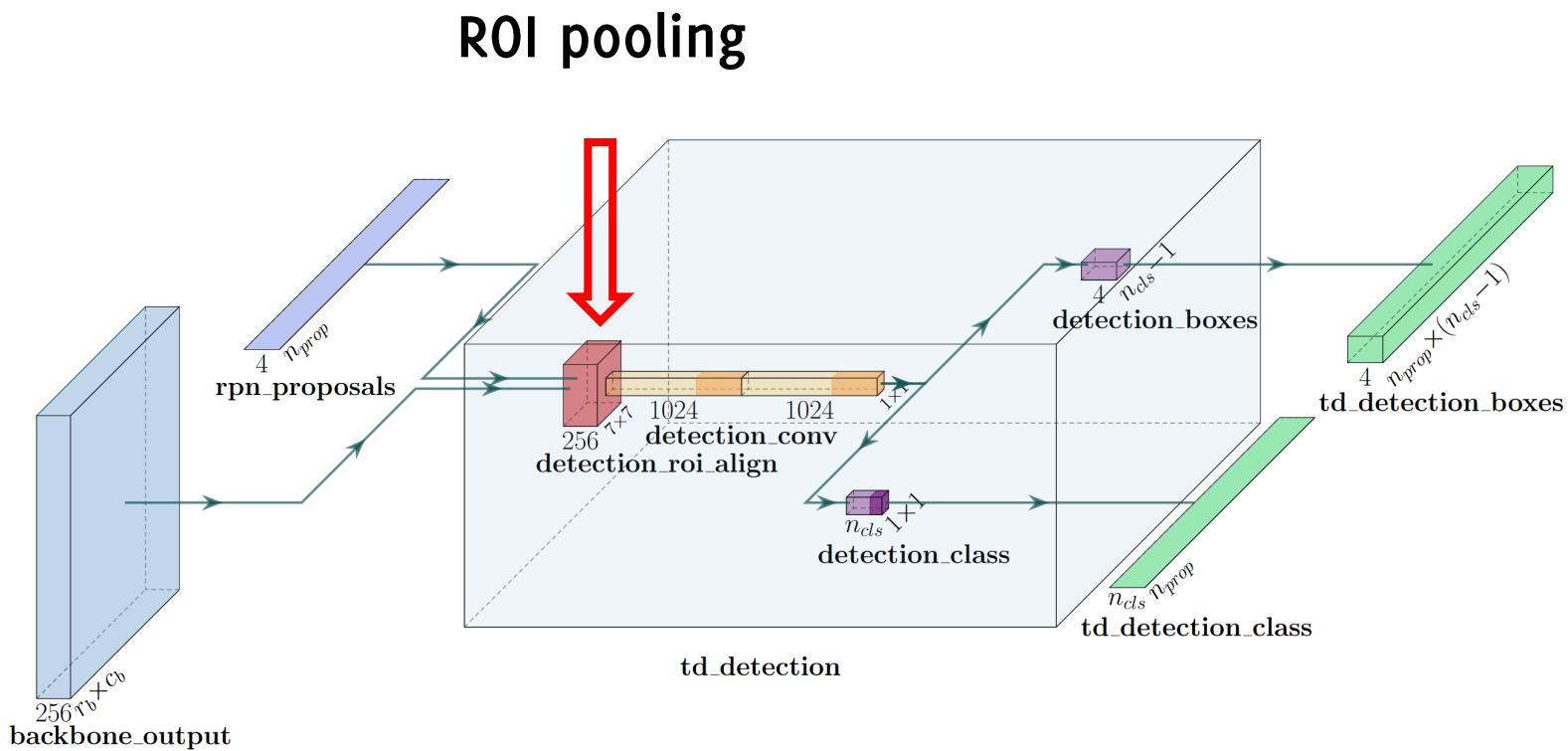
Processing repeated for
each of the n_{prop} proposal



From the CNN
backbone

Faster R-CNN

RPN returns n_{prop} region proposals, thus replaces the region proposal algorithms (selective search) that are fed to the ROI pooling and then classified by the standard Fast-RCNN architecture.



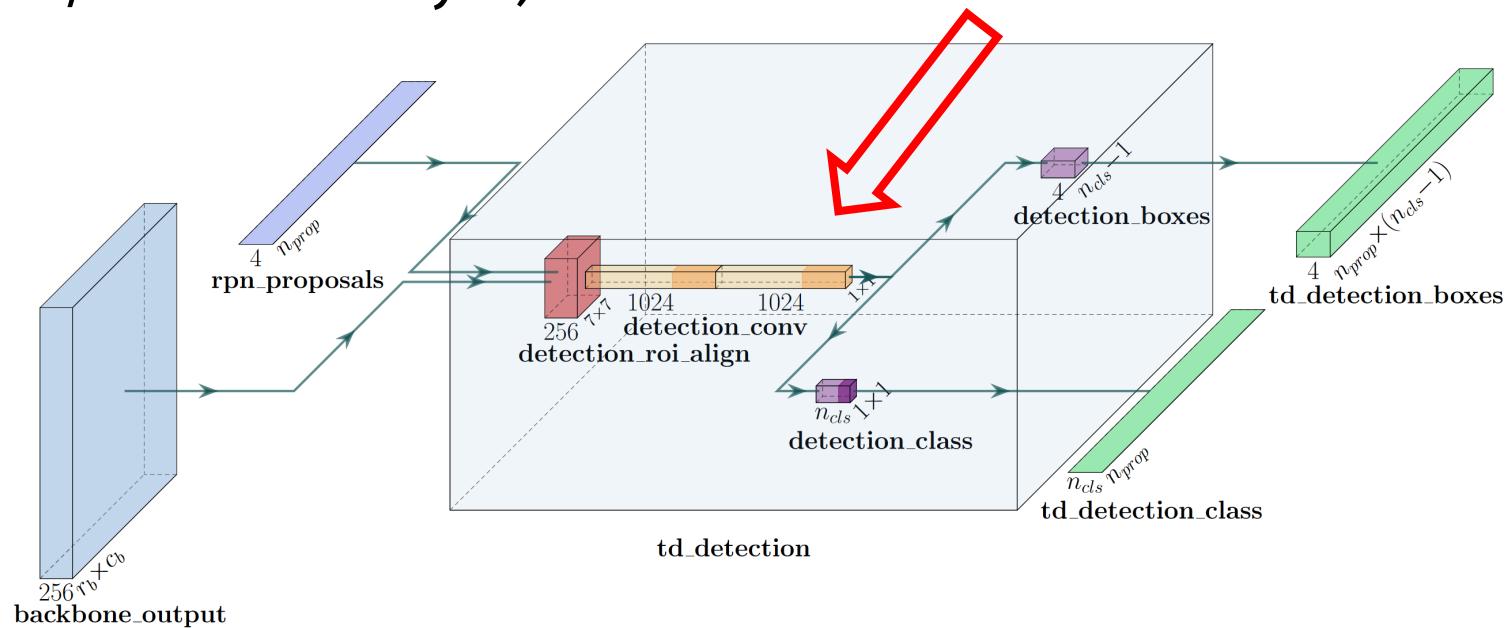
ROI Align and Pooling

This is meant to extract regions in the latent space referring to proposals

The features are then subsampled to a fixed size.

For each proposal there is a latent vector that is further processed by the detection head.

After some MLP/1x1 conv layer, it returns **ROI feature vector**



RPN and the Anchor Shapes

- If you want to train the network to predict the objectiveness score and the offsets for different anchors, there is no need to design different RPN, but just to **define different targets (associated to different anchor shape and size) when training the RPN**.
- You can in principle adopt **anchors having non-rectangular shapes** (e.g. ellipses, or simply rotated BBs)....
- However it can be difficult to compute IoU, which is instead extremely easy and fast in case of BBs... there are shapes that cannot be easily intersected (no closed form expression) to define the loss

Detection Head

DETECTION HEAD

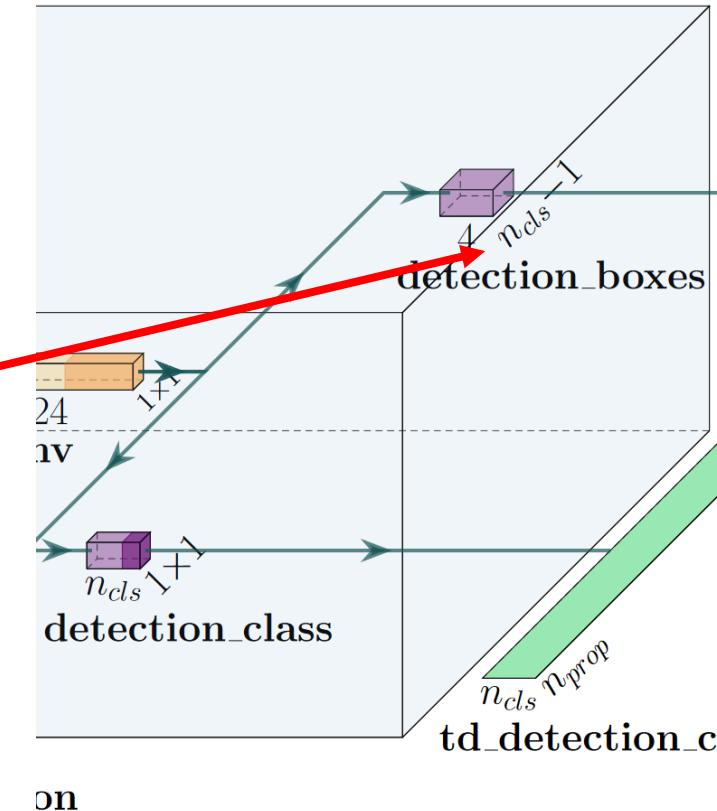
The “detection head has two heads” and takes as input a single ROI feature vector.

It returns two outputs:

- L classification scores (class probabilities)
- $4 \times (L - 1)$ bounding box corrections

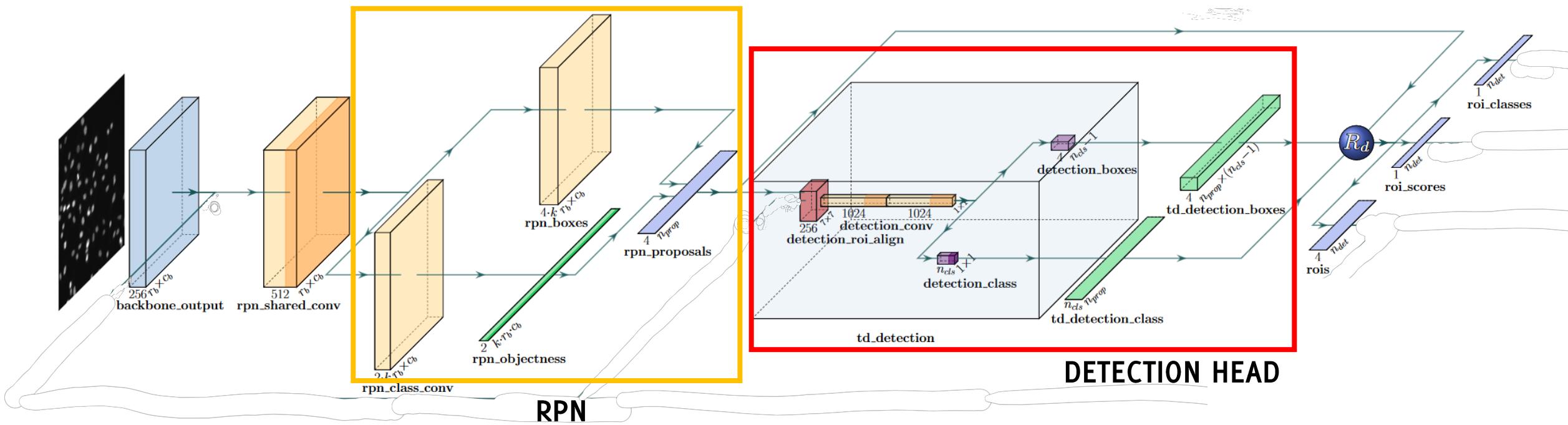
The classes include the “**background**” class for which we do not compute the bounding box corrections

The Faster output corresponds to the predicted objects with large confidence (fewer than the n proposals)



Faster R-CNN Training

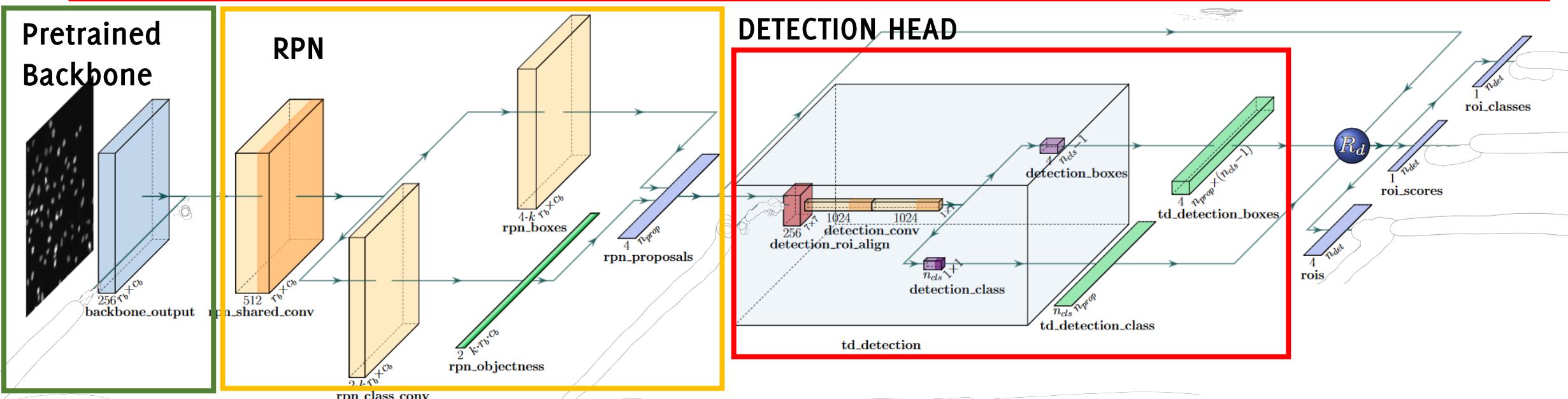
- Training now involves 4 losses, in a weighted sum:
 - RPN classify object/non object
 - RPN regression coordinates
 - Final classification score
 - Final BB coordinates
- During training, object/non object ground truth is defined by measuring the overlap with annotated BB
- The loss include a term of final BB coordinates computed form detection head



Faster R-CNN Training

Training procedure

1. Train RPN keeping backbone network frozen and training only RPN layers. This ignores object classes but just bounding box locations (Multi-task loss cls + reg)
2. Train Fast-RCNN using proposals from RPN trained before. Fine tune the whole Fast-RCNN including the backbone
3. Fine tune the RPN in cascade of the new backbone
4. Freeze backbone and RPN and fine tune only the last layers of the Faster R-CNN



Faster R-CNN

At test time,

- Take the top ~ 300 anchors according to their object scores
- Consider the refined bounding box location of these 300 anchors
- These are the ROI to be fed to a Fast R-CNN
- Classify each ROI and provide the non-background ones as output

Faster R-CNN provides as output to each image a set of BB with their classifier posterior

The network becomes much faster (0.2s test time per image)

Faster R-CNN



It's still a two stage detector

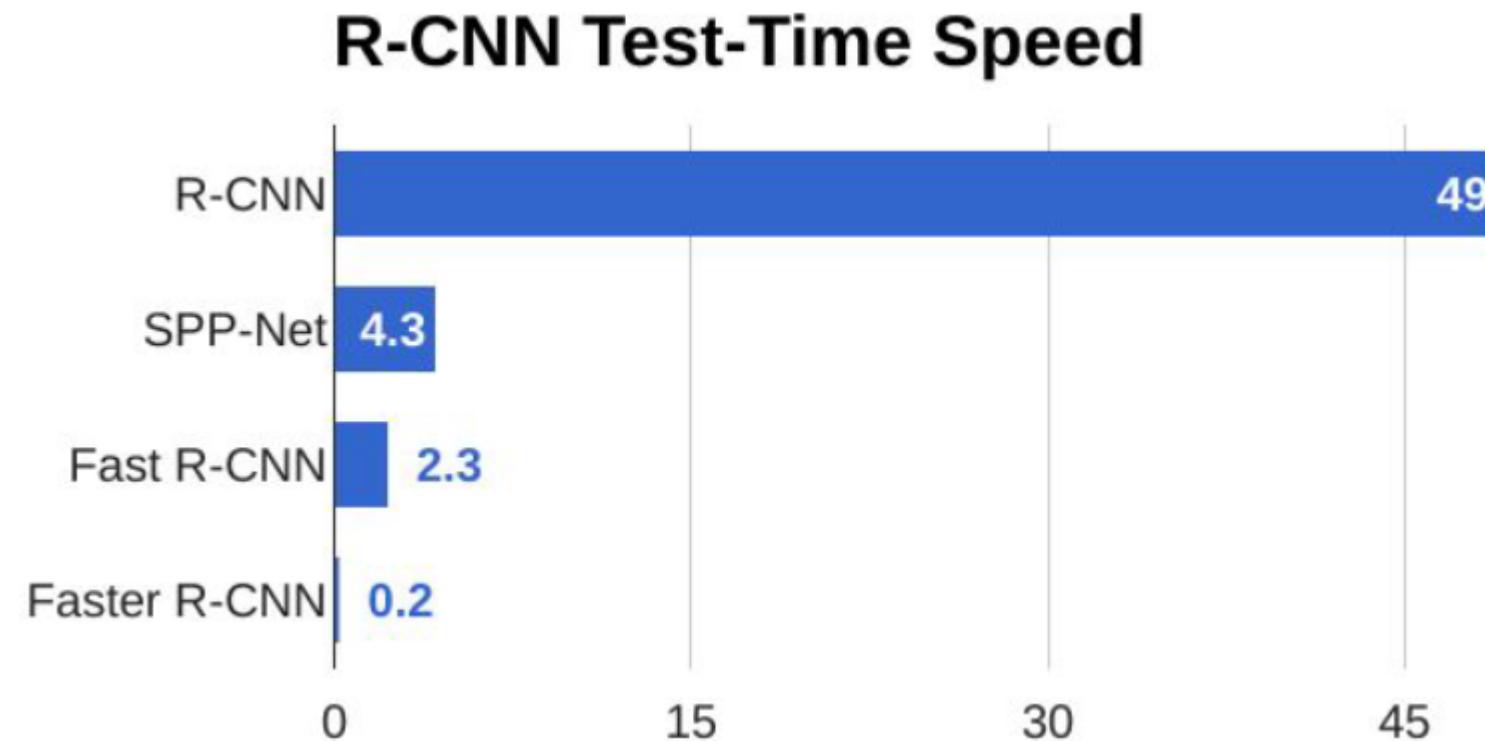
First stage: *SELECTING PROPOSALS*.

- run a backbone network (e.g. VGG16) to extract features
- run the Region Proposal Network to estimate ~ 300 ROI

Second stage (the same as in Fast R-CNN): *CLASSIFYING + BB*.

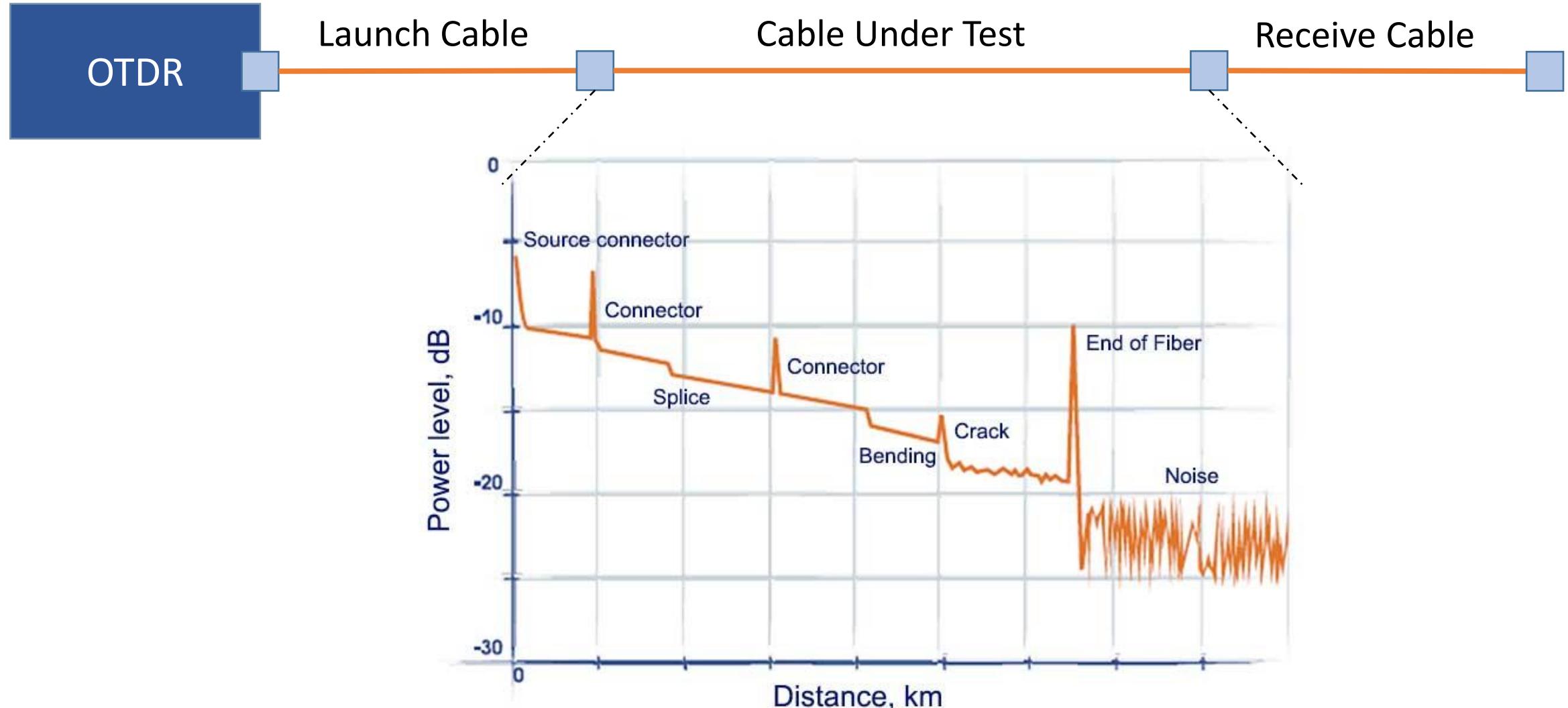
- Crop Features through ROI pooling (with alignment)
- Predict object class using FC + softmax
- Predict bounding box offset to improve localization using FC + softmax

