

Beyond classification: Object detection, Segmentation, Human pose estimation

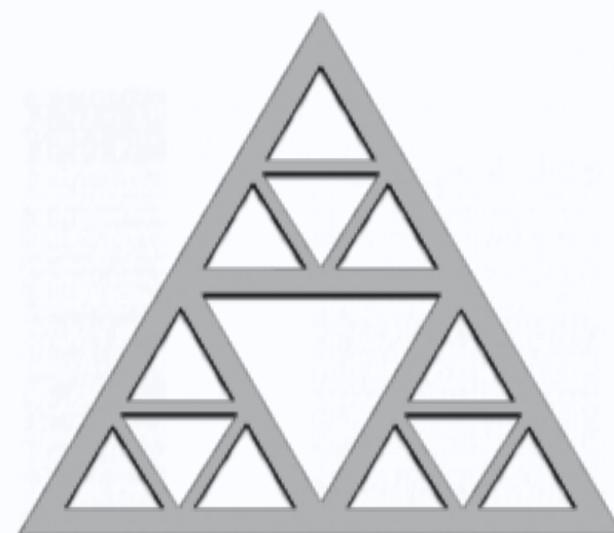
Gül Varol

IMAGINE team, École des Ponts ParisTech

gul.varol@enpc.fr

<http://imagine.enpc.fr/~varolg/>

@RecVis, 12.11.2024



École des Ponts
ParisTech

Announcements

- A2 due today.
- FP proposal due next week.
- A3 due Nov 26.

Neural Networks

Previous week: Introduction to neural networks

Last week: Neural networks for **visual recognition**:

Convolutional Neural Networks (CNNs) for image classification

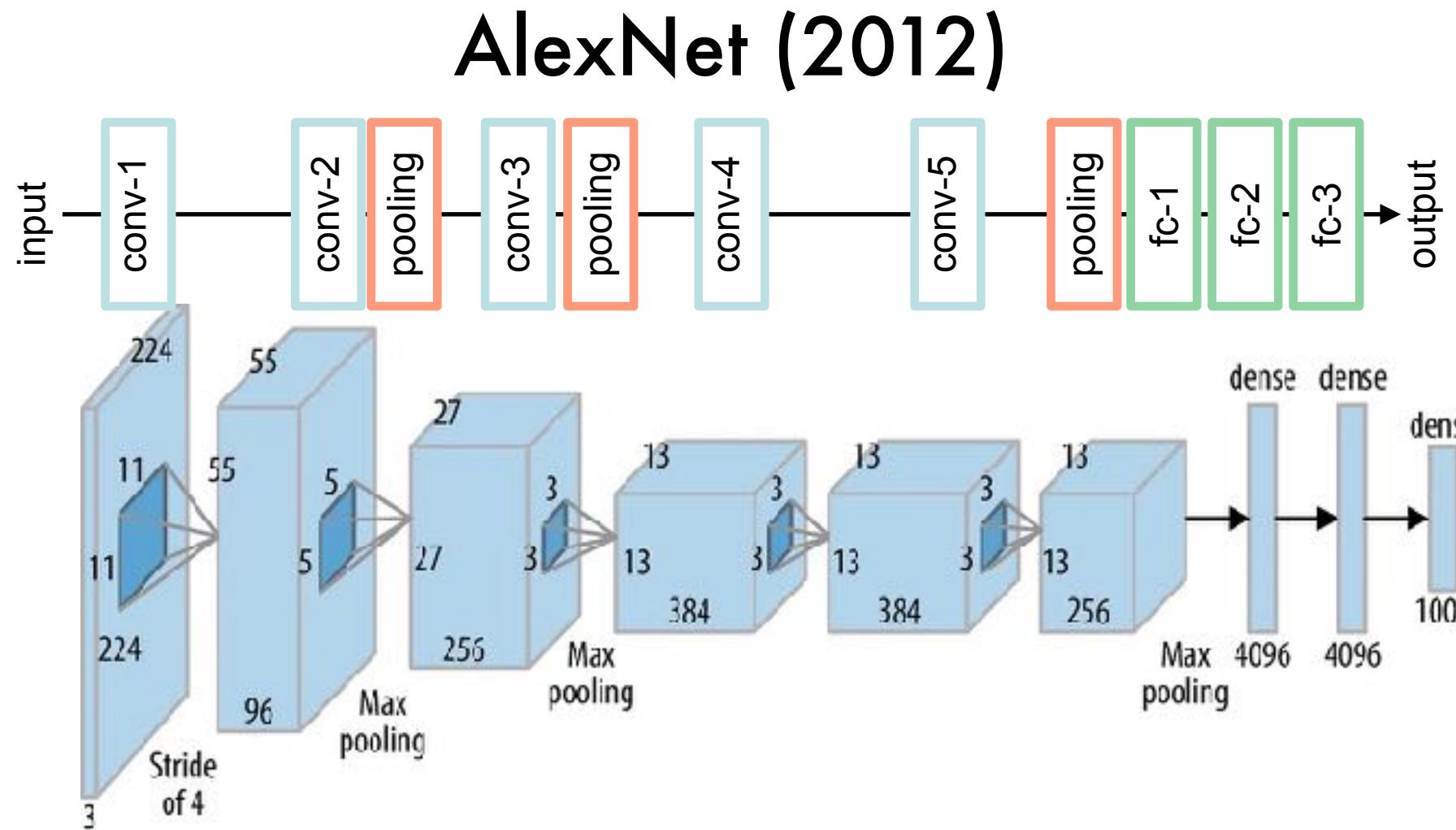
This week: Beyond CNNs: Transformers

Beyond classification: Object detection, Segmentation, Human pose;

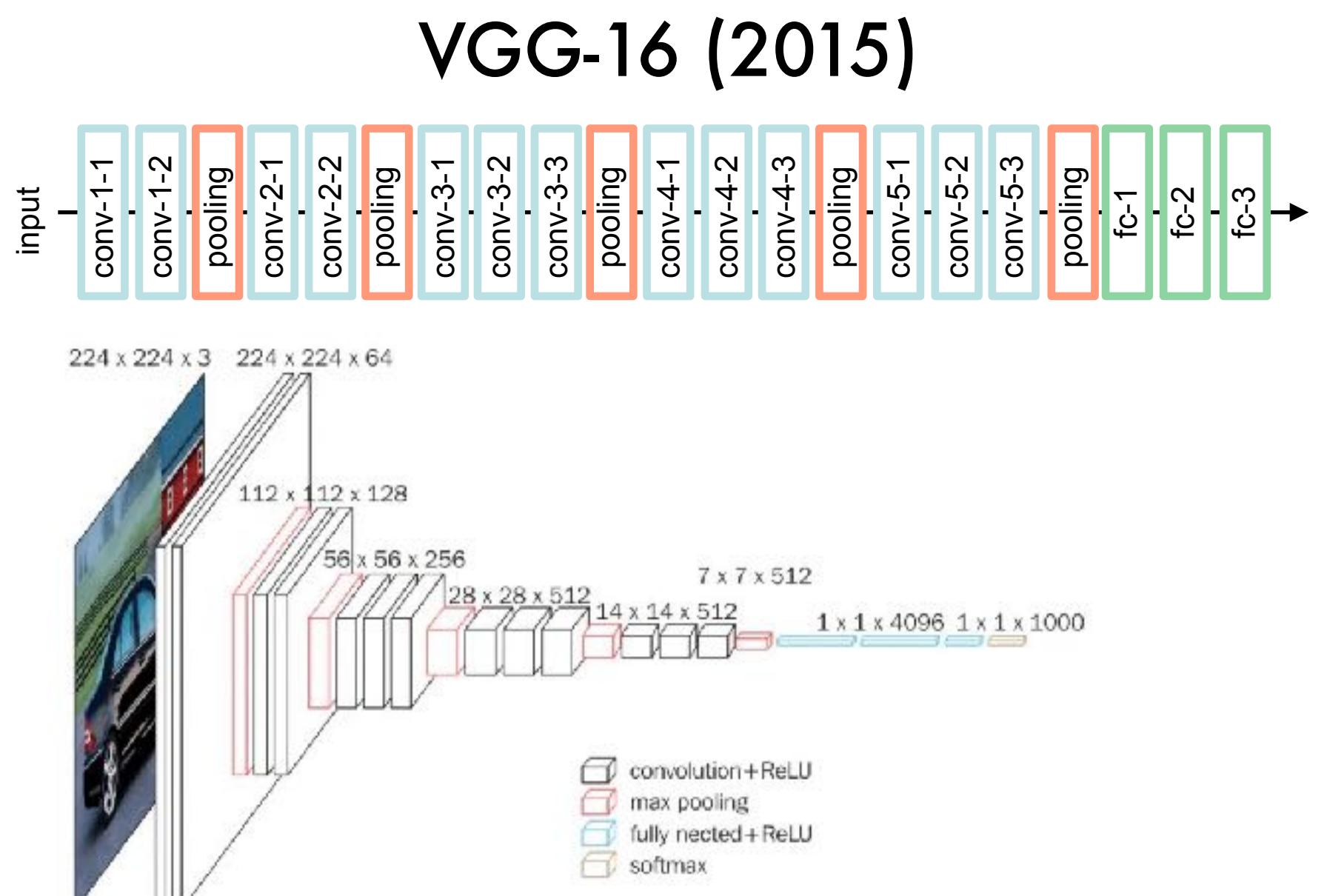
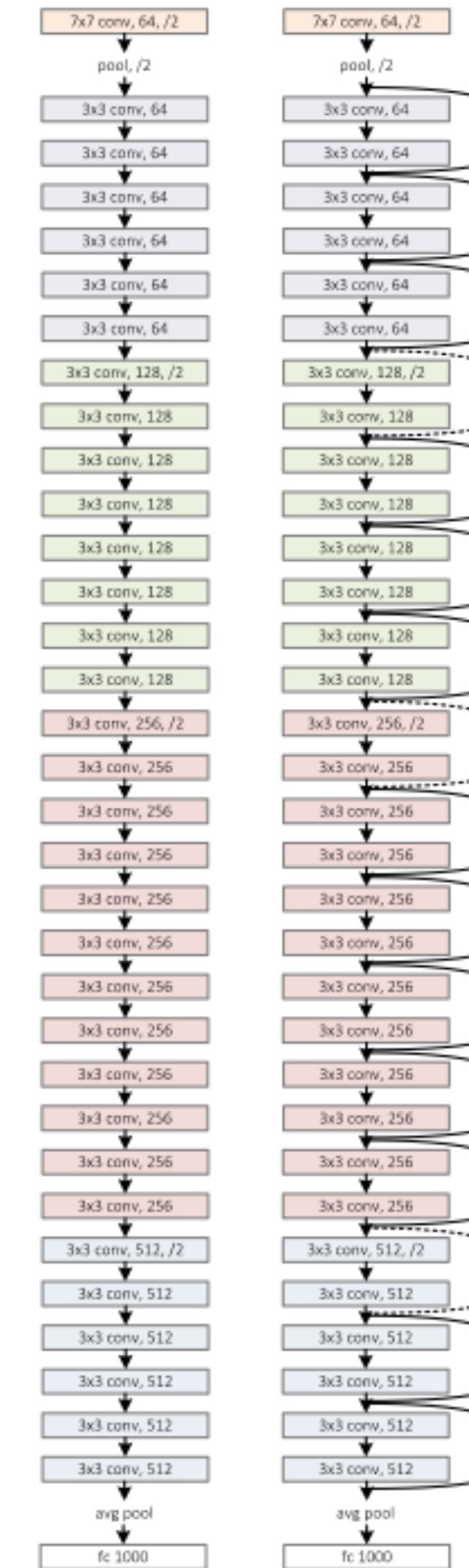
Next week: Generative models; Vision & language

Recap: Neural networks for image classification

Image source: oreilly.com



ResNet (2016)

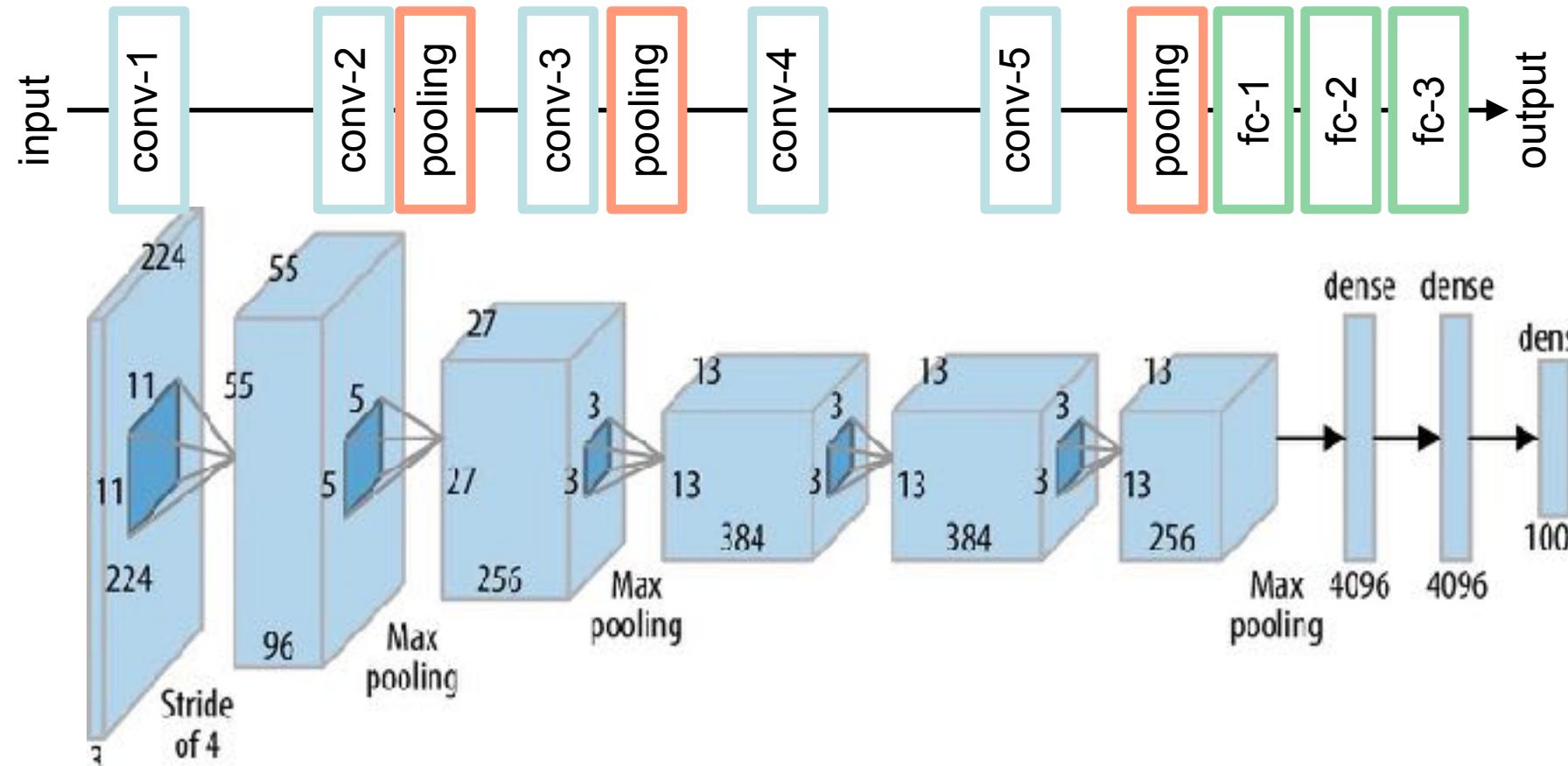


VGG-19

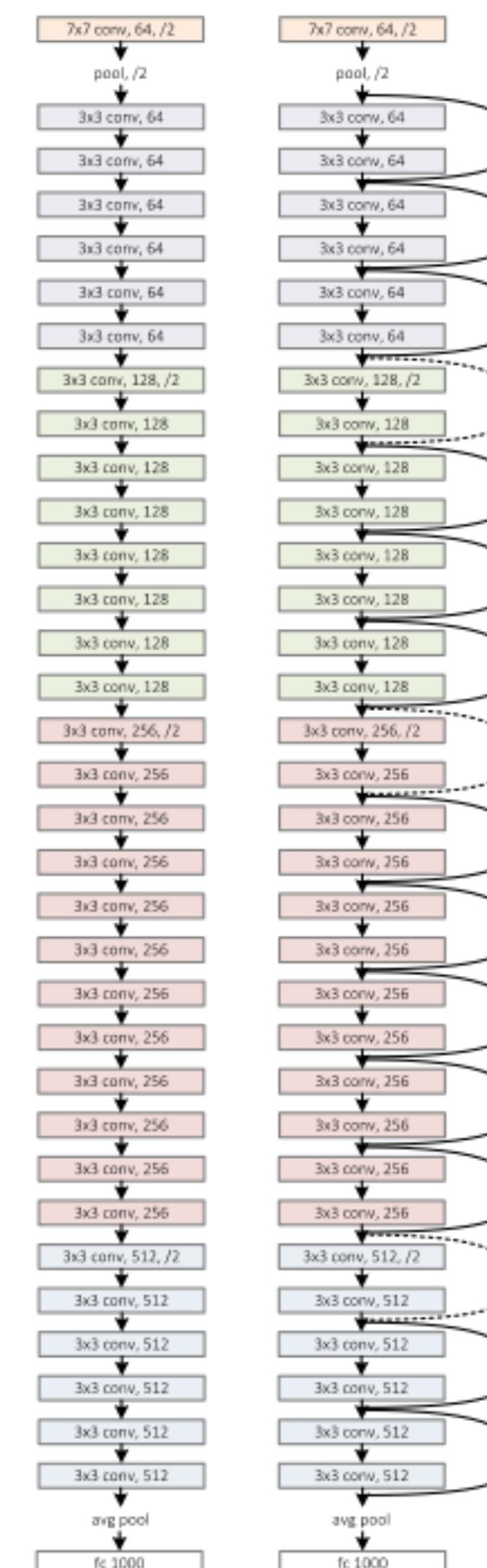
ResNet-34

Recap: Neural networks for image classification

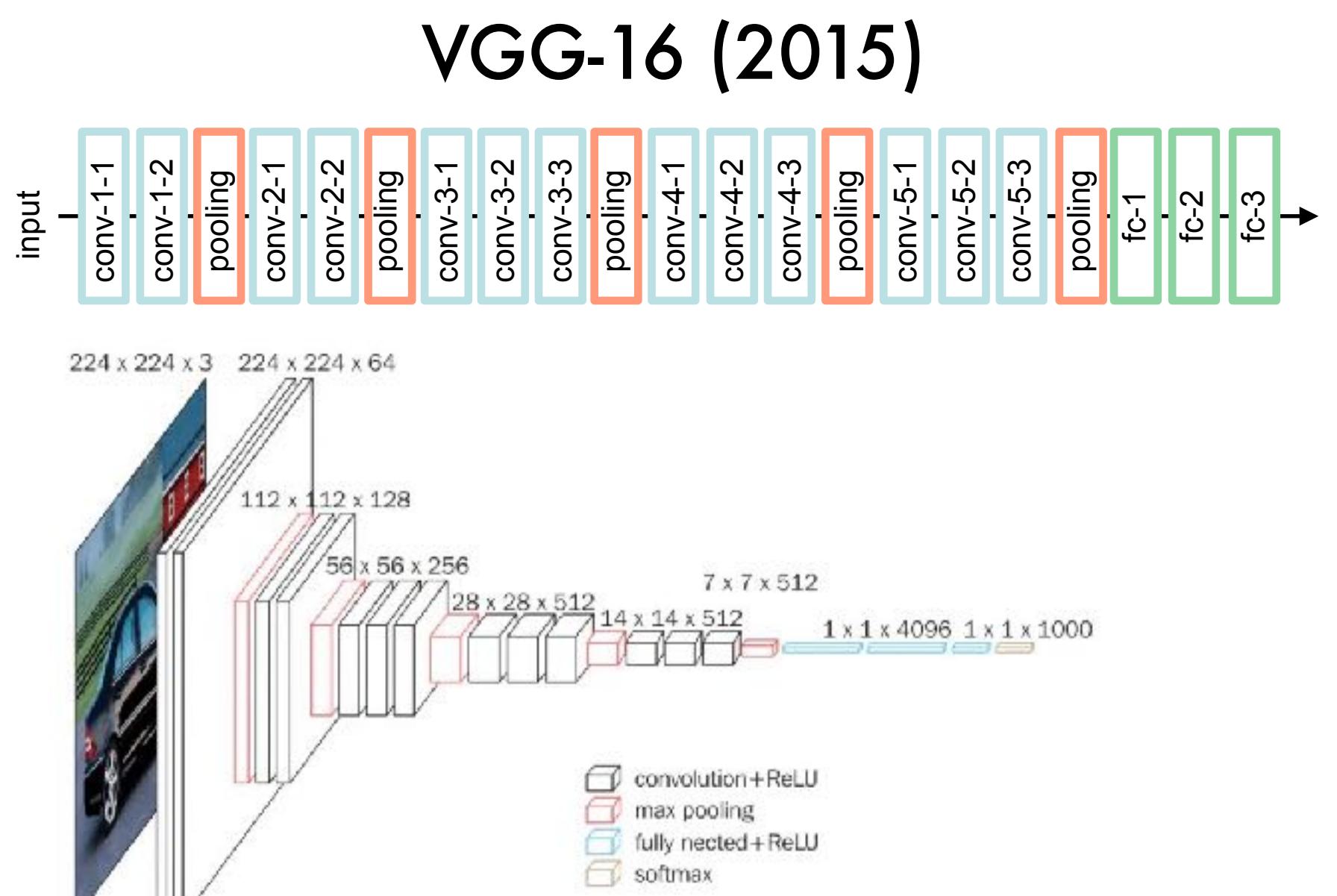
AlexNet (2012)



ResNet (2016)

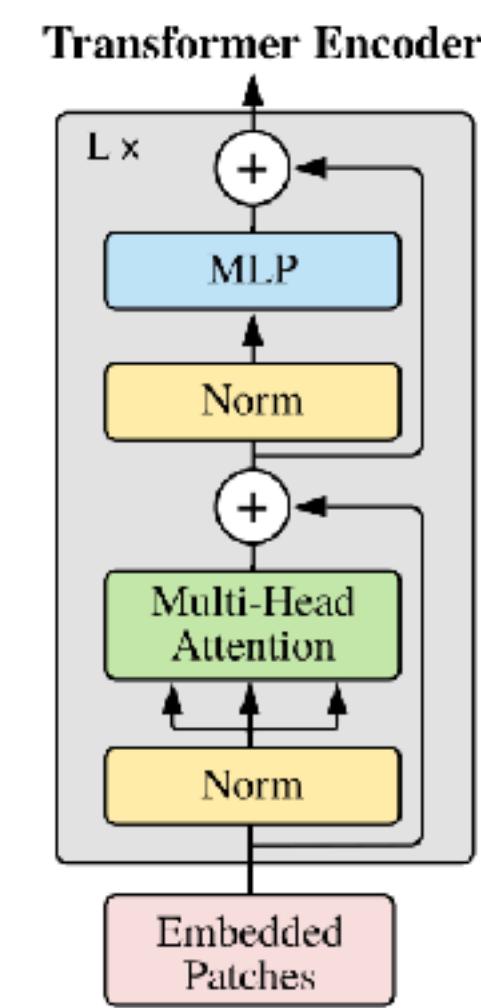
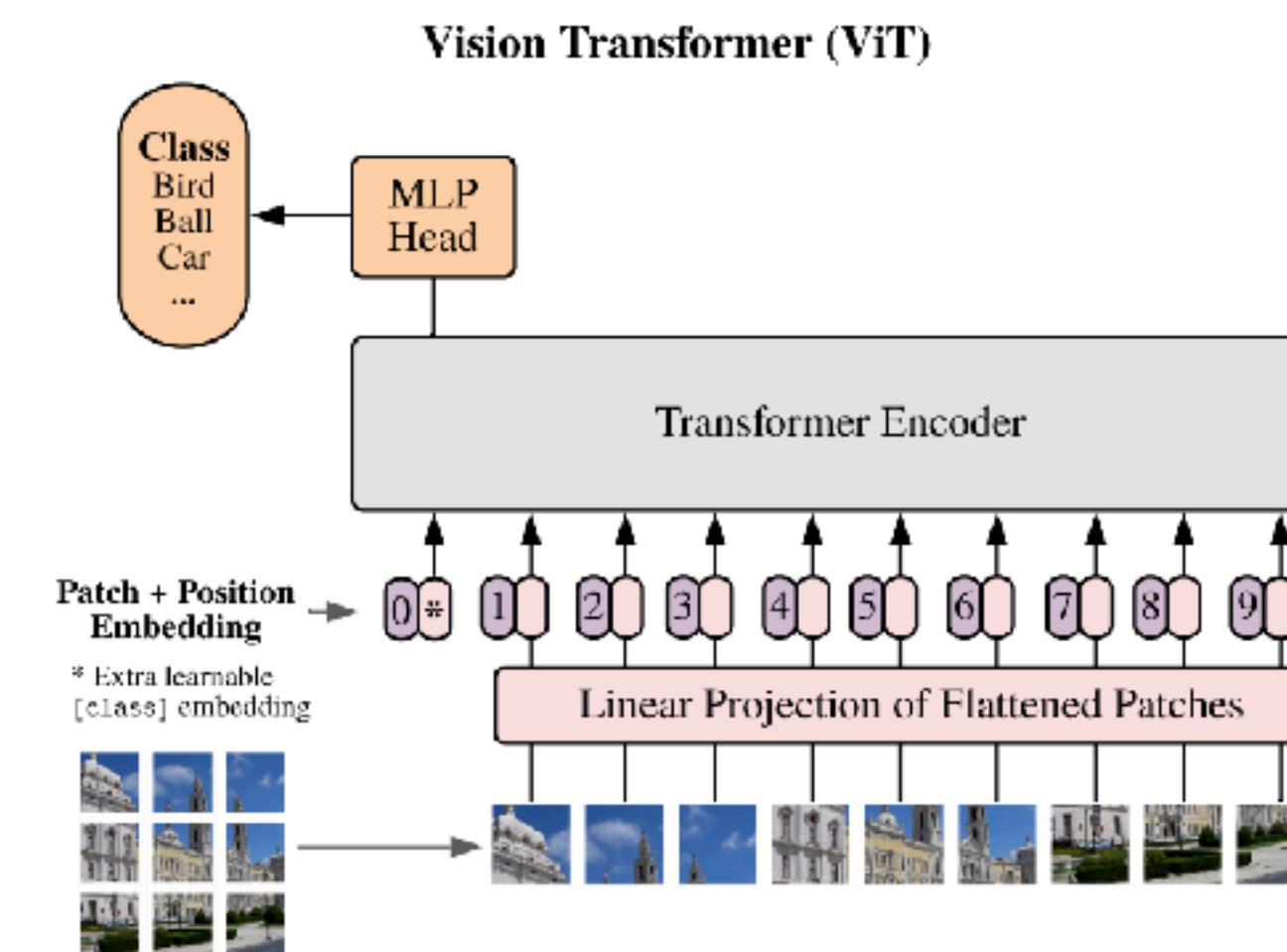


VGG-19



ResNet-34

ViT (2021)



1. Beyond CNNs

- **Attention & Transformer**
- **Vision Transformers**

2. Beyond classification

- **a. Intro to structured outputs**
- **b. Object detection (localization)**
- **c. Segmentation**
- **d. Human pose estimation**

1. Beyond CNNs

- **Attention & Transformer**
- **Vision Transformers**

1. Beyond CNNs

- **Attention & Transformer**
- **Vision Transformers**

Do we need convolutions?