Faster R-CNN

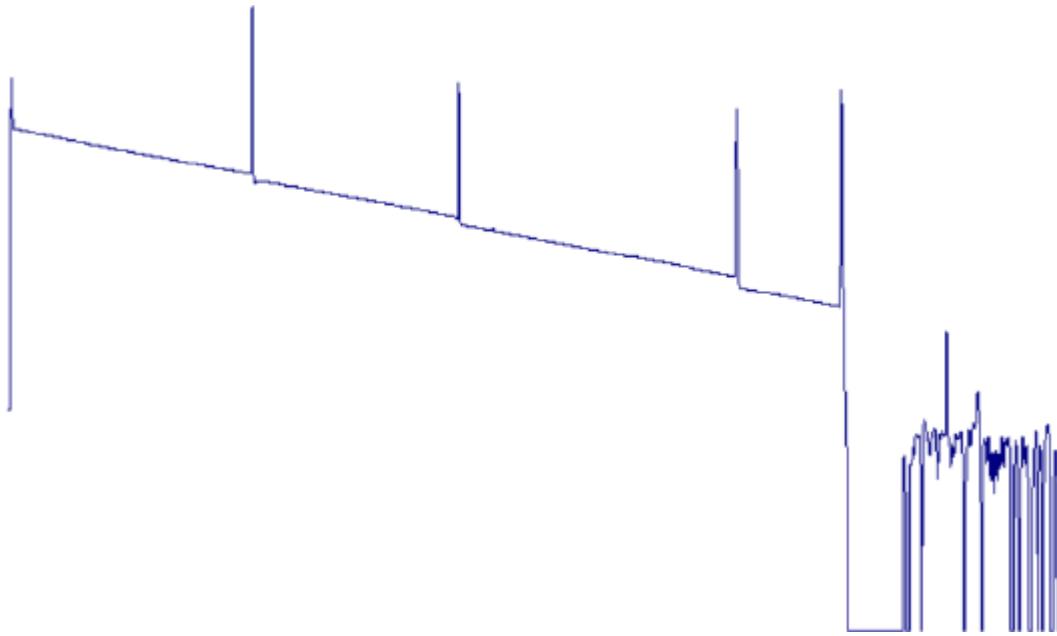


An industrial application

Event Detection in Optical Signals



Event Detection in Optical Signals



OTDR trace: events indicate problems over a span of optical fiber

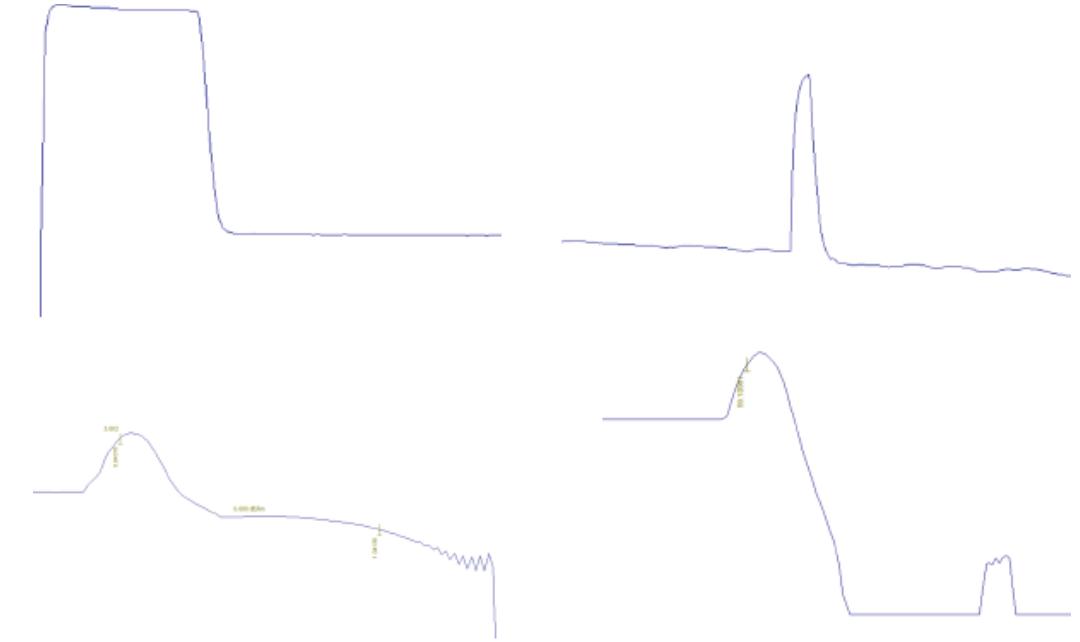
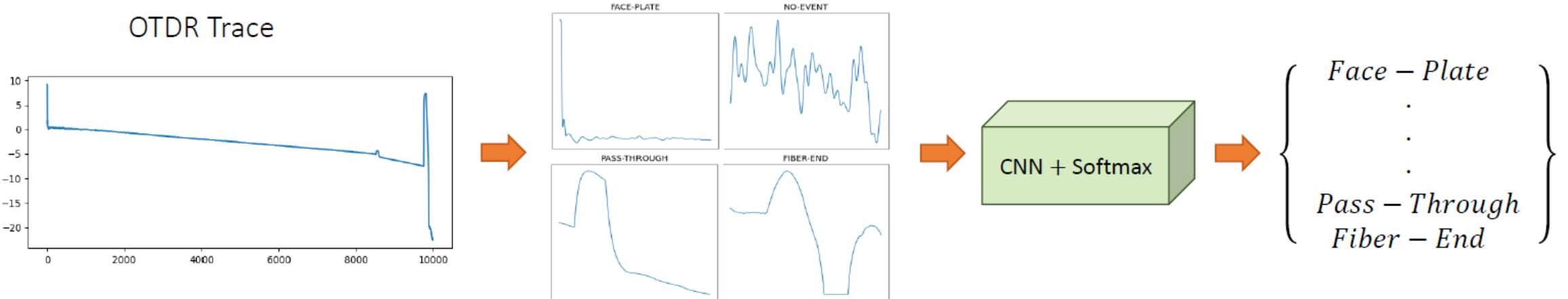


Fig. 1 (b): OTDR events examples: face-plate, pass-through, fiber-cut, fiber-end

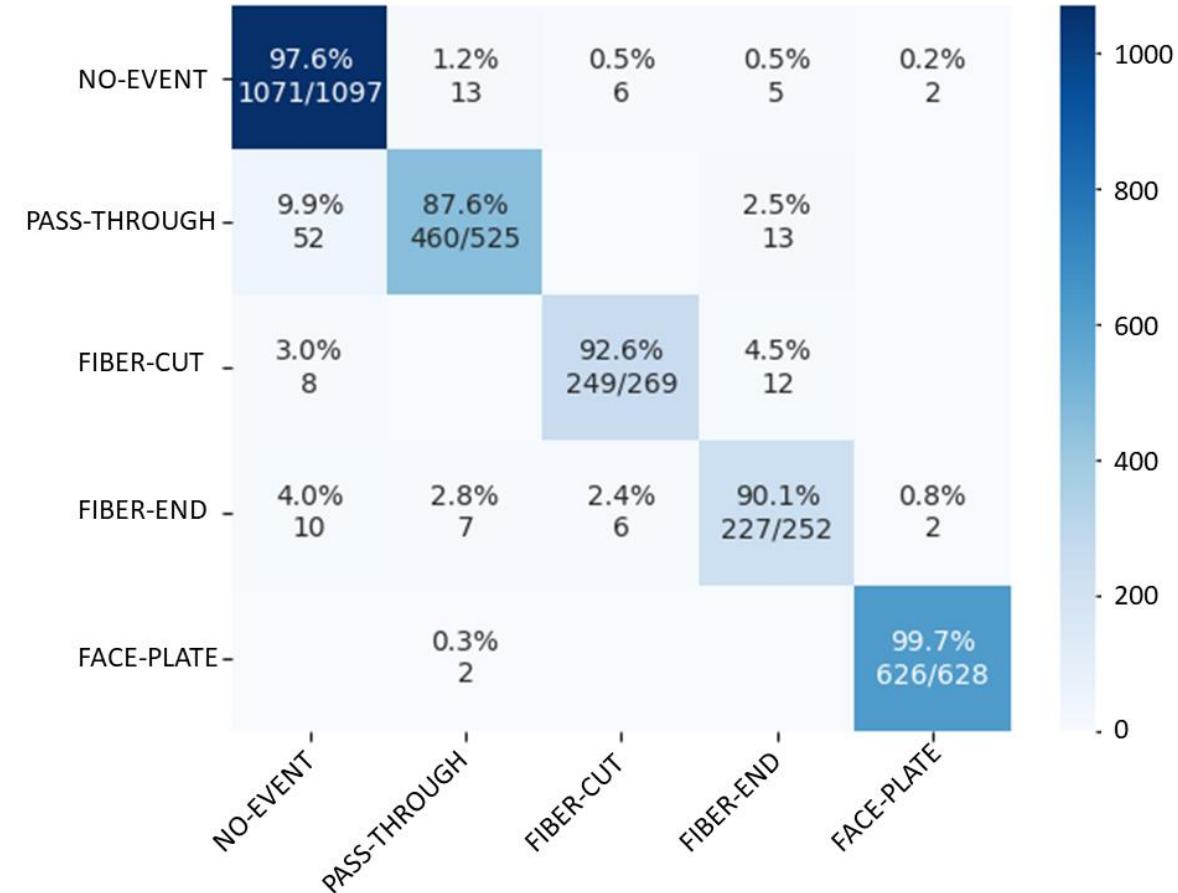
Event Detection in Optical Signals

Train a classification network over fixed windows (300 samples)

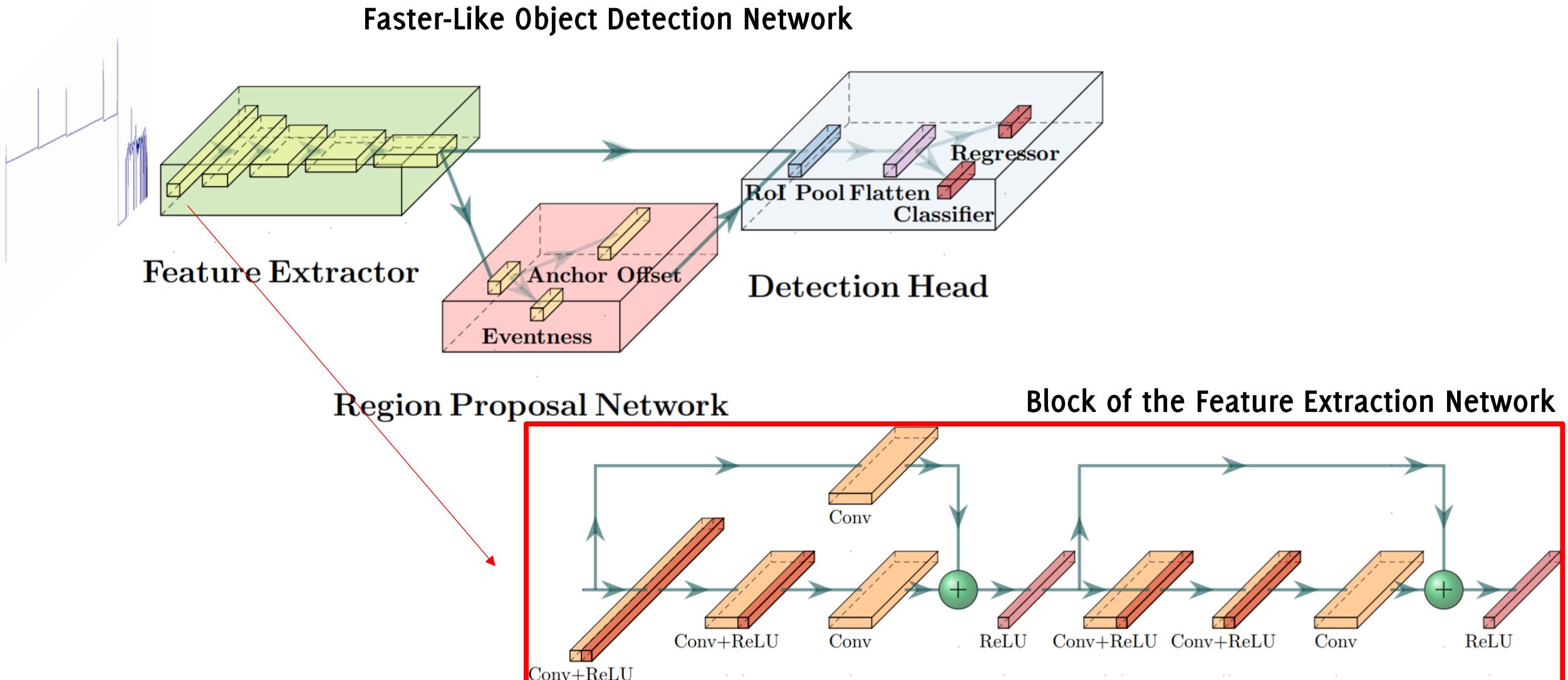


Event Detection in Optical Signals

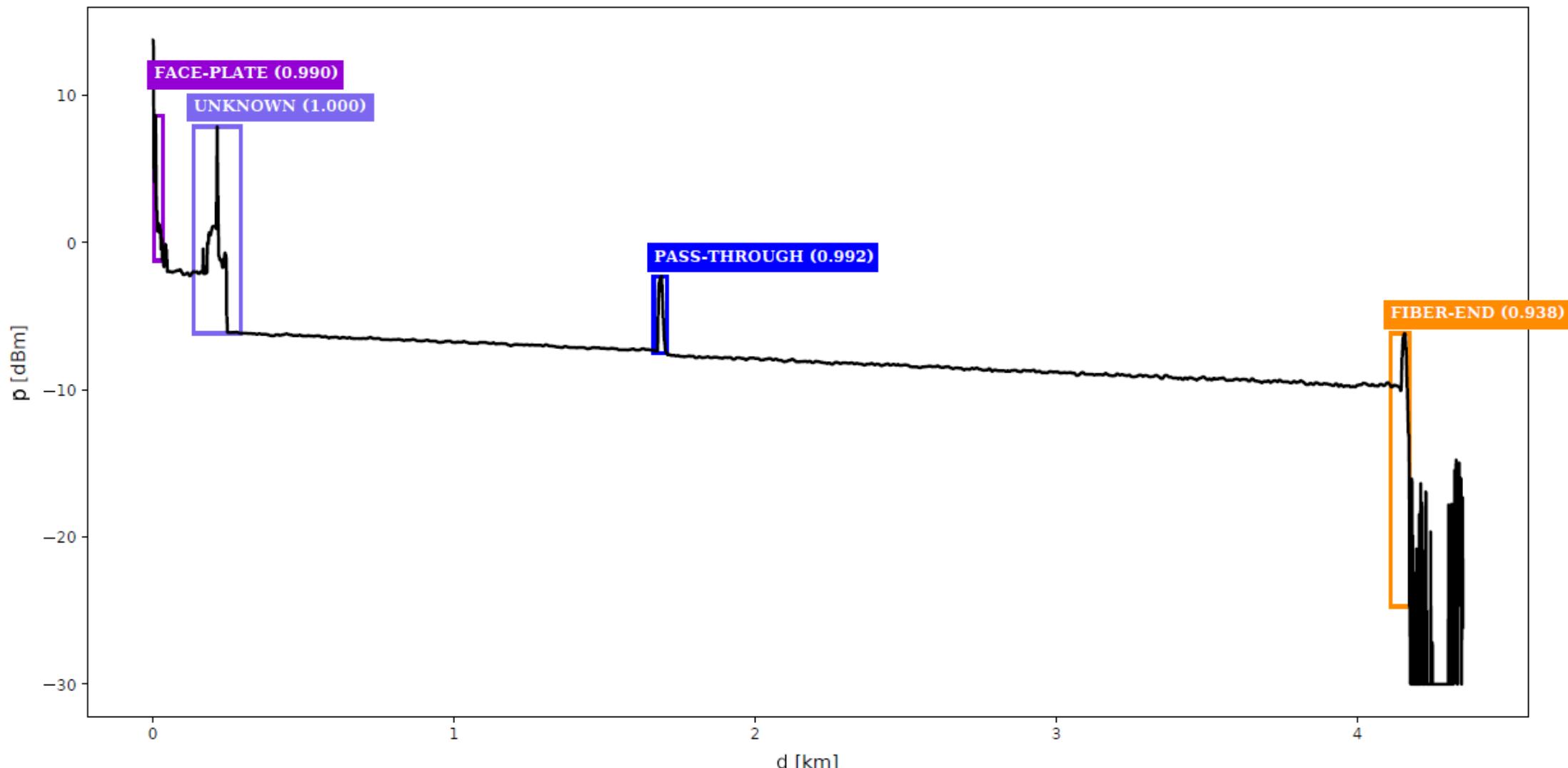
- Requires an additional step to split OTDR traces into windows
- Limited to single-scale events resolution
- Works under the assumption that each window includes at most one event
- Does not consider events position
- Does not share computations



Event Detection in Optical Signals



Able to detect both known and unknown events





This CVPR paper is the Open Access version, provided by the Computer Vision Foundation.
Except for this watermark, it is identical to the version available on IEEE Xplore.

You Only Look Once: Unified, Real-Time Object Detection

Joseph Redmon*, Santosh Divvala*[†], Ross Girshick[¶], Ali Farhadi*[†]

University of Washington*, Allen Institute for AI[†], Facebook AI Research[¶]

<http://pjreddie.com/yolo/>

YOLO/SSD

R-CNN methods are based on region proposals

There are also region-free methods, like:

YOLO: You Only Look Once

SSD: Single Shot Detectors

The rationale

Detection networks are indeed a pipeline of multiple steps.

In particular, **region-based methods make it necessary to have two steps** during inference

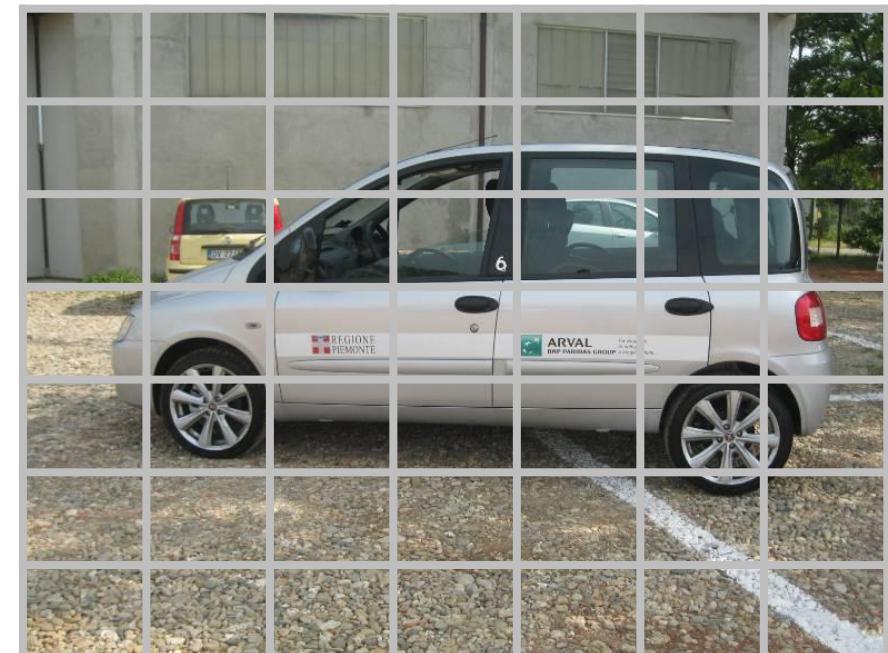
This can be slow to run and hard to optimize, because each individual component must be trained separately.

In Yolo "*we reframe the object detection as a single regression problem, straight from image pixels to bounding box coordinates and class probabilities*"

And solve these regression problems **all at once**, with a large CNN

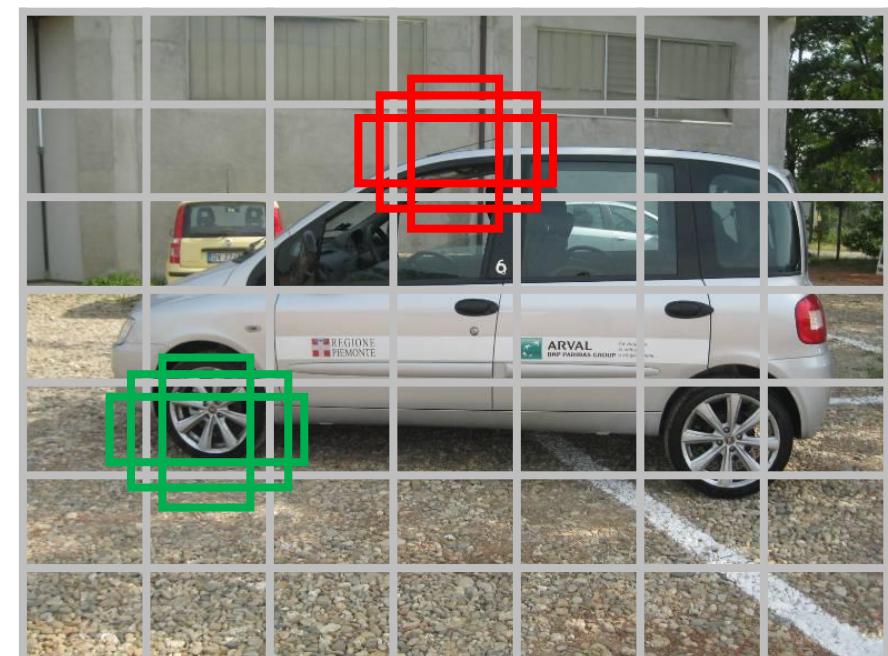
YOLO

1. divide the image in a coarse grid (e.g. 7×7)



YOLO

1. divide the image in a coarse grid (e.g. 7×7)
2. each grid cell contains *B anchors* (*base bounding box*) associated



YOLO

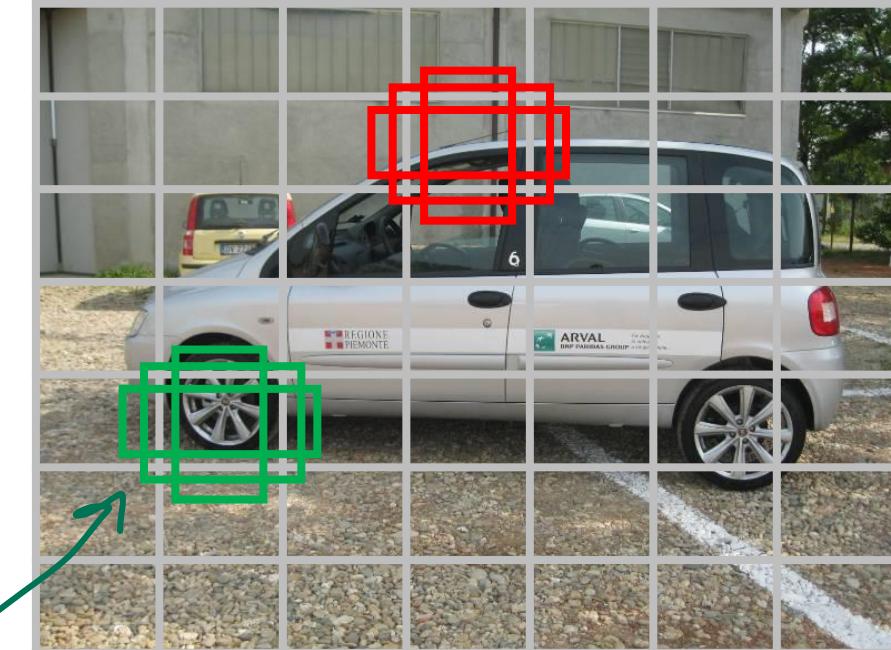
3. For each cell and anchor we predict:

- The **offset of the base bounding box**, to better match the object:
 $(dx, dy, dh, dw, objectness_score)$
- The **classification score** of the base-bounding box over the C considered categories (including background)

So, the output of the network has dimension

$$7 \times 7 \times B \times (5 + C)$$

" B basic bounding boxes"



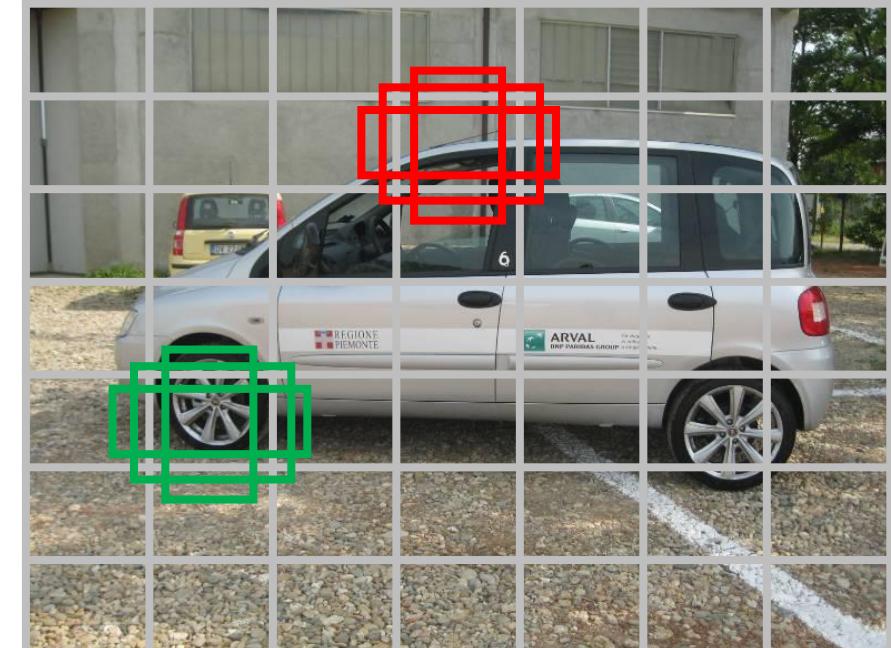
YOLO

The whole prediction is performed in a single forward pass over the image, by a single convolutional network

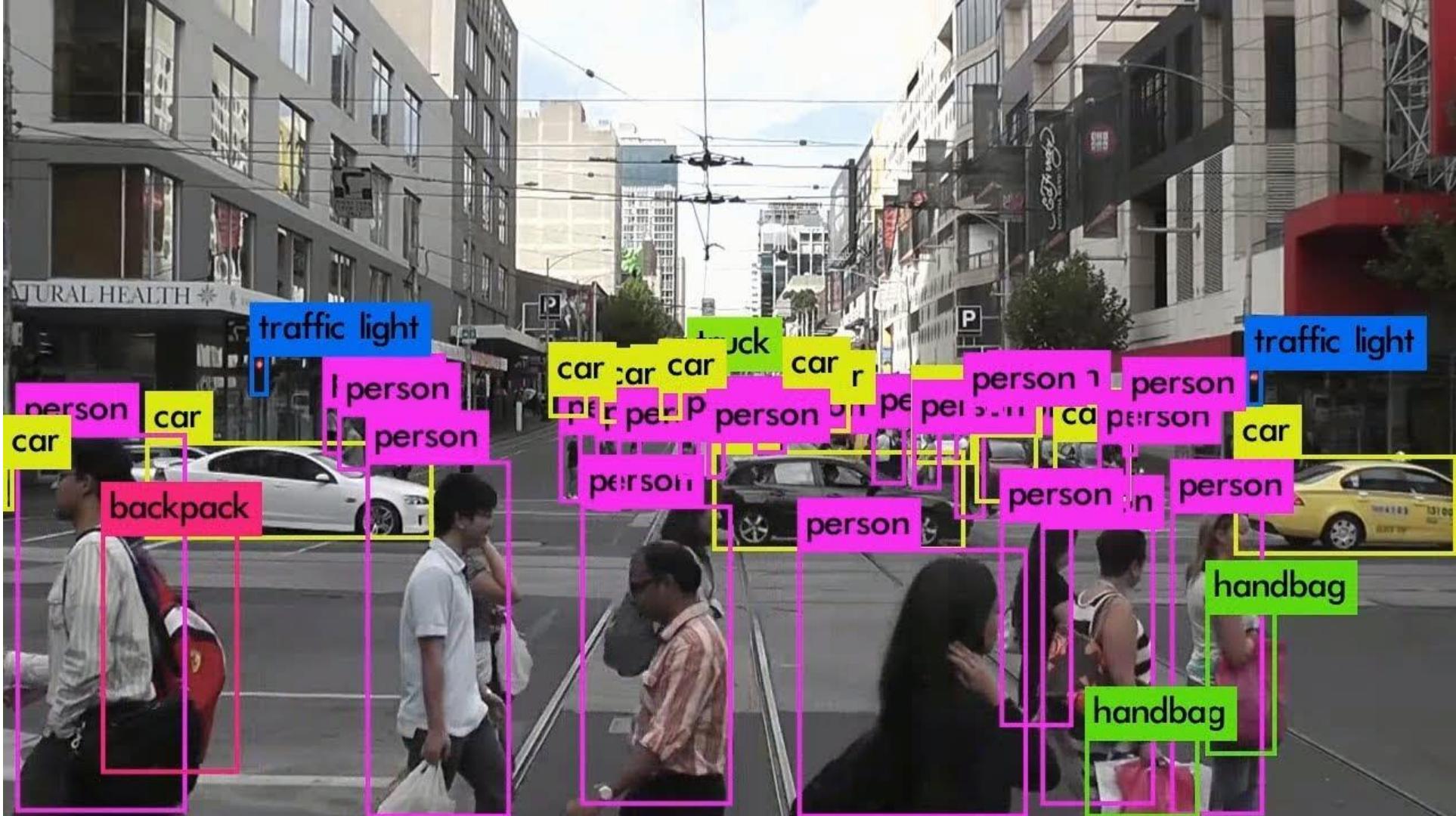
Training this network is sort of tricky to assess the loss (matched / not matched)

YOLO/SSD shares a similar ground of the RPN used in Faster R-CCN

Typically, networks based on region-proposals are more accurate, single shot detectors are faster but less accurate



Object Detection



Do it yourself!

https://colab.research.google.com/drive/1chjXM6ckE1s-xoZZjGG_Q0dmlCeBPzPR

Perform inference using a pre-trained
object detection network

Instance Segmentation

Object Detection vs Instance Segmentation

**Object
Detection**



**Instance
Segmentation**



Instance Segmentation, the problem

Assign to an input image I :

- multiple labels $\{l_i\}$ from a fixed set of categories $\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$, each corresponding to an instance of that object
- the coordinates $\{(x, y, h, w)_i\}$ of the bounding box enclosing each object
- the set of pixels S in each bounding box corresponding to that label

$$I \rightarrow \{(x, y, h, w, l, S)_1, \dots, (x, y, h, w, l, S)_N\}$$



This ICCV paper is the Open Access version, provided by the Computer Vision Foundation.
Except for this watermark, it is identical to the version available on IEEE Xplore.

Mask R-CNN

Kaiming He Georgia Gkioxari Piotr Dollár Ross Girshick
Facebook AI Research (FAIR)

Instance Segmentation

It combines the challenges of:

- **Object detection** (multiple instances present in the image)
- **Semantic segmentation** (associate a label to each pixel) separating each object instance

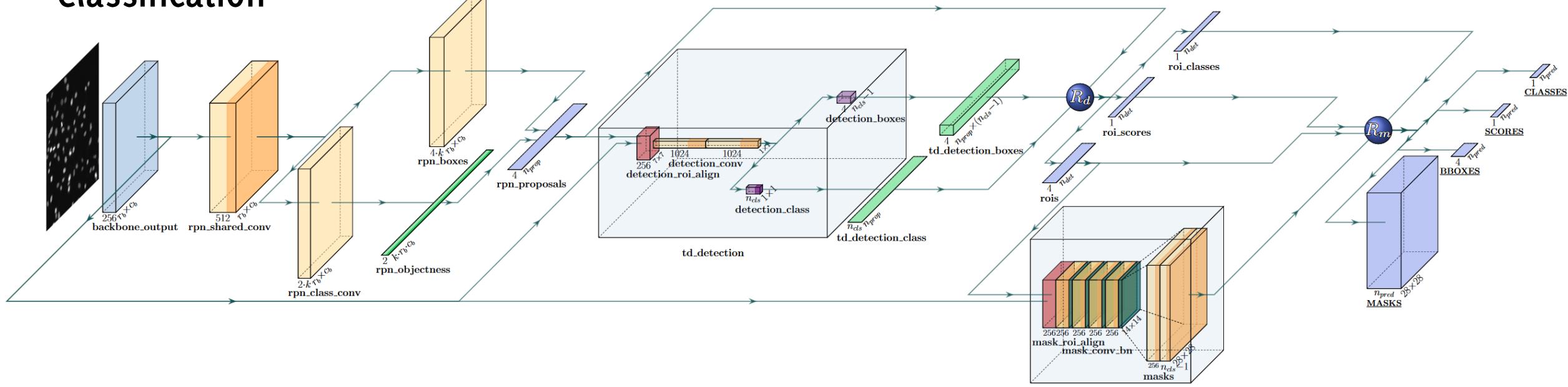
Color the exact pixels
→ mask.

Mask R-CNN

As in Faster R-CNN, each ROI is classified and the bounding boxes are estimated

Semantic segmentation inside each ROI

Conv backbone Classification



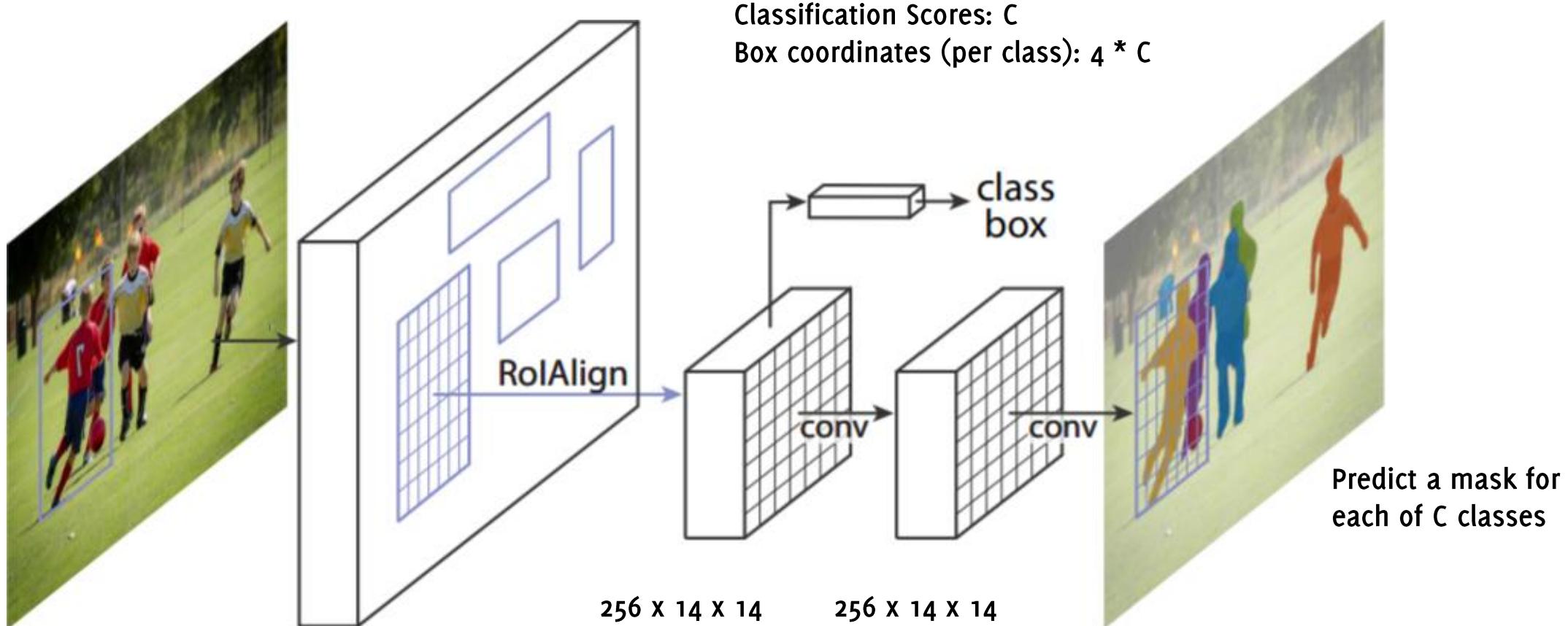
Predict a mask for each of C classes

Mask head

Mask R-CNN

As in Faster R-CNN, each ROI is classified and the bounding boxes are estimated

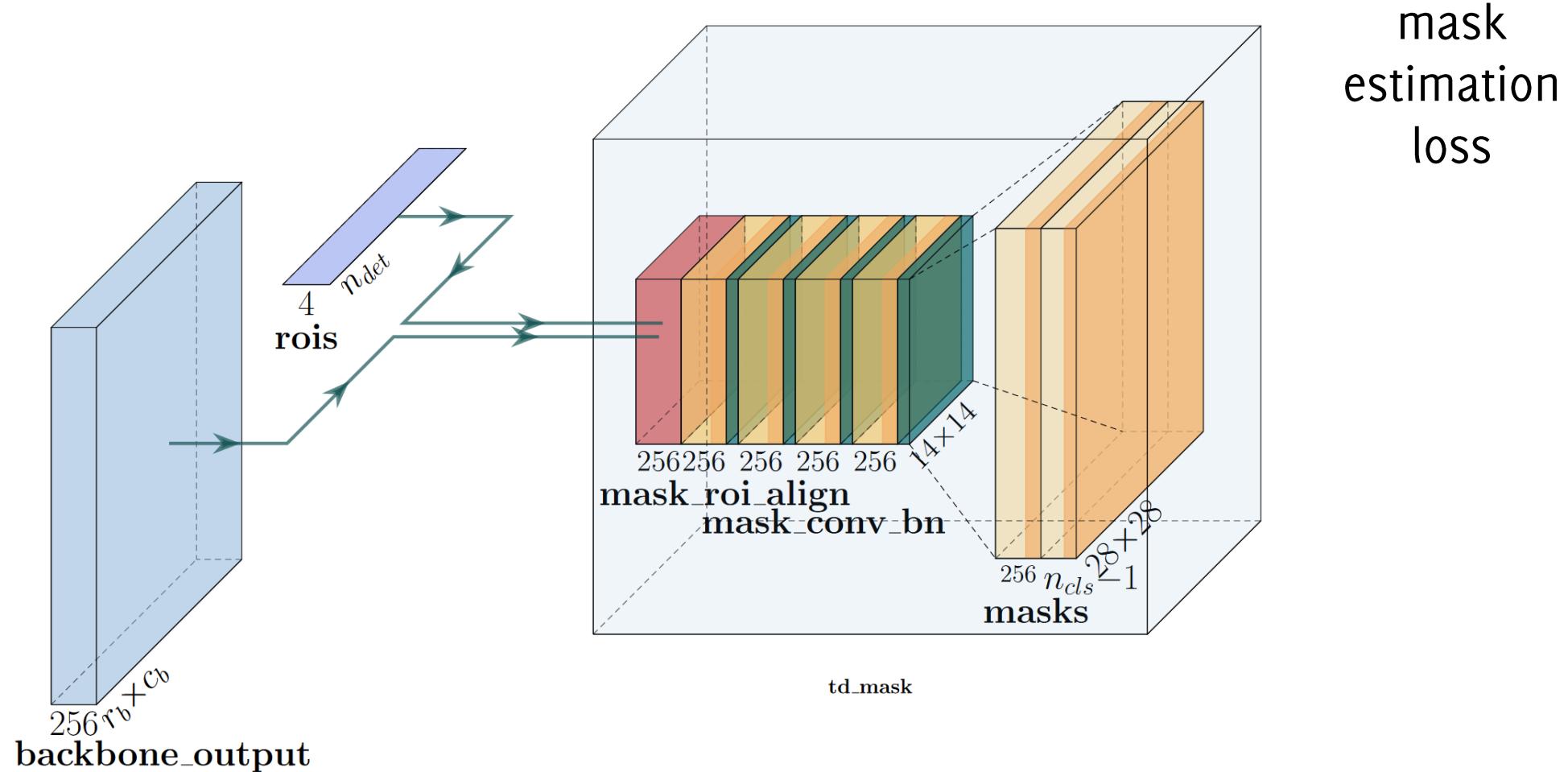
Semantic segmentation inside each ROI



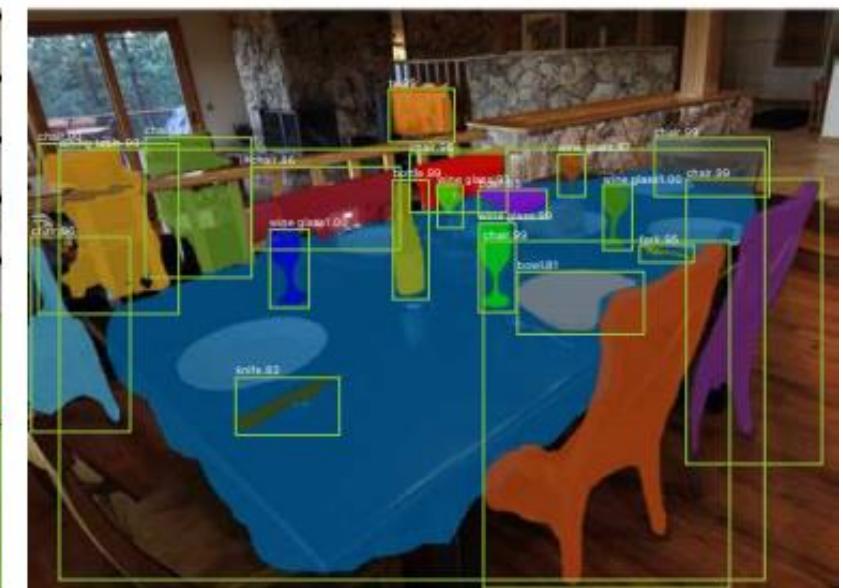
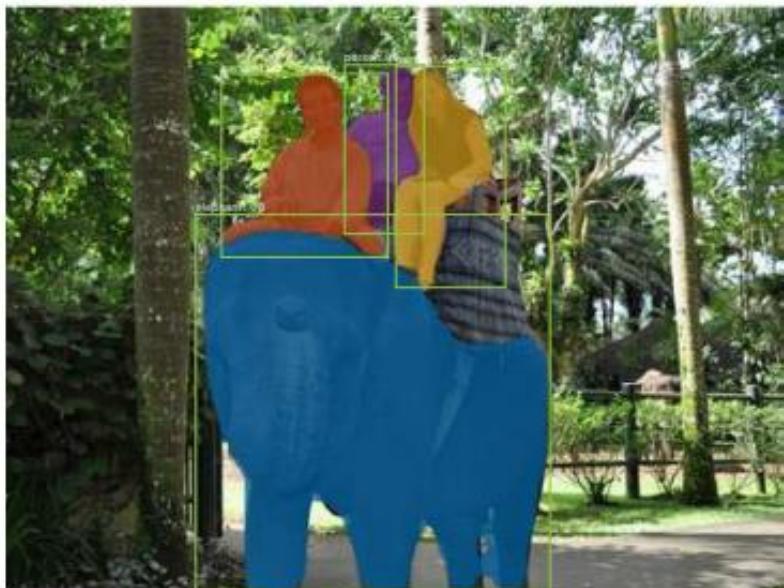
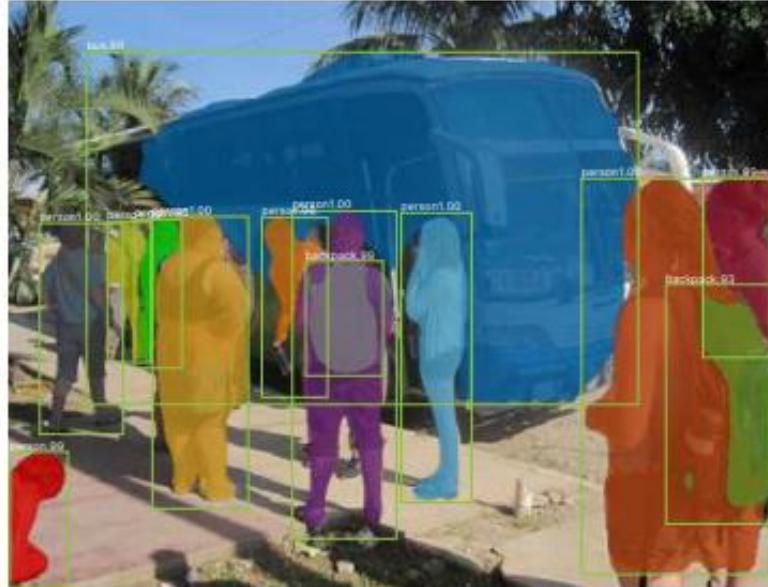
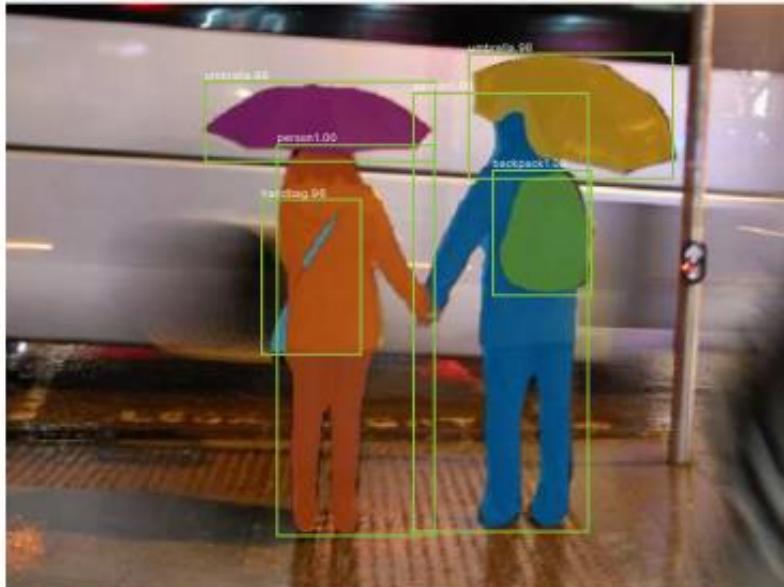
Mask R-CNN

Mask is estimated for each ROI and each class

Extends Faster R-CNN by adding a branch for predicting an object mask in parallel with the existing branch for bounding box recognition



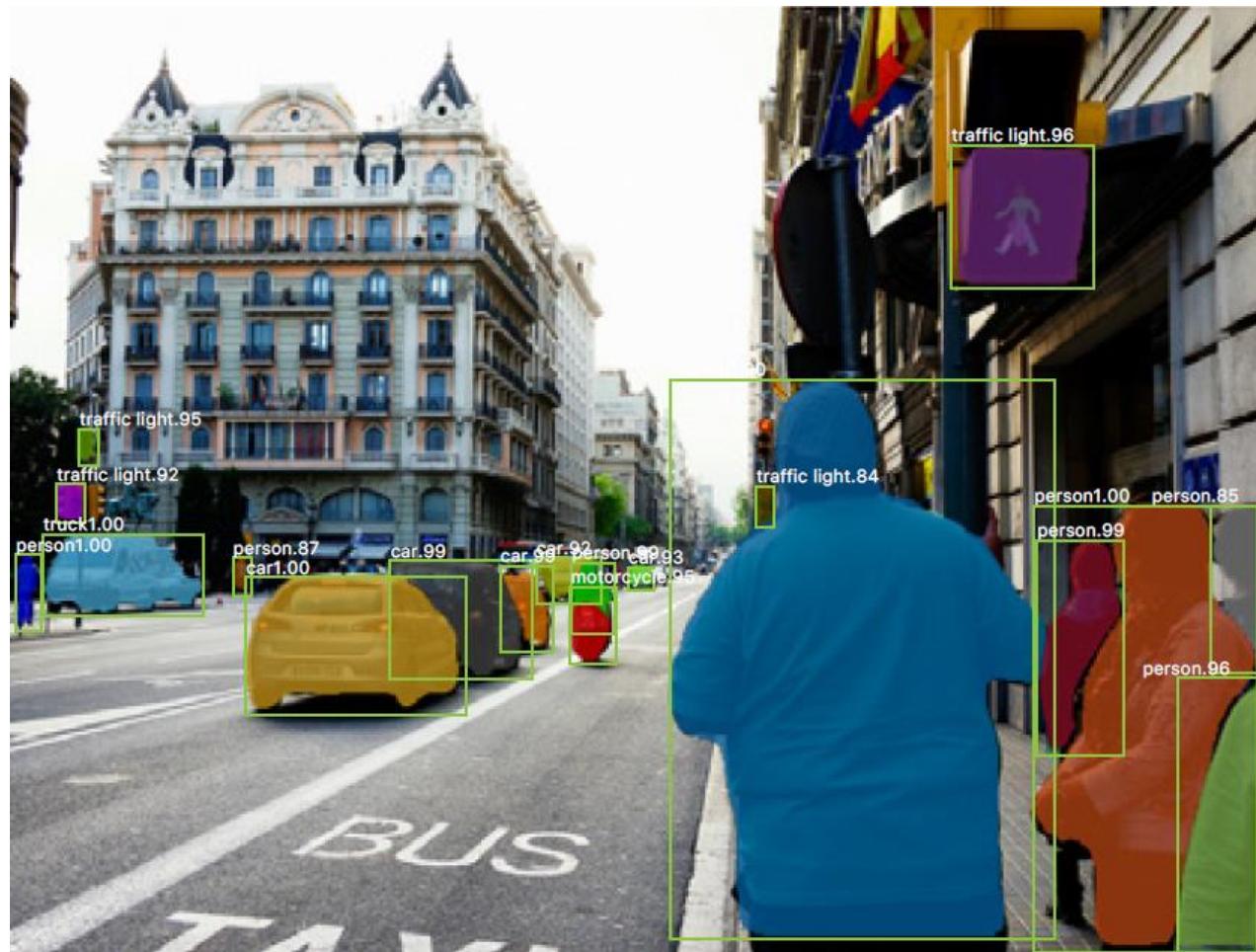
Mask R-CNN (end-to-end training)