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Published as a conference paper at ICLR 2021

AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy*,†, Lucas Beyer*, Alexander Kolesnikov*, Dirk Weissenborn*,
Xiaohua Zhai*, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby*,†

*equal technical contribution, †equal advising
Google Research, Brain Team
{adosovitskiy, neilhoulsby}@google.com

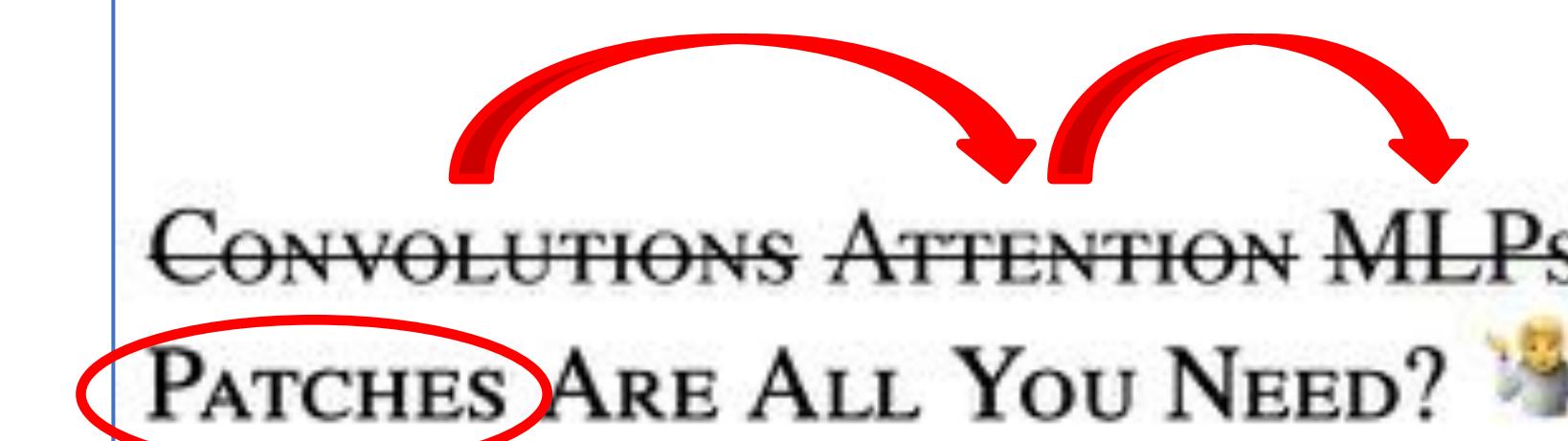
ICLR 2021

MLP-Mixer: An all-MLP Architecture for Vision

Ilya Tolstikhin*, Neil Houlsby*, Alexander Kolesnikov*, Lucas Beyer*,
Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner,
Daniel Keysers, Jakob Uszkoreit, Mario Lucic, Alexey Dosovitskiy

*equal contribution
Google Research, Brain Team
{tolstikhin, neilhoulsby, akolesnikov, lbeyer,
xzhai, unterthiner, jessicayung†, andstein,
keysers, usz, lucic, adosovitskiy}@google.com

Under review as a conference paper at ICLR 2022



CONVOLUTIONS ATTENTION MLPs PATCHES ARE ALL YOU NEED? 🧩

Anonymous authors
Paper under double-blind review

arXiv 2022

<https://github.com/KentoNishi/awesome-all-you-need-papers>

Recent Hype#1: Transformers

- Transformers = neural network architectures stacking "attention" layers¹
- Initially successful for natural language processing
- Then applied to computer vision². Better performance than CNNs given enough data.
- The hype still continues today.
- What is attention?

¹ Vaswani et al. "Attention is all you need", NeurIPS 2017.

² Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale", ICLR 2021.

Attention

- Motivation: sequence-to-sequence models

$$\mathbf{Y} = \text{AttentionLayer}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$$

$$\mathbf{Y} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right) \mathbf{V}$$

$$\mathbf{Y} = \mathbf{A}\mathbf{V}$$

$$Y_i = \sum_j A_{i,j} V_j$$

$$A_{i,j} = \text{softmax} \left(\frac{Q_i \cdot K_j}{\sqrt{d}} \right)$$

Attention

- Used to capture **context** *within the sequence*

The animal
didn't cross
the street because **it** was too tired .

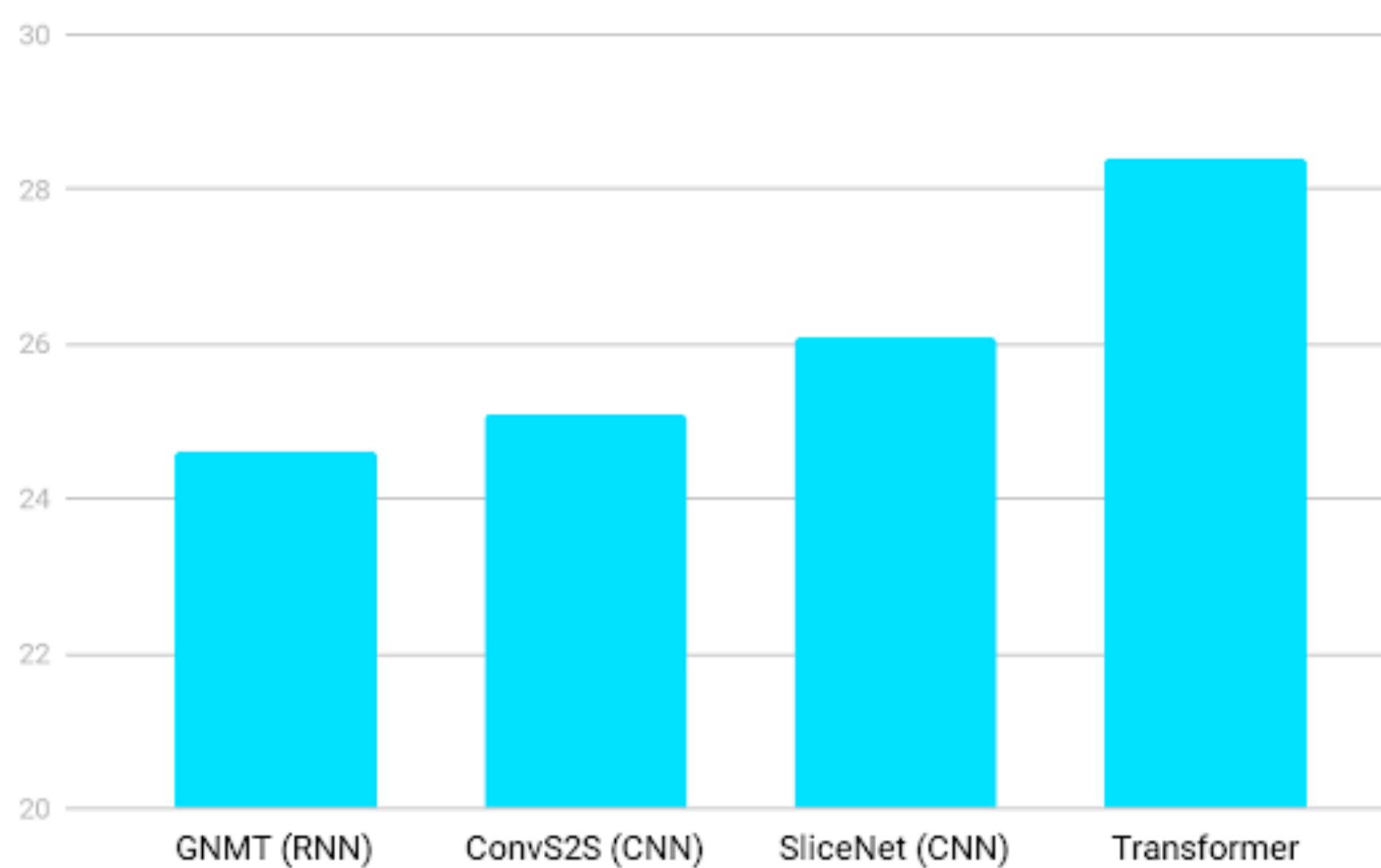
The animal
didn't cross
the street because **it** was too wide .

As we are encoding “it”, we should focus on “the animal”

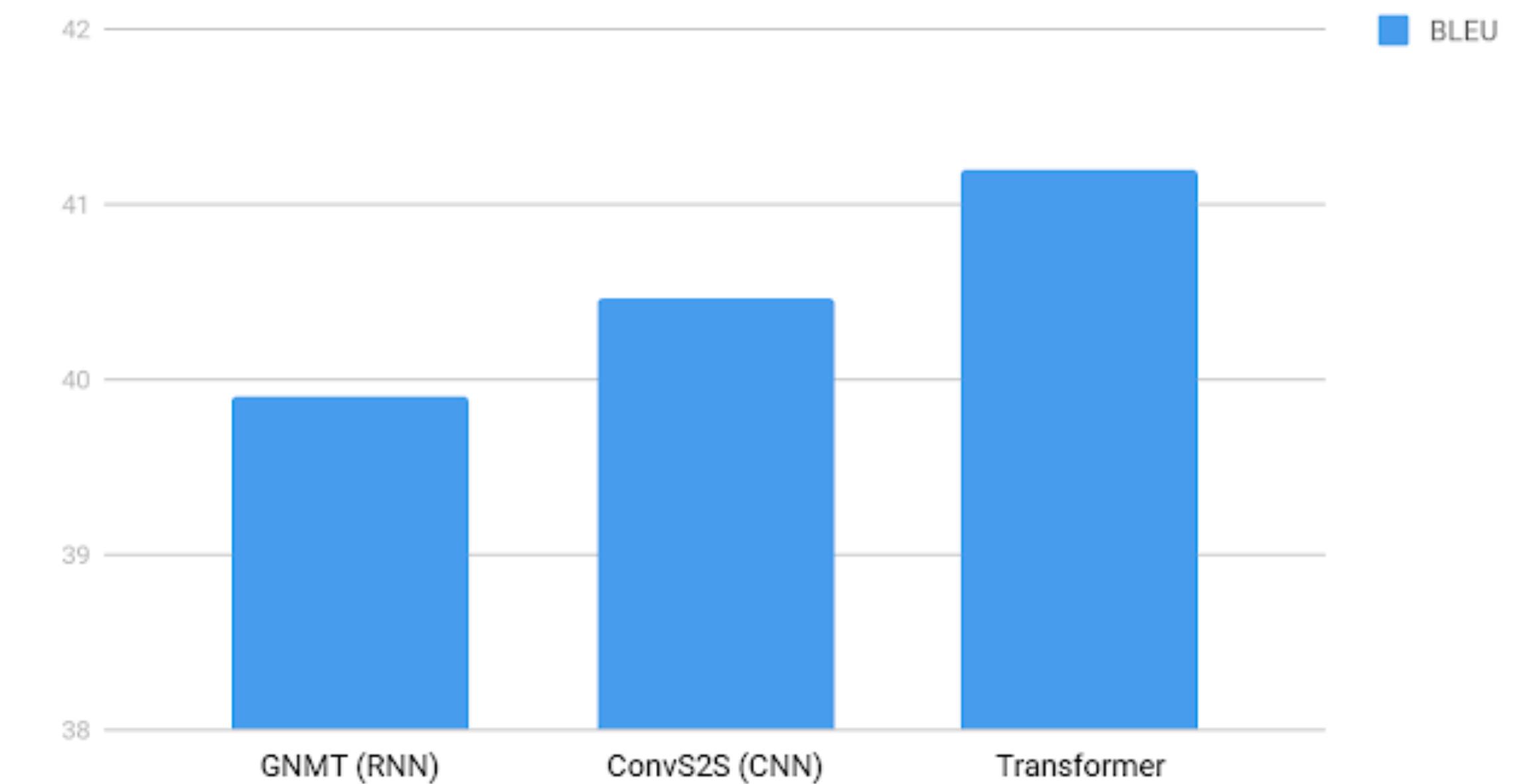
As we are encoding “it”, we should focus on “the street”

Original transformer results on machine translation

English German Translation quality



English French Translation Quality



<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

Attention

$$\mathbf{Y} = \text{AttentionLayer}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$$

$$\mathbf{Y} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

$$\mathbf{Y} = \mathbf{A}\mathbf{V}$$

$$Y_i = \sum_j A_{i,j} V_j$$
$$A_{i,j} = \text{softmax}\left(\frac{Q_i \cdot K_j}{\sqrt{d}}\right)$$

Attention

Query Key Value

$$Y = \text{AttentionLayer}(Q, K, V)$$

$$Y = \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V$$

?

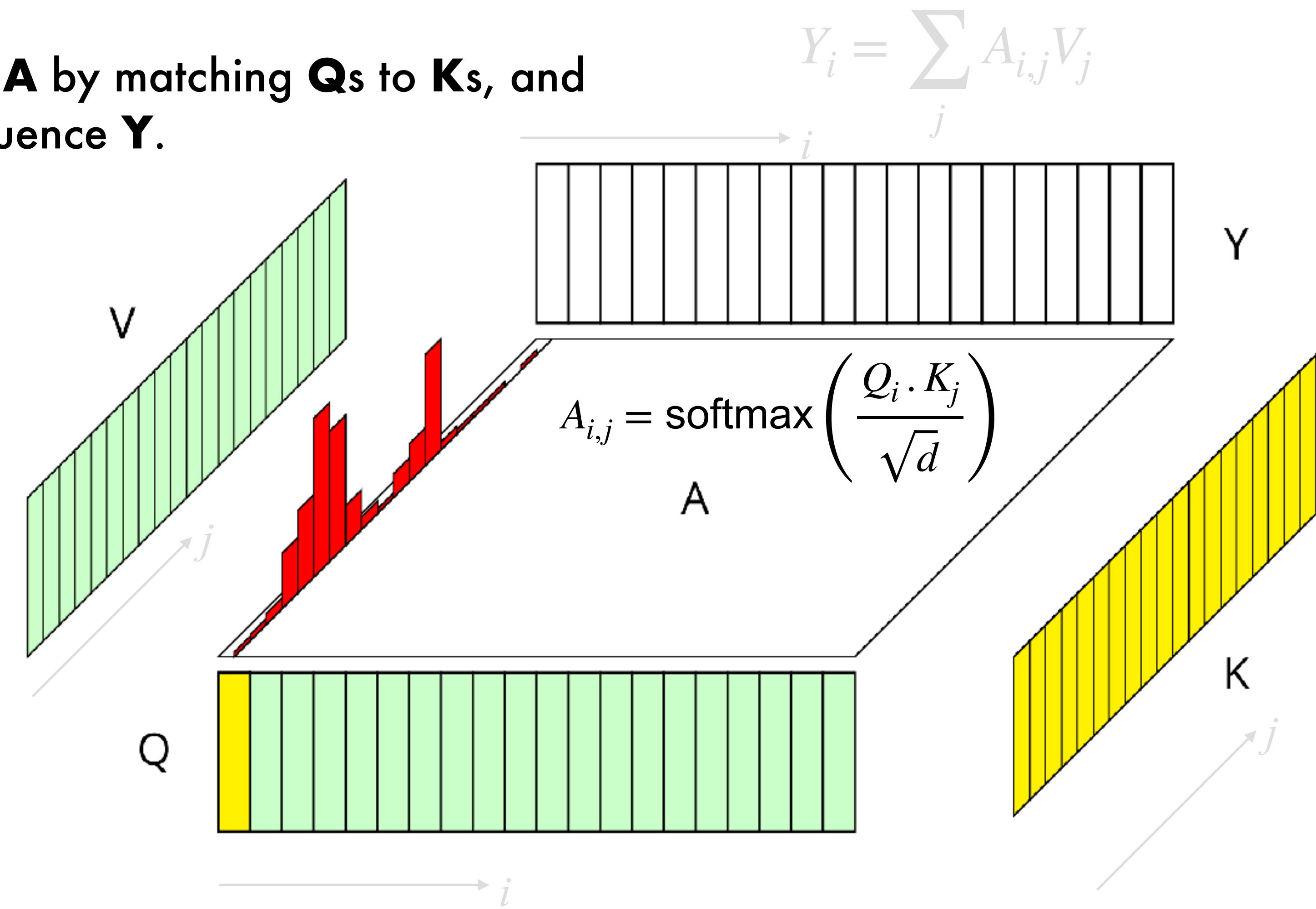
Output
Attention

$$Y = A V$$

$$Y_i = \sum_j A_{i,j} V_j$$
$$A_{i,j} = \text{softmax} \left(\frac{Q_i \cdot K_j}{\sqrt{d}} \right)$$

Attention mechanism

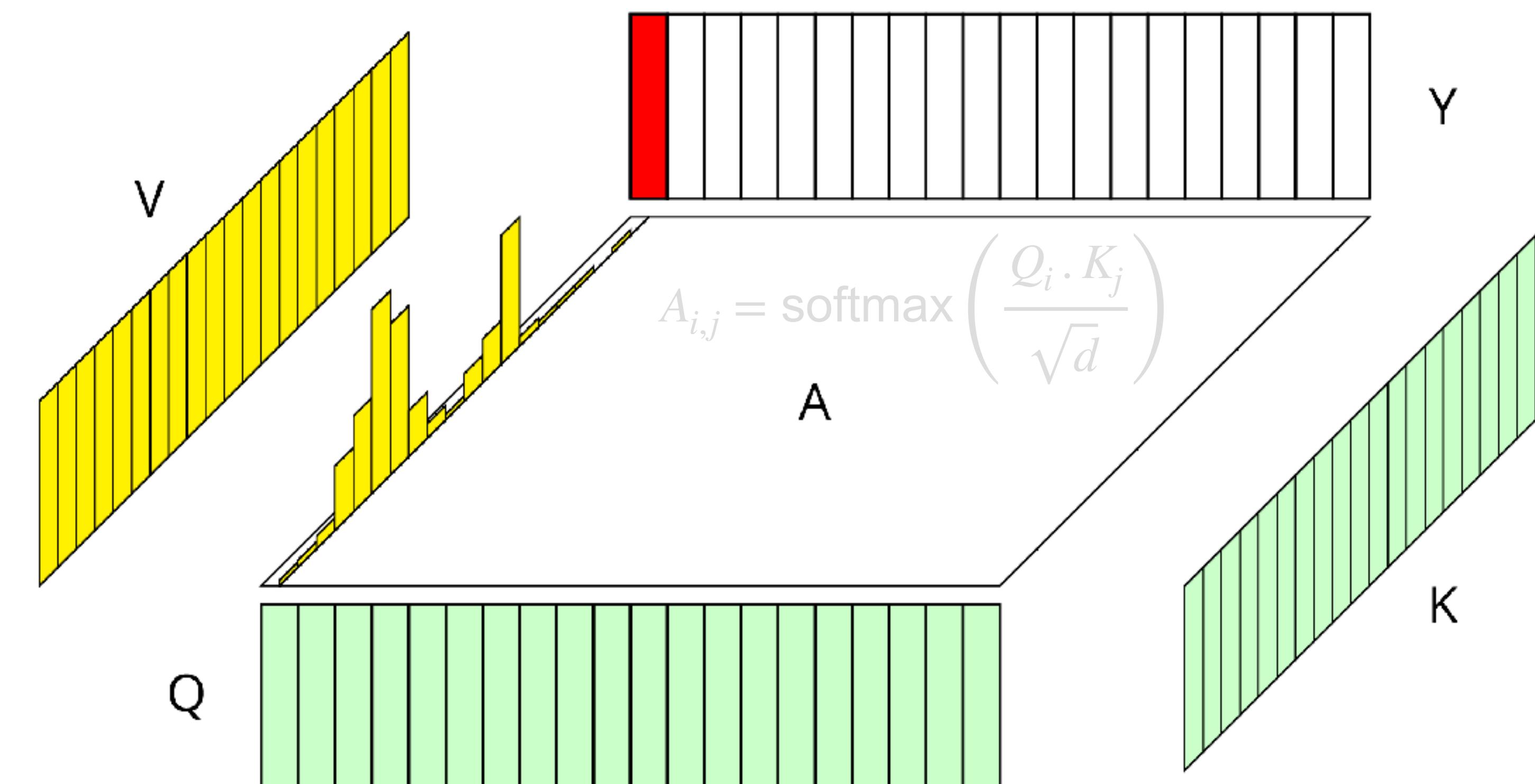
Given a **query** sequence **Q**,
a **key** sequence **K**, and
a **value** sequence **V**,
compute an attention matrix **A** by matching **Qs** to **Ks**, and
weight **V** with to get the sequence **Y**.



Attention mechanism

Given a **query** sequence **Q**,
a **key** sequence **K**, and
a **value** sequence **V**,
compute an attention matrix **A** by matching **Qs** to **Ks**, and
weight **V** with to get the sequence **Y**.

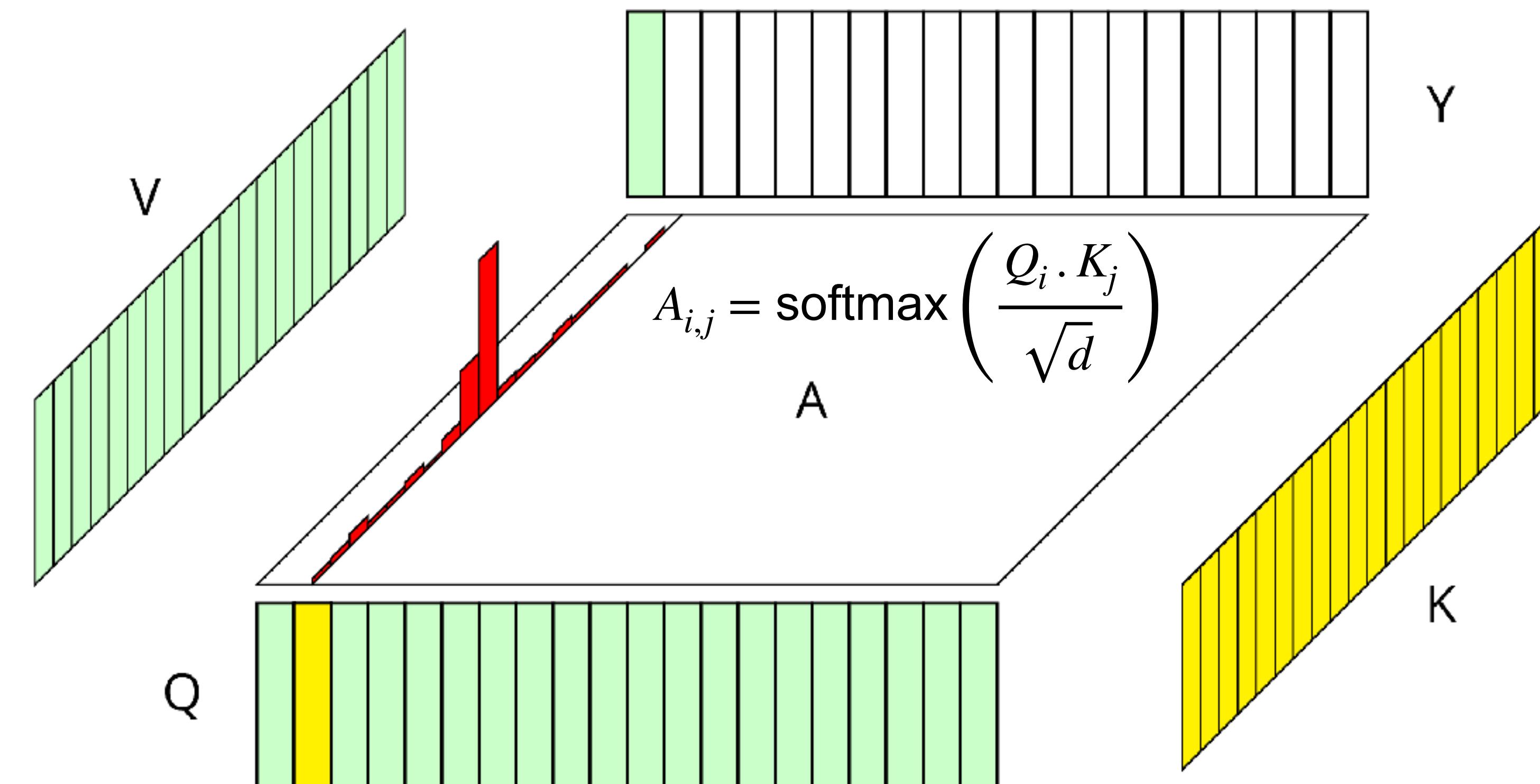
$$Y_i = \sum_j A_{i,j} V_j$$



Attention mechanism

Given a **query** sequence **Q**,
a **key** sequence **K**, and
a **value** sequence **V**,
compute an attention matrix **A** by matching **Qs** to **Ks**, and
weight **V** with to get the sequence **Y**.