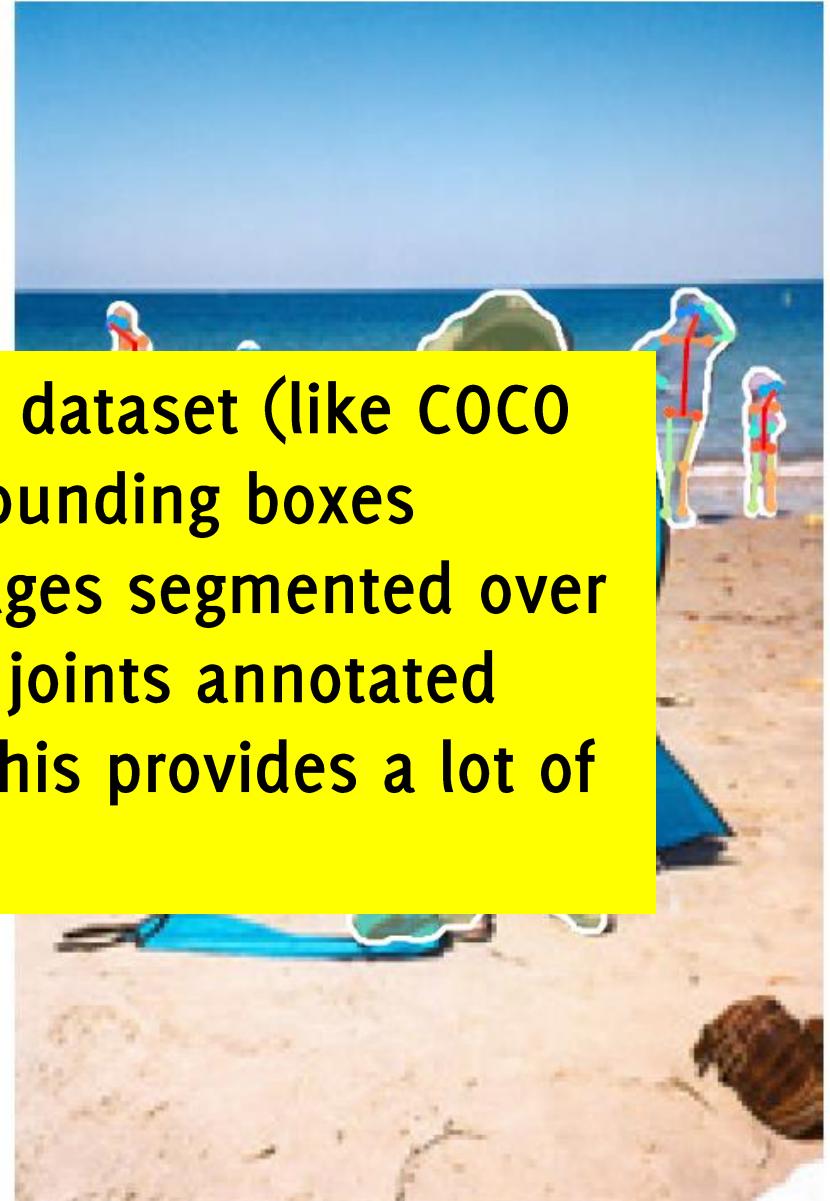
Mask R-CNN (including Pose estimation)



Instance Segmentation / Human Pose Estimation



Instance Segmentation / Human Pose Estimation



This can be also trained from segmentation dataset (like COCO dataset), from where you can infer bounding boxes Microsoft COCO dataset contains 200.000 images segmented over 80 categories. Persons are provided with joints annotated Since there are many instances per image, this provides a lot of training data

Feature Pyramid Network

Slide credits: Daniele Casciani, Luca Alessandrini

Goal

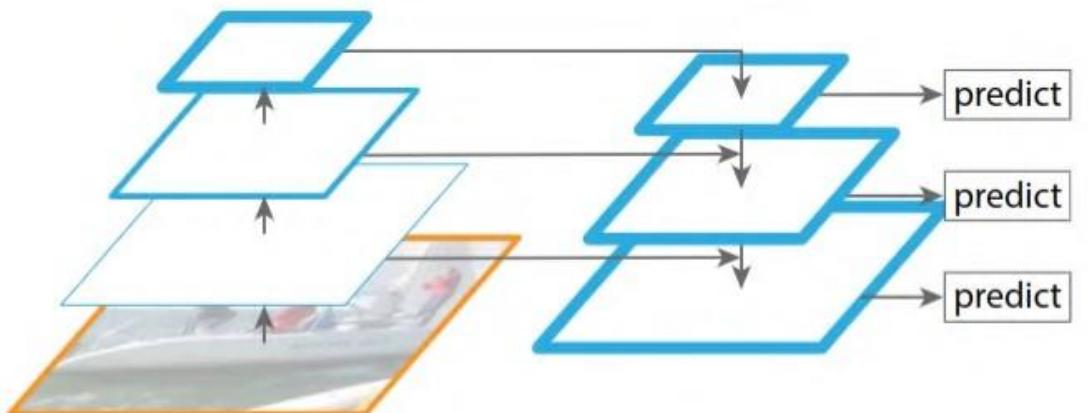
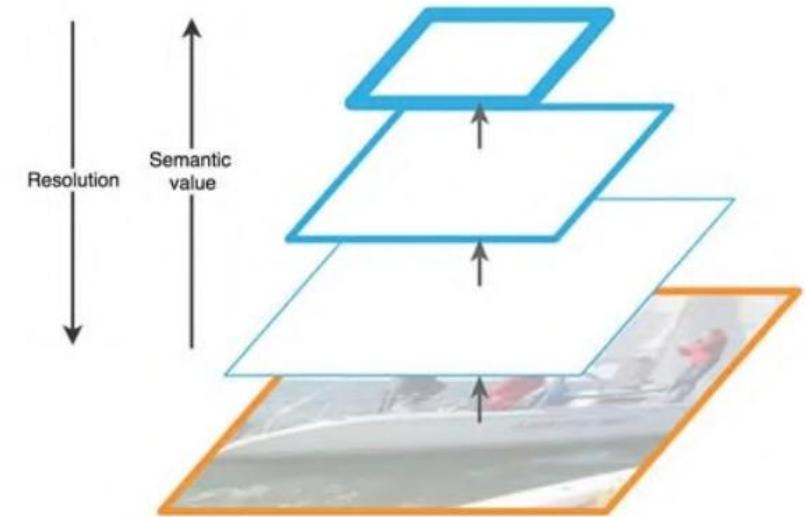
The goal is to naturally leverage the pyramidal shape of a ConvNet's feature hierarchy while creating a feature pyramid that has **strong semantics at all scales**.

Combines

- semantically strong features (i.e low resolution)
- semantically weak features (i.e high-resolution)

Using

- top-down pathway
- lateral connections



Motivations

Needing of recognizing objects at vastly different scales. Feature pyramid built upon image pyramids (*featurized image pyramids*) were the baseline

Problem

At the time, state-of-the-art methods performed multi-scale testing on featurized image pyramids, but ... inference time and memory increased a lot! Moreover, training in an end-to-end manner with image pyramids is infeasible. Many used them only at inference time: this leads to inconsistency w.r.t. training.

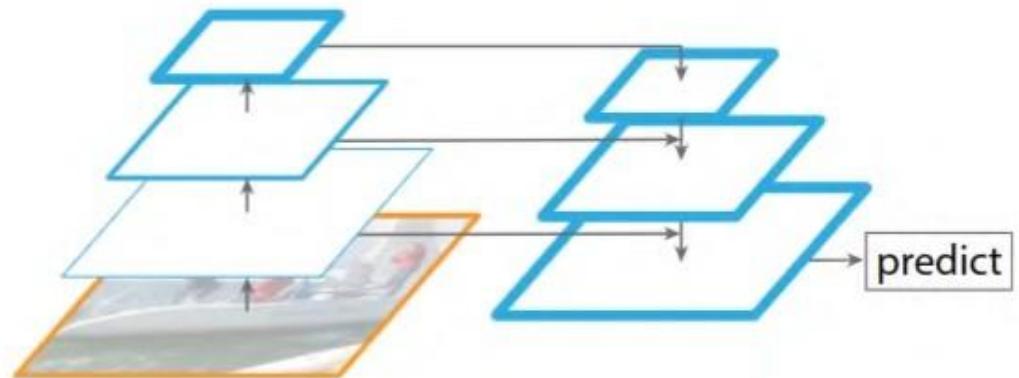
For this reason, Fast and Fast R-CNN opt to not use them:

"Recent and more accurate detection methods like Fast R-CNN [11] and Faster R-CNN [29] advocate using features computed from a single scale, because it offers a good trade-off between accuracy and speed. Multi-scale detection, however, still performs better, especially for small objects."

FPN vs Others

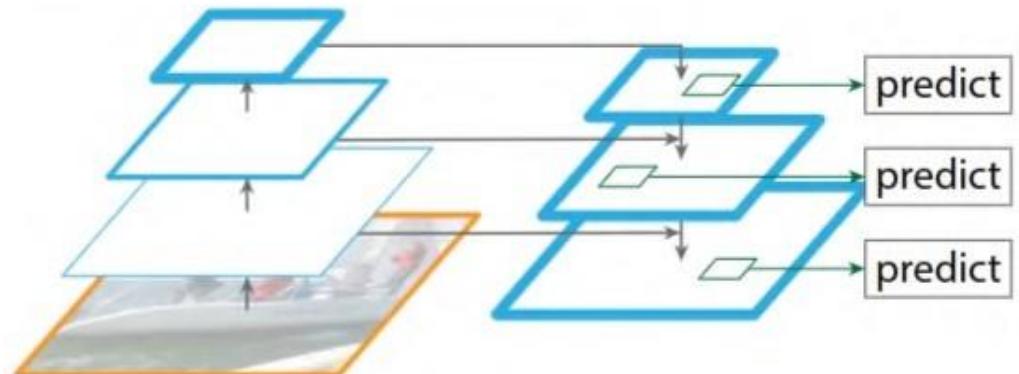
Others (e.g U-Net)

- Architectures also adopting top-down and skip connections
- goal is to produce a **single high-level feature map** on which make predictions
- What is missing: needed to recognize objects across multiple scales



Feature Pyramid Network

- predictions (e.g., object detections) are independently made on each level



Method

Input

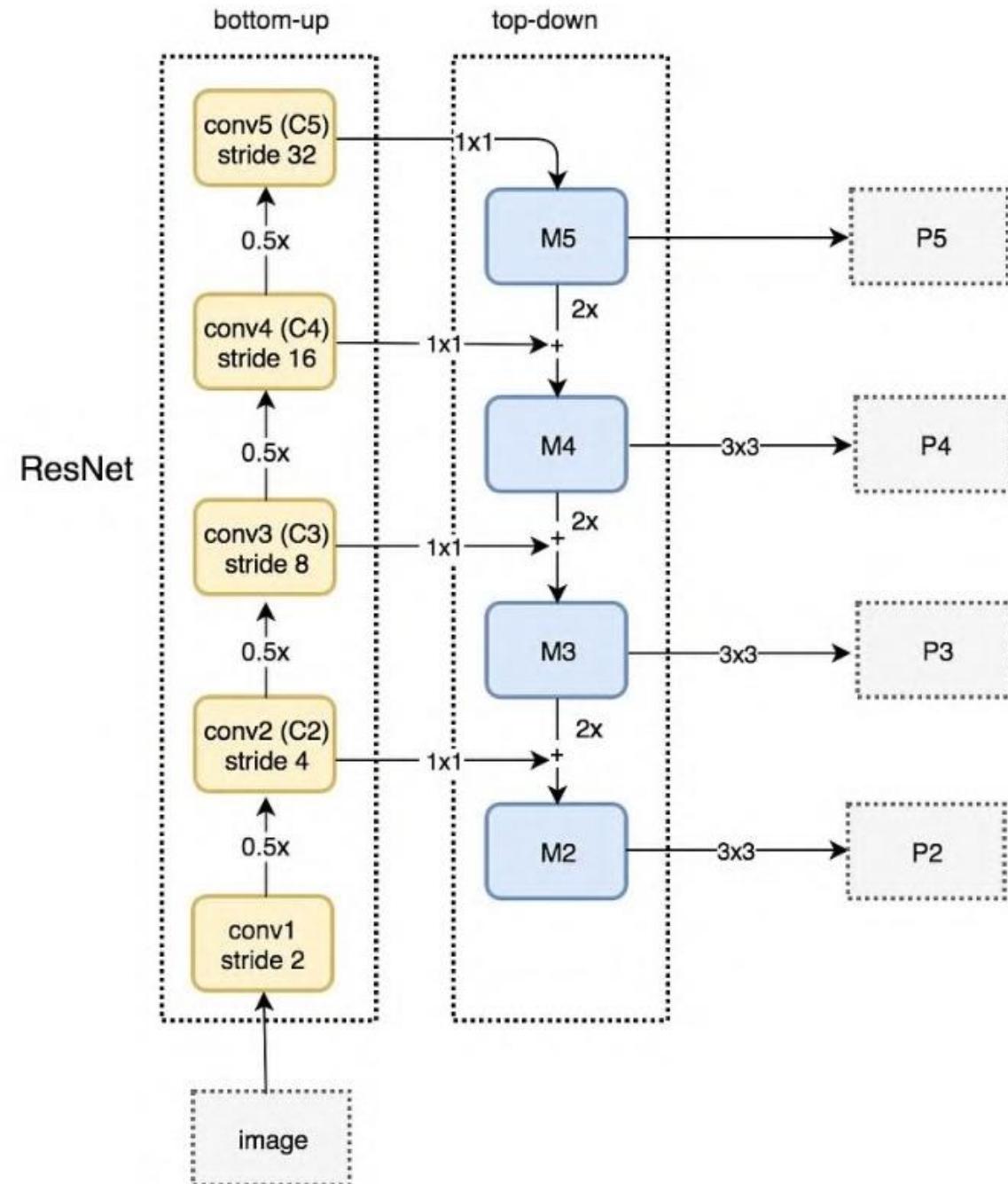
- A single-scale image of an arbitrary size

Output

- proportionally sized feature maps

Components

- Bottom-up pathway
- Top-down pathway
- Lateral-connections



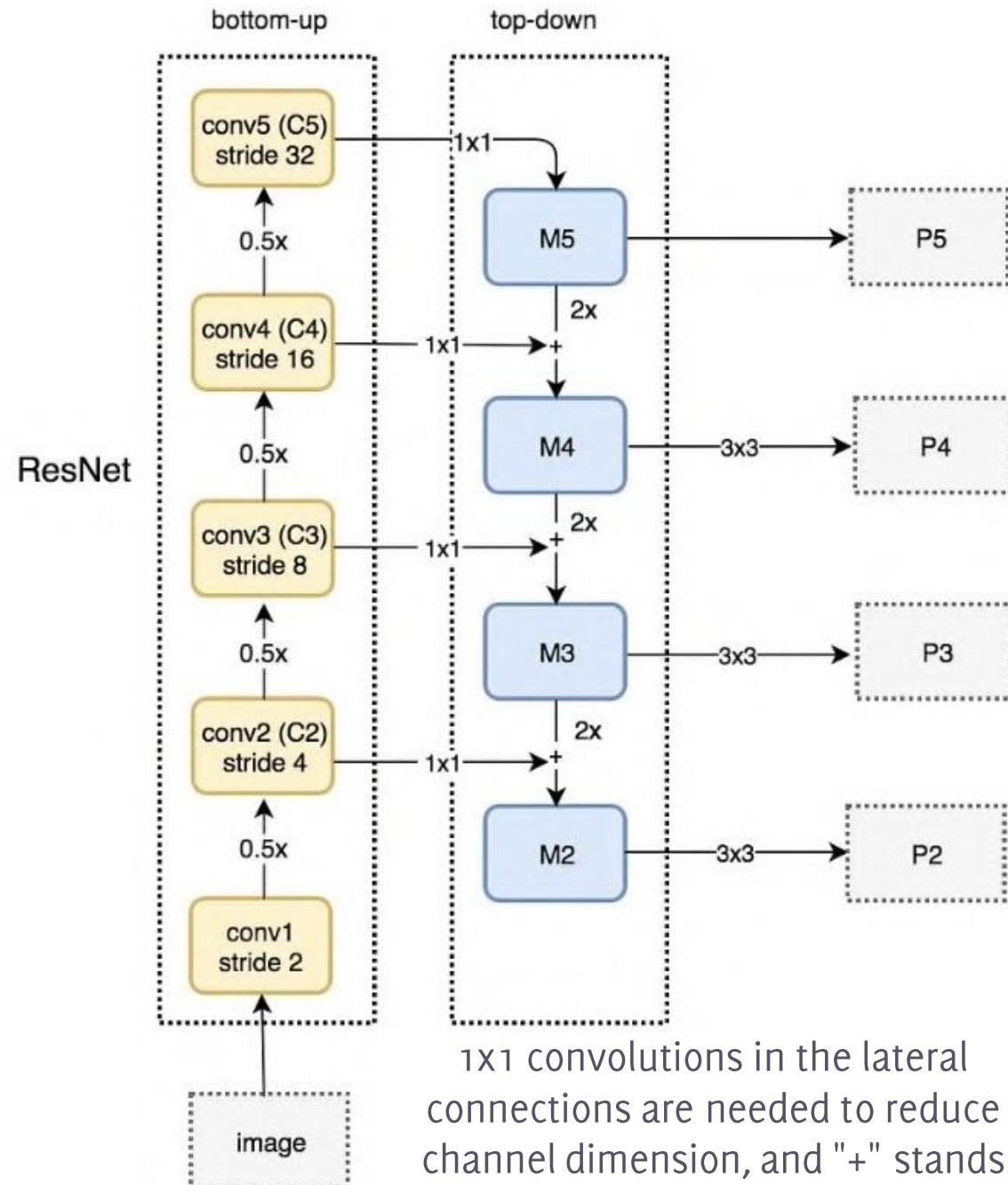
Method

Bottom-up

- Proportionally sized feature maps
- Moving up, the spatial dimension is reduced by 1/2
- Output of each convolution is used in the top-down pathway

Top-down with connections

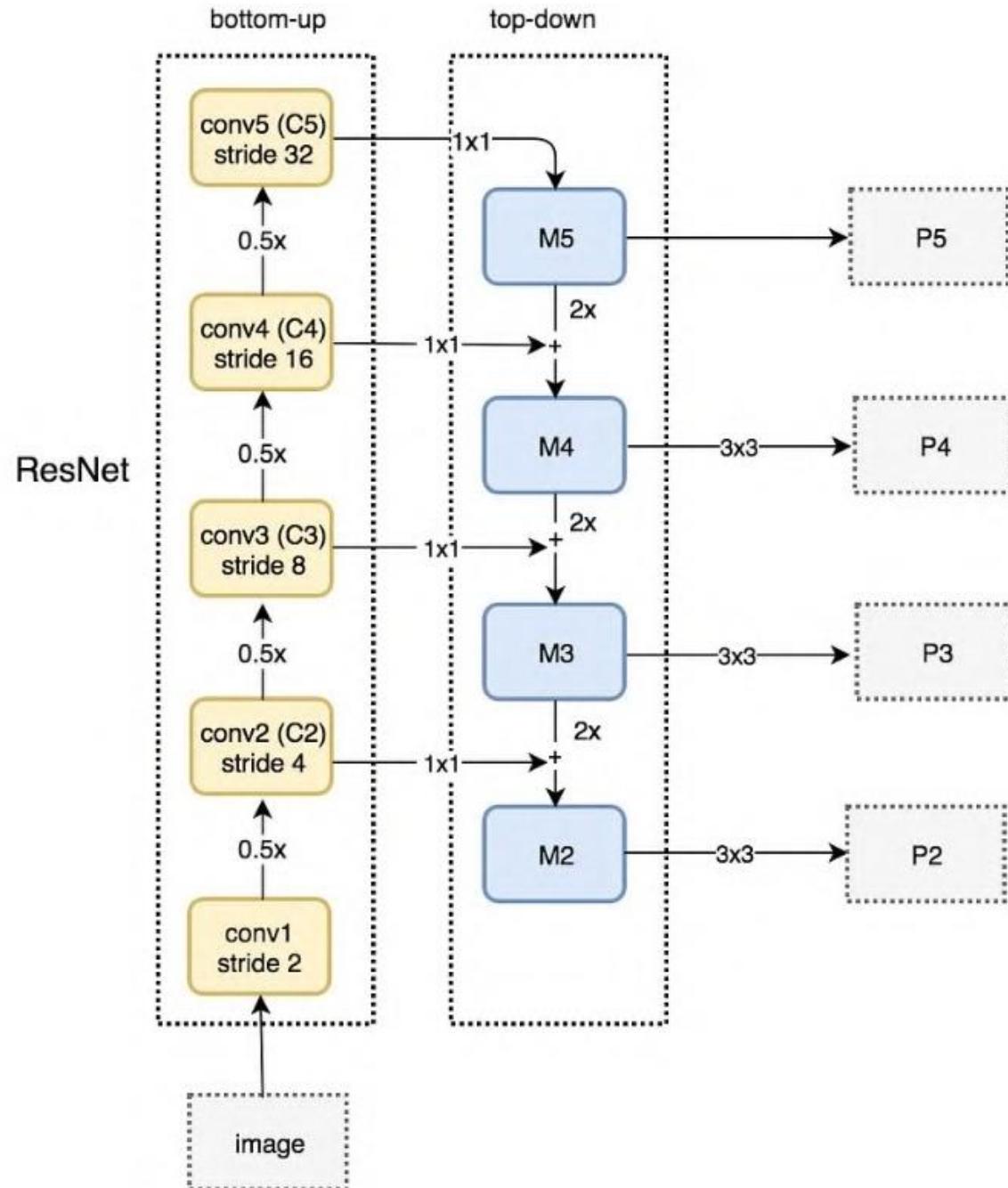
- Derives higher resolution features by upsampling .
- Enhanced with features from the bottom up part through lateral connections (reduce channel dimensions).
- 3×3 convolution to merged layers, that reduces the aliasing effect.

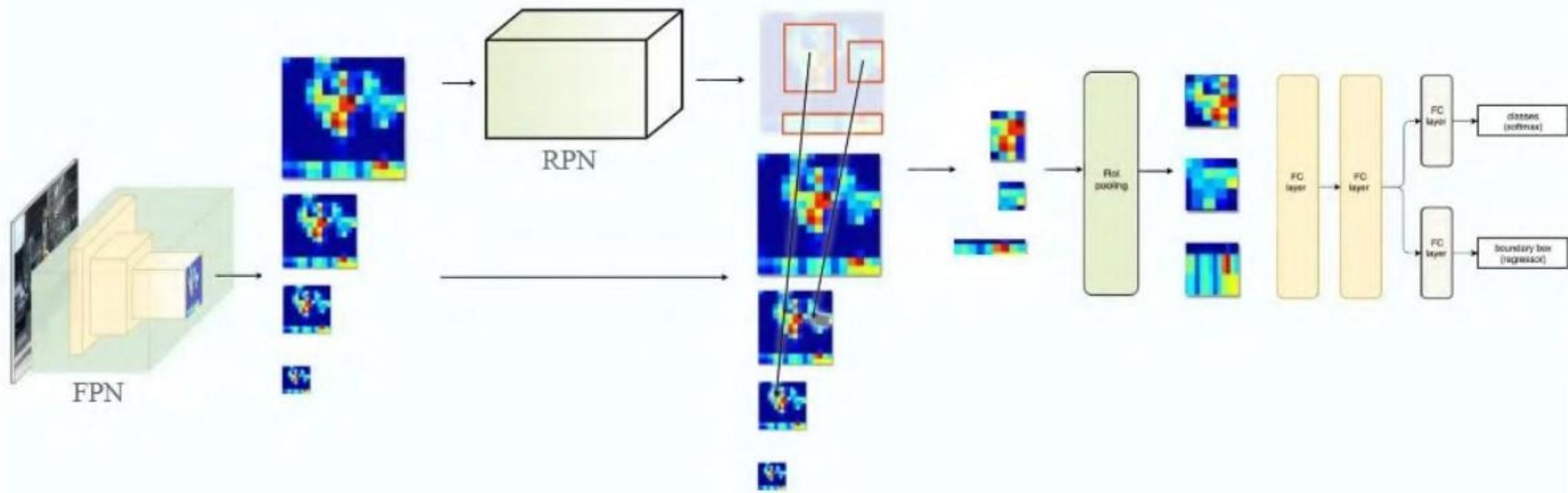


FPN with RPN

FPN is not an object detector by itself. It is a feature extractor that works with object detectors

RPN was adapted by replacing the single-scale feature map with FPN by attaching a **predictor head** (3×3 conv and two sibling 1×1 convs) to each level on the feature pyramid. Since the head slides over all locations in all pyramid levels, it is not necessary to have multi-scale anchors on a specific level.





FPN with Fast R-CNN or Faster R-CNN

We apply the RPN to generate RoIs. Based on the size of the RoI, we select the feature map layer in the proper scale to extract the feature patches. Formally, a RoI of width w and height h is assigned to the level k of the feature pyramid:

$$k = \lfloor k_0 + \log_2(\sqrt{wh}/224) \rfloor.$$

where: $k_0 = 4$
 k is the P_k layer in the FPN used to generate the feature patch.

So if $k = 3$, we select P_3 as our feature maps. We apply the ROI pooling and feed the result to the Fast R-CNN head to finish the prediction.