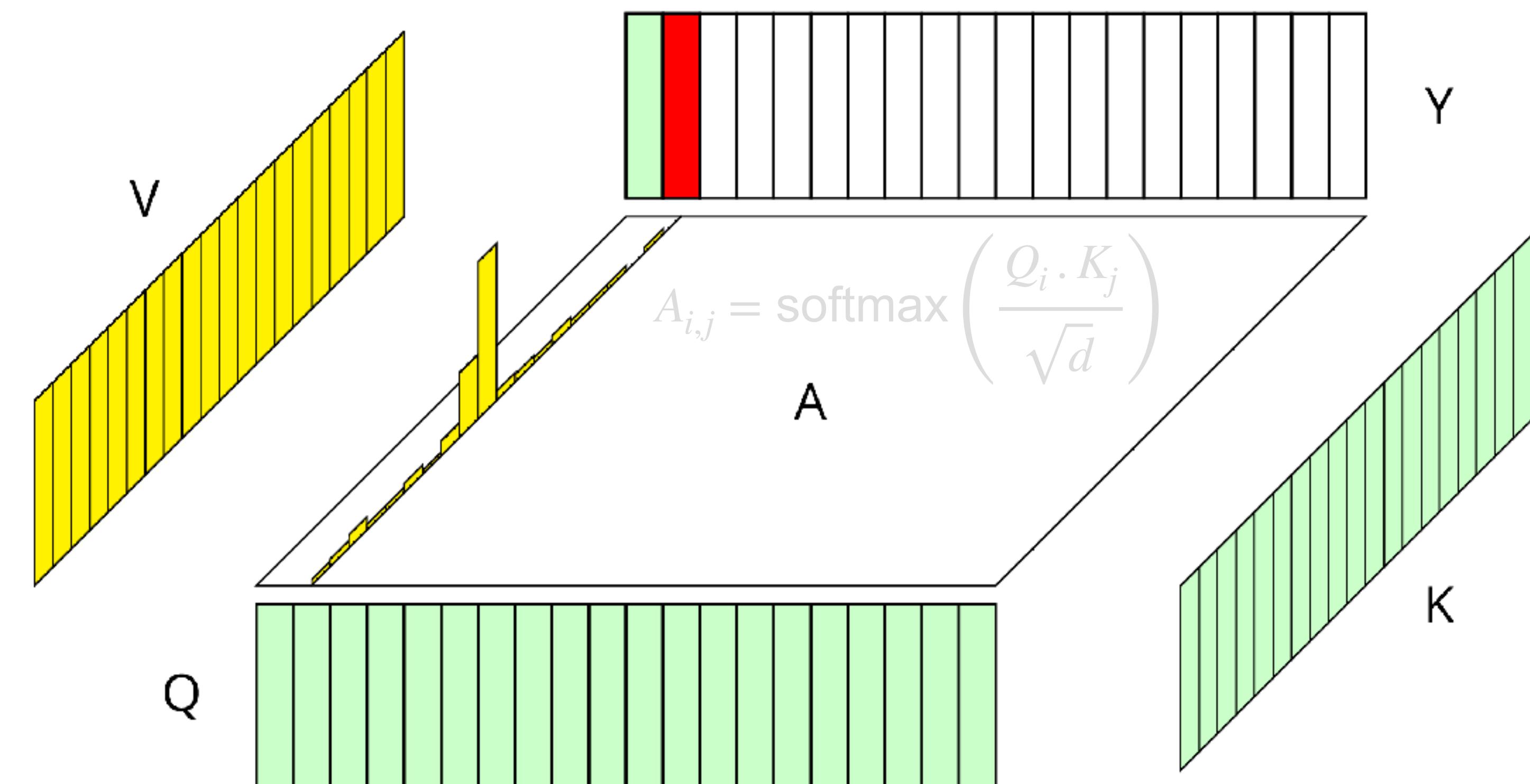
$$Y_i = \sum_j A_{i,j} V_j$$



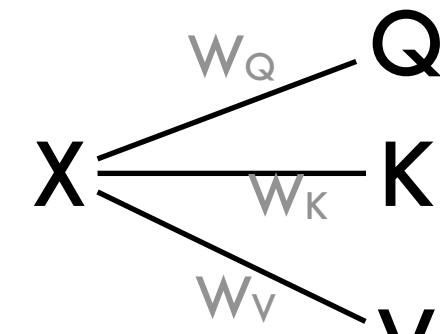
Attention mechanism

Given a **query** sequence **Q**,
a **key** sequence **K**, and
a **value** sequence **V**,
compute an attention matrix **A** by matching **Qs** to **Ks**, and
weight **V** with to get the sequence **Y**.

$$Y_i = \sum_j A_{i,j} V_j$$

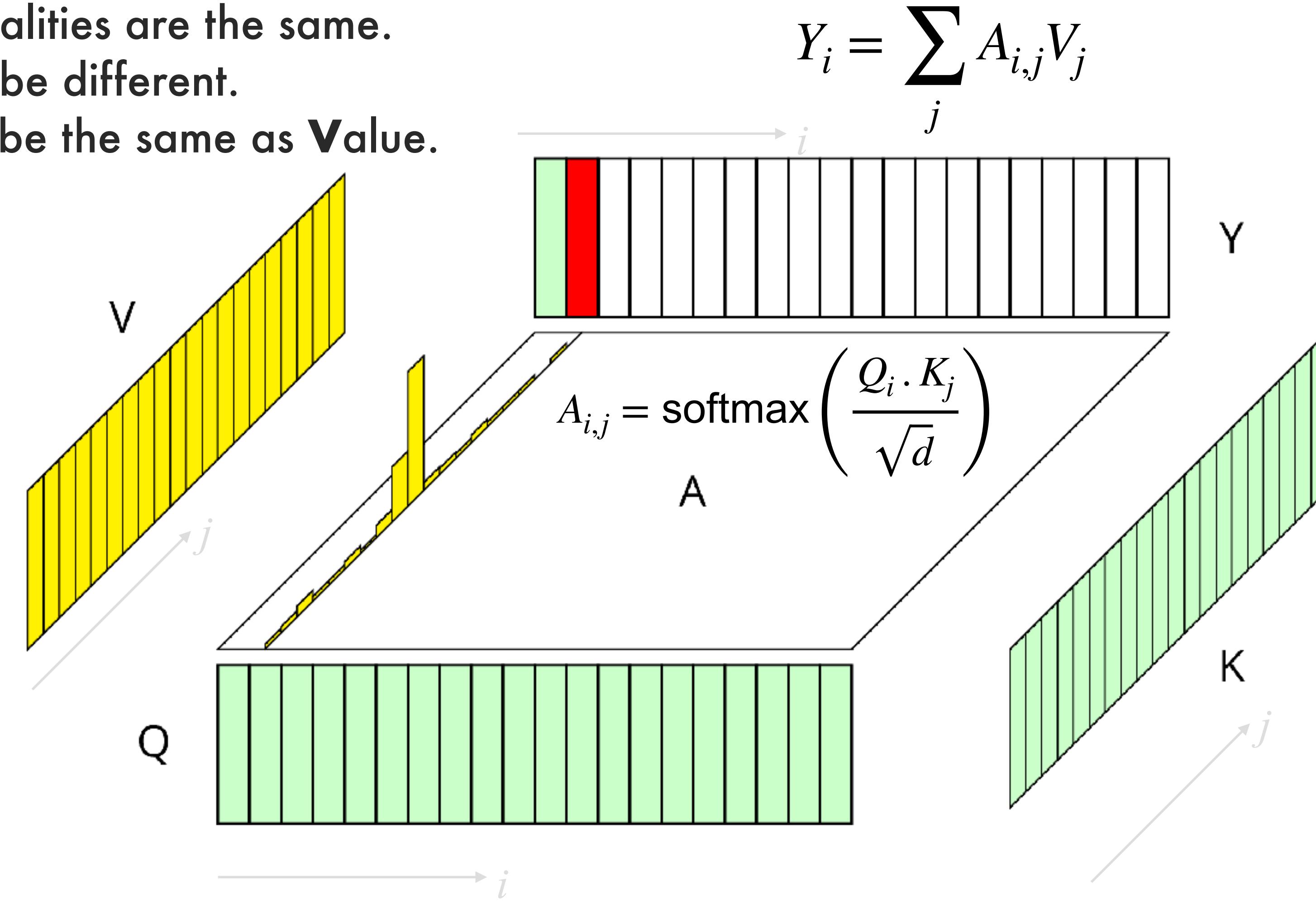


Attention mechanism

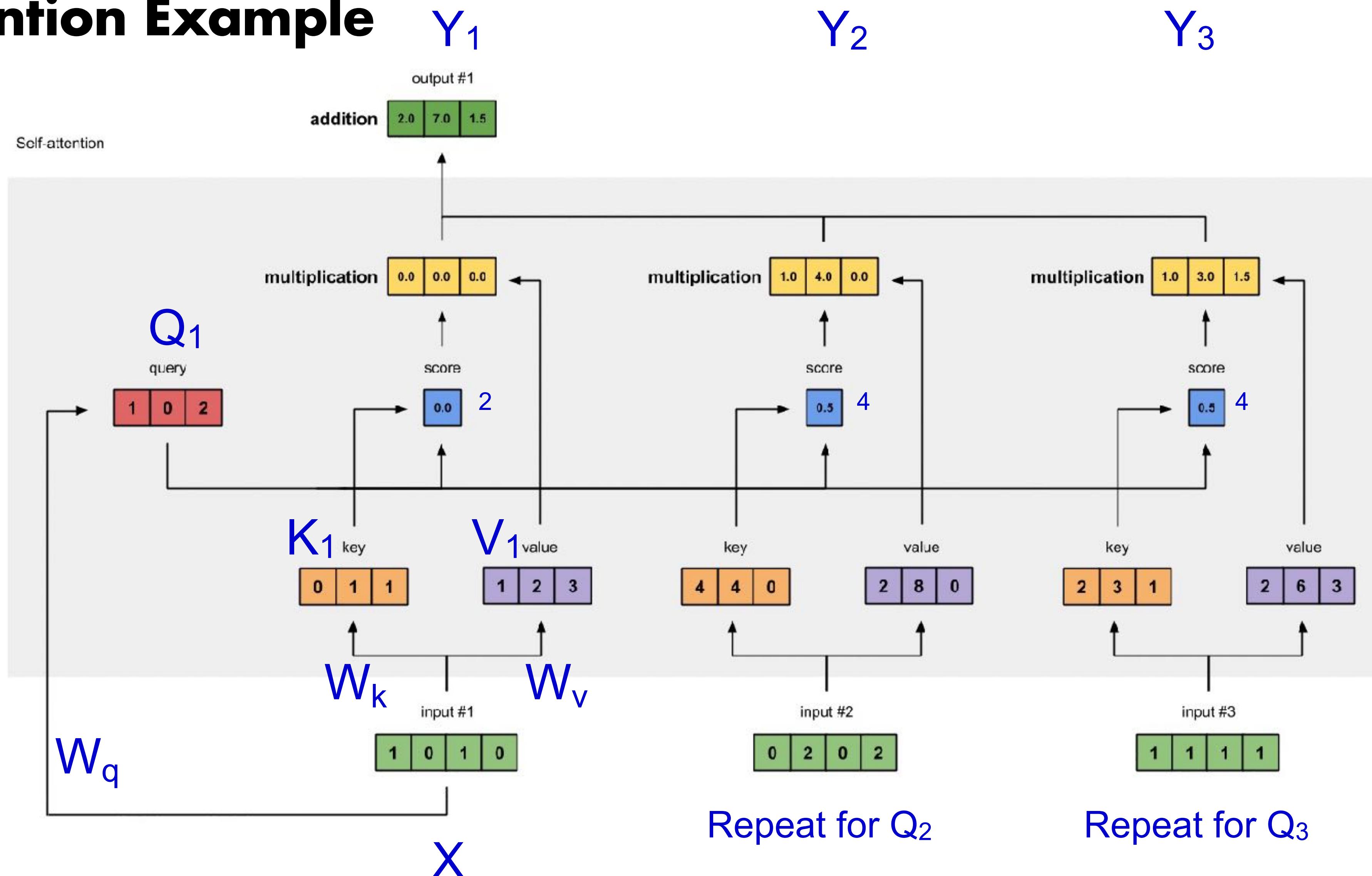


In "self-attention", (Q, K, V) obtained from the same input, linearly projected three times.

- **Query and Key** dimensionalities are the same.
- **Value** dimensionality may be different.
- Output dimensionality will be the same as **Value**.



Self-Attention Example



$$A_{i,j} = \text{softmax}\left(\frac{Q_i \cdot K_j}{\sqrt{d}}\right)$$

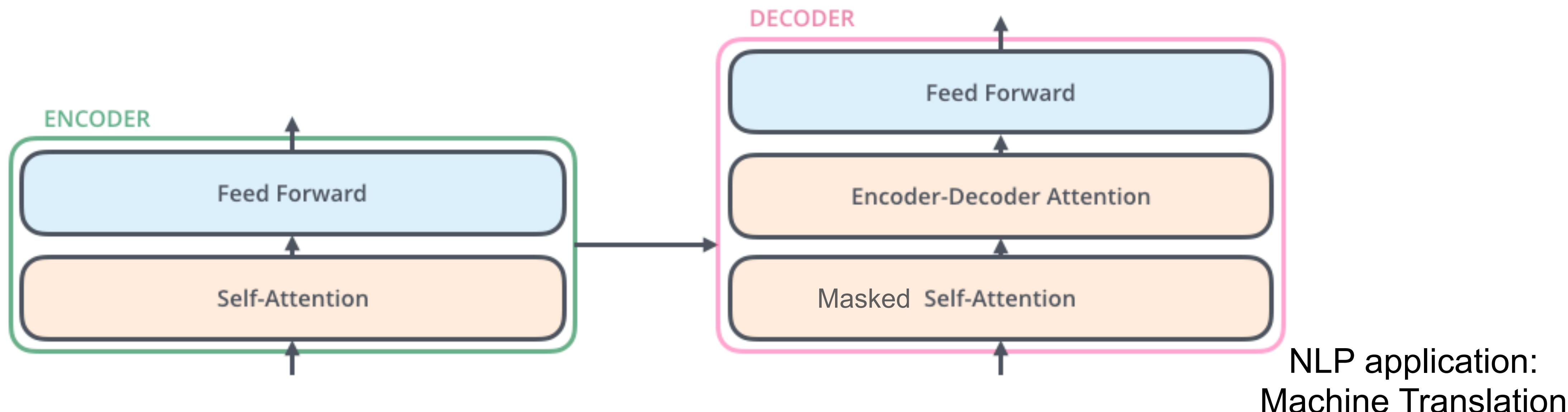
$$Y_i = \sum_j A_{i,j} V_j$$

Basic transformer model

- Sequence-to-sequence architecture using only point-wise processing and attention (no recurrent units or convolutions)

Encoder: receives entire input sequence and outputs encoded sequence of the same length

Decoder: predicts next token conditioned on encoder output and previously predicted tokens



A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, I. Polosukhin,
[Attention is all you need](#), NeurIPS 2017

Key-Value-Query attention model

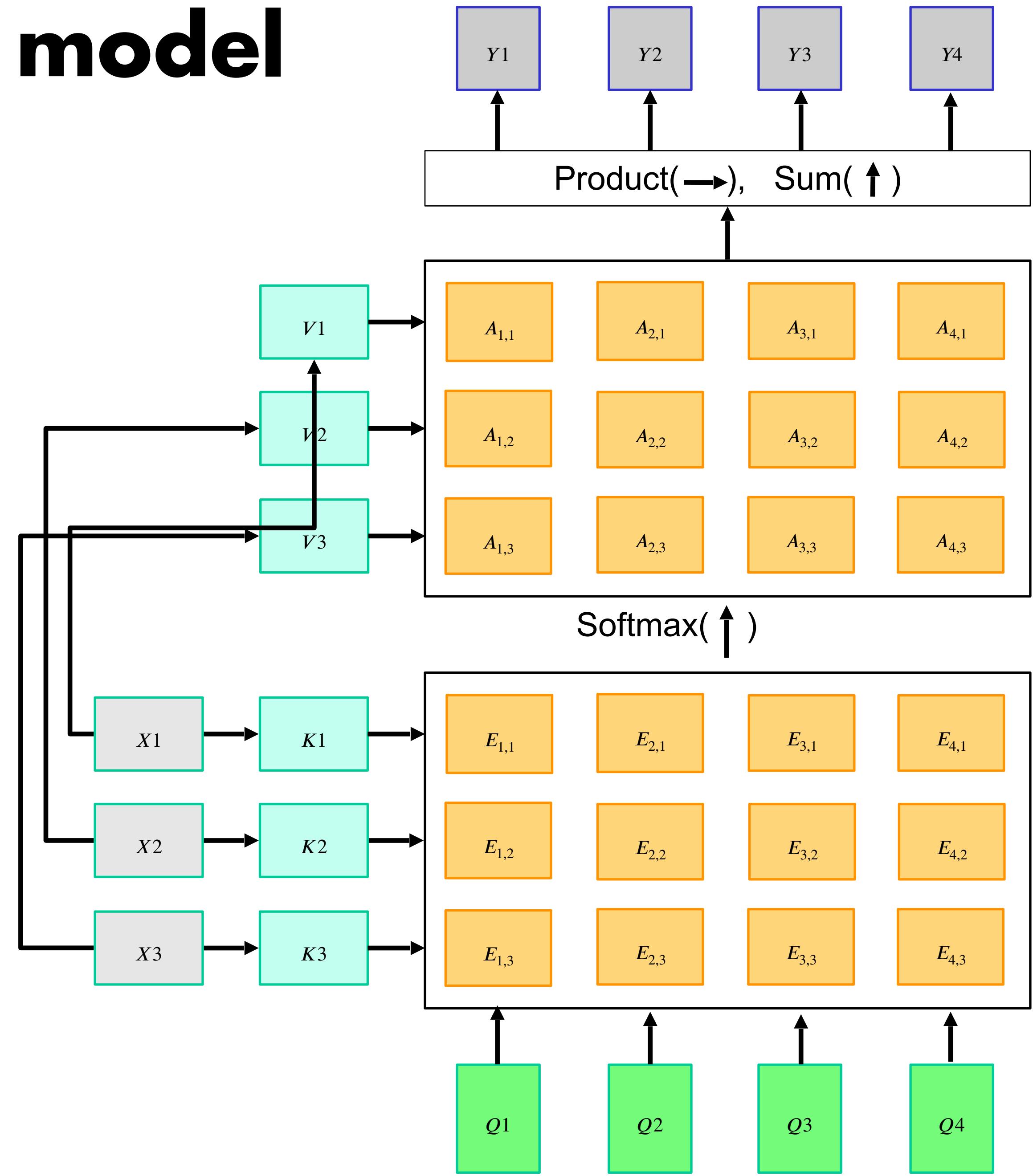
- Key vectors: $K = XW_K$
- Value Vectors: $V = XW_V$
- Query vectors
- Similarities: scaled dot-product attention

$$E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}} \text{ or } E = QK^T / \sqrt{D}$$

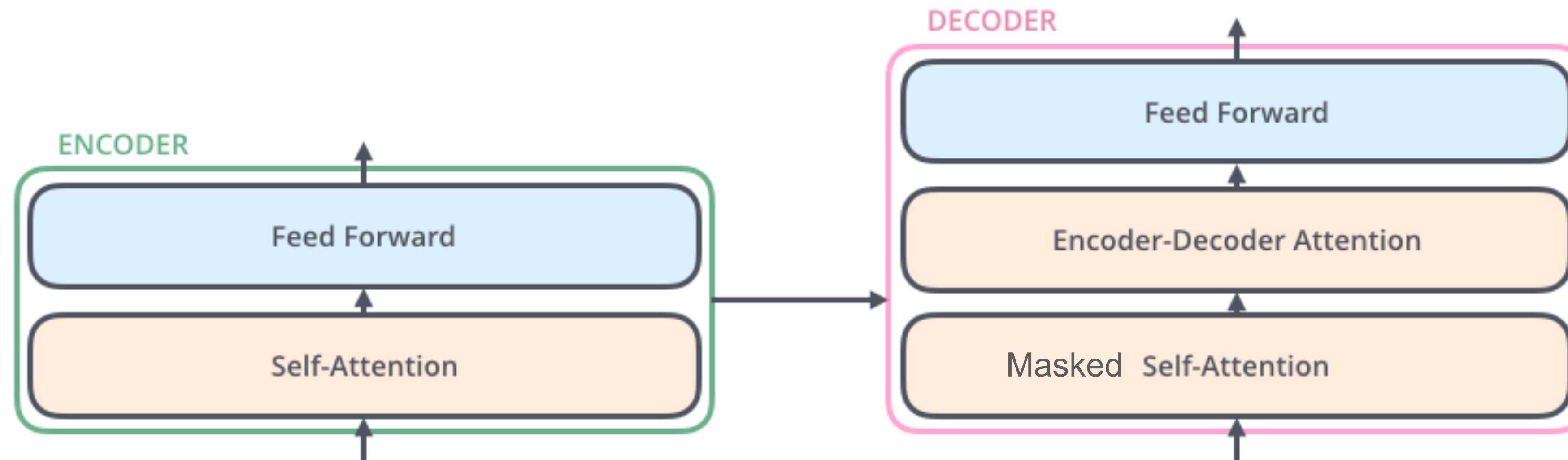
(D is the dimensionality of the keys)

- Attn. weights: $A = \text{softmax}(E, \text{ dim} = 1)$
- Output vectors:

$$Y_i = \sum_j A_{i,j} V_j \text{ or } Y = AV$$



Attention mechanisms



- **Encoder self-attention:** queries, keys, and values come from previous layer of encoder
- **Decoder self-attention:** values corresponding to future decoder outputs are masked out
- **Encoder-decoder attention:** queries come from previous decoder layer, keys and values come from output of encoder

Self-attention

- Used to capture context *within the sequence*

The animal didn't cross the street because it was too tired .

The animal didn't cross the street because it was too wide .

As we are encoding “it”, we should focus on “the animal”

As we are encoding “it”, we should focus on “the street”

Self-attention layer

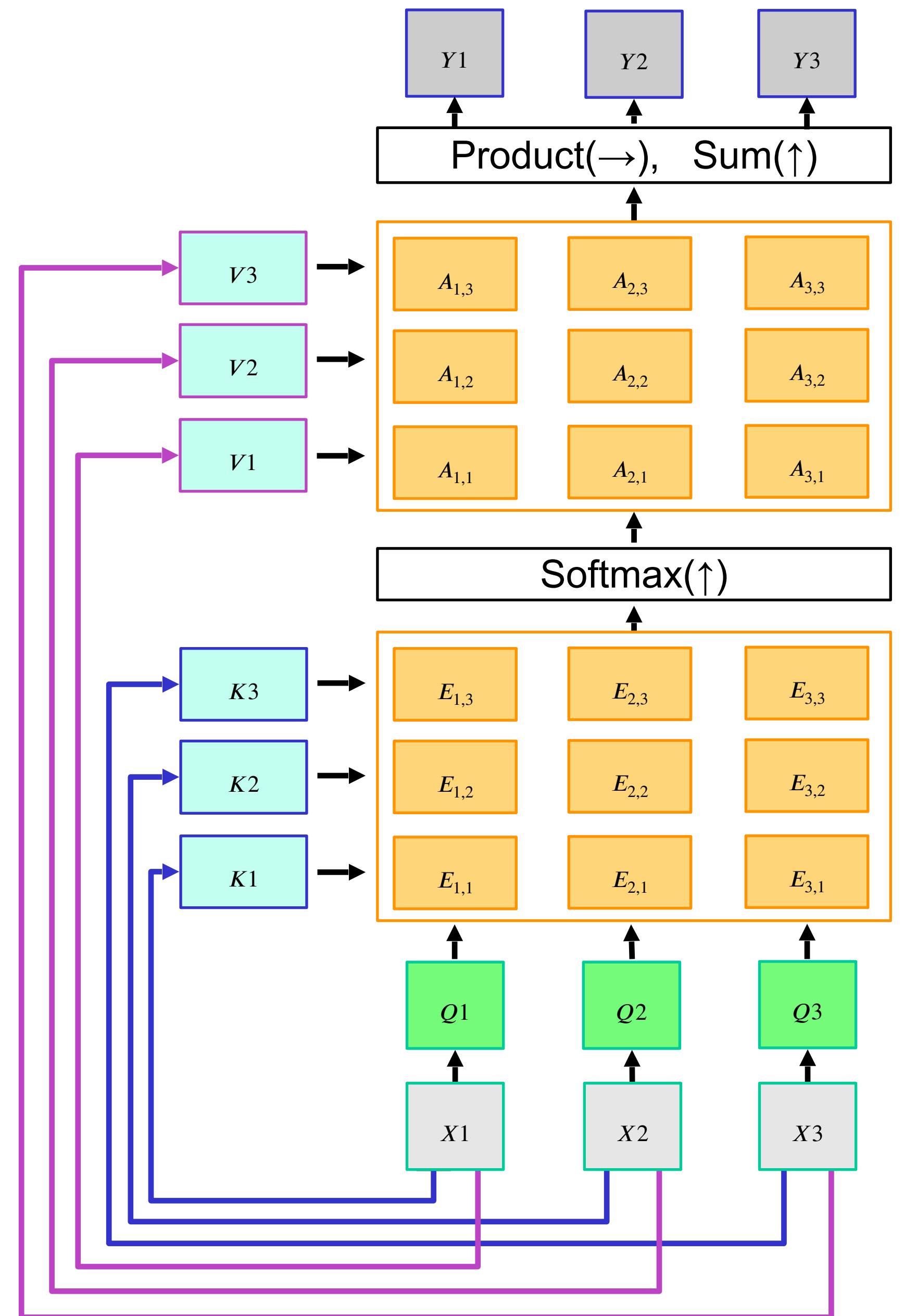
- **Query vectors:** $Q = XW_Q$
- **Key vectors:** $K = XW_K$
- **Value vectors:** $V = XW_V$
- **Similarities:** scaled dot-product attention

$$E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}} \text{ or } E = QK^T / \sqrt{D}$$

(D is the dimensionality of the keys)

- **Attn. weights:** $A = \text{softmax}(E, \text{ dim} = 1)$
- **Output vectors:**

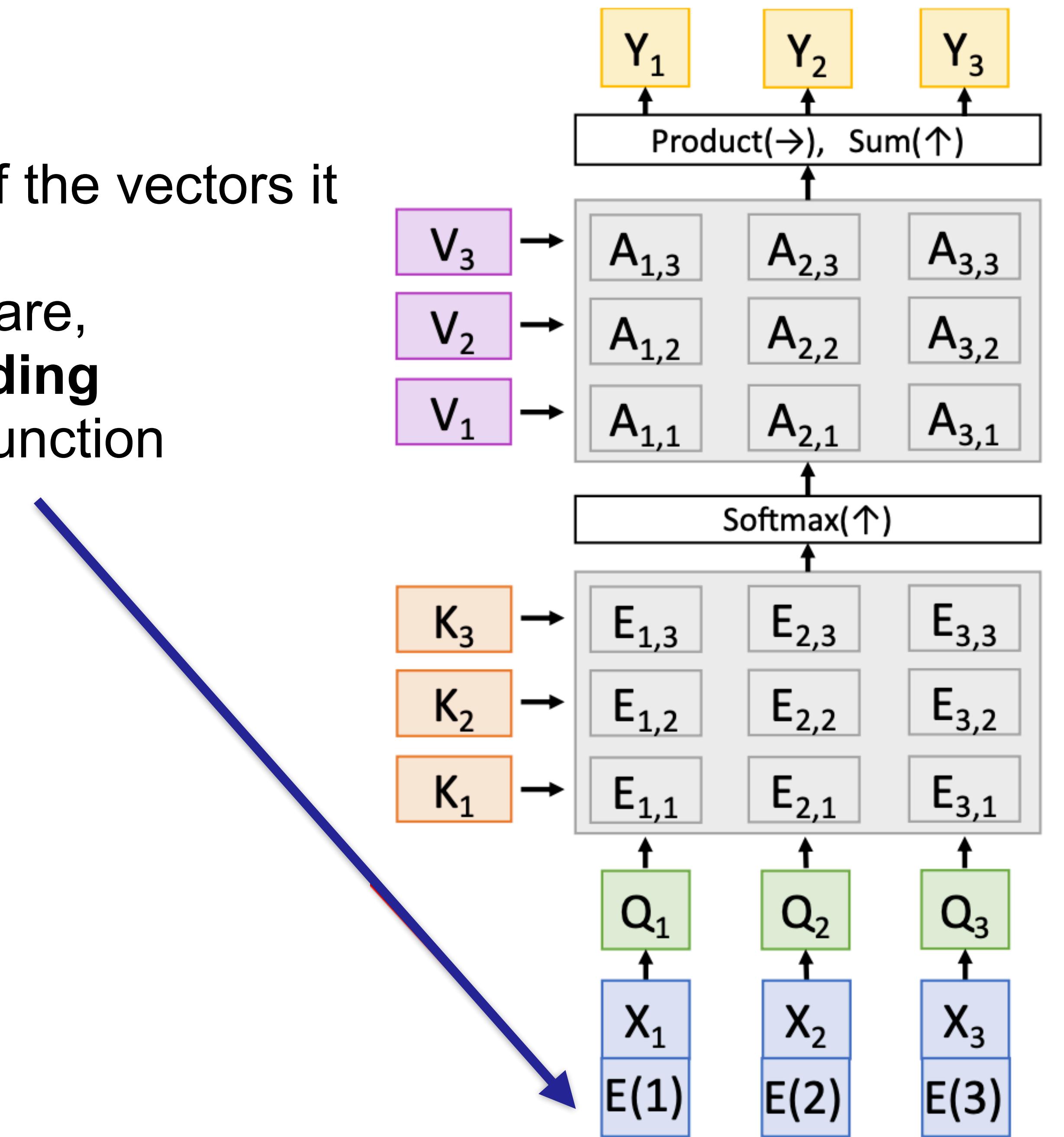
$$Y_i = \sum_j A_{i,j} V_j \text{ or } Y = AV$$



One query per input vector

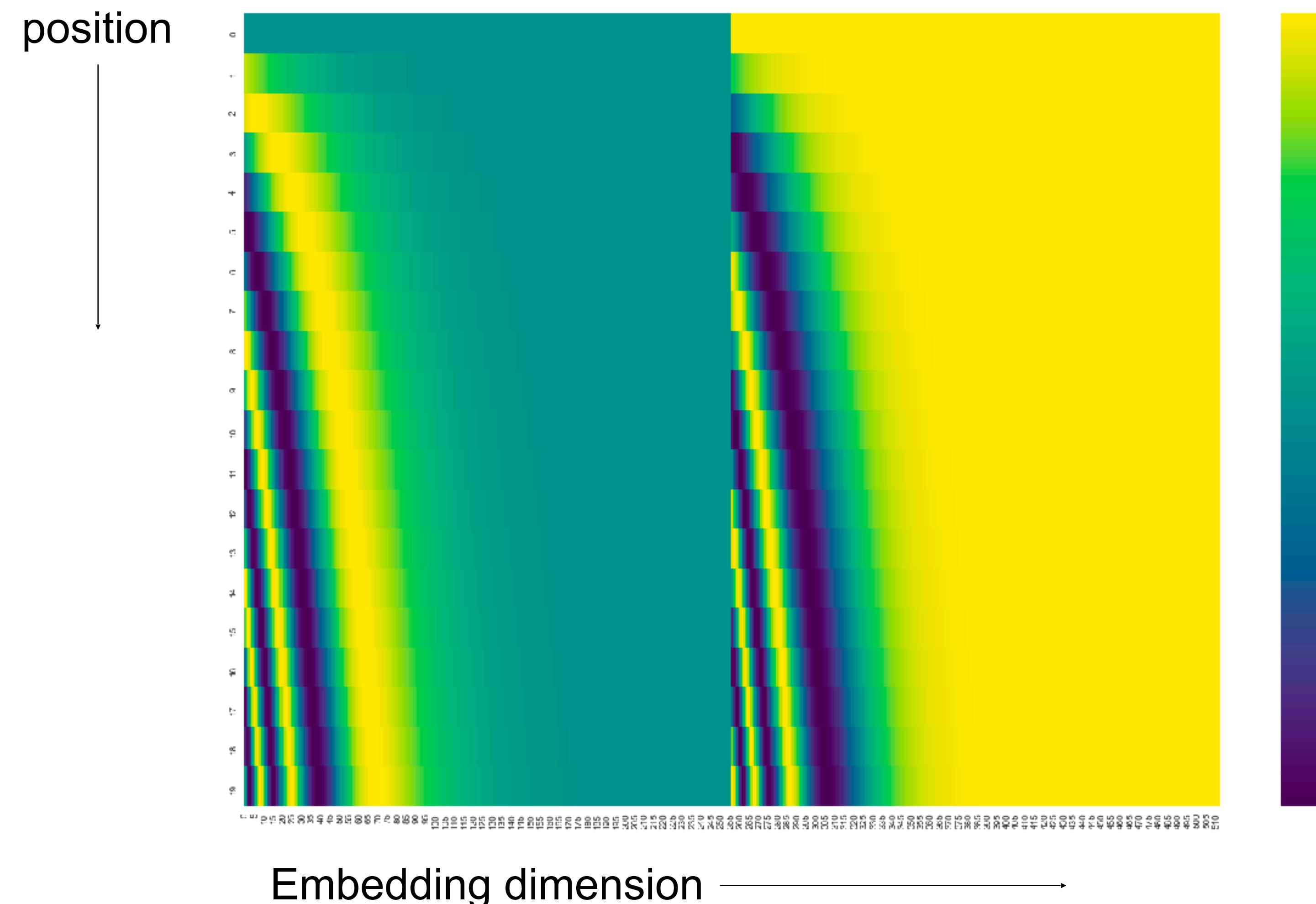
Positional encoding

- Self attention doesn't "know" the order of the vectors it is processing!
- In order to make processing position-aware, concatenate input with **positional encoding**
- E can be learned lookup table, or fixed function

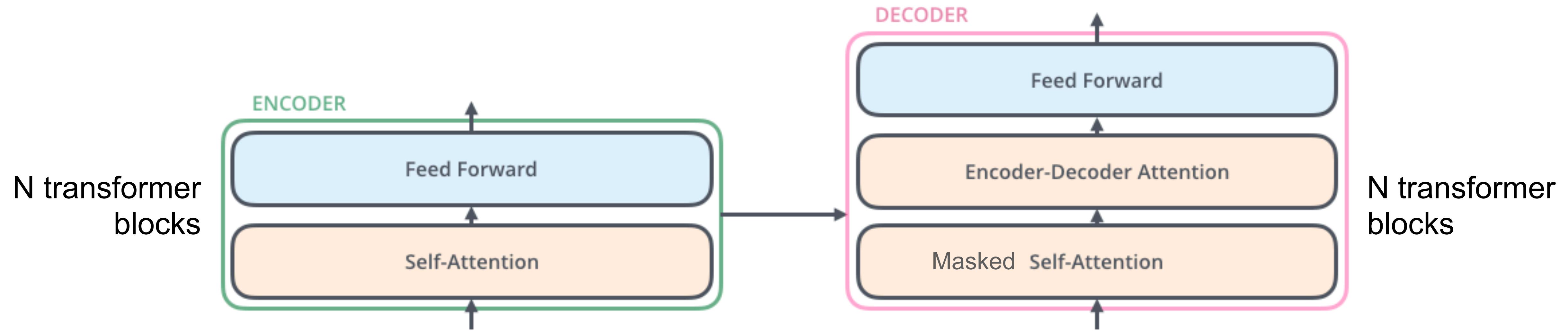


Positional encoding

- To give transformer information about ordering of tokens, add function of position (based on sines and cosines) to every input



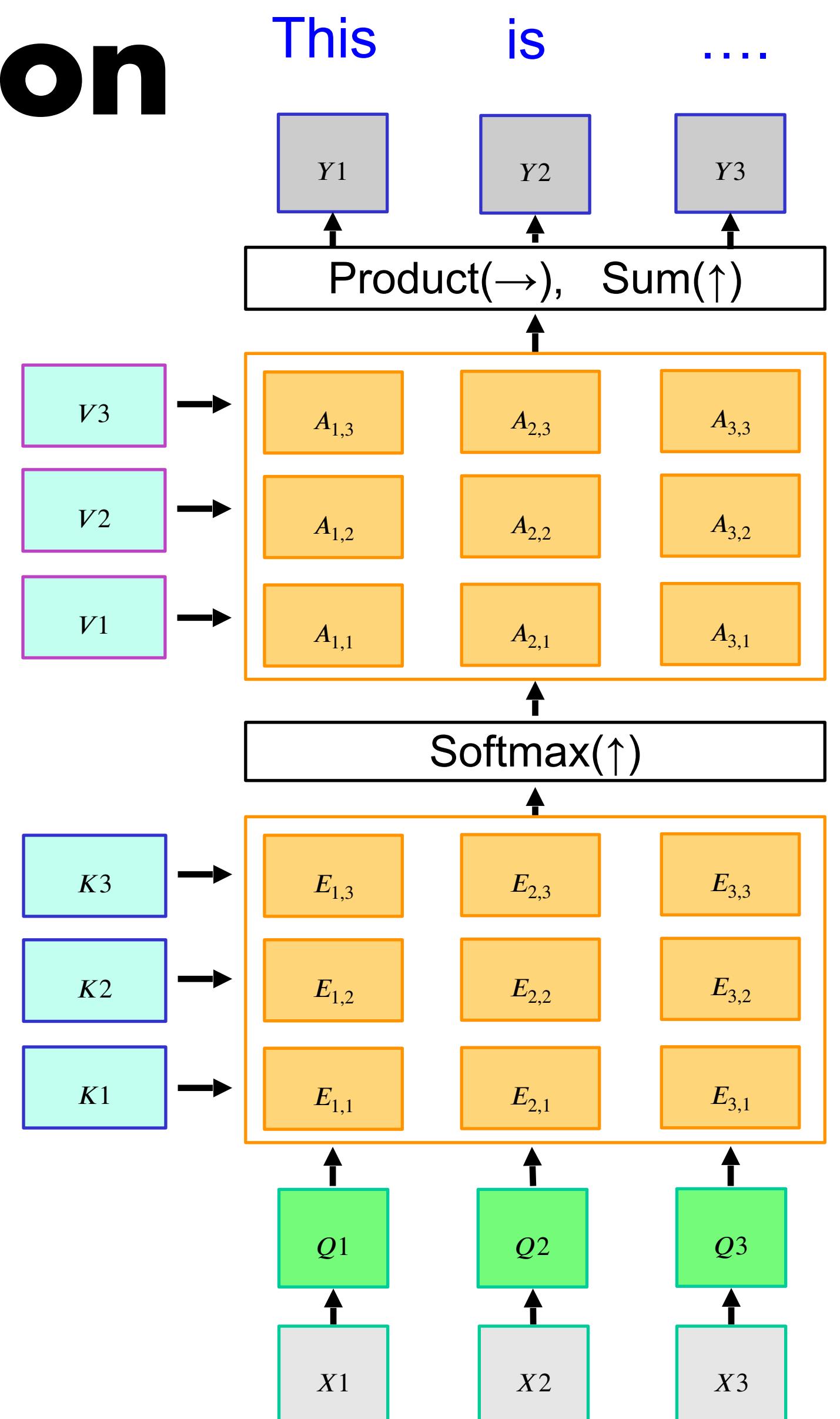
Attention mechanisms: Overview



- **Encoder self-attention:** queries, keys, and values come from previous layer of encoder
- **Decoder self-attention:** values corresponding to future decoder outputs are masked out
- **Encoder-decoder attention:** queries come from previous decoder layer, keys and values come from output of encoder

Decoder: Masked self-attention

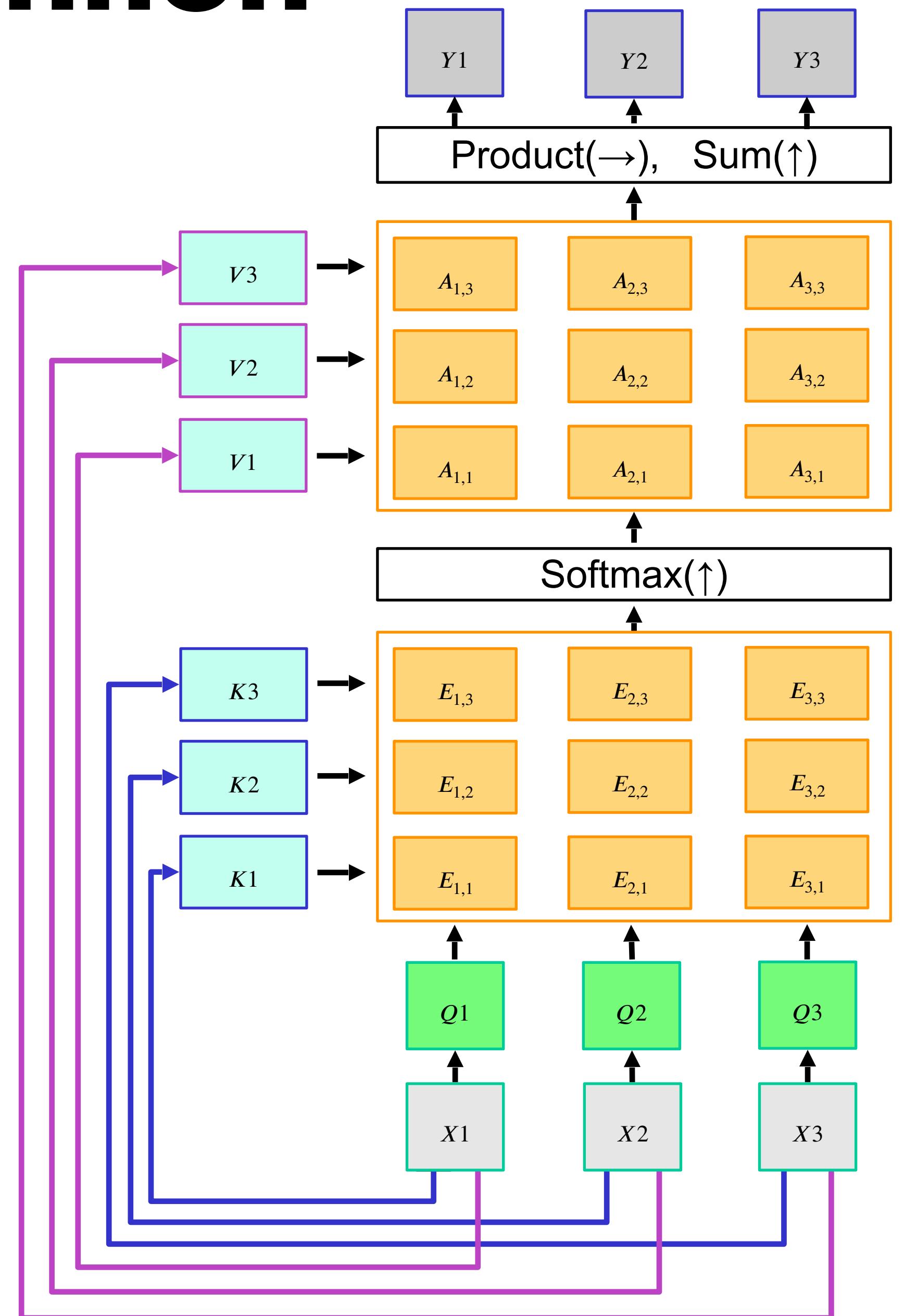
- The decoder should not “look ahead” in the output sequence



$\langle \text{START} \rangle \text{ This } \text{ is }$

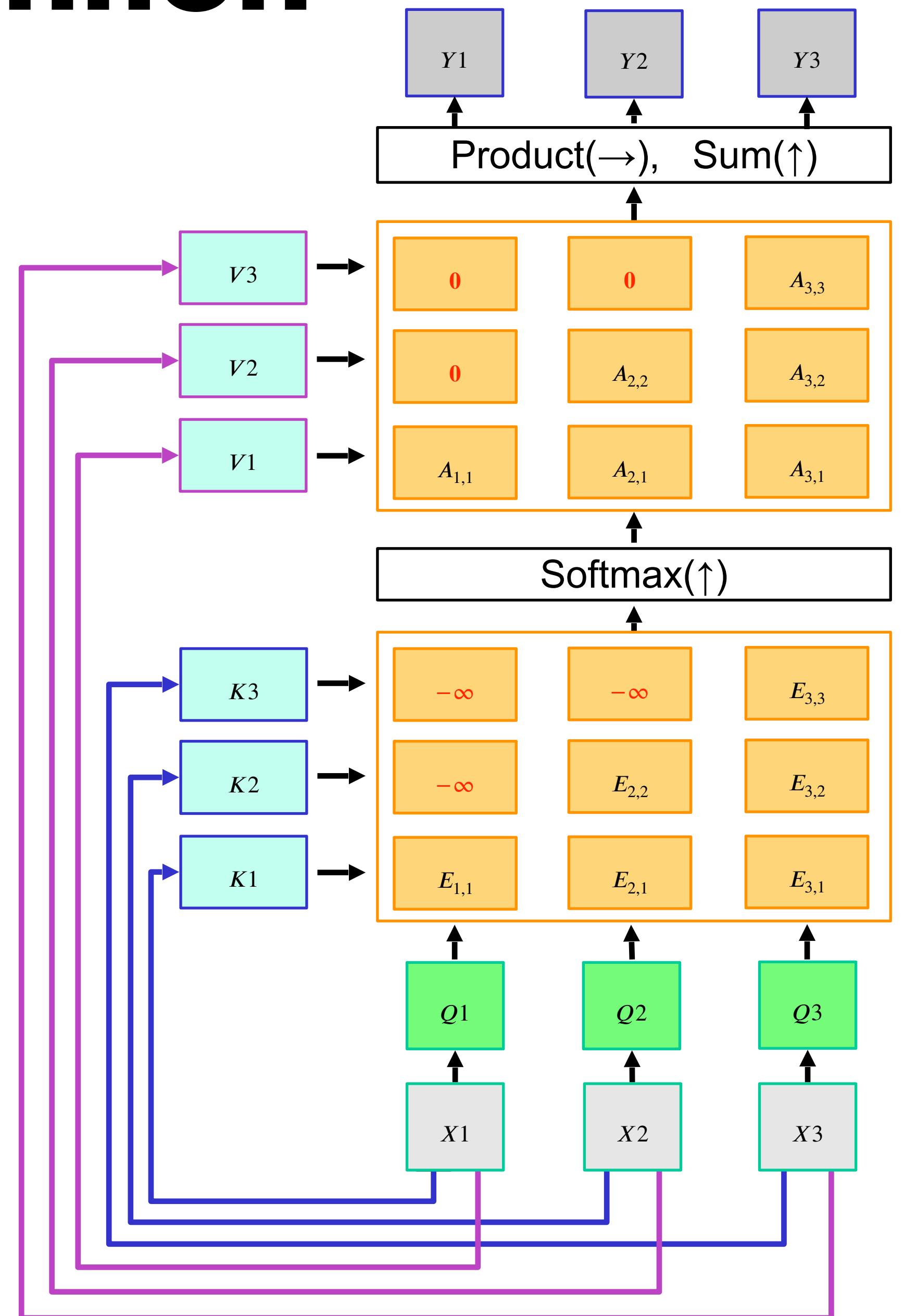
Decoder: Masked self-attention

- The decoder should not “look ahead” in the output sequence

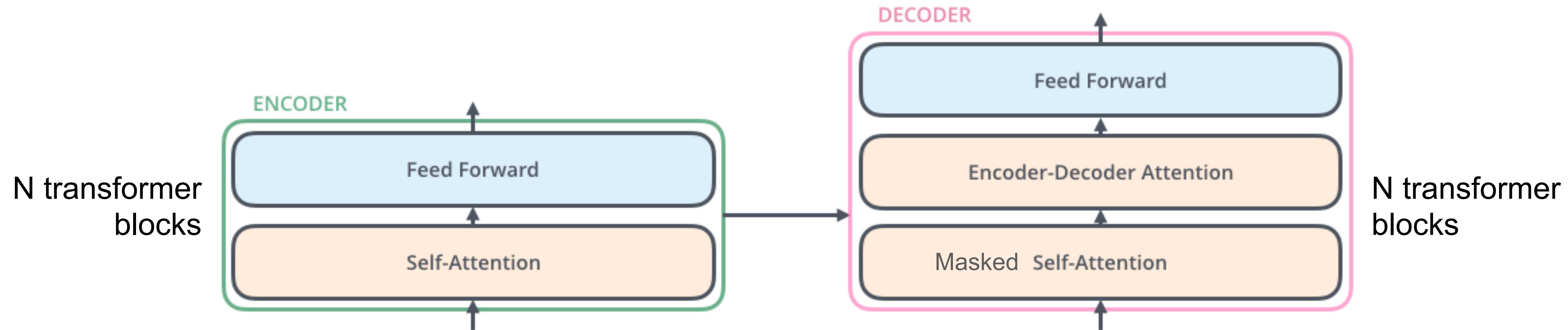


Decoder: Masked self-attention

- The decoder should not “look ahead” in the output sequence

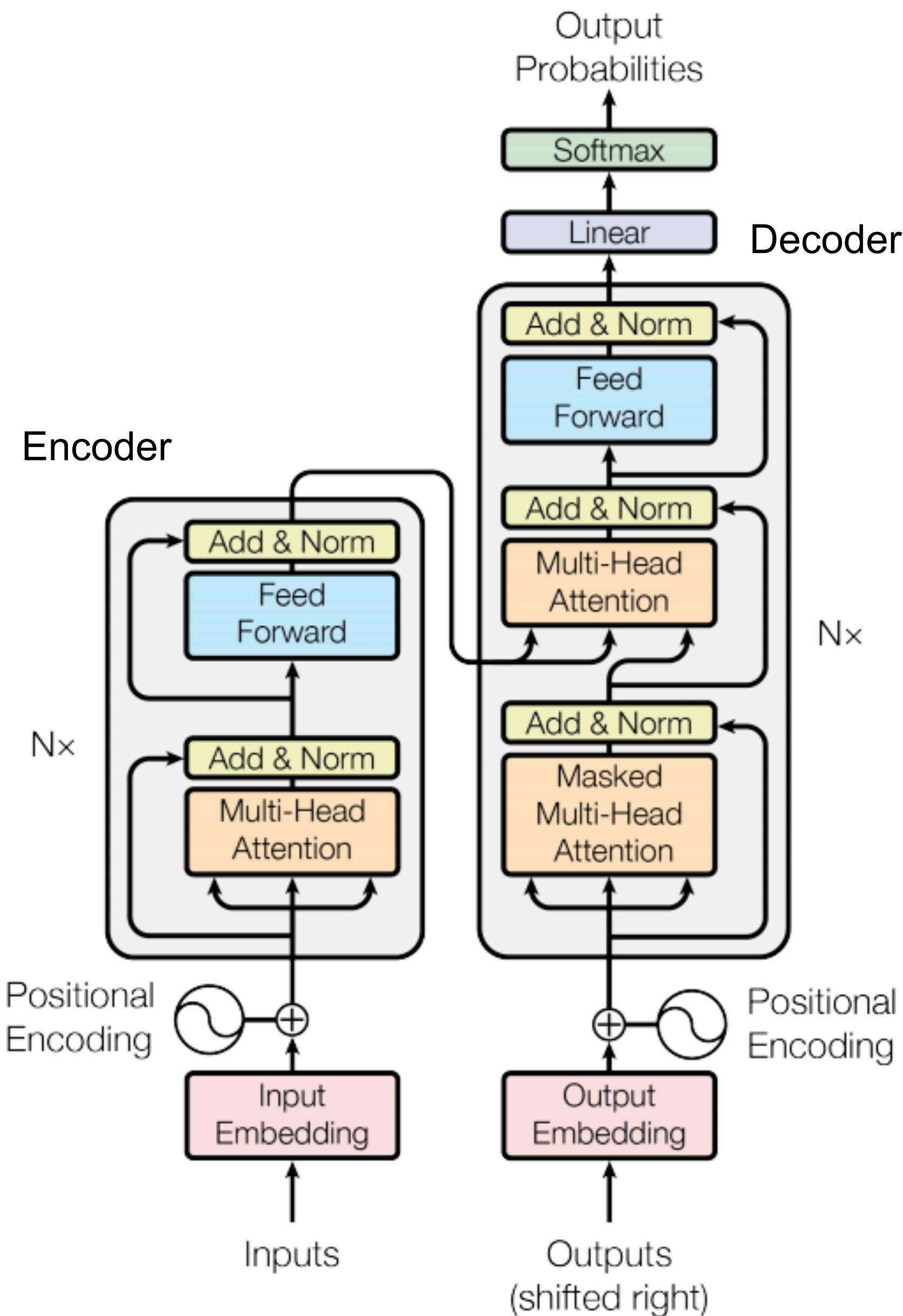


Attention mechanisms: Overview



- **Encoder self-attention:** queries, keys, and values come from previous layer of encoder
- **Decoder self-attention:** values corresponding to future decoder outputs are masked out
- **Encoder-decoder attention:** queries come from previous decoder layer, keys and values come from output of encoder