References

1. <https://jonathan-hui.medium.com/understanding-feature-pyramid-networks-for-object-detection-fpn-45b227b9106c>
2. https://openaccess.thecvf.com/content_cvpr_2017/html/Lin_Feature_Pyramid_Networks_CVPR_2017_paper.html

In collaboration with



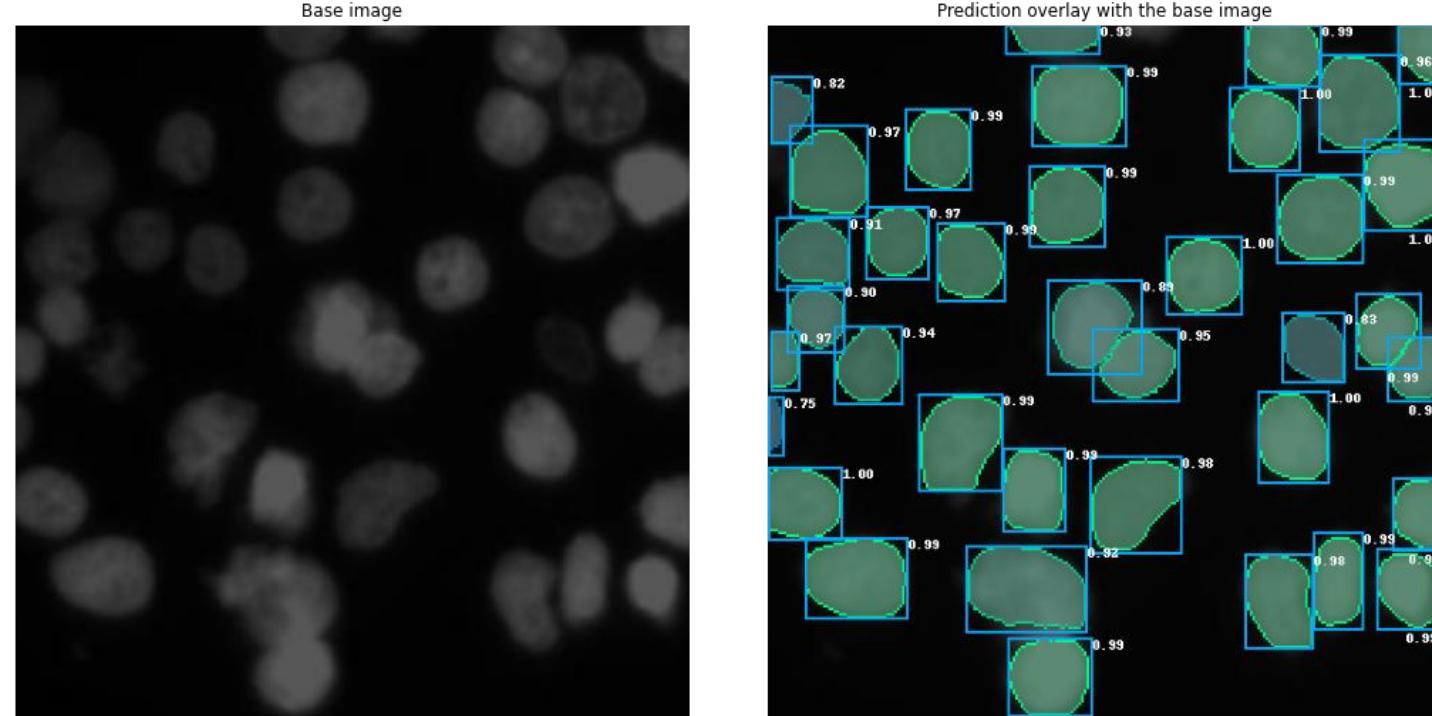
Instance Segmentation

Research Project with Ikonisys

Credits Roberto Basla

Instance Segmentation Network

The instance segmentation network is able to identify nuclei in fluoroscopy images.

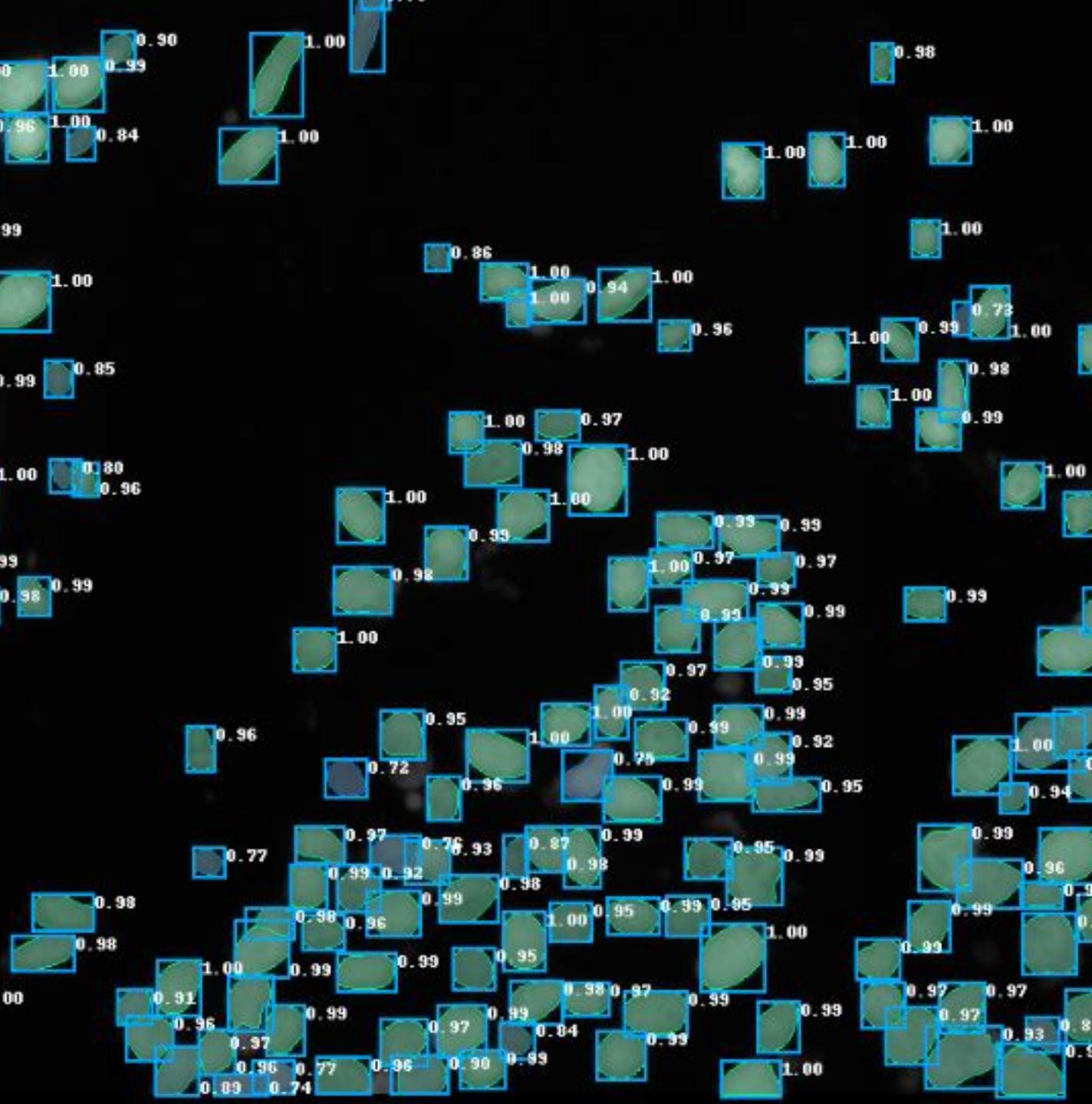


The network provides, for each identified nucleus, its ^{Image}segmentation mask, ^{Prediction}bounding box, and prediction confidence.

Model Performance

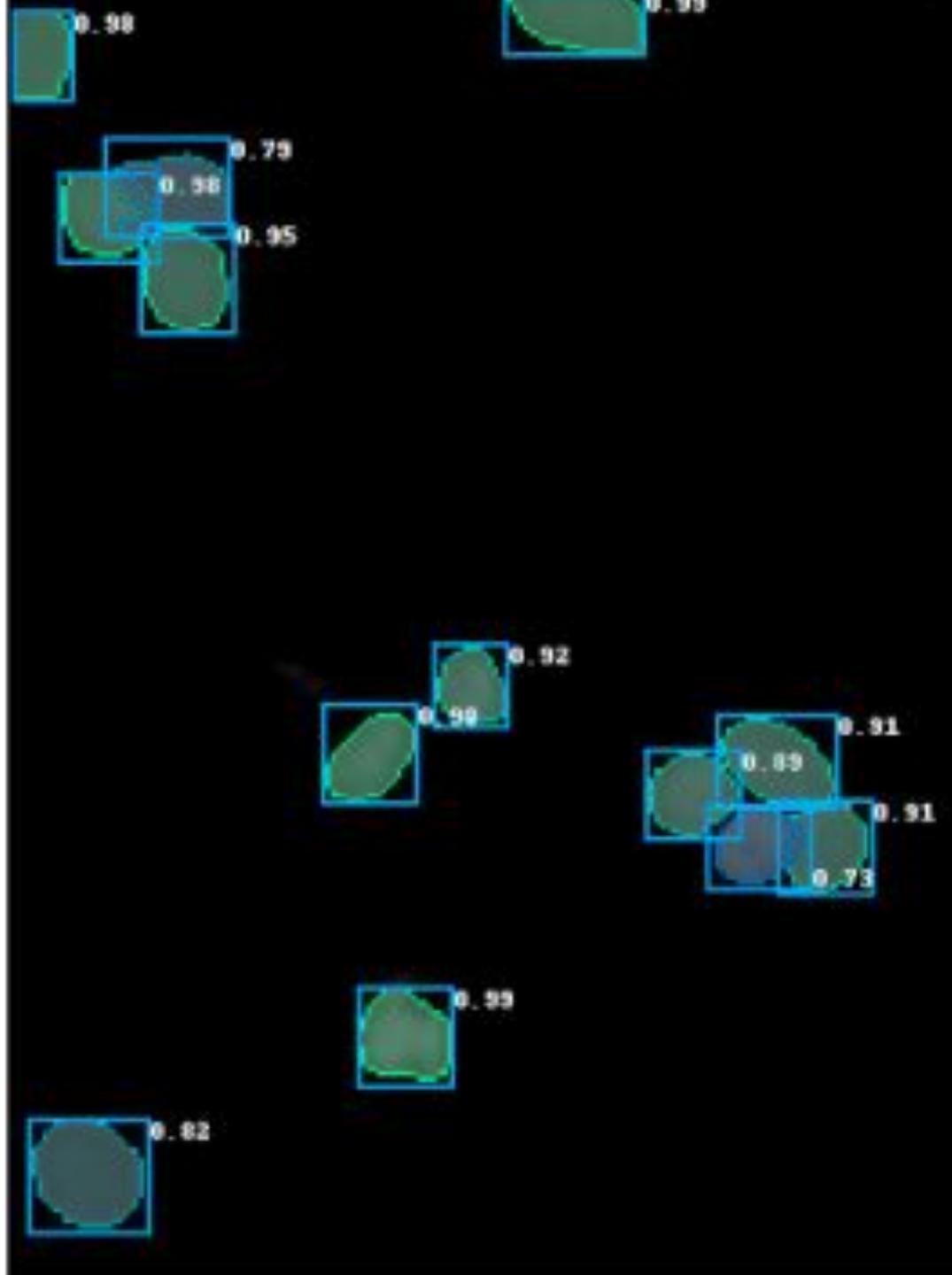
The model achieves great results even when trained with **very few images from the target domain**, by making use of transfer learning from these datasets:

- DSB2018
- DSB2019
- BBBCo39
- NSDE



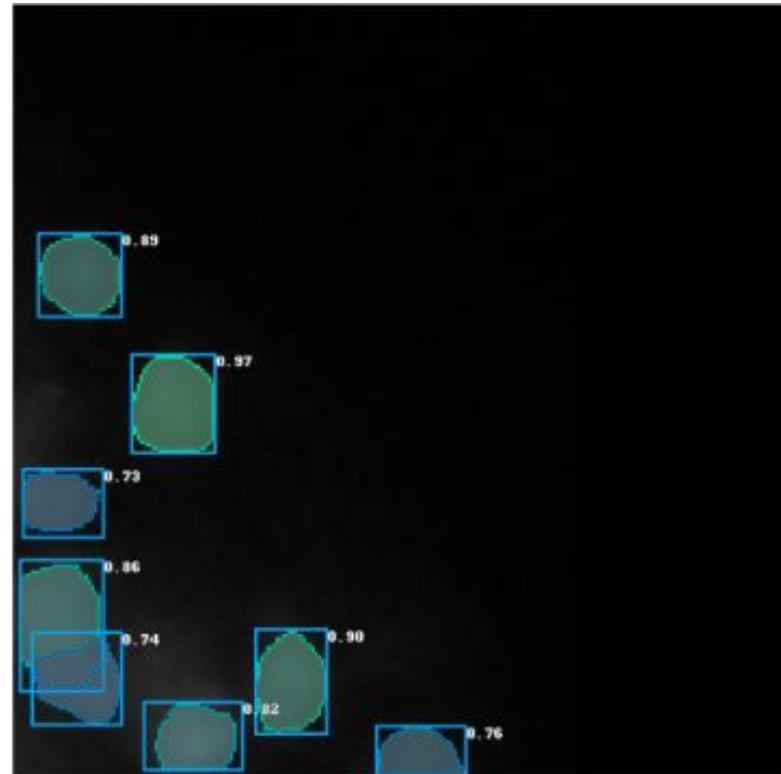
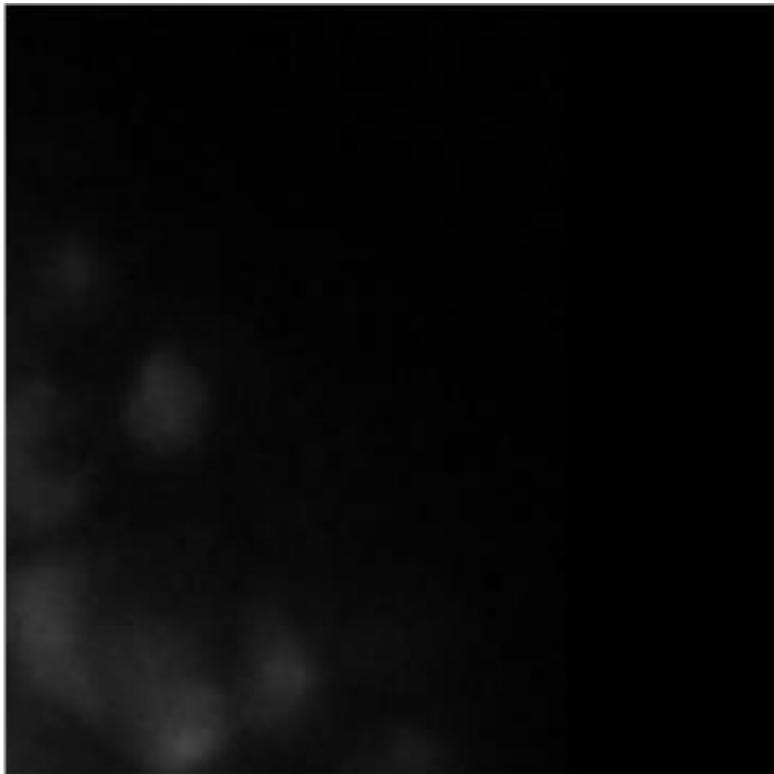
Overlapping segmentation

Overlapping cells can be identified and masked independently.



Robustness

The model can identify even the faintest of nuclei in the image, further proving the performance of the system.



Ongoing Research Activities

... good for thesis projects

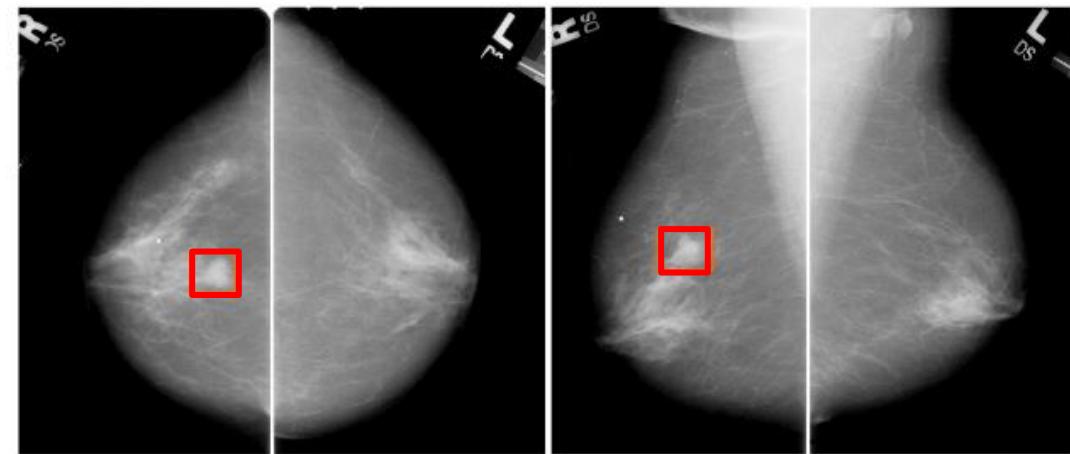
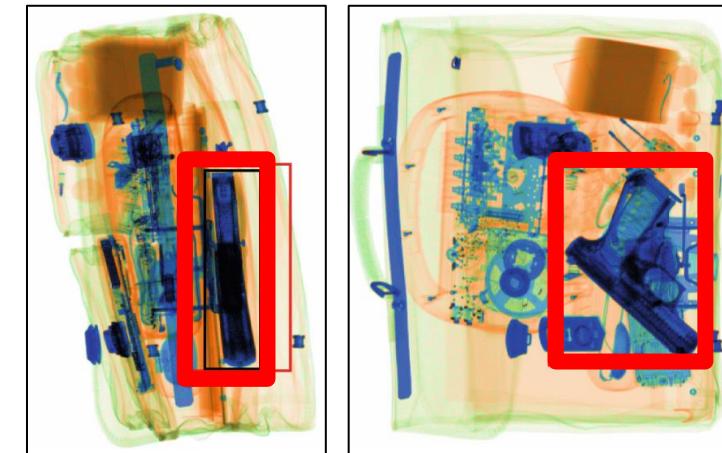
Multi-view object detection CNN

Detection systems usually take only a single input image. In some applications, such as

- Baggage inspection in airports
- Tumor detection in X-ray images

you have multiple views of the same objects.

Exploiting multiple-views can increase detection power.



[1] Steitz et al., "Multi-view x-ray r-cnn." (2018)

[2] Wimmer et al. "Multi-task fusion for improving mammography screening data classification" (2021)

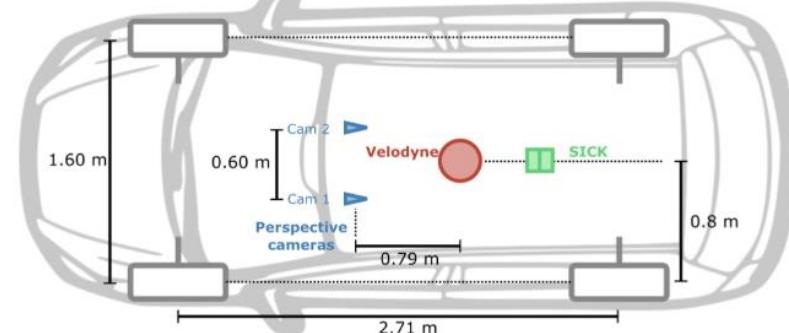
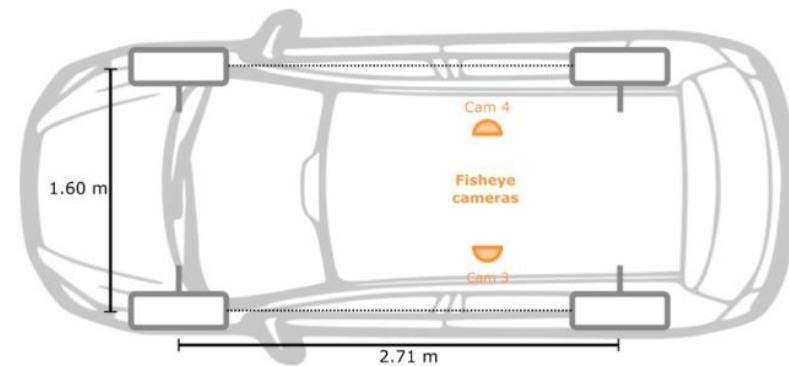
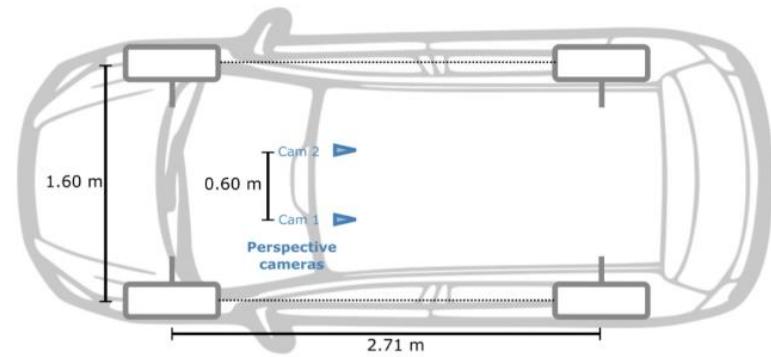
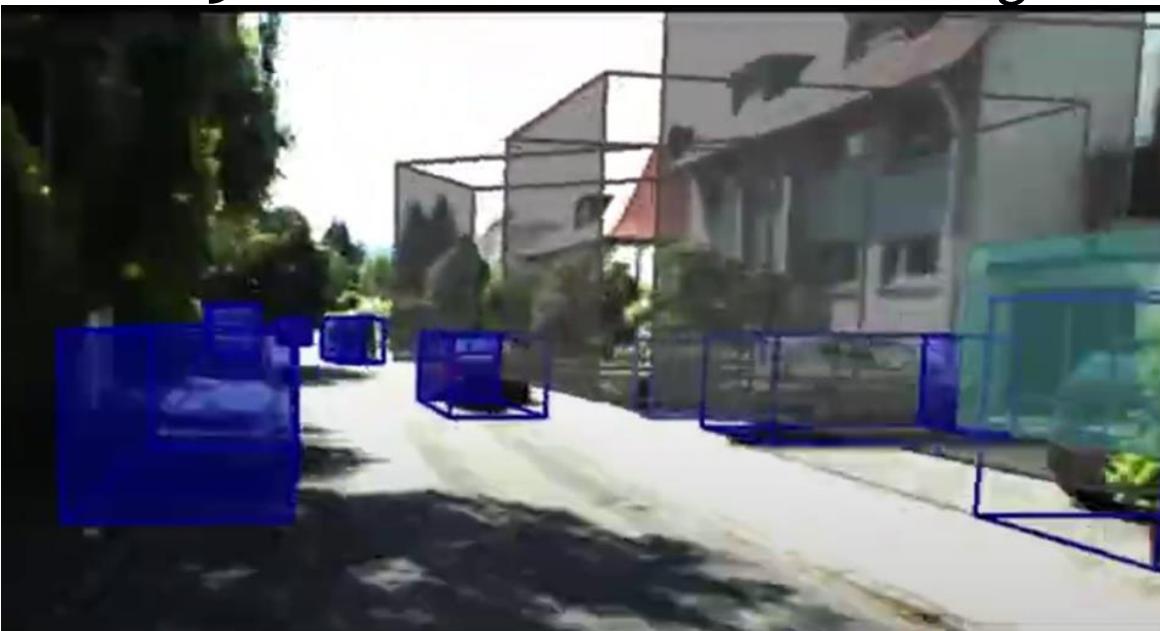
Multimodal object detection CNN