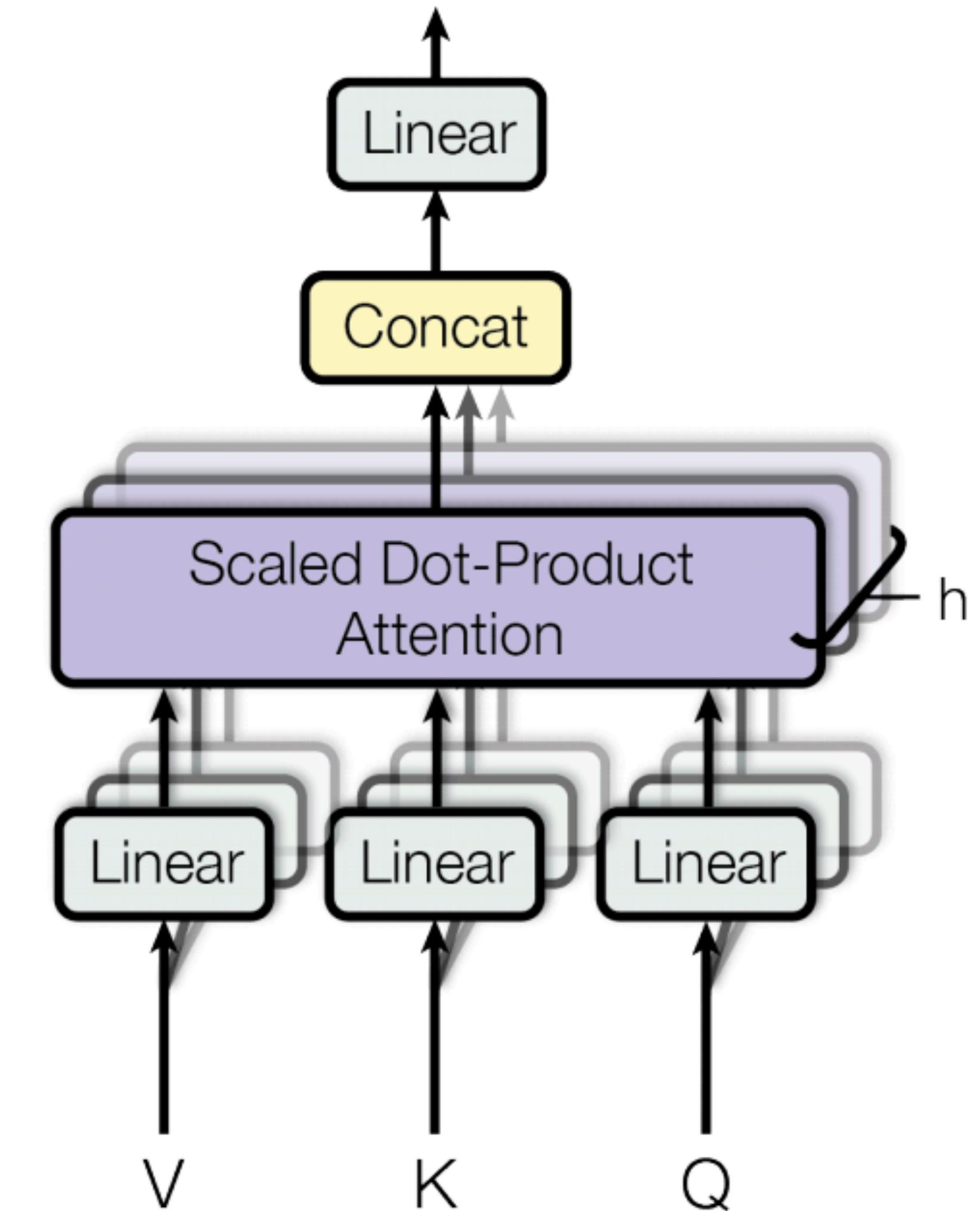
Transformer architecture: Details



A. Vaswani et al., [Attention is all you need](#), NeurIPS 2017

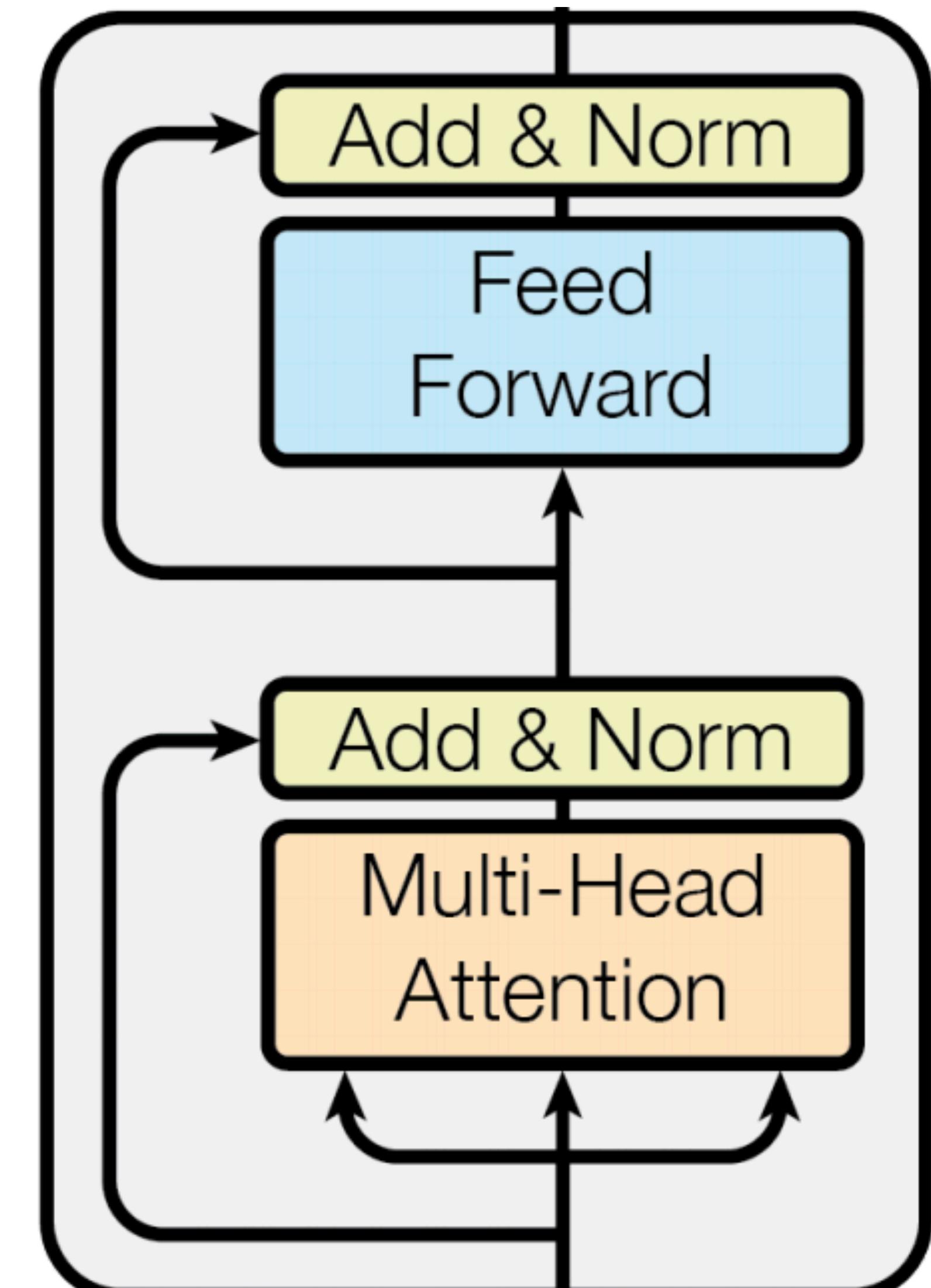
Multi-head attention

- Run h attention models in parallel on top of different linearly projected versions of Q , K , V ; concatenate and linearly project the results
- Intuition: enables model to attend to different kinds of information at different positions (see [visualization tool](#))

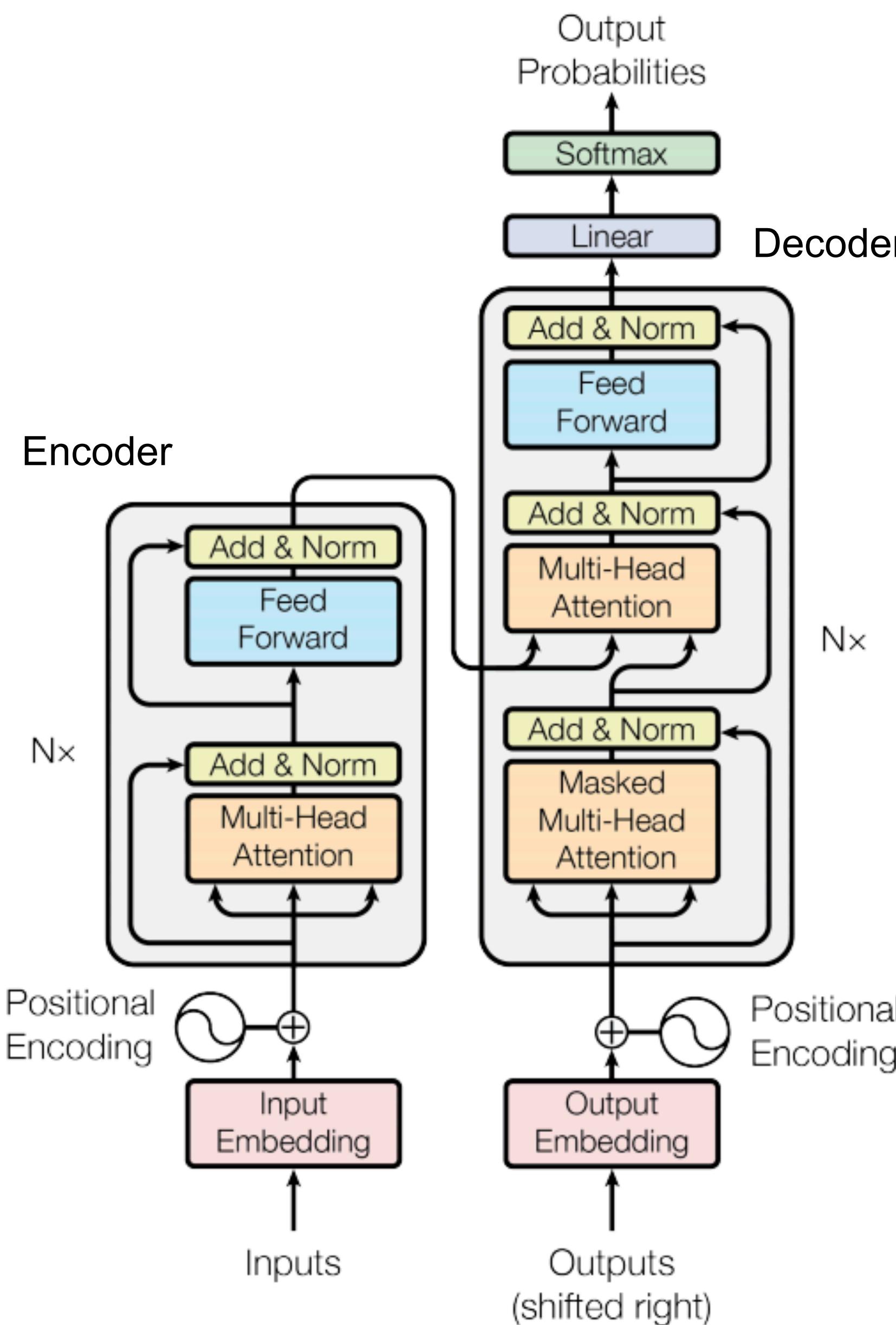


Transformer blocks

- A Transformer is a sequence of transformer blocks
 - Vaswani et al.: $N=12$ blocks, embedding dimension = 512, 6 attention heads
 - **Add & Norm:** residual connection followed by layer normalization
 - **Feedforward:** two linear layers with ReLUs in between, applied *independently* to each vector
- **Attention is the only interaction between inputs!**



Transformer architecture: Zooming back out



A. Vaswani et al., [Attention is all you need](#), NeurIPS 2017

Transformer implementation

```
class MultiHeadedAttention(nn.Module):
    """
    Multi-Head Attention module from "Attention is All You Need"
    Implementation modified from OpenNMT-py.
    https://github.com/OpenNMT/OpenNMT-py
    """

    def __init__(self, num_heads: int, size: int, dropout: float = 0.1):
        """
        Create a multi-headed attention layer.
        :param num_heads: the number of heads
        :param size: model size (must be divisible by num_heads)
        :param dropout: probability of dropping a unit
        """
        super(MultiHeadedAttention, self).__init__()

        assert size % num_heads == 0

        self.head_size = head_size = size // num_heads
        self.model_size = size
        self.num_heads = num_heads

        self.k_layer = nn.Linear(size, num_heads * head_size)
        self.v_layer = nn.Linear(size, num_heads * head_size)
        self.q_layer = nn.Linear(size, num_heads * head_size)

        self.output_layer = nn.Linear(size, size)
        self.softmax = nn.Softmax(dim=-1)
        self.dropout = nn.Dropout(dropout)
```

nn.Linear:
Learnable params

Transformer implementation

```
def forward(  
    self, k: Tensor, v: Tensor, q: Tensor, rel_mouth_times=None, mask: Tensor = None  
):  
    """  
    Computes multi-headed attention.  
    :param k: keys [B, K, D] with K being the sentence length.  
    :param v: values [B, K, D]  
    :param q: query [B, Q, D] with Q being the sentence length.  
    :param mask: optional mask [B, 1, K]  
    :return:  
    """  
  
    batch_size = k.size(0) # B  
    num_heads = self.num_heads # H  
    # project the queries (q), keys (k), and values (v)  
    → k = self.k_layer(k)  
    → v = self.v_layer(v)  
    → q = self.q_layer(q)  
  
    # reshape q, k, v for our computation to [B, H, ..., D/H]  
    k = k.view(batch_size, -1, num_heads, self.head_size).transpose(1, 2)  
    v = v.view(batch_size, -1, num_heads, self.head_size).transpose(1, 2)  
    q = q.view(batch_size, -1, num_heads, self.head_size).transpose(1, 2)  
  
    # compute scores  
    q = q / math.sqrt(self.head_size)
```

$$\begin{aligned}K &= XW_K \\V &= XW_V \\Q &= XW_Q\end{aligned}$$

Transformer implementation

```
# [B, H, Q, K]
scores = torch.matmul(q, k.transpose(2, 3))

# apply the mask (if we have one)
# we add a dimension for the heads to it below: [B, 1, 1, K]
if mask is not None:
    scores = scores.masked_fill(~mask.unsqueeze(1), float("-inf"))

# apply attention dropout and compute context vectors.
# [B, H, Q, K]
attention_map = self.softmax(scores)
attention = self.dropout(attention_map)

# get context vector (select values with attention)
# [B, H, Q, D/H]
context = torch.matmul(attention, v)
# reshape back to [B, Q, D]
context = (
    context.transpose(1, 2)
    .contiguous()
    .view(batch_size, -1, num_heads * self.head_size)
)

# [B, Q, D]
output = self.output_layer(context)

return output, attention_map
```

$$E = QK^T / \sqrt{D}$$

$$A = \text{softmax}(E, \text{dim} = 1)$$

$$Y = AV$$

```
class TransformerEncoderLayer(nn.Module):
    """
    One Transformer encoder layer has a Multi-head attention layer plus
    a position-wise feed-forward layer.
    """

    def __init__(self, size: int = 0, ff_size: int = 0, num_heads: int = 0, dropout: float = 0.1):
        """
        A single Transformer layer.
        :param size:
        :param ff_size:
        :param num_heads:
        :param dropout:
        """
        super(TransformerEncoderLayer, self).__init__()

        self.layer_norm = nn.LayerNorm(size, eps=1e-6)
        self.src_src_att = MultiHeadedAttention(num_heads, size, dropout=dropout)
        self.feed_forward = PositionwiseFeedForward(
            size, ff_size=ff_size, dropout=dropout
        )
        self.dropout = nn.Dropout(dropout)
        self.size = size

    # pylint: disable=arguments-differ
    def forward(self, x: Tensor, mask: Tensor) -> Tensor:
        """
        Forward pass for a single transformer encoder layer.
        First applies layer norm, then self attention,
        then dropout with residual connection (adding the input to the result),
        and then a position-wise feed-forward layer.
        :param x: layer input
        :param mask: input mask
        :return: output tensor
        """

        x_norm = self.layer_norm(x)
        h, att_map_src_src = self.src_src_att(k=x_norm, v=x_norm, q=x_norm, mask=mask)
        h = self.dropout(h) + x
        o = self.feed_forward(h)
        return o
```

1. Beyond CNNs

- Attention & Transformer

- Vision Transformers

Image transformer – Google

Self-attention only locally

- Image generation and super-resolution with 32x32 output, attention restricted to local neighborhoods

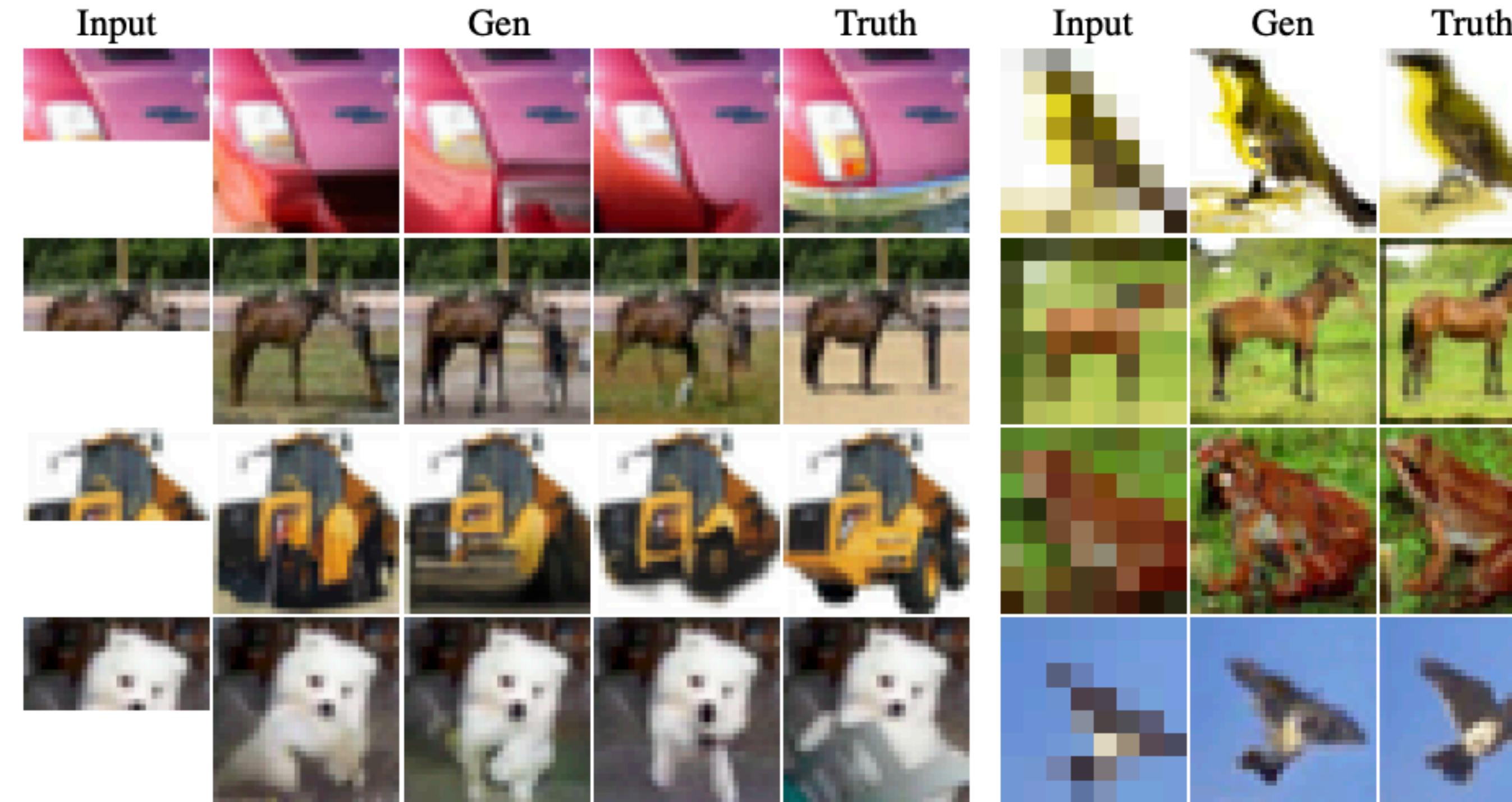


Table 2. On the left are image completions from our best conditional generation model, where we sample the second half. On the right are samples from our four-fold super-resolution model trained on CIFAR-10. Our images look realistic and plausible, show good diversity among the completion samples and observe the outputs carry surprising details for coarse inputs in super-resolution.

Sparse transformers – OpenAI

scalable approximations to global self-attention in
order to be applicable to images

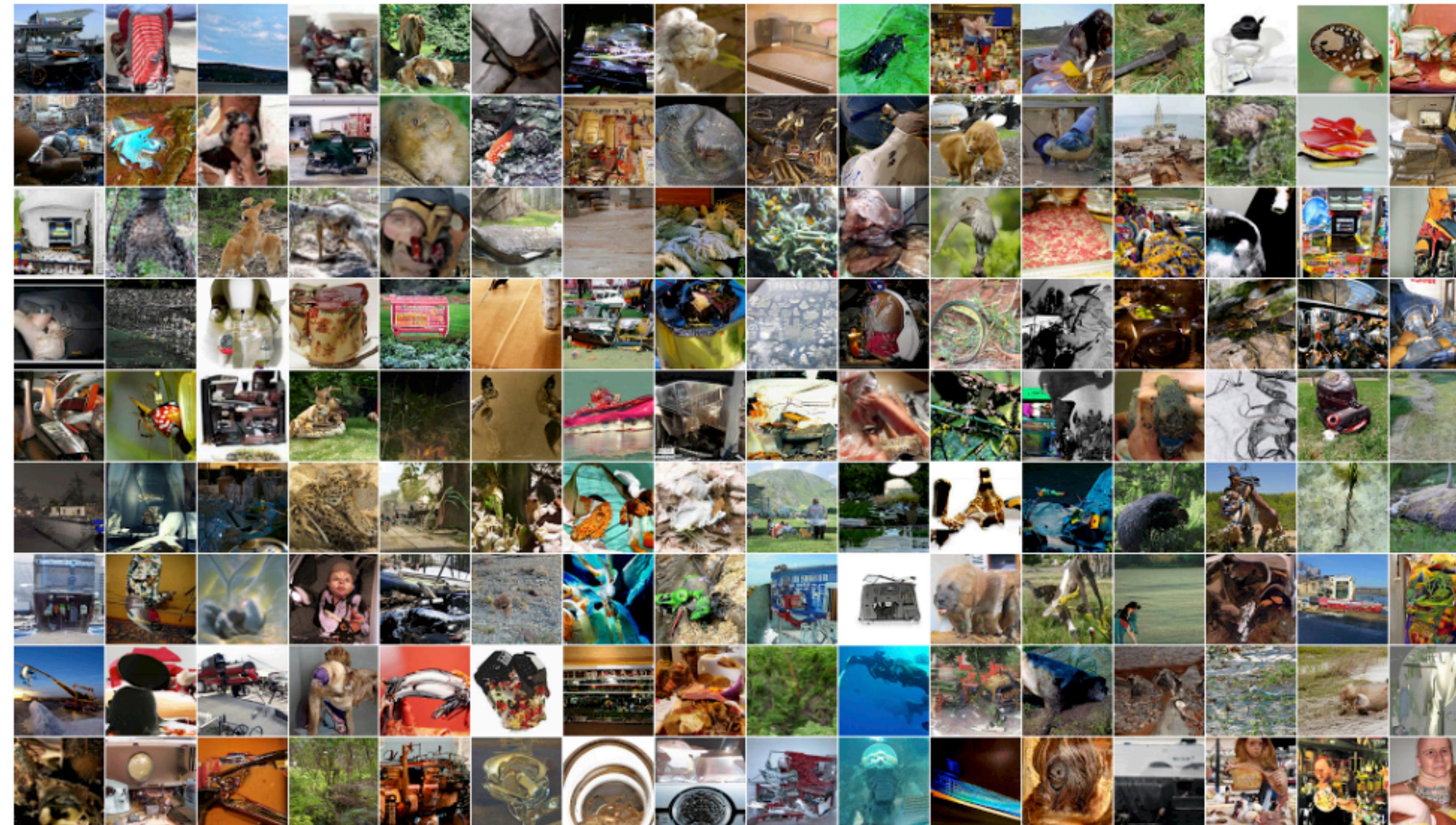
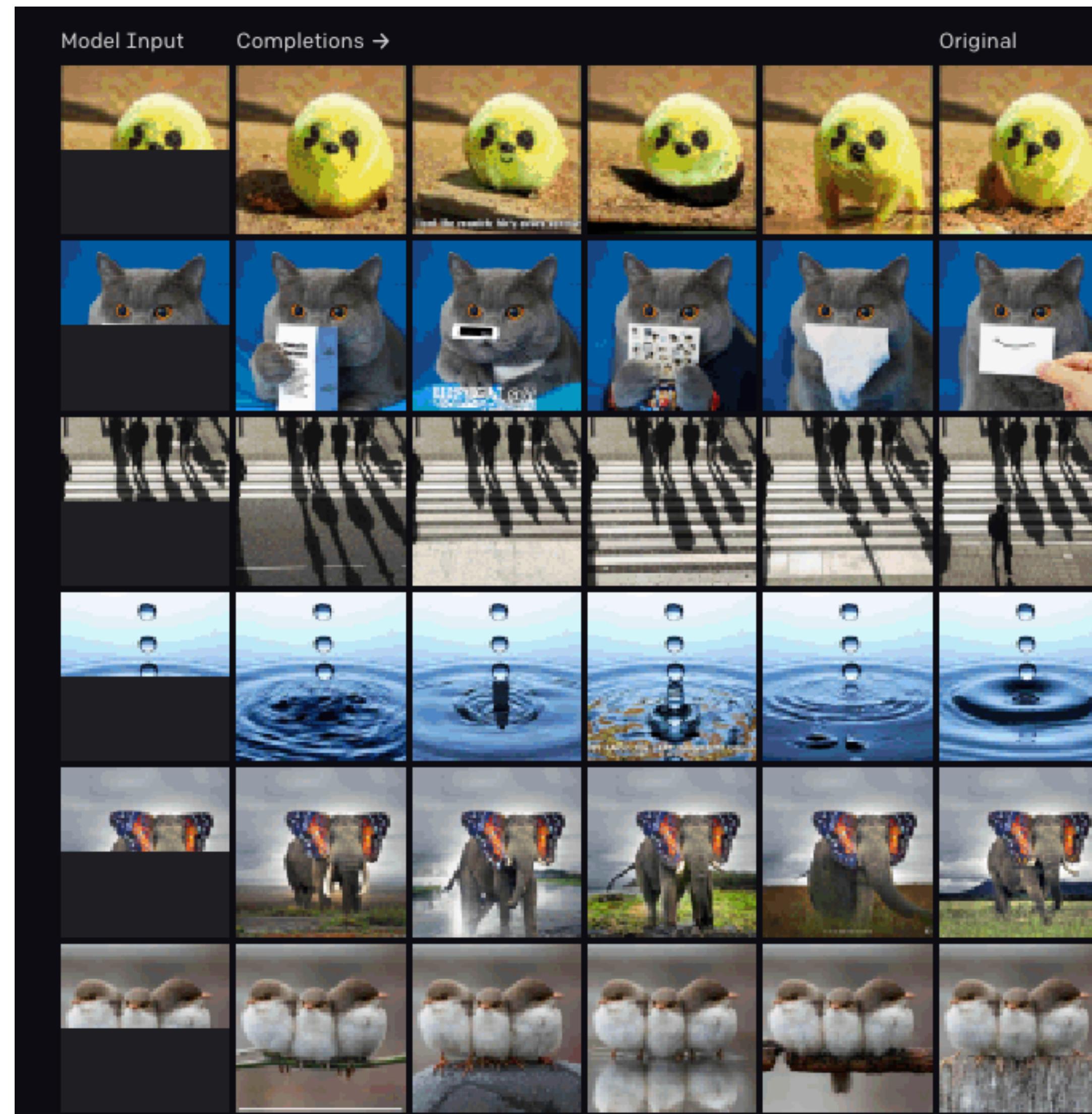


Figure 5. Unconditional samples from ImageNet 64x64, generated with an unmodified softmax temperature of 1.0. We are able to learn long-range dependencies directly from pixels without using a multi-scale architecture.

Image GPT* – OpenAI

works on reduced resolutions



<https://openai.com/blog/image-gpt/>

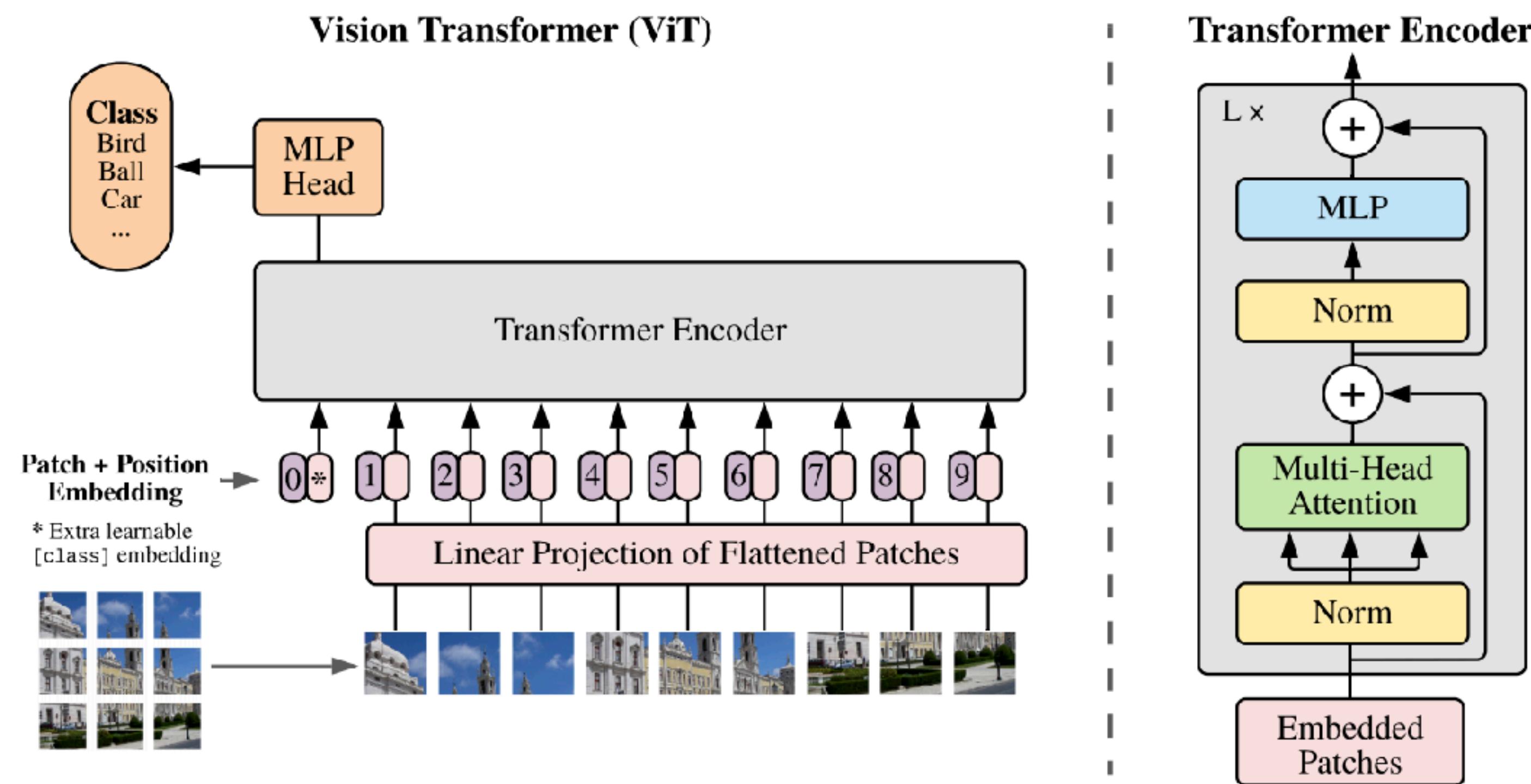
M. Chen et al., [Generative pretraining from pixels](#), ICML 2020

*GPT: Generative pre-trained Transformer

Vision transformer (ViT) - Google

Full resolution

- Split an image into patches, feed linearly projected patches into standard transformer encoder
- With patches of 14×14 pixels, you need $16 \times 16 = 256$ patches to represent 224×224 images



Vision transformer (ViT)

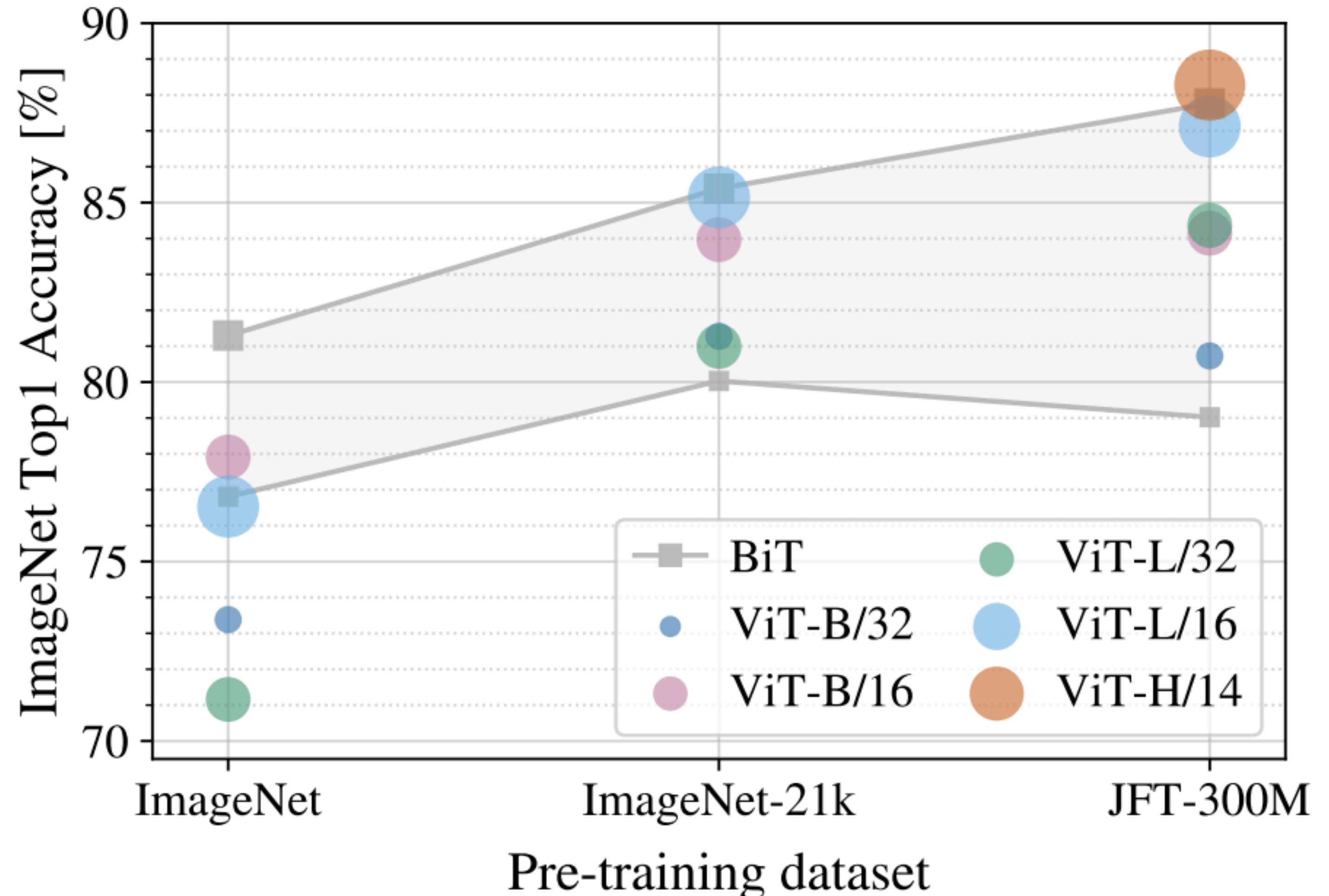


Figure 3: Transfer to ImageNet. While large ViT models perform worse than BiT ResNets (shaded area) when pre-trained on small datasets, they shine when pre-trained on larger datasets. Similarly, larger ViT variants overtake smaller ones as the dataset grows.

BiT: [Big Transfer](#) (ResNet)
ViT: Vision Transformer (Base/Large/Huge,
patch size of 14x14, 16x16, or 32x32)

[Internal Google dataset](#) (not public)

Masked autoencoders are scalable vision learners

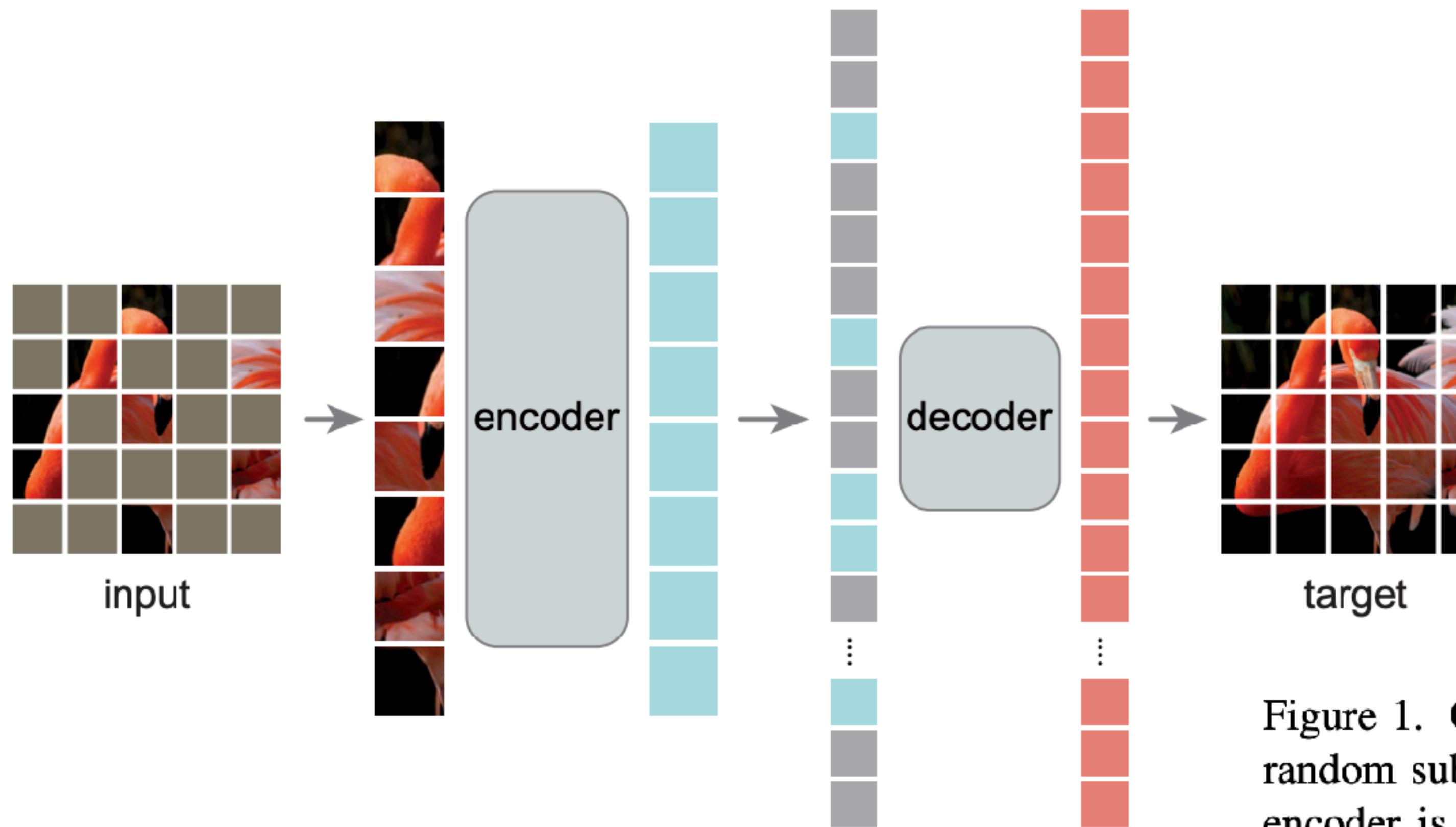


Figure 1. **Our MAE architecture.** During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images to produce representations for recognition tasks.

Masked autoencoders are scalable vision learners

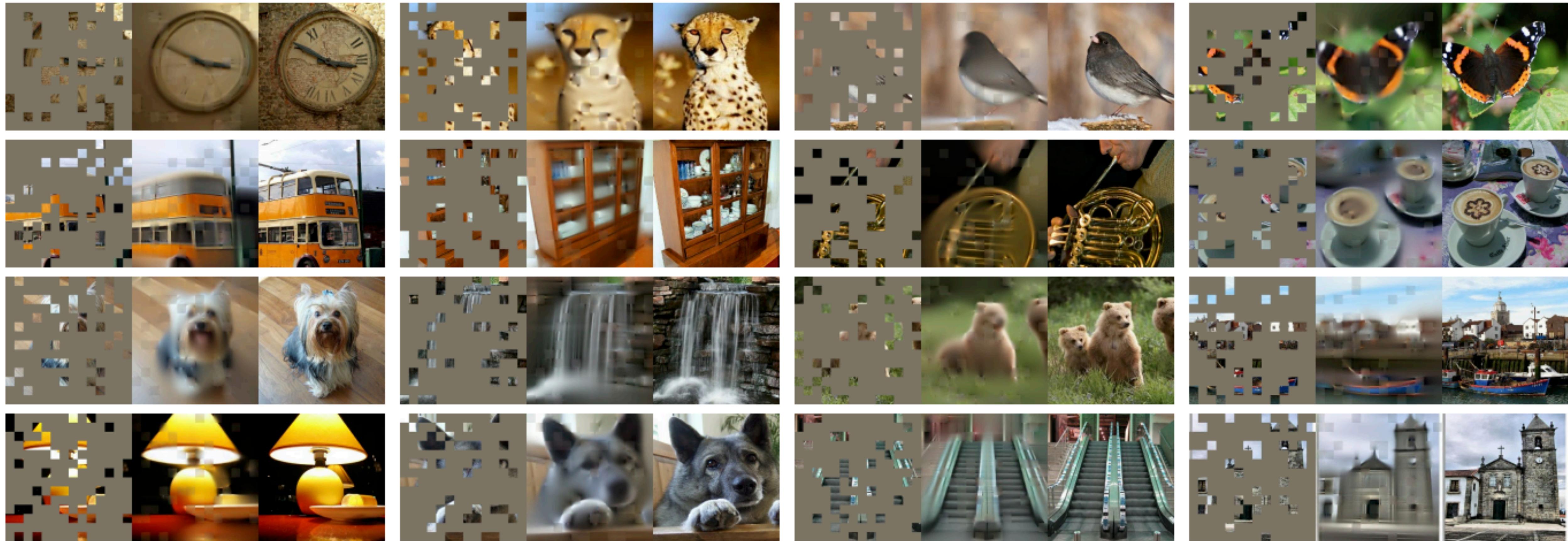
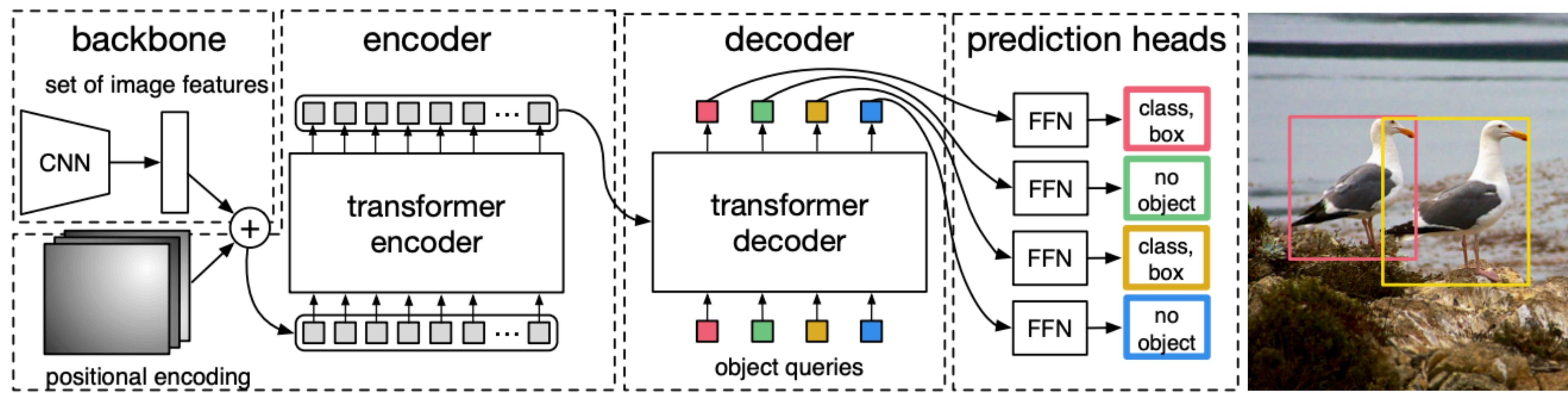
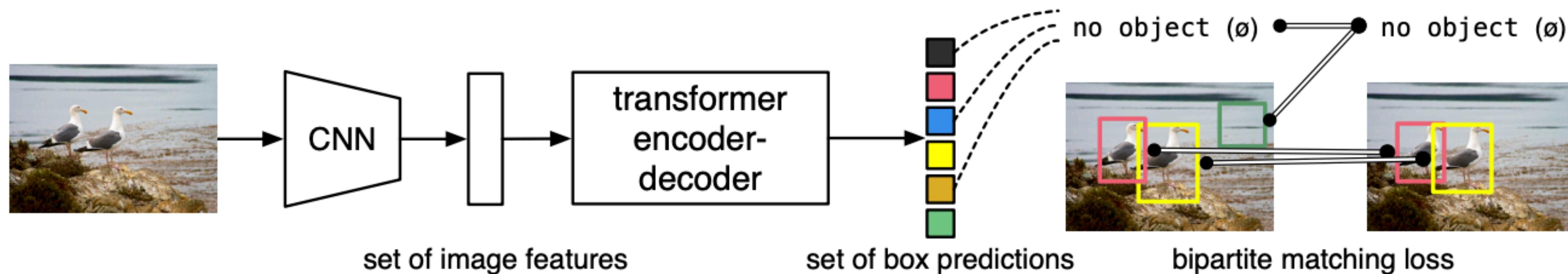


Figure 2. Example results on ImageNet *validation* images. For each triplet, we show the masked image (left), our MAE reconstruction[†] (middle), and the ground-truth (right). The masking ratio is 80%, leaving only 39 out of 196 patches. More examples are in the appendix.

[†]*As no loss is computed on visible patches, the model output on visible patches is qualitatively worse. One can simply overlay the output with the visible patches to improve visual quality. We intentionally opt not to do this, so we can more comprehensively demonstrate the method’s behavior.*

Detection Transformer (DETR)

- Hybrid of CNN and transformer, aimed at standard recognition task



Do we need attention?