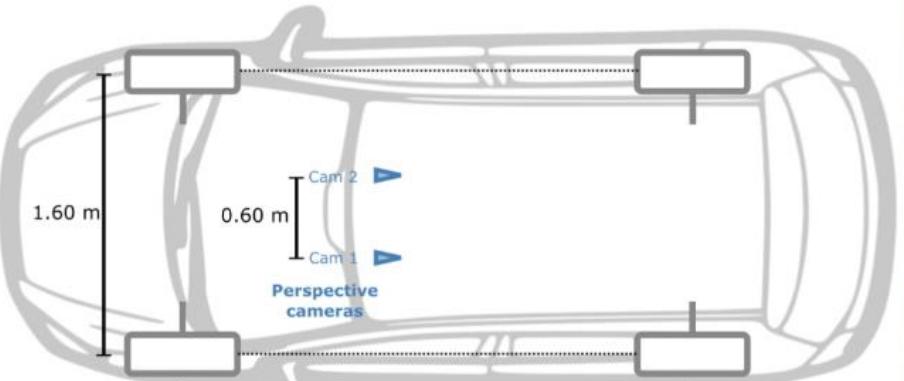
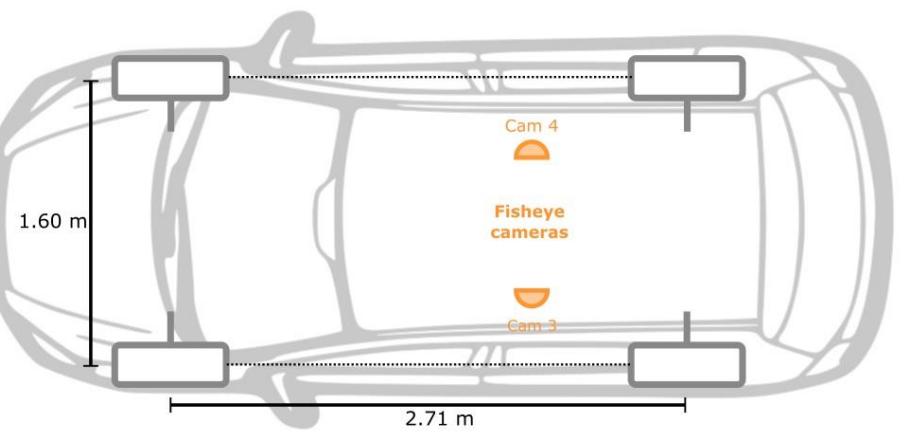
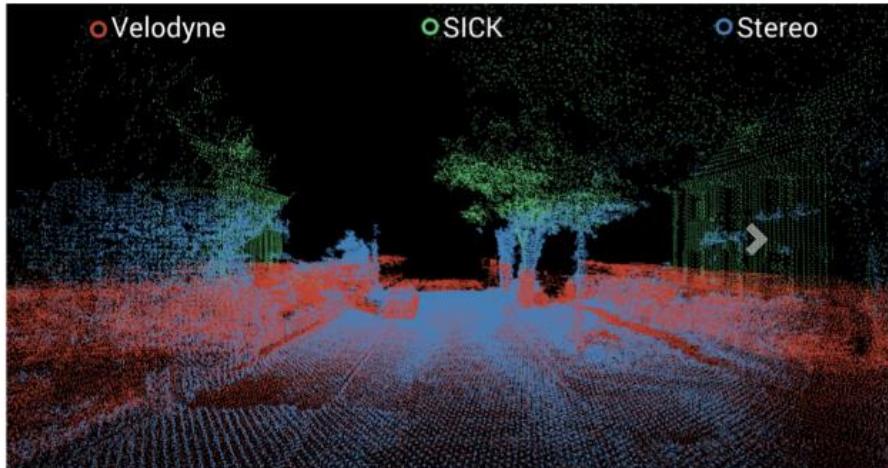
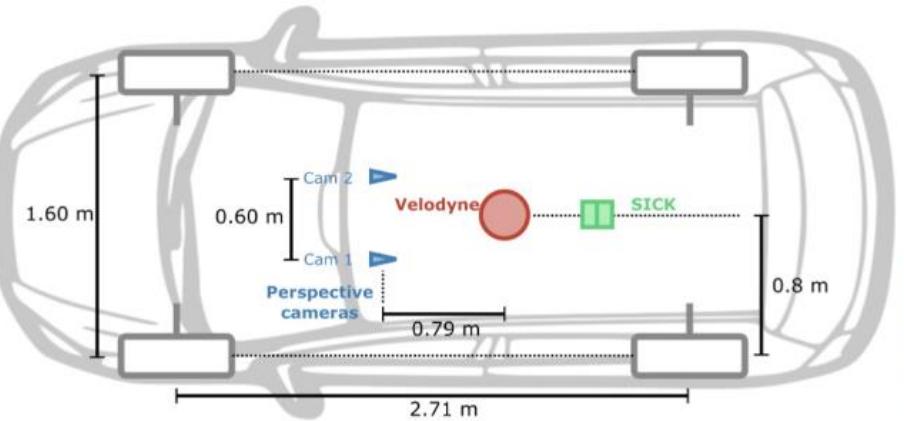
In collaboration with Prof.
Savarese, MOVE Research Group

Object detection can leverage multiple streams:

- RGB images (possibly stereo pairs)
- Point Clouds from lasers scanners
- event cameras..

To infer 3D coordinates of bounding boxes







AN2DL, Boracchi

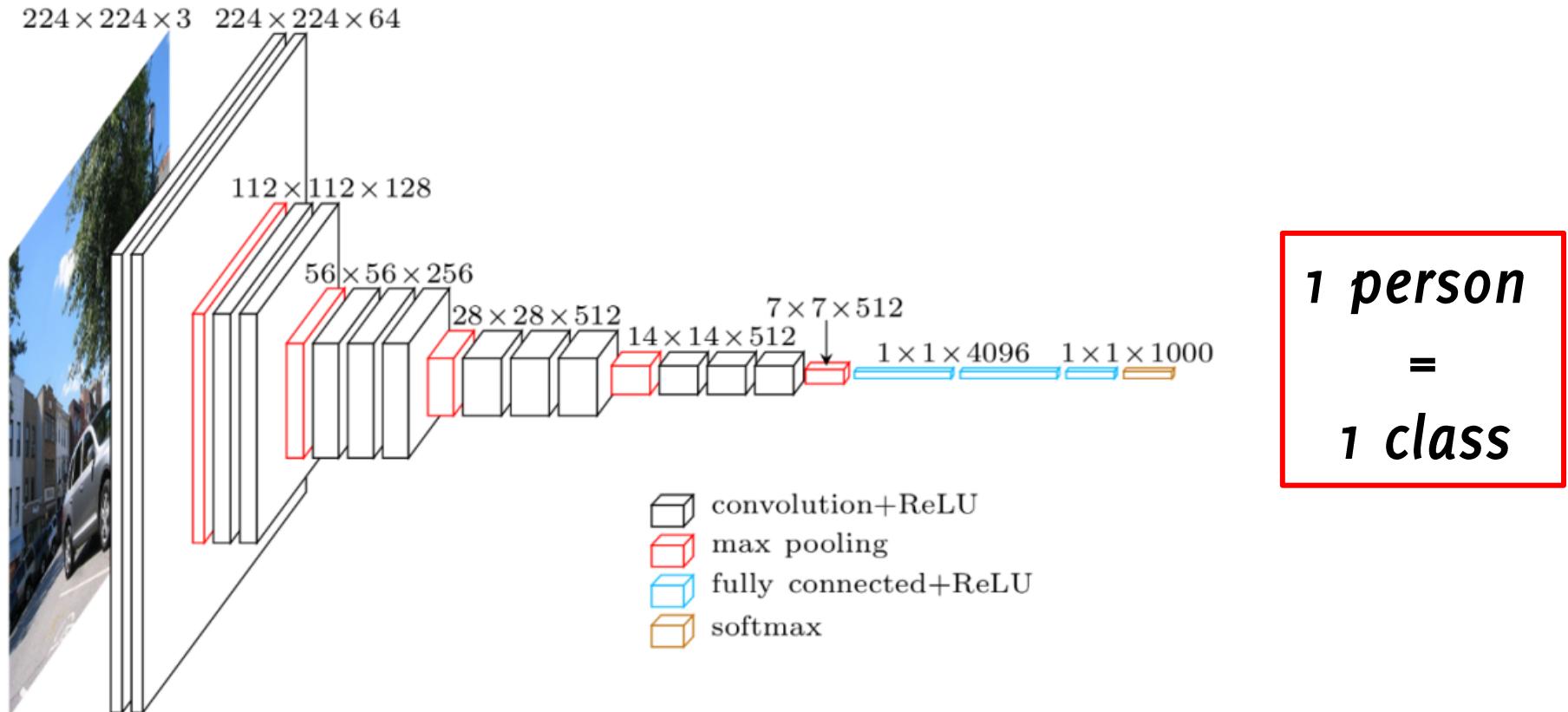
Metric Learning

Say you are asked to implement a face
identification system
to open Bocconi door!

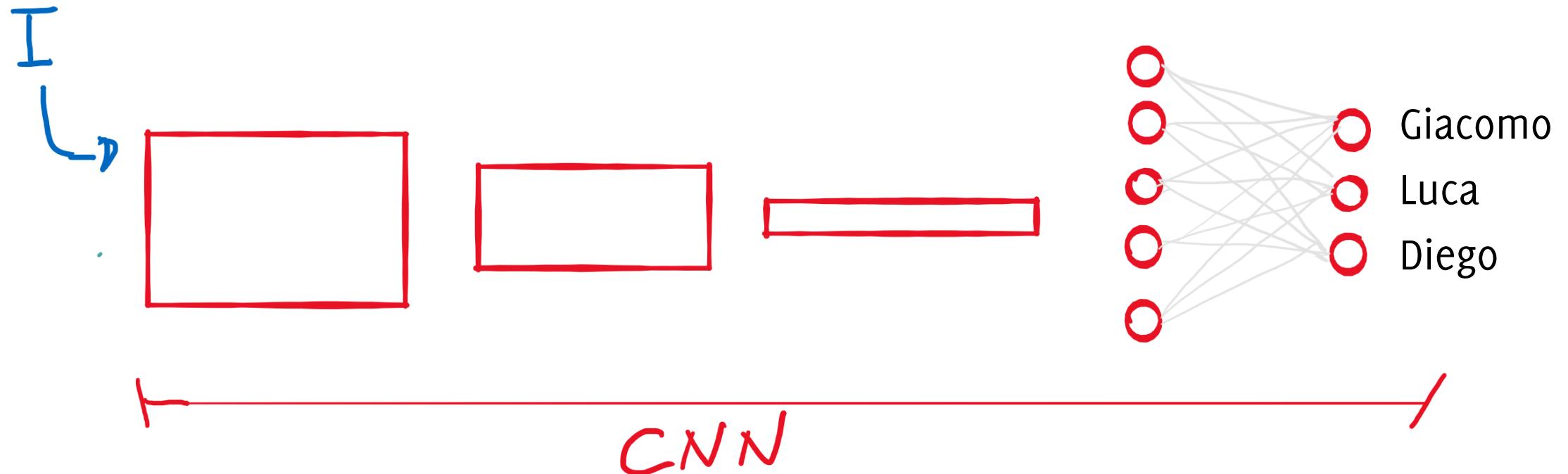


Classification? Detection? Segmentation?

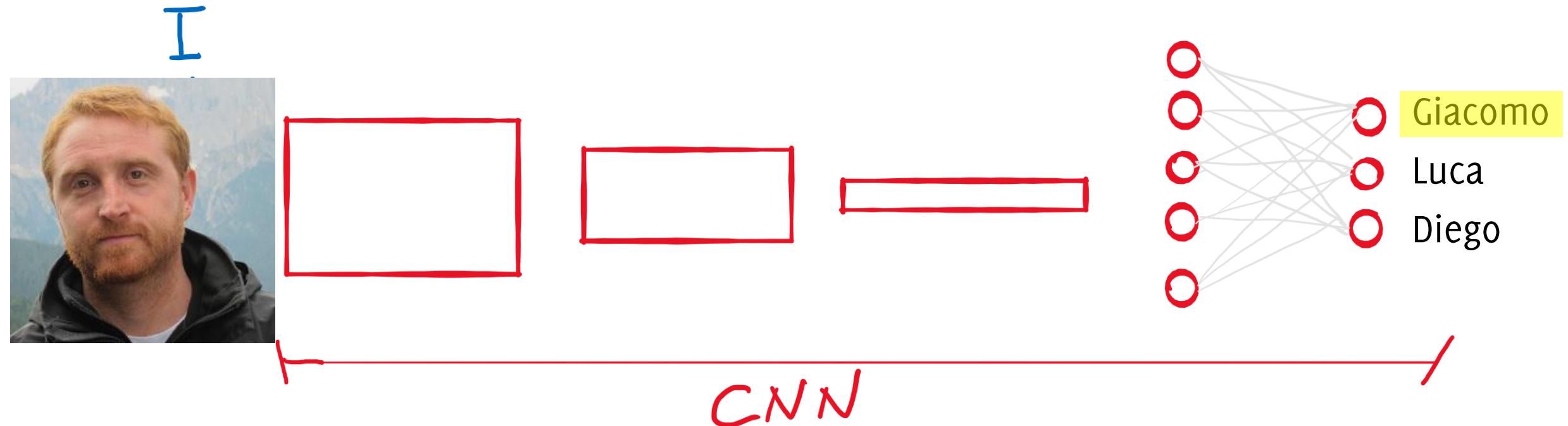
Let's start fresh from a Convolutional Neural Network for classification



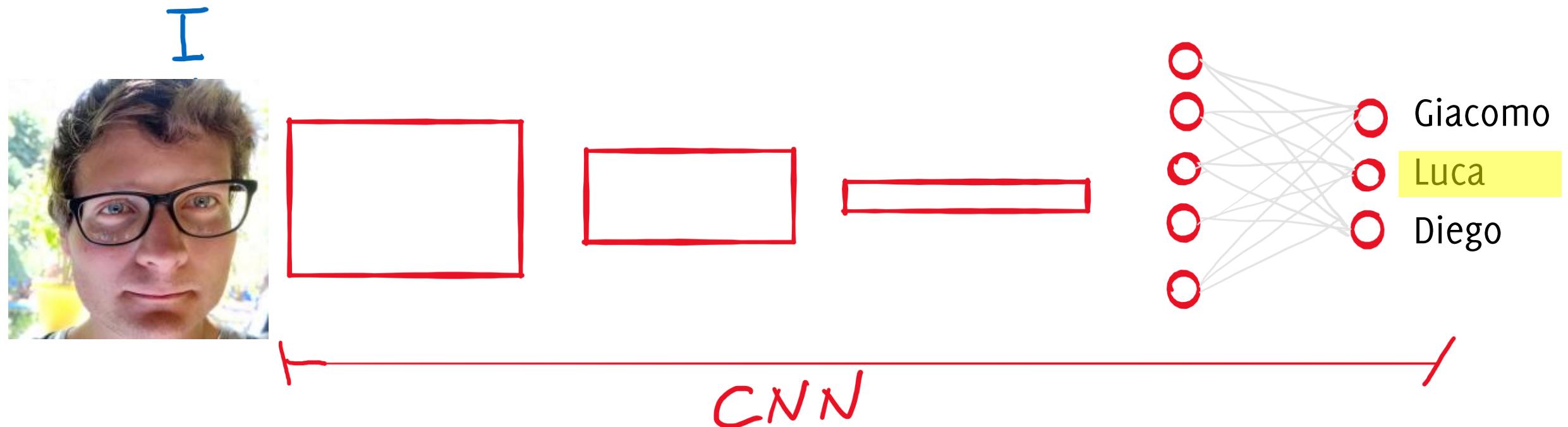
Face identification by Classification



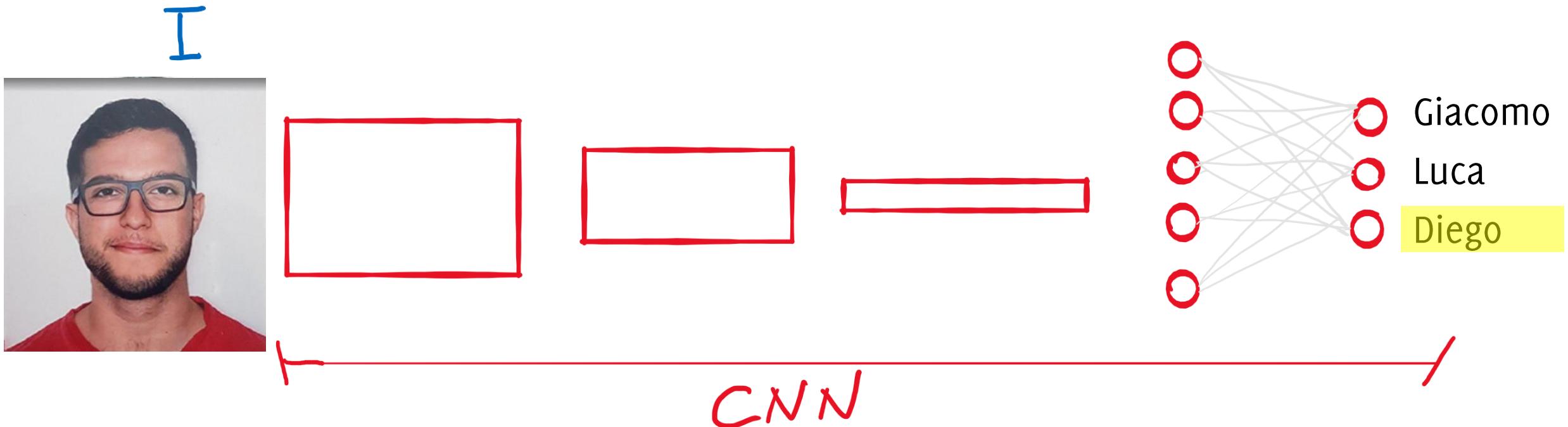
Face identification by Classification



Face identification by Classification



Face identification by Classification



What do we need?

A training set

- A few images per class / person
- Possibly Images in different conditions (position, light, facial expression, clothes ...)

A few Py snippets and a GPU...

That's easy...

What do we need?

A training set

- A few images per class / person
- Possibly Images in different conditions (position, light, facial expression, clothes ...)

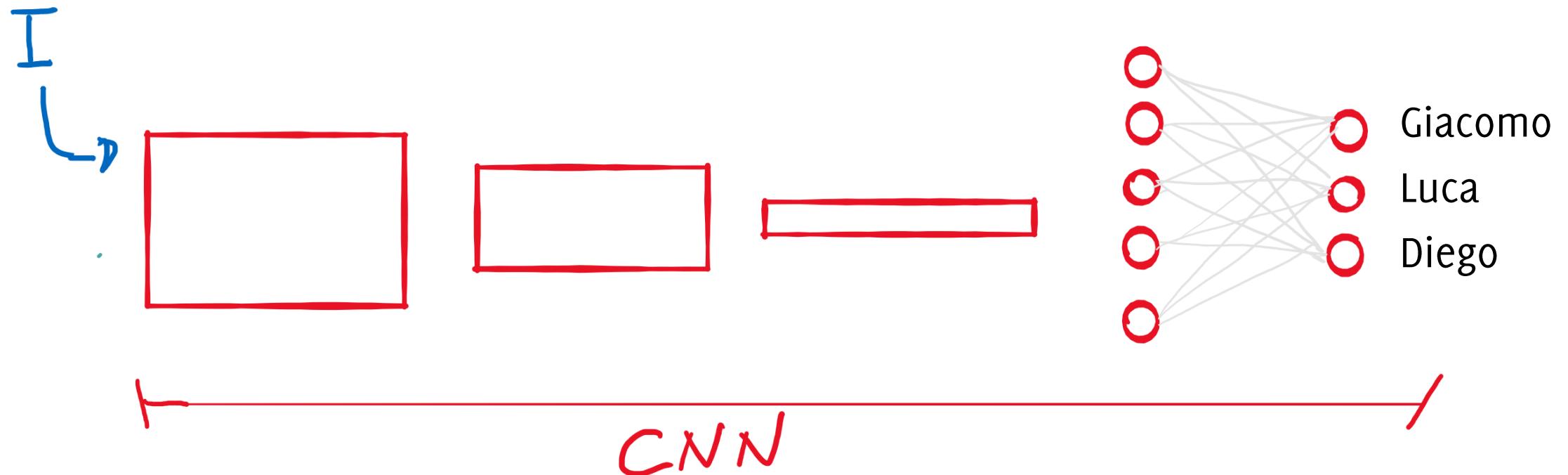
A few Py snippets and a GPU...

That's easy...

Are we happy with this solution?

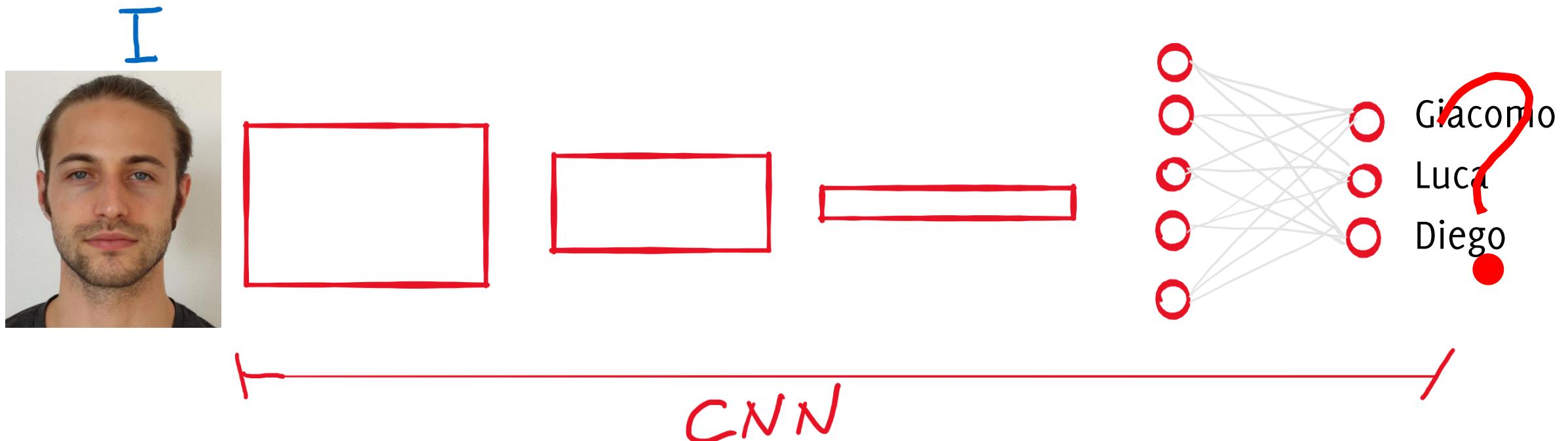
... Not quite

What happens when you need to enroll a new employee?



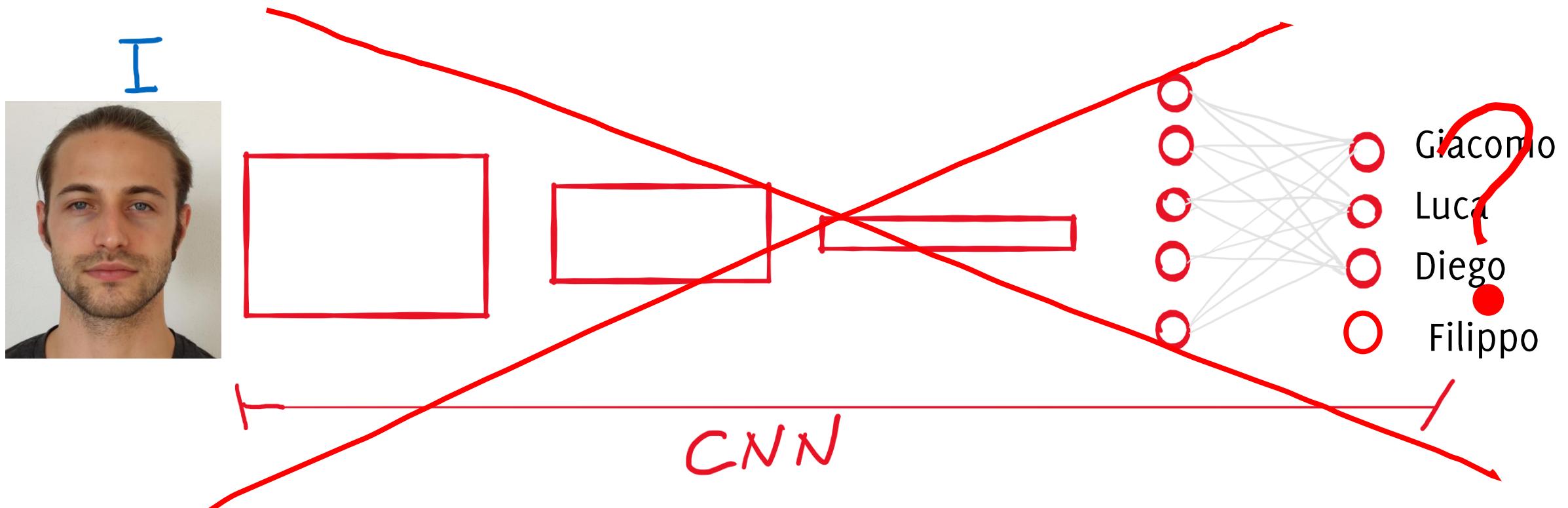
... Not quite

What happens when you need to enroll a new employee?