Posted by u/L-MK 7 months ago 3 4 4



576 [R] Do You Even Need Attention? A Stack of Feed-Forward Layers Does Surprisingly Well on ImageNet

Research

Do we need attention?

MLP-Mixer: An all-MLP Architecture for Vision

Ilya Tolstikhin*, Neil Houlsby*, Alexander Kolesnikov*, Lucas Beyer*,
Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner,
Daniel Keysers, Jakob Uszkoreit, Mario Lucic, Alexey Dosovitskiy

*equal contribution

Google Research, Brain Team

{tolstikhin, neilhoulsby, akolesnikov, lbeyer,
xzhai, unterthiner, jessicayung[†], andstein,
keysers, usz, lucic, adosovitskiy}@google.com

4 May 2021

Do You Even Need Attention? A Stack of Feed-Forward Layers Does Surprisingly Well on ImageNet

Luke Melas-Kyriazi
Oxford University

lukemk@robots.ox.ac.uk

6 May 2021

Pay Attention to MLPs

Hanxiao Liu, Zihang Dai, David R. So, Quoc V. Le
Google Research, Brain Team
{hanxiao.liu, zihang.dai, david.so, qvl}@google.com

17 May 2021

Recent Hype#2: MLPs (!)

- Back to basics
- MLPs perform similar to Transformers while being more efficient
- CNNs and MLPs complexity linear with the number of input pixels, Transformers quadratic

MLP-Mixer: An all-MLP Architecture for Vision

Ilya Tolstikhin*, Neil Houlsby*, Alexander Kolesnikov*, Lucas Beyer*,
Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner,
Daniel Keysers, Jakob Uszkoreit, Mario Lucic, Alexey Dosovitskiy

*equal contribution

Google Research, Brain Team

{tolstikhin, neilhoulsby, akolesnikov, lbeyer,
xzhai, unterthiner, jessicayung[†], andstein,
keysers, usz, lucic, adosovitskiy}@google.com

4 May 2021

Do You Even Need Attention? A Stack of Feed-Forward Layers Does Surprisingly Well on ImageNet

Luke Melas-Kyriazi
Oxford University
lukemk@robots.ox.ac.uk

6 May 2021

Pay Attention to MLPs

Hanxiao Liu, Zihang Dai, David R. So, Quoc V. Le
Google Research, Brain Team
{hanxiao1, zihangd, davidso, qvl}@google.com

17 May 2021

Summary: Beyond CNNs

- CNNs (convolution), Transformers (attention), MLPs (fully connected)
- There is no answer to which architecture is better.
- Often depends on the data.
- If you have infinite data, more complex can be better (e.g., MLP ~ Transformers > CNN).
- Similar performance can be obtained with more efficient models (e.g., MLP ~ Transformers)
- It is possible there will be newer/better architectures/hypes before you graduate. Stay tuned.

Summary: Beyond CNNs

- Some Advantages of Transformers:
 - Flexibility to handle variable input sequences
 - Flexibility to combine multimodal data [more on this next week, in Vision & Language]

2. Beyond classification

- **a. Intro to structured outputs**
- **b. Object detection (localization)**
- **c. Segmentation**
- **d. Human pose estimation**

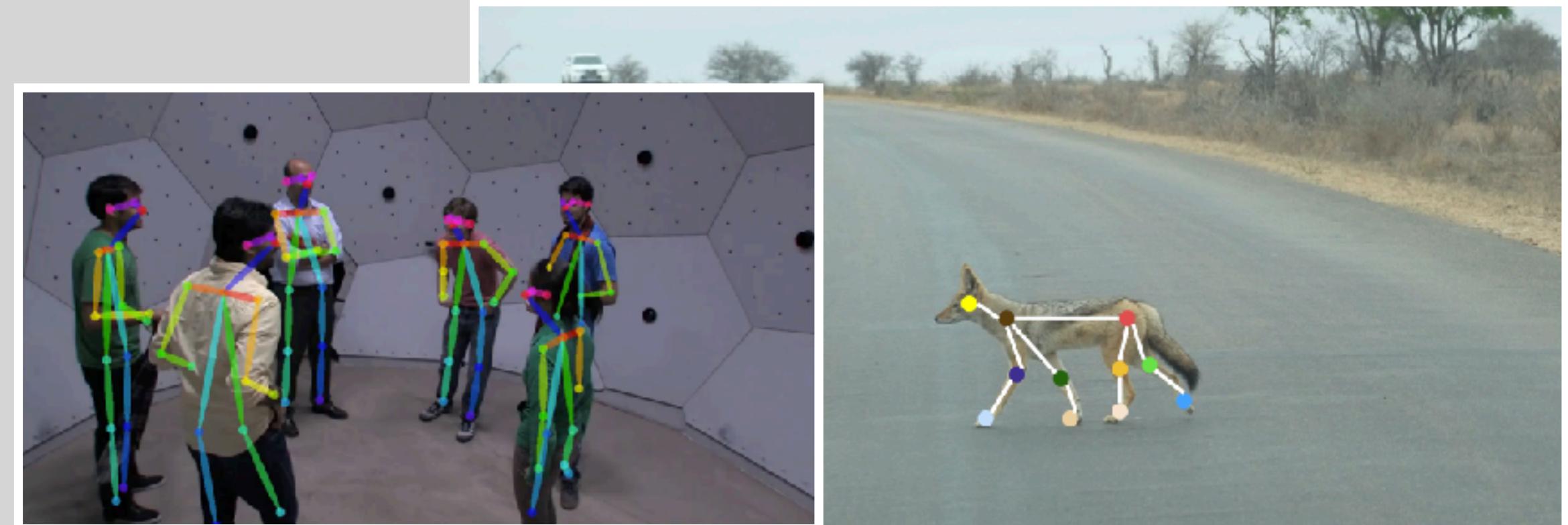
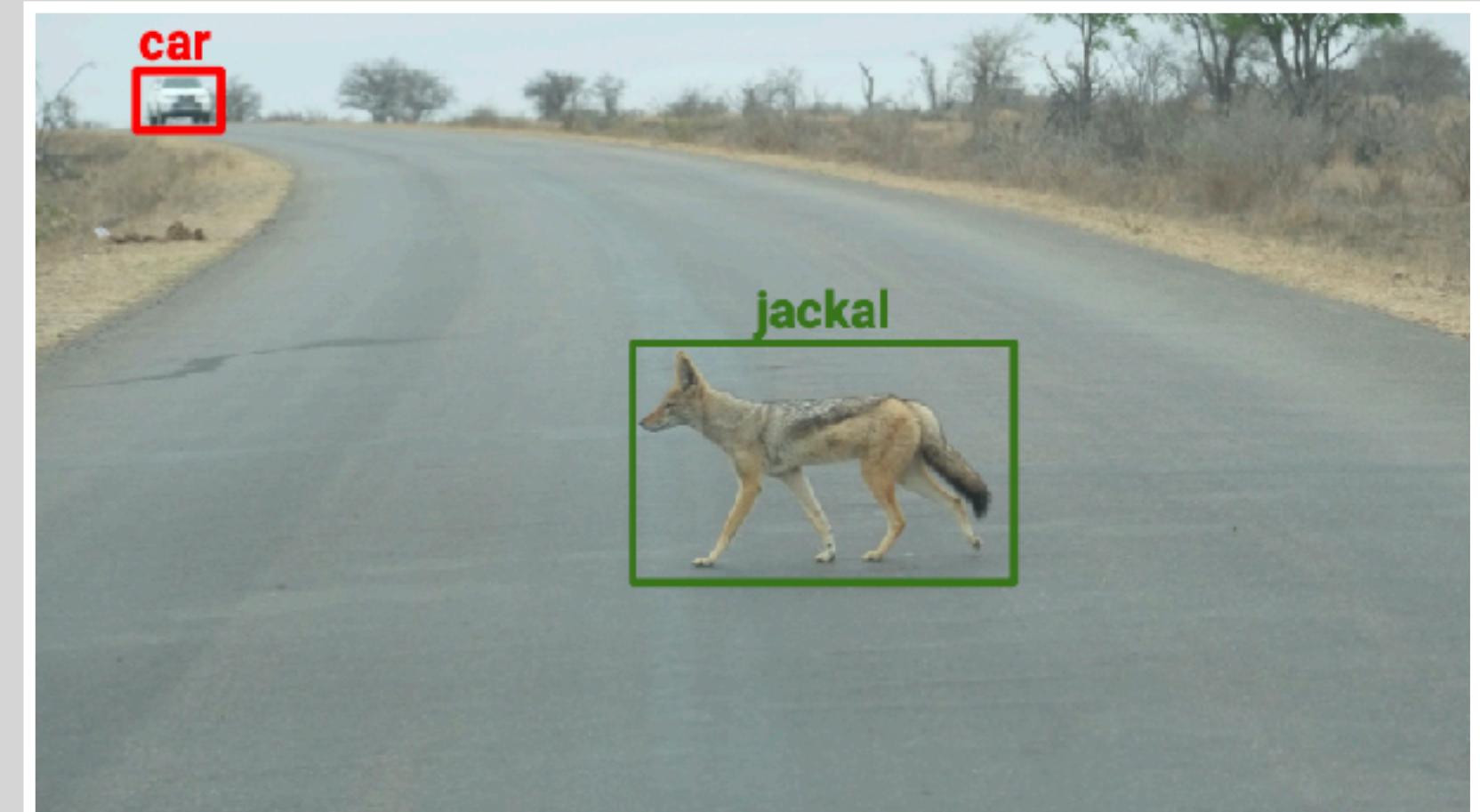
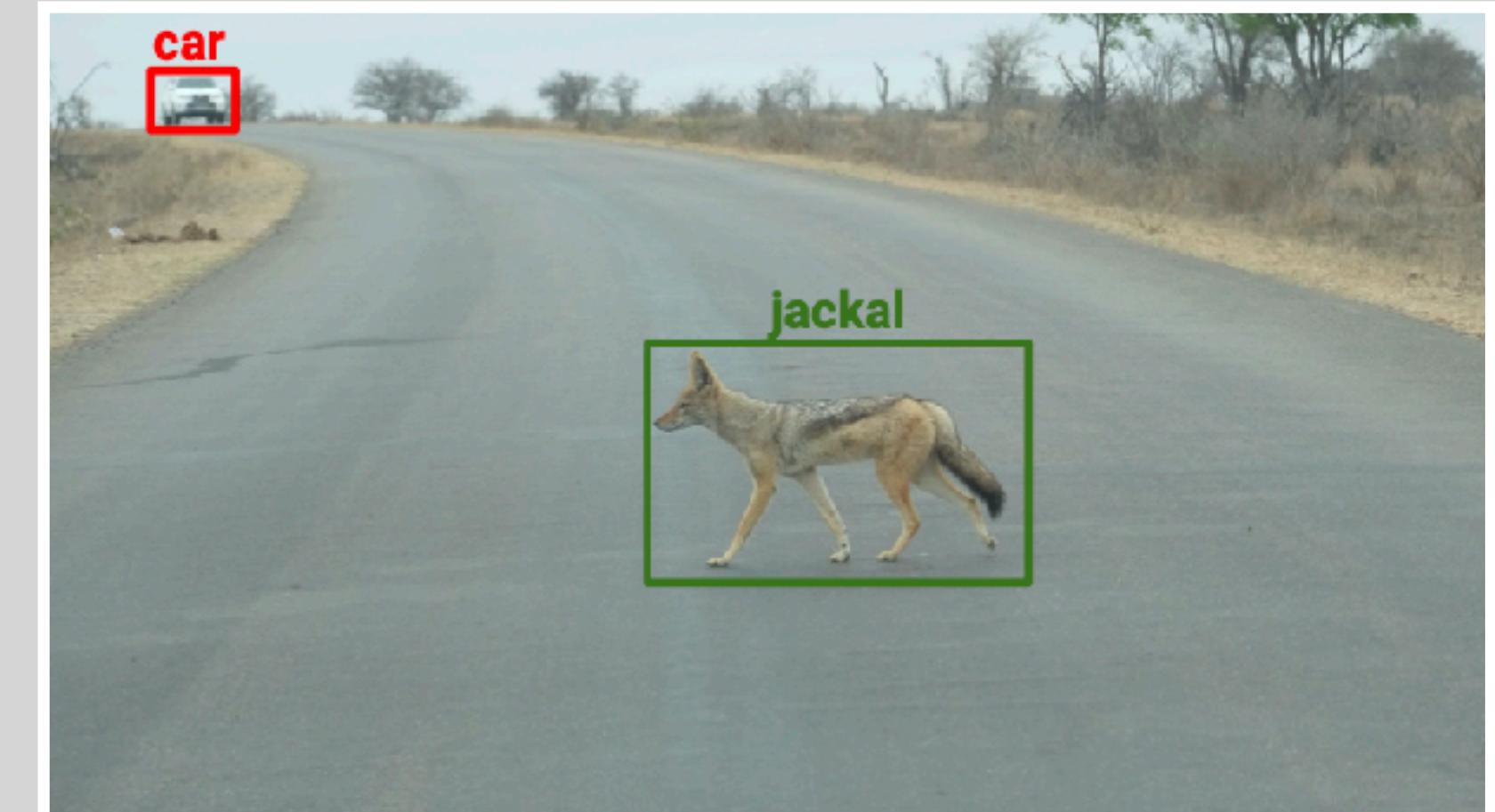


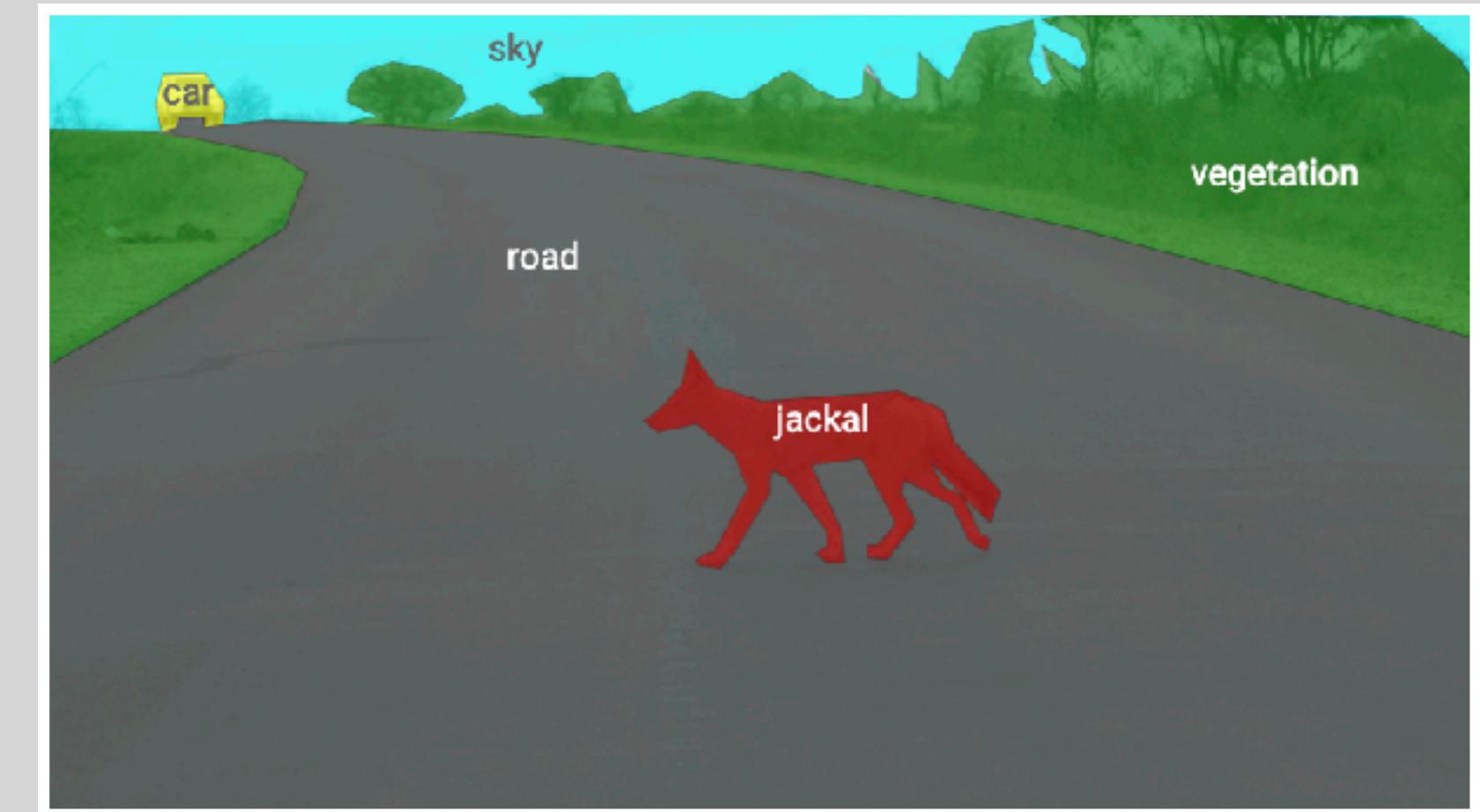
Image credits: Naila Murray

2. Beyond classification

- **a. Intro to structured outputs**



- **b. Object detection (localization)**



- **c. Segmentation**



- **d. Human pose estimation**



Image credits: Naila Murray

So far:

- Image in



- Category out

→ Motorbike

What we would like to do...

- Visual scene understanding
- What is in the image and where



- Object categories, identities, properties, activities, relations, ...

(Some) Fundamental Tasks in Computer Vision

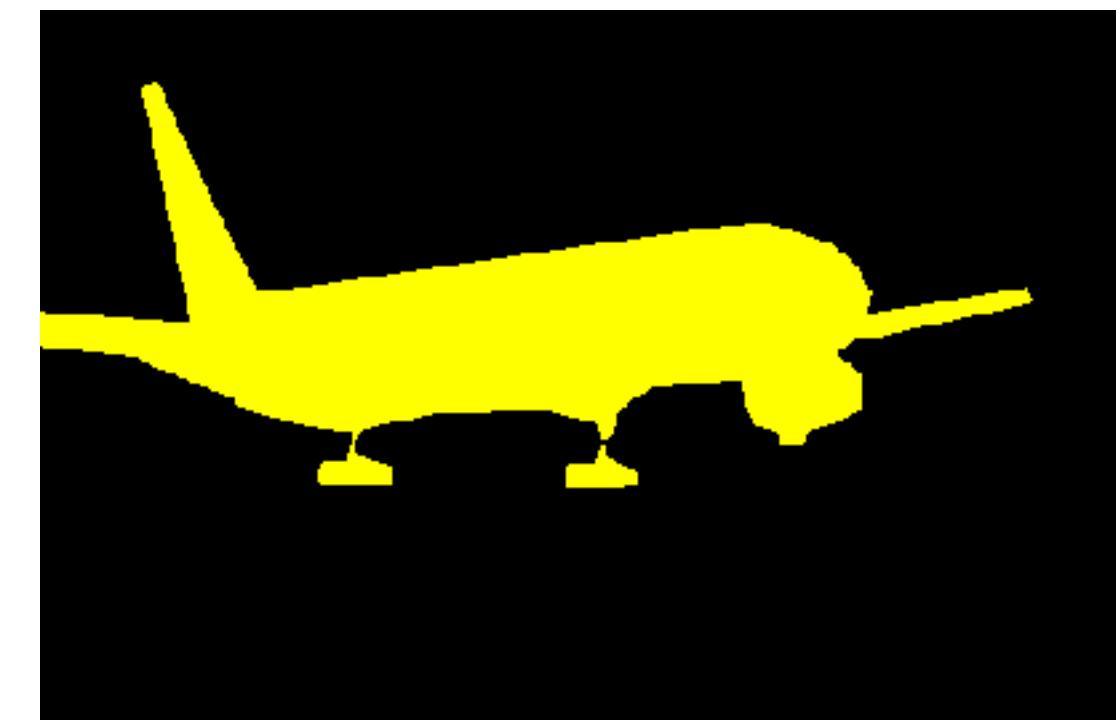
- **Image Classification**
 - Does the image contain an aeroplane?
(last lecture)



- **Object Class Detection/Localization**
 - Where are the aeroplanes (if any)?



- **Object Class Segmentation**
 - Which pixels are part of an aeroplane (if any)?

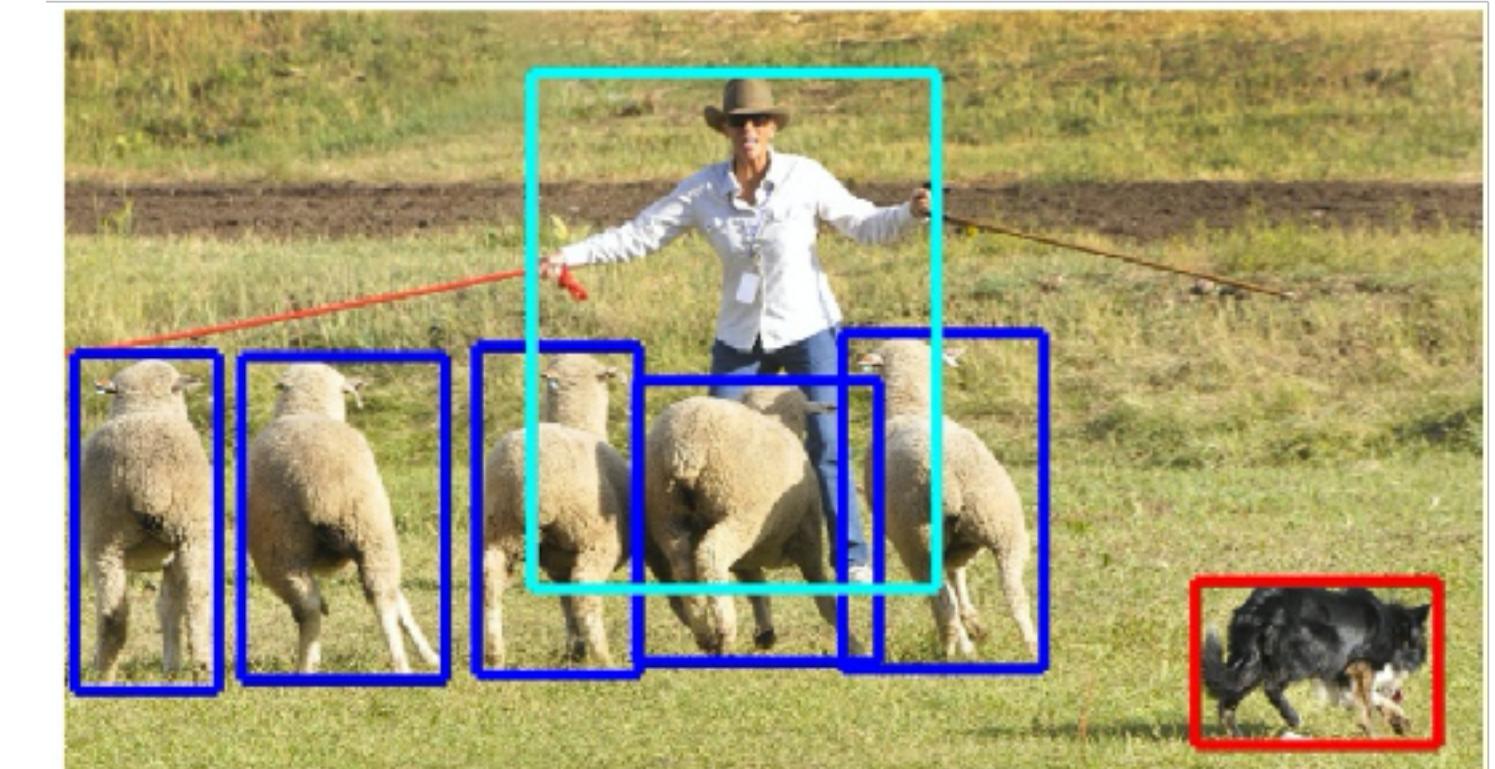


(Some) Fundamental Tasks in Computer Vision

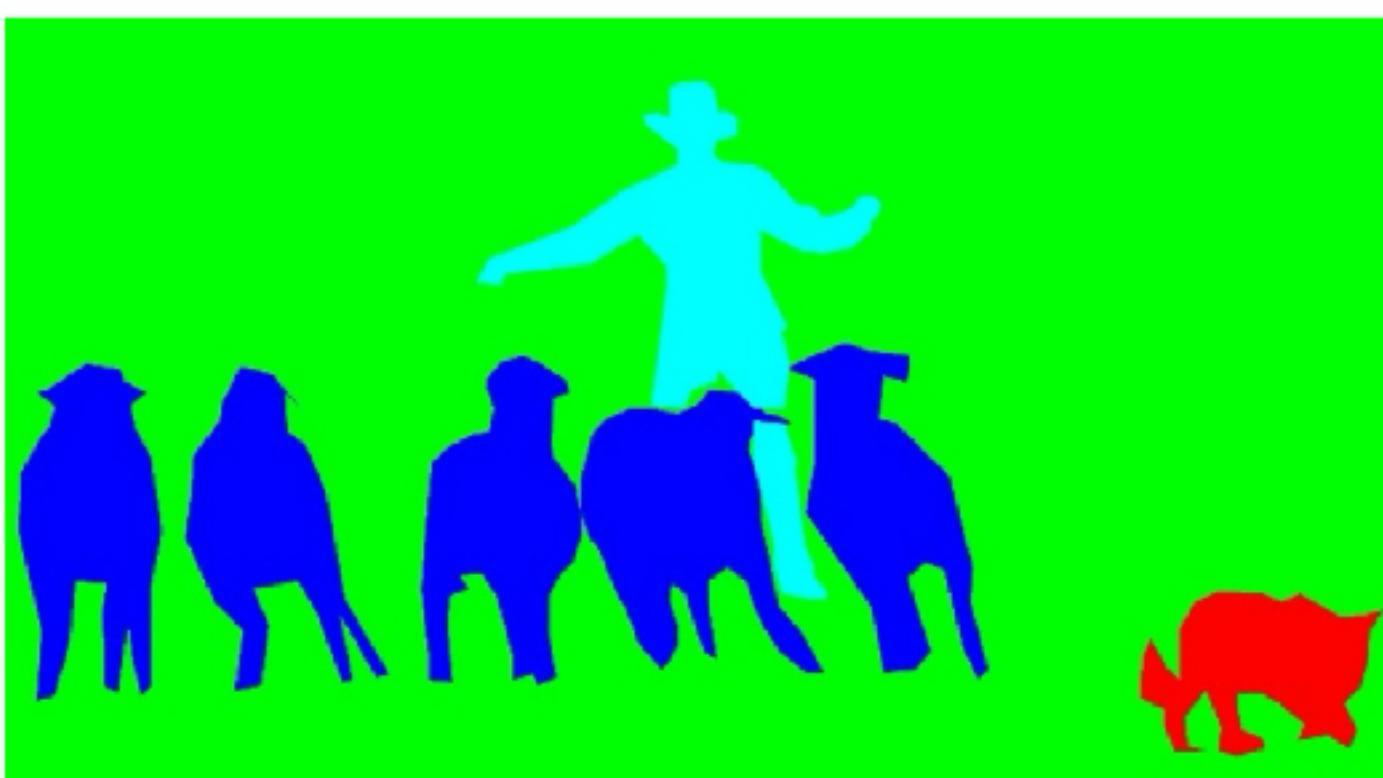
- Representing objects in the image:
 - Class labels
 - Bounding box
 - Semantic pixel-wise labels
 - Instance pixel-wise labels