... Not quite

What happens when you need to enroll a new employee?



The whole network has to be retrained for each new person to be identified

A Different Approach is Needed!

What about image comparison?

Why don't we store a picture for each employee (**template**), and then perform identification by pairing the input to its closest template?

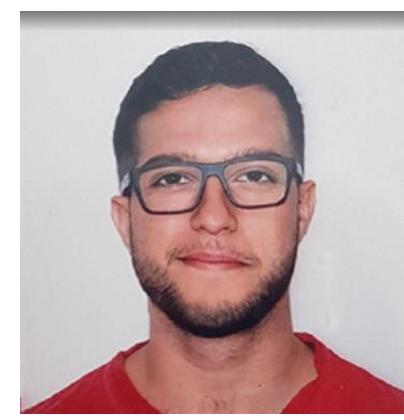
T_1



T_2



T_3



What about image comparison?

So, identification becomes:

$$\hat{i} = \operatorname{argmin}_{j=1,\dots,3} \|I - T_j\|_2$$

I



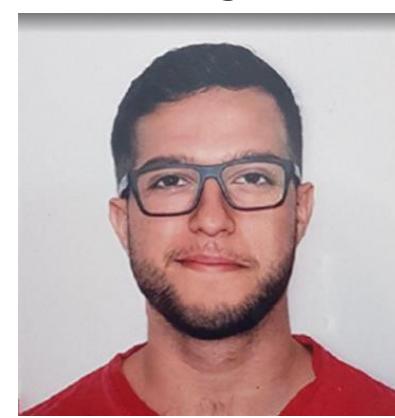
T_1



T_2



T_3



What about image comparison?

Enrolling a new individual would be straightforward



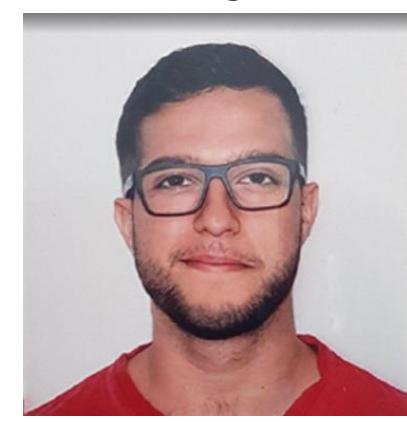
T_1



T_2



T_3



What about image comparison?

Enrolling a new individual would be straightforward... it is enough to add a template to the database!

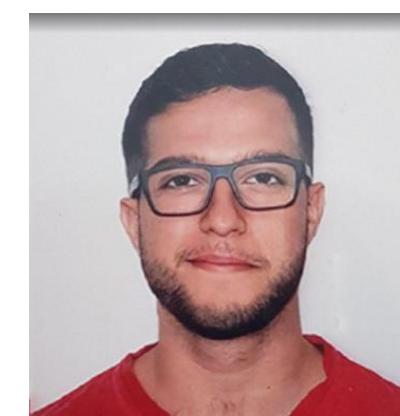
T_1



T_2



T_3



T_4



What about image comparison?

... but how to perform identification?

$$\hat{i} = \operatorname{argmin}_{j=1,\dots,4} \|I - T_j\|_2$$

I



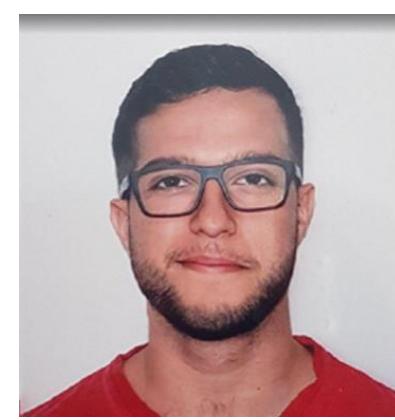
T_1



T_2



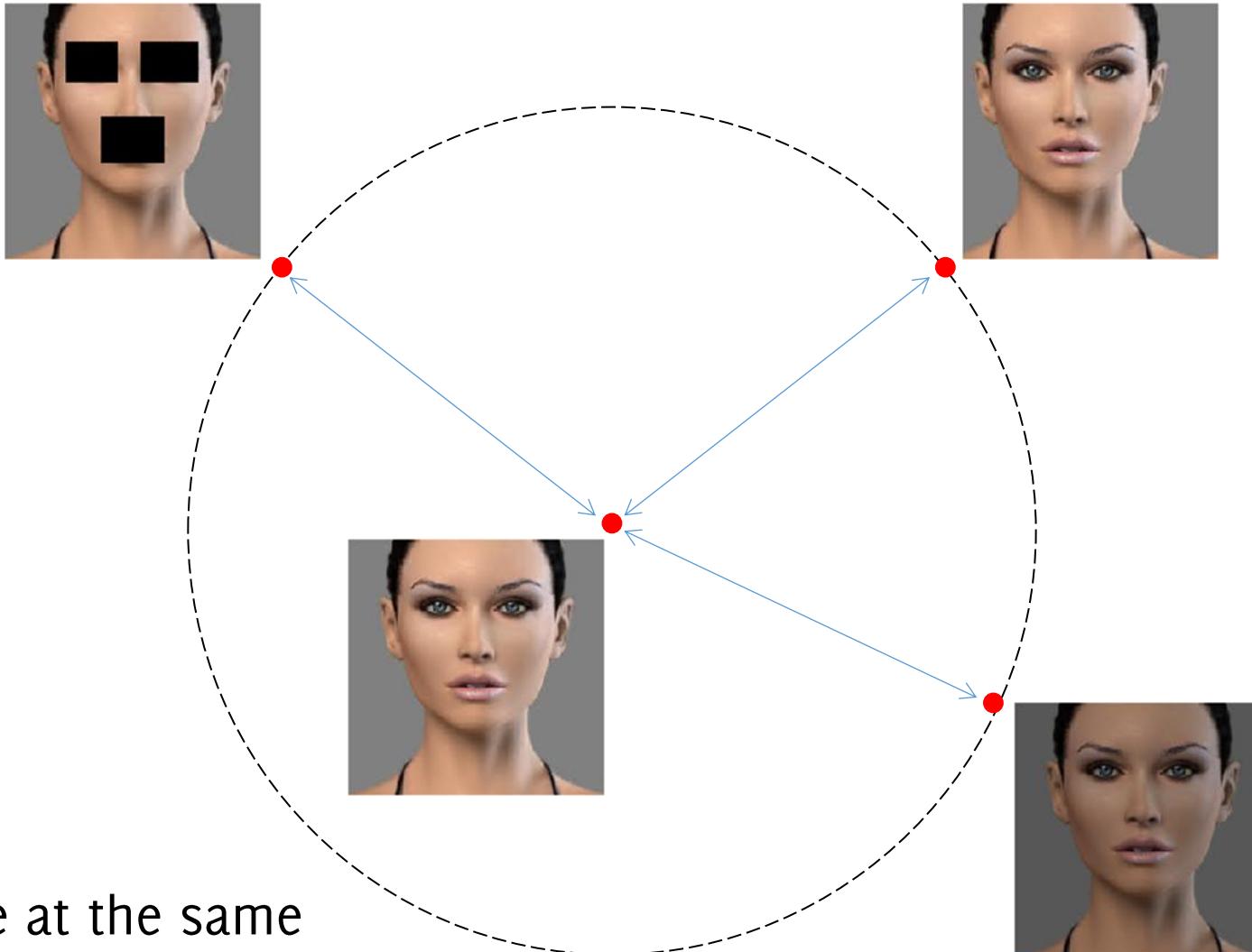
T_3



T_4



Bad News...



The three images are at the same distance from the reference at the center

What about image comparison?

... but how to perform identification?

$$\hat{i} = \operatorname{argmin}_{j=1,\dots,4} \|I - T_j\|_2$$

I



T_1



T_2



T_3



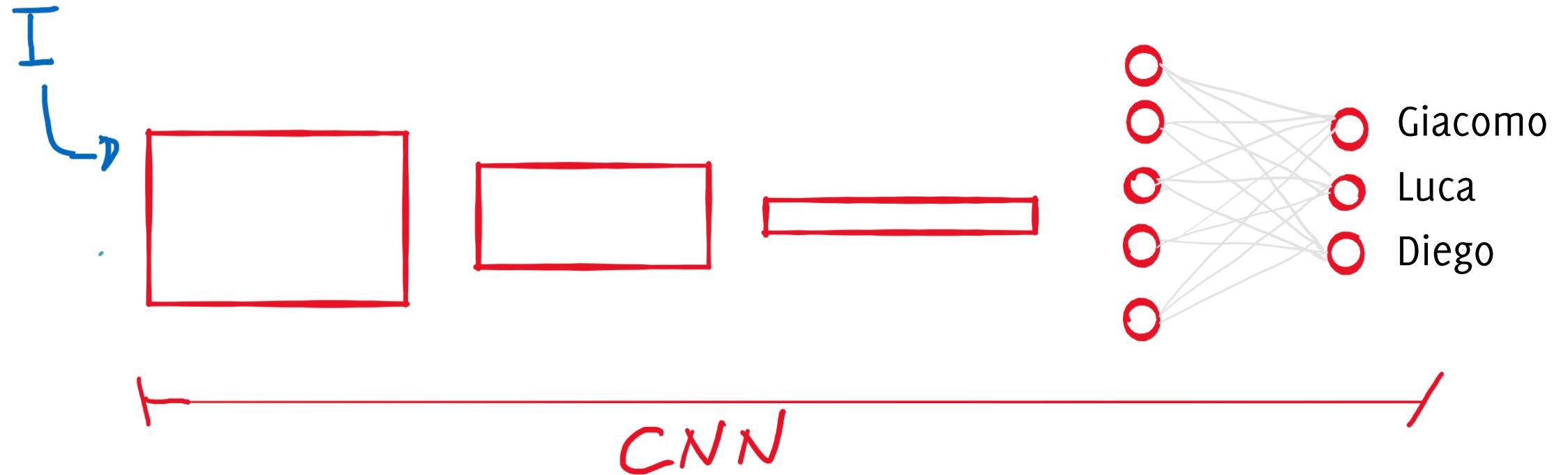
T_4



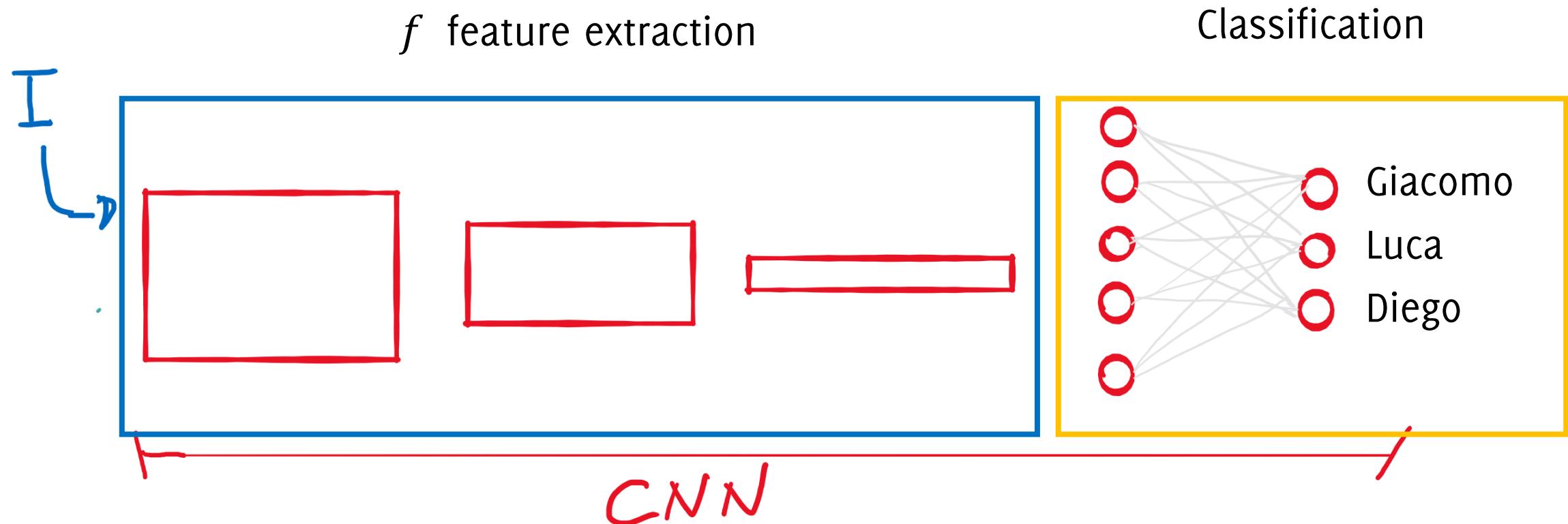
Is it possible to move to a learned distance,
and to train this relying only on these examples?

...we need a better distance measure for
face identification!

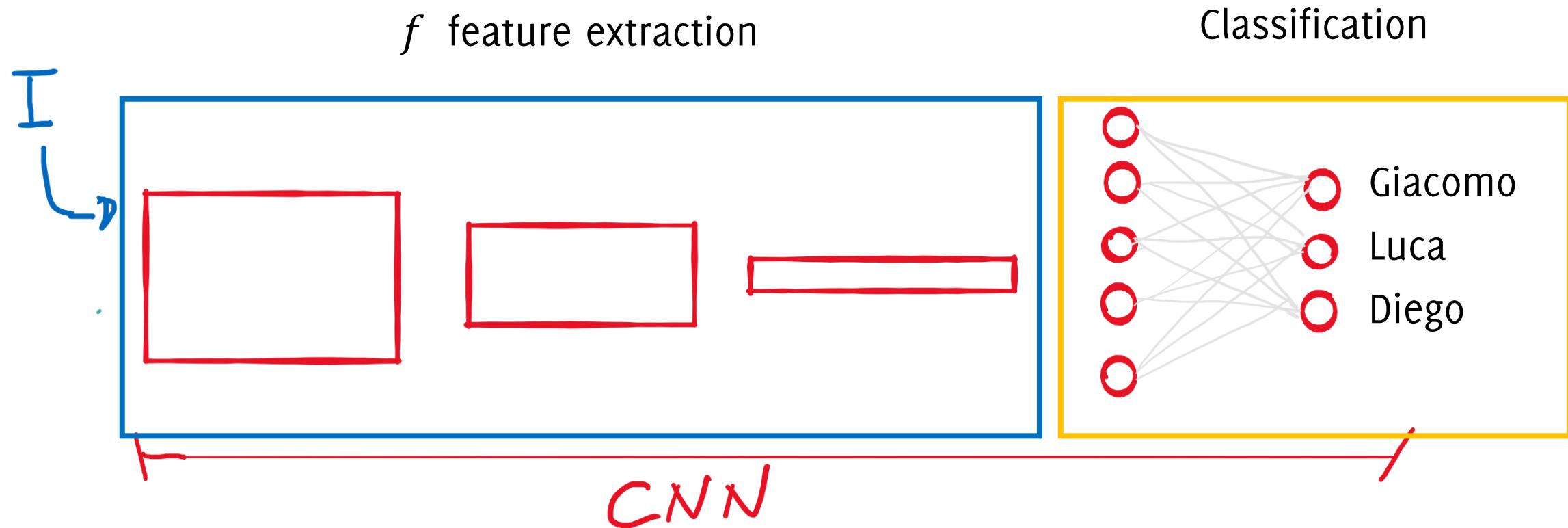
Dissecting CNN



Dissecting CNN

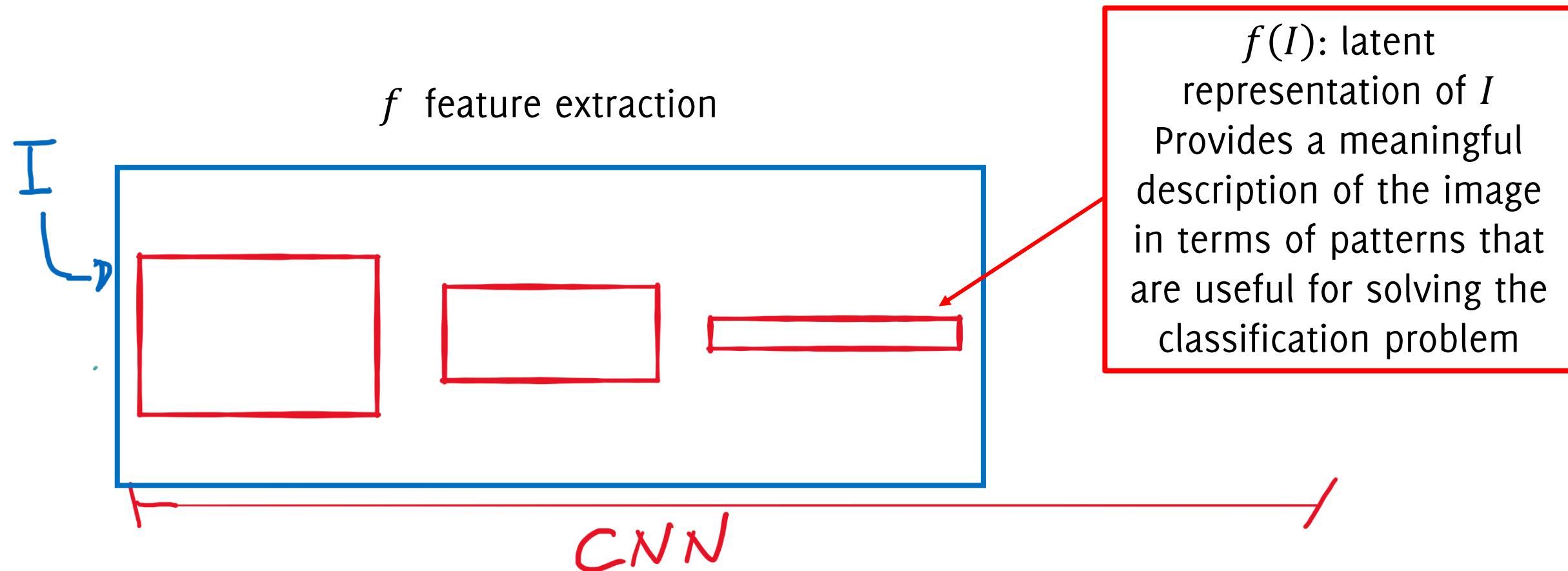


Dissecting CNN

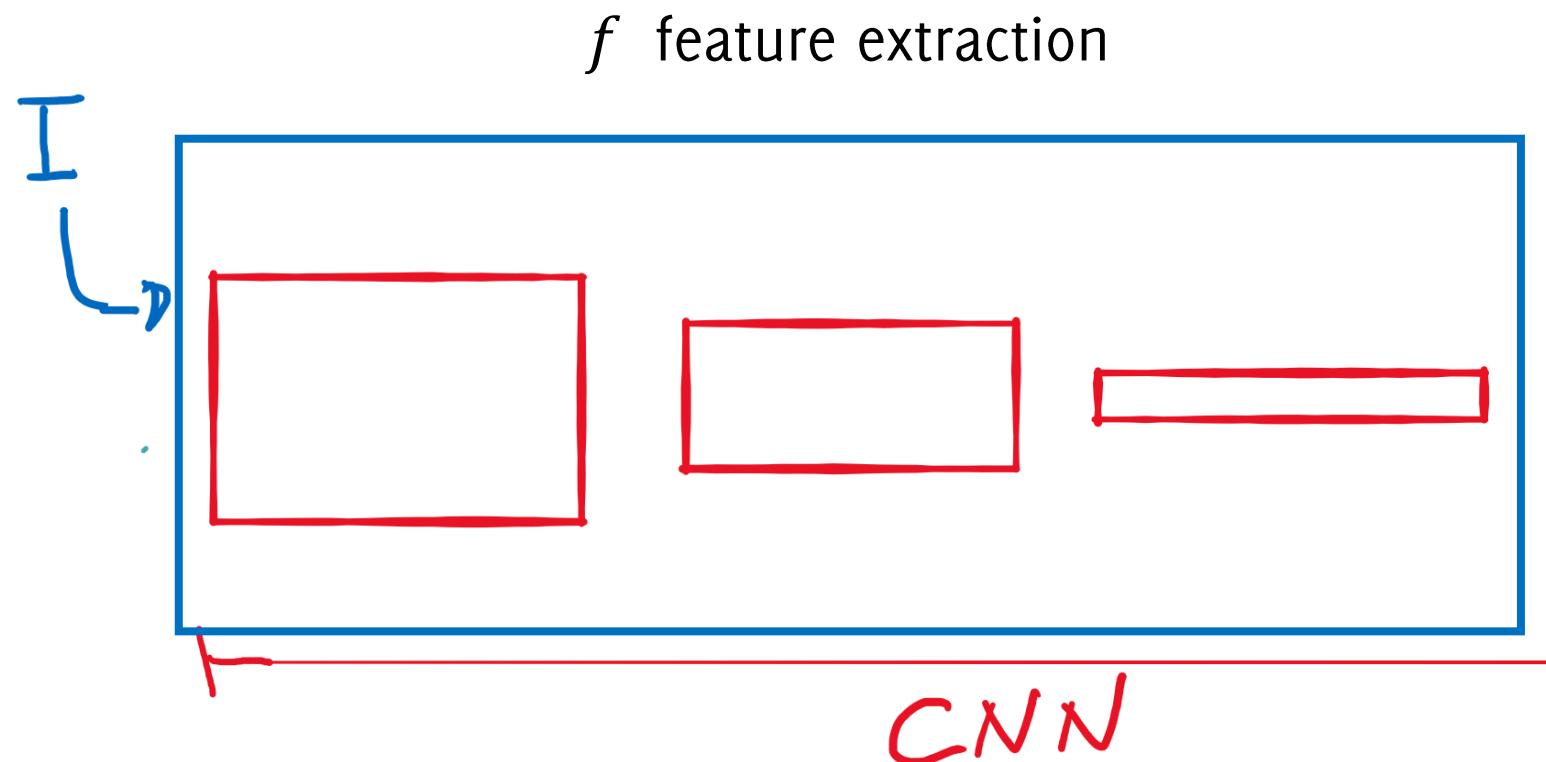


The feature extraction part is typically more general-purpose than the classification part (see Transfer Learning / Domain Adaptation)

Latent representation



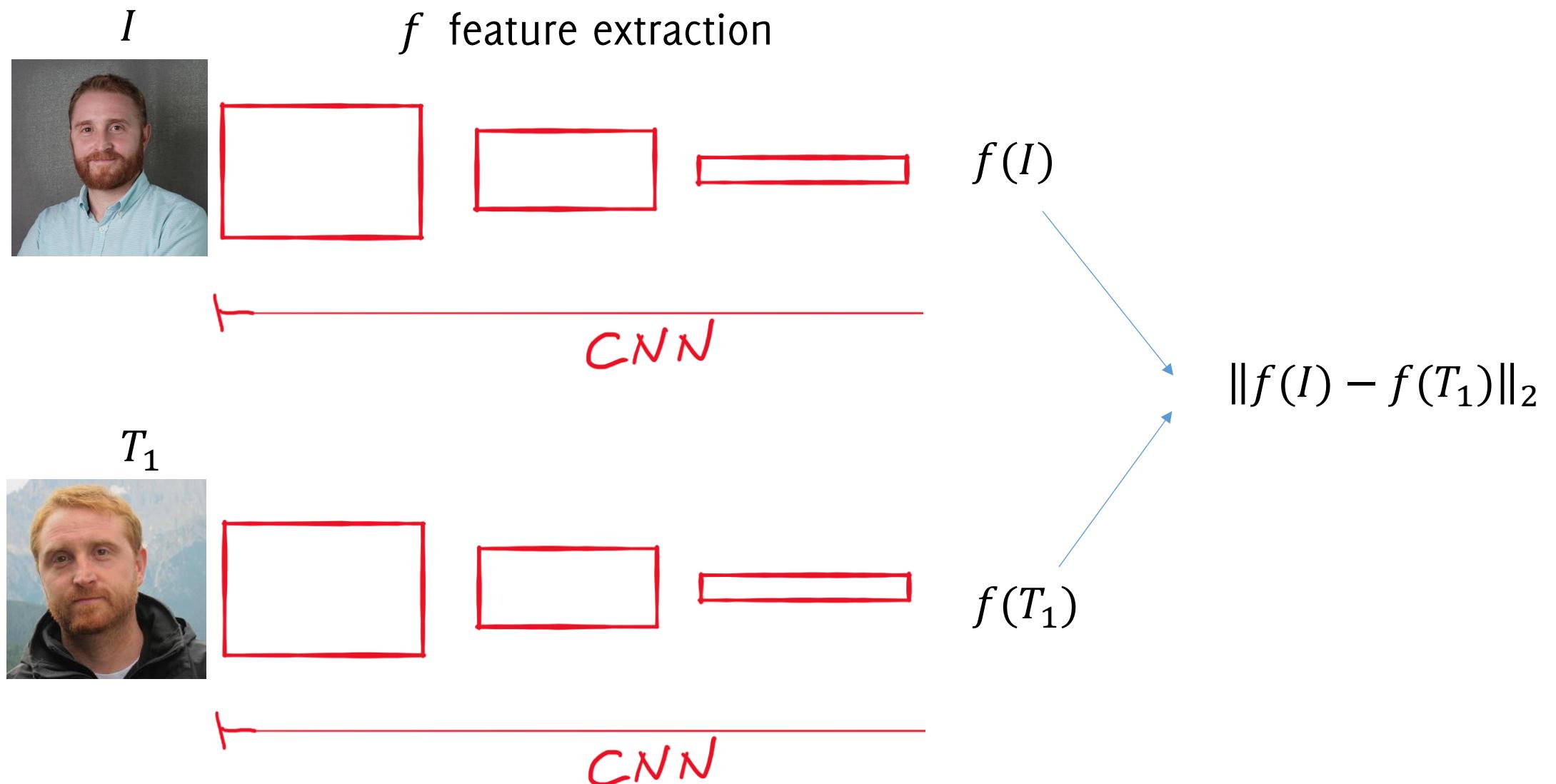
Distance among latent representations



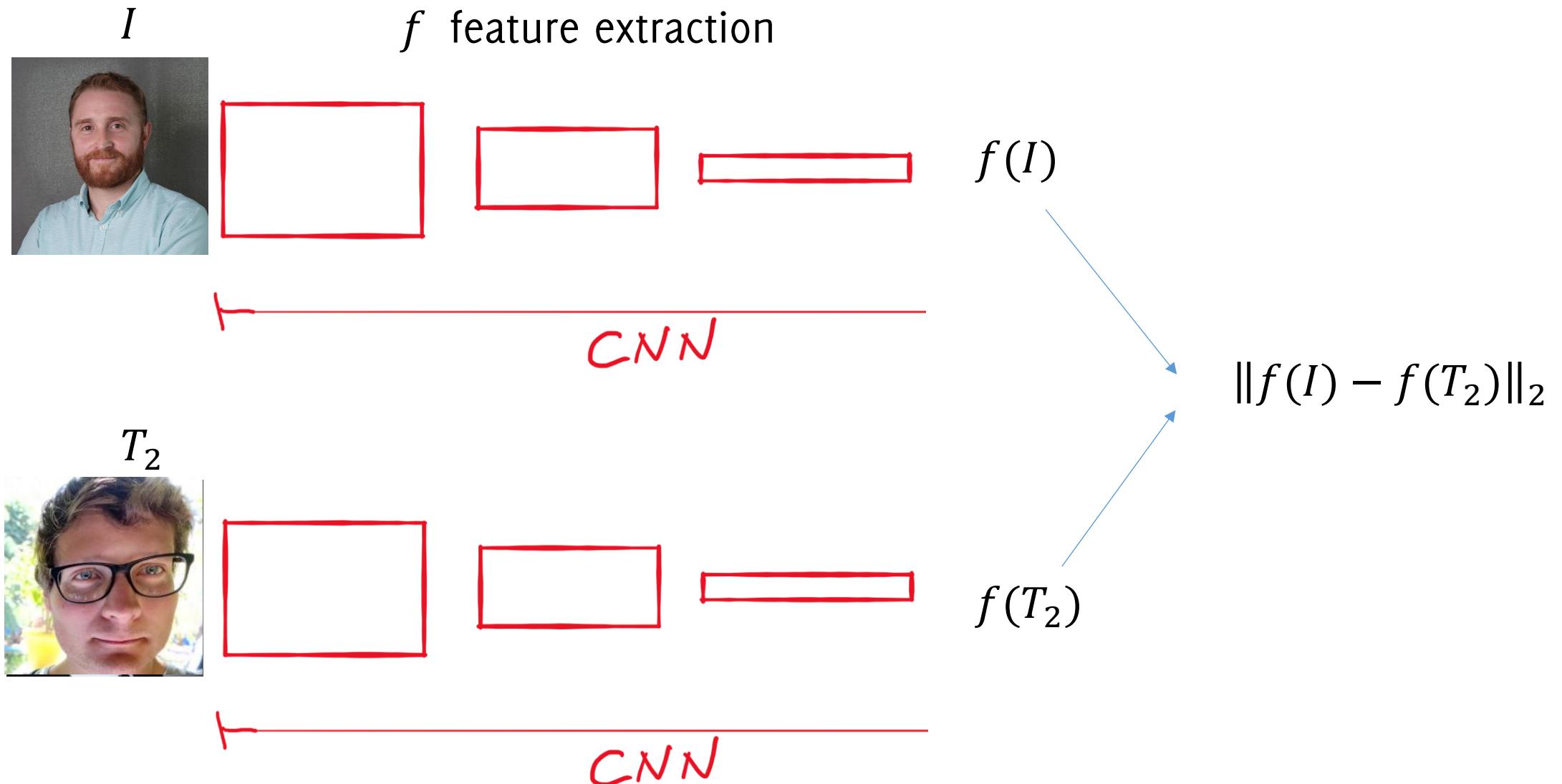
What about using?

$$\hat{i} = \operatorname{argmin}_{j=1,\dots,4} \|f(I) - f(T_j)\|_2$$

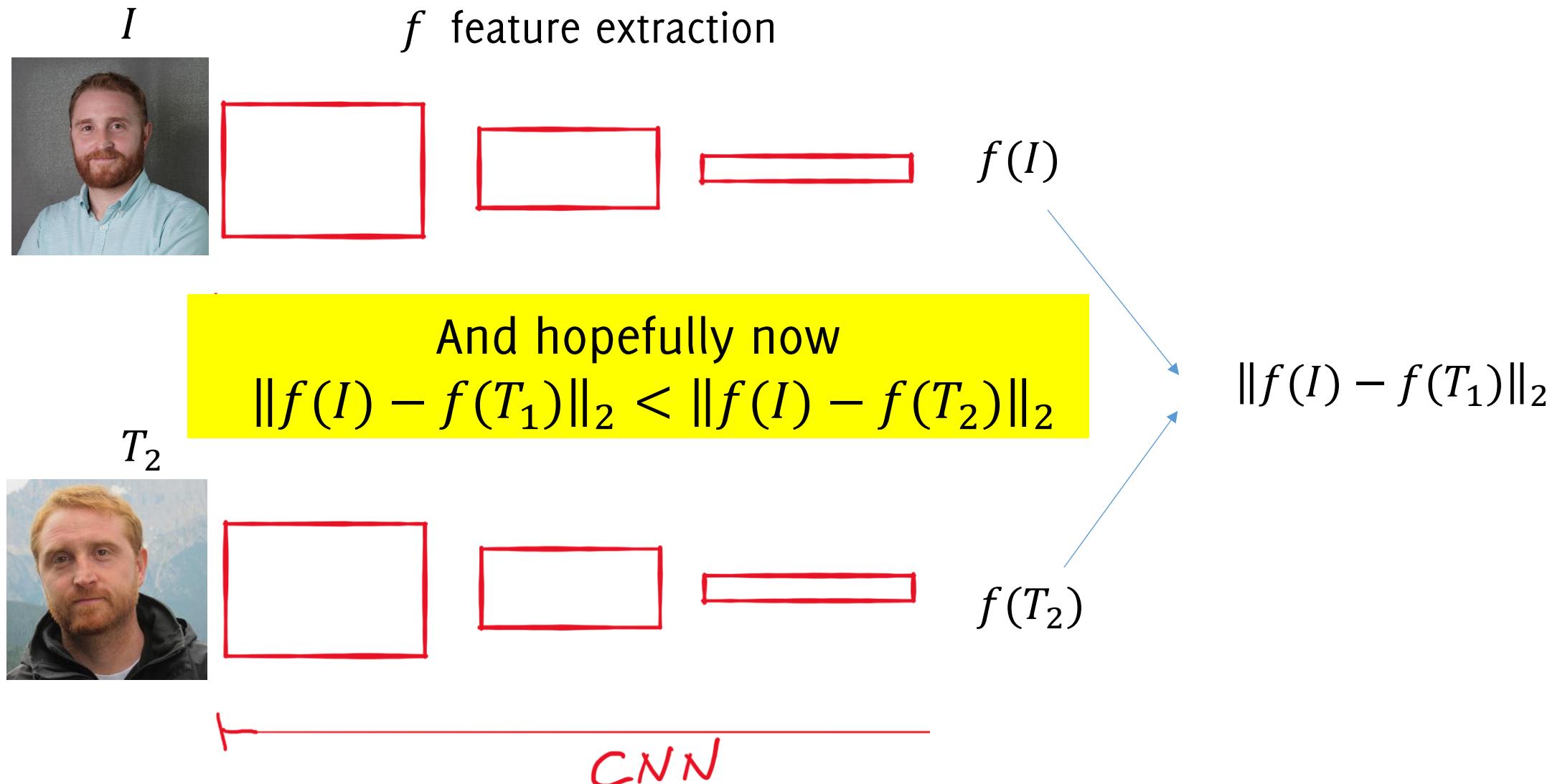
Distance among latent representations



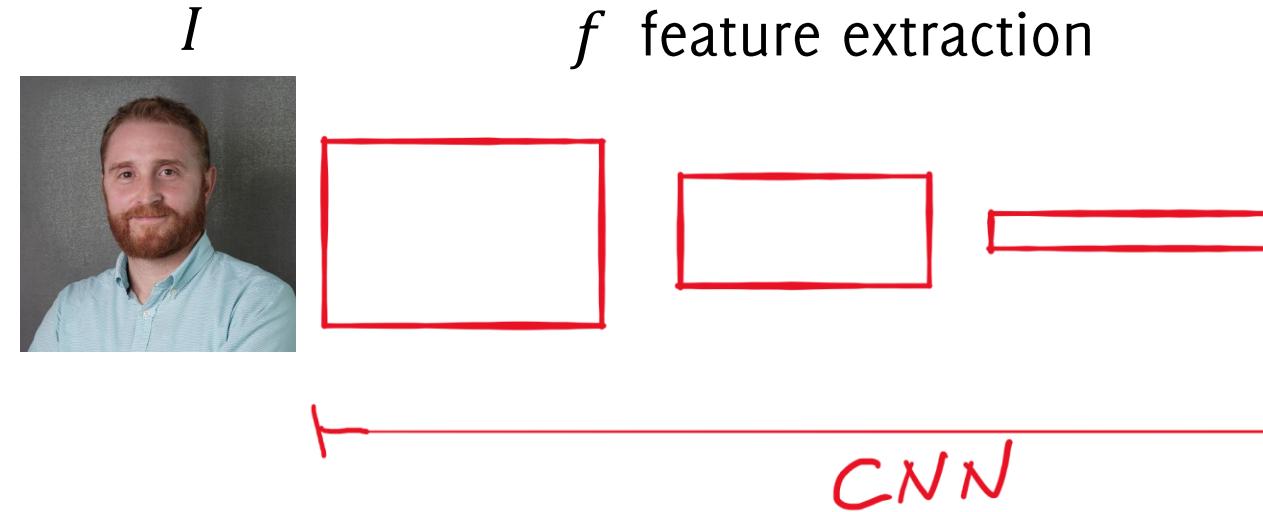
Distance among latent representations



Distance among latent representations



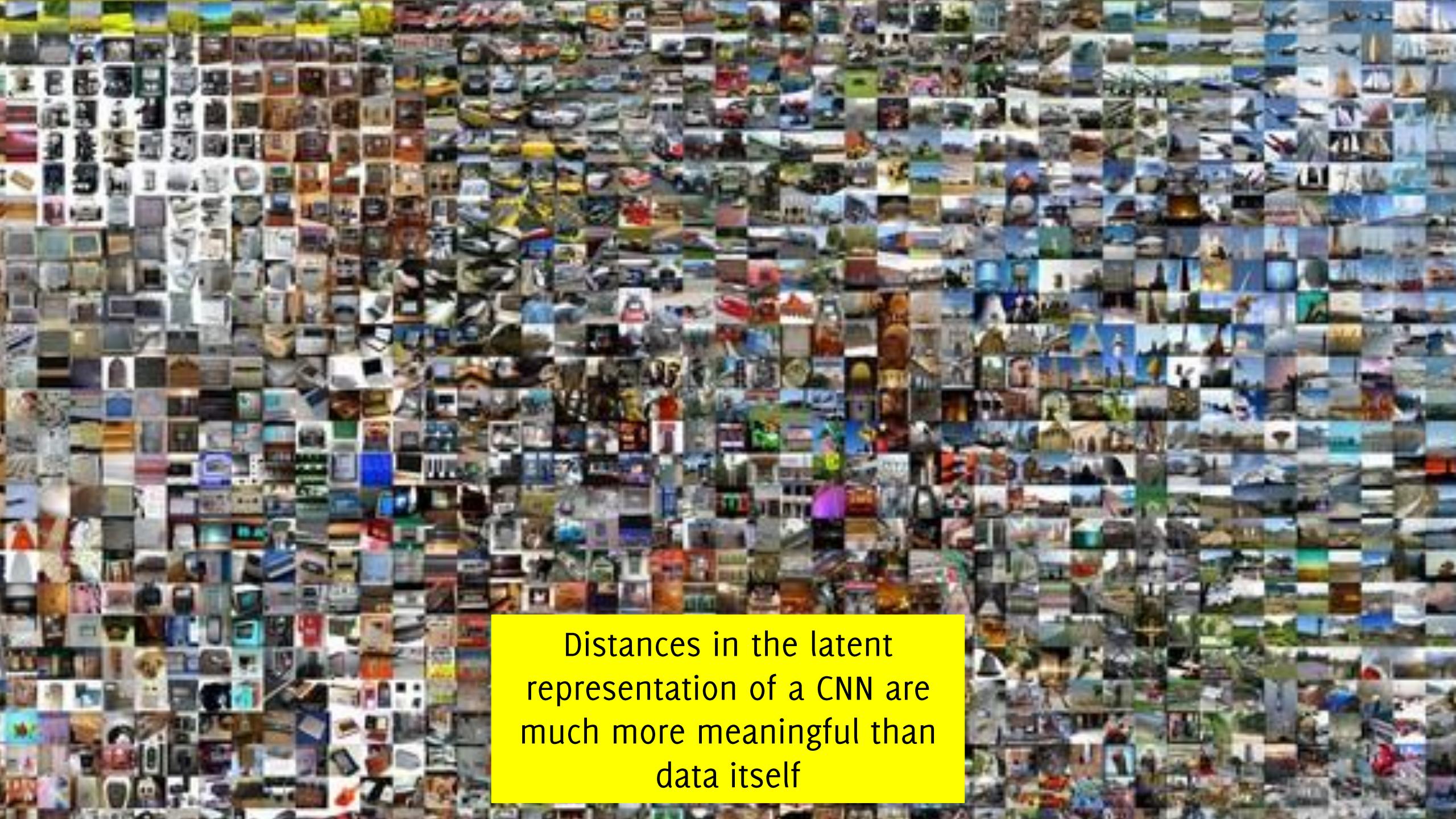
In practice



1. Extract features from the first image $f(I)$
2. Perform identification as $\hat{i} = \operatorname{argmin}_{j=1,\dots,4} \|f(I) - f(T_j)\|_2$

No need to compute $f(T_j) \forall j$, these can be directly stored

This is equivalent to perform image retrieval in the latent space... we know this can work!



Distances in the latent representation of a CNN are much more meaningful than data itself

However... can we do any better?
.... After all, the network $f(\cdot)$ was not
trained for this purpose...

Metric Learning

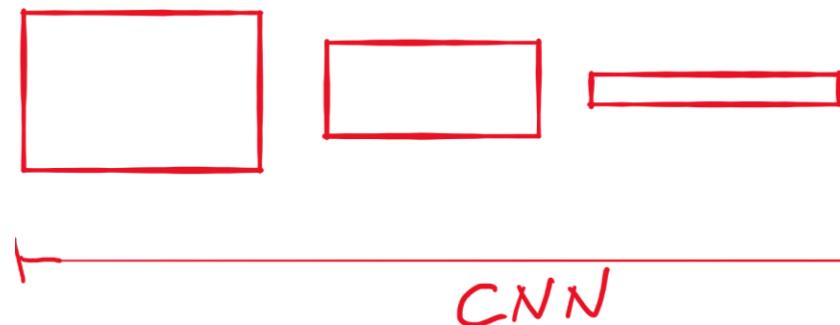
We are still comparing latent representations trained for classification (over a few persons?), while not for comparing images.

A more appealing perspective would be to train the network to measure distances between images.

We would like to train the weights W of our CNN such that

$$\|f_W(I) - f_W(T_i)\|_2 < \|f_W(I) - f_W(T_j)\|_2 \quad \forall j \neq i$$

When I belongs to class i



Siamese Networks

I



$$f_W(I)$$

T_2

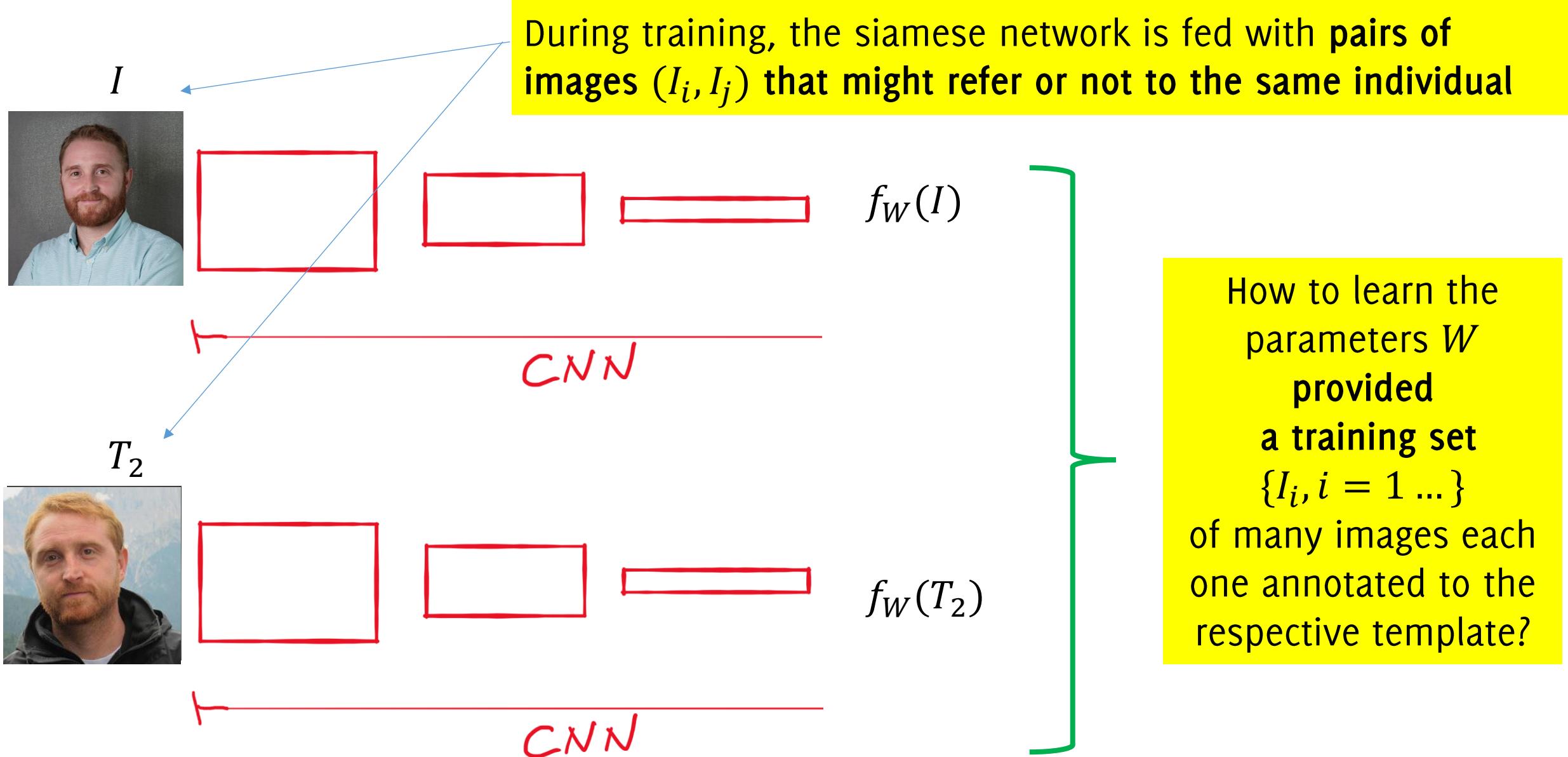


$$f_W(T_2)$$



The «two» networks have to **perform the same operations**, using the same weights W . Hence the term **Siamese**

How to Train the Siamese Network?



Contrastive Loss

The *contrastive loss* function is defined as follow:

$$W = \operatorname{argmin}_{\omega} \sum_{i,j} \mathcal{L}_{\omega}(I_i, I_j, y_{i,j})$$

Where:

$$\mathcal{L}_{\omega}(I_i, I_j, y_{i,j}) = \frac{(1 - y_{i,j})}{2} \|f_{\omega}(I_i) - f_{\omega}(I_j)\|_2 + \frac{y_{i,j}}{2} \max(0, m - \|f_{\omega}(I_i) - f_{\omega}(I_j)\|_2)$$

- $y_{i,j} \in \{0,1\}$ is the label associate with the input pair (I_i, I_j) :
 - 0 when (I_i, I_j) refers to the same person
 - 1 otherwise
- m is a hyperparameter indicating the margin we want (like in Hinge Loss)
- $\|f_{\omega}(I_i) - f_{\omega}(I_j)\|_2$ is the distance in the latent space.

Triplet Loss

A loss function such that a training sample I is compared against

- P a positive input, referring to the same person
- N a negative input, referring to a different person

We train the network to minimize the distance from the positive samples and maximize the distance from the negative ones

$$\mathcal{L}_\omega(I, P, N) = \max(0, m + (\|f_\omega(I) - f_\omega(P)\|_2 - \|f_\omega(I) - f_\omega(N)\|_2))$$

Triplet loss forces that a pair of samples from the same individual are smaller in distance than those with different ones.

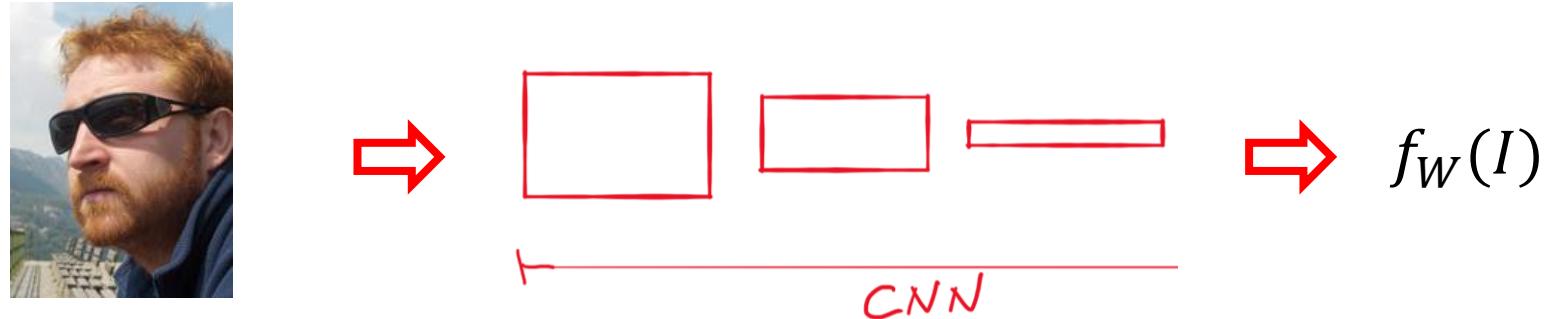
m always play the role of the margin.

The selection of triplets for training is a matter of study

Just to Recap

When a new image has to be verified

1. We feed the image I to the trained network, thus compute $f_W(I)$



2. Identify the person having average minimum distance from templates (in case there are many associated to the same individual)

$$\hat{u} = \operatorname{argmin}_u \frac{\sum_{T_{u,j}} \|f_W(I) - f_W(T_{u,j})\|_2}{\#\{T_u\}}$$

3. Assess whether

$$\frac{\sum_{T_{i,j}} \|f_W(I) - f_W(T_{i,j})\|_2}{\#\{T_{i,j}\}} < \gamma$$

is sufficiently small, otherwise **no identification**

Other decision rules can be adopted (e.g. searching for the person giving the a template with a minimum distance)