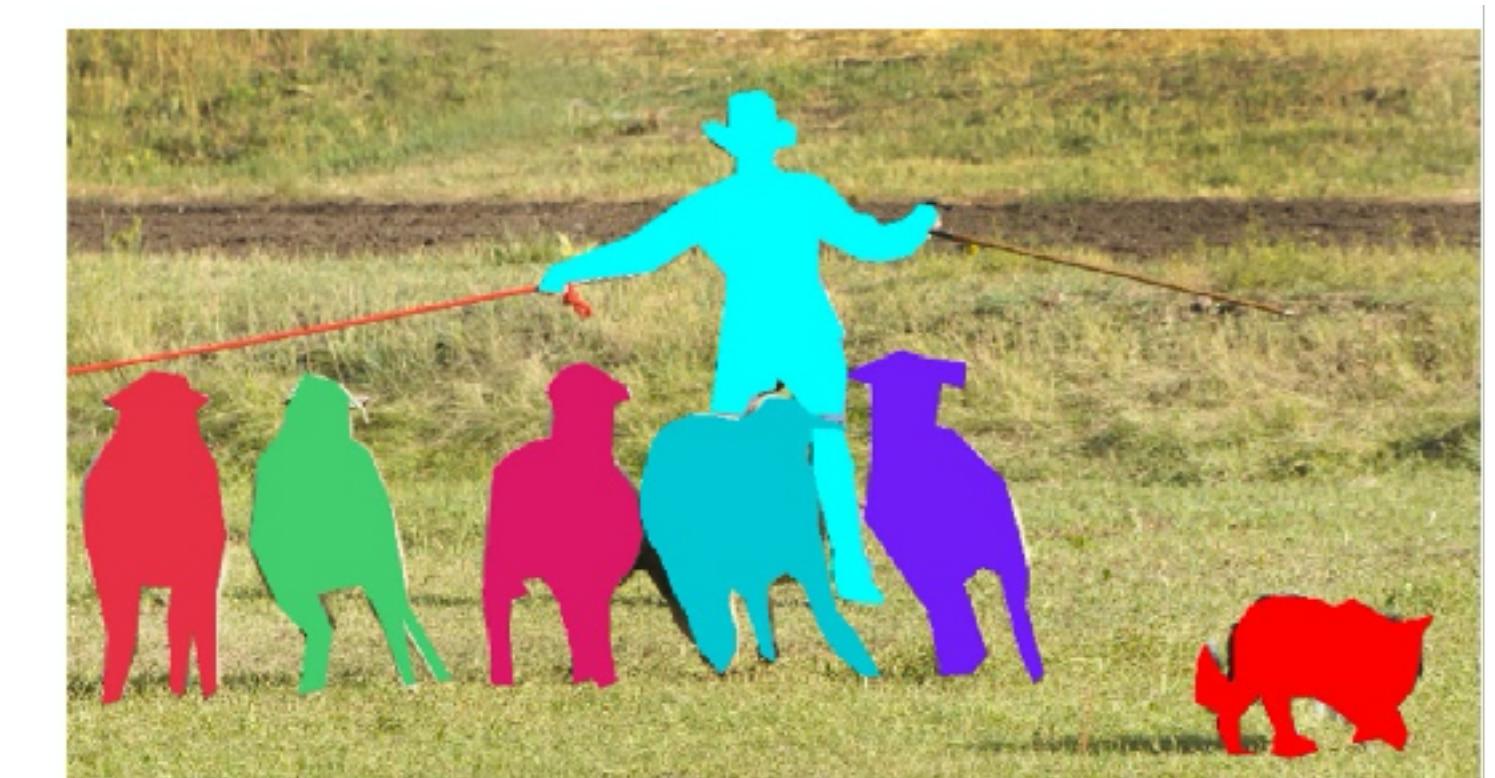
Image Classification



Object Detection

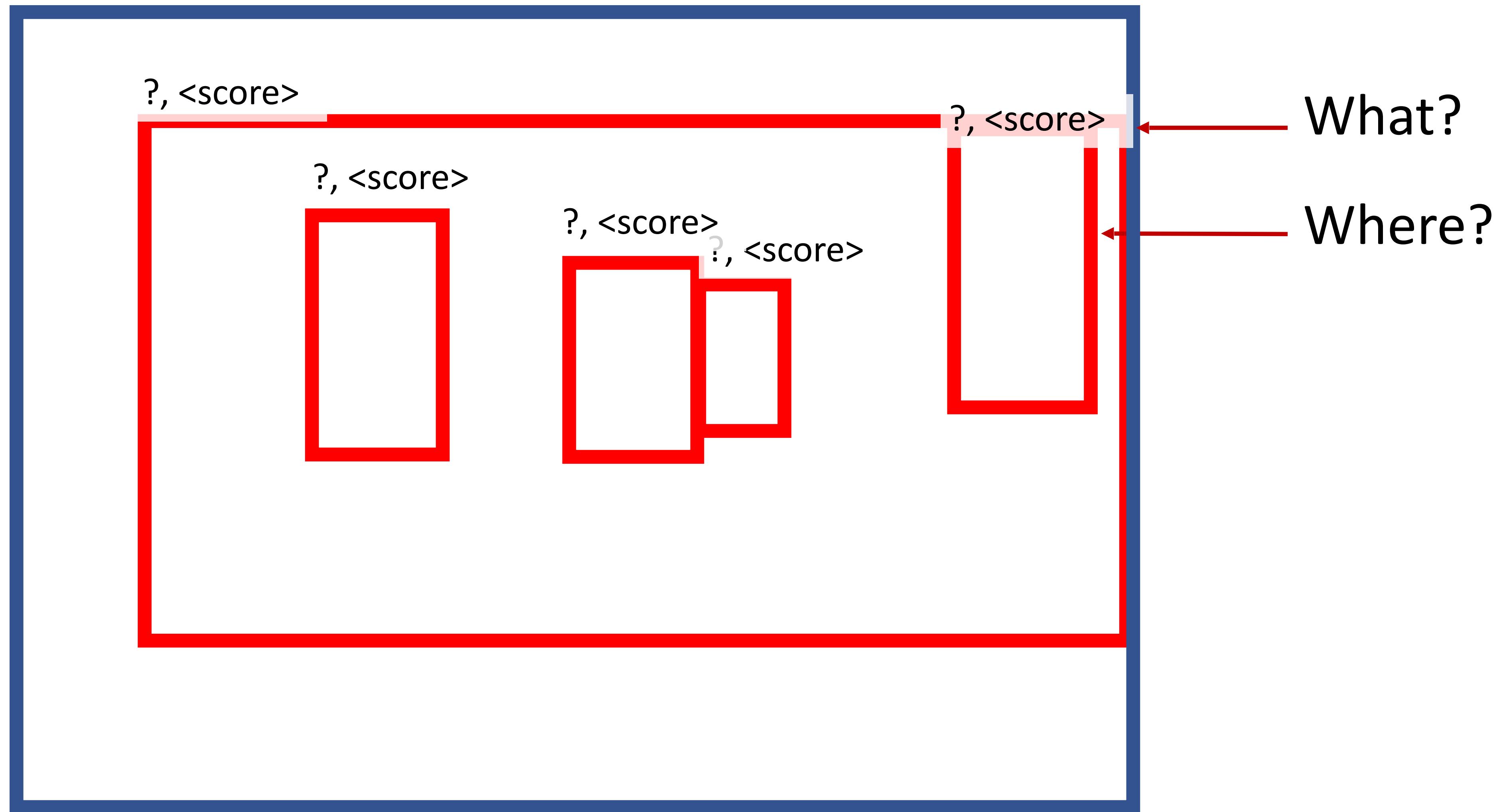


Semantic Segmentation



Instance Segmentation

Object Detection



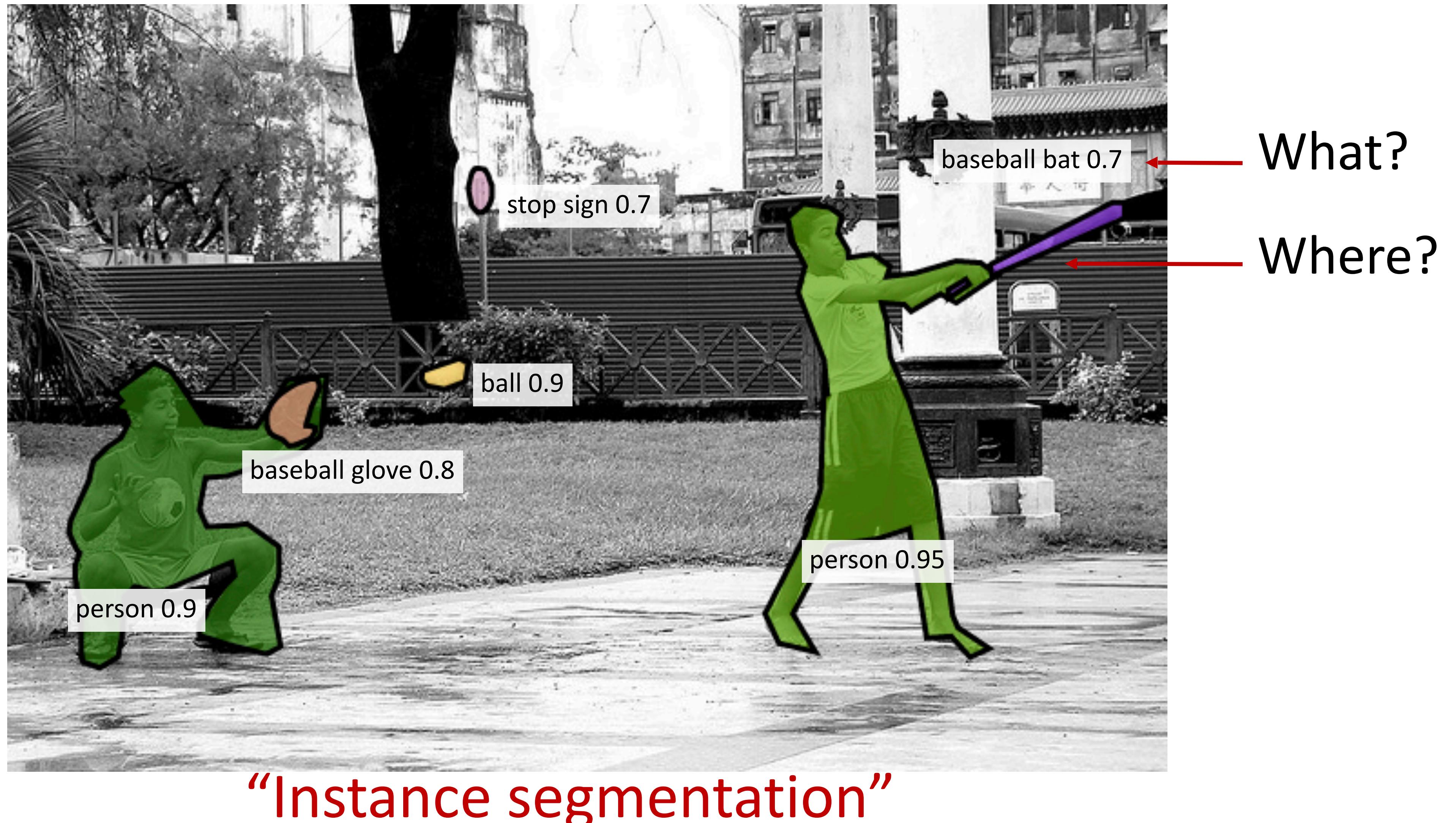
Object Detection with Bounding Boxes



What?
Where?

“Object detection”

Object Detection with Segmentation Masks



Classification vs. Detection

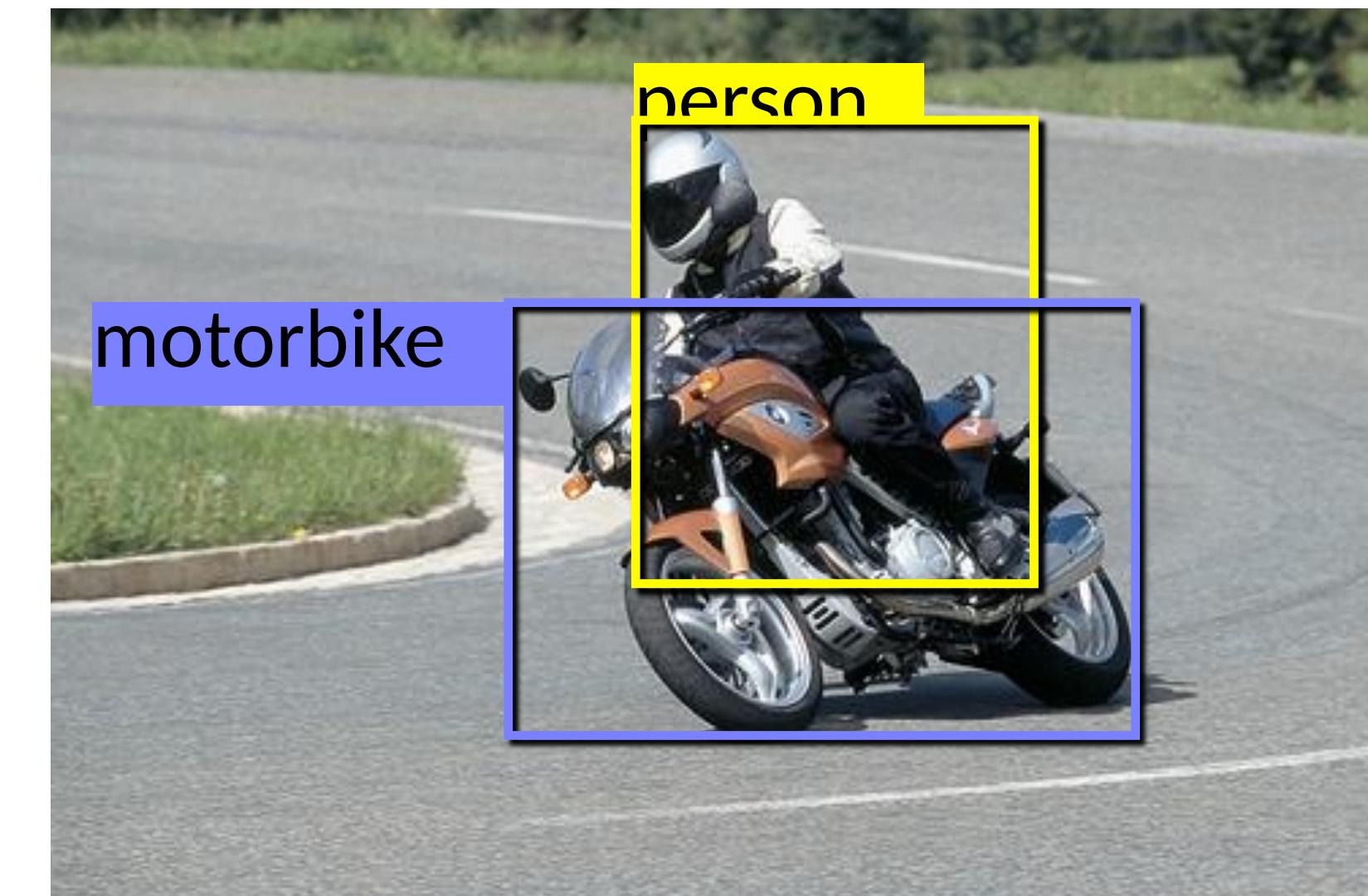


Problem formulation

{ airplane, bird, motorbike, person, sofa }



Input



Desired output

Things vs. Stuff

Ted Adelson, Forsyth et al. 1996.

Thing (n): An object with a specific size and shape.



Stuff (n): Material defined by a homogeneous or repetitive pattern of fine-scale properties, but has no specific or distinctive spatial extent or shape.



Panoptic segmentation



(a) image



(b) semantic segmentation



(c) instance segmentation

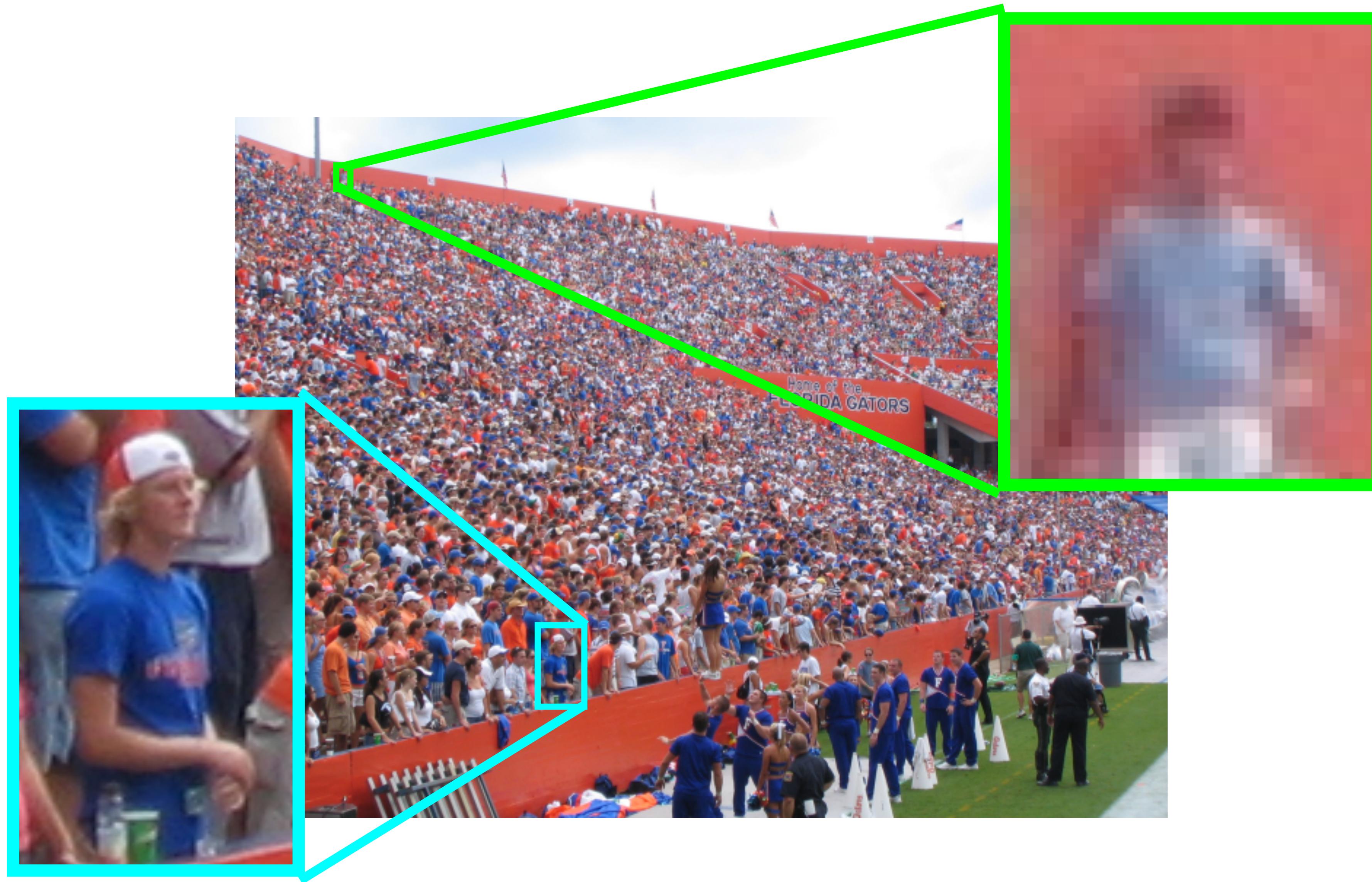


(d) panoptic segmentation

things – countable objects such as
people, animals, tools

stuff – amorphous regions of similar
texture or material such as **grass, sky, road**

Challenges: Scale



Challenges: Occlusion and truncation



Challenges: Background Clutter

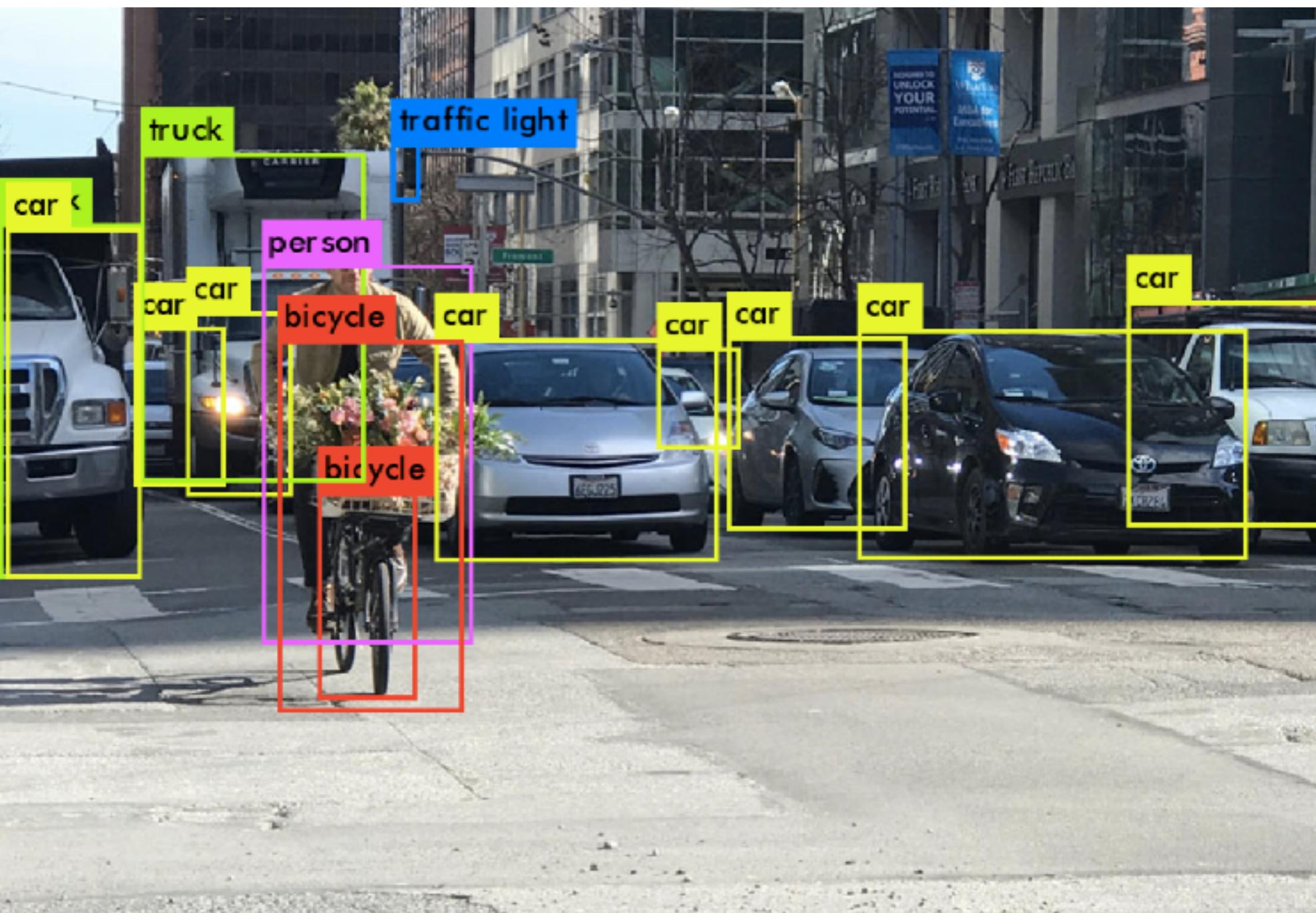


Challenges: Intra-class variation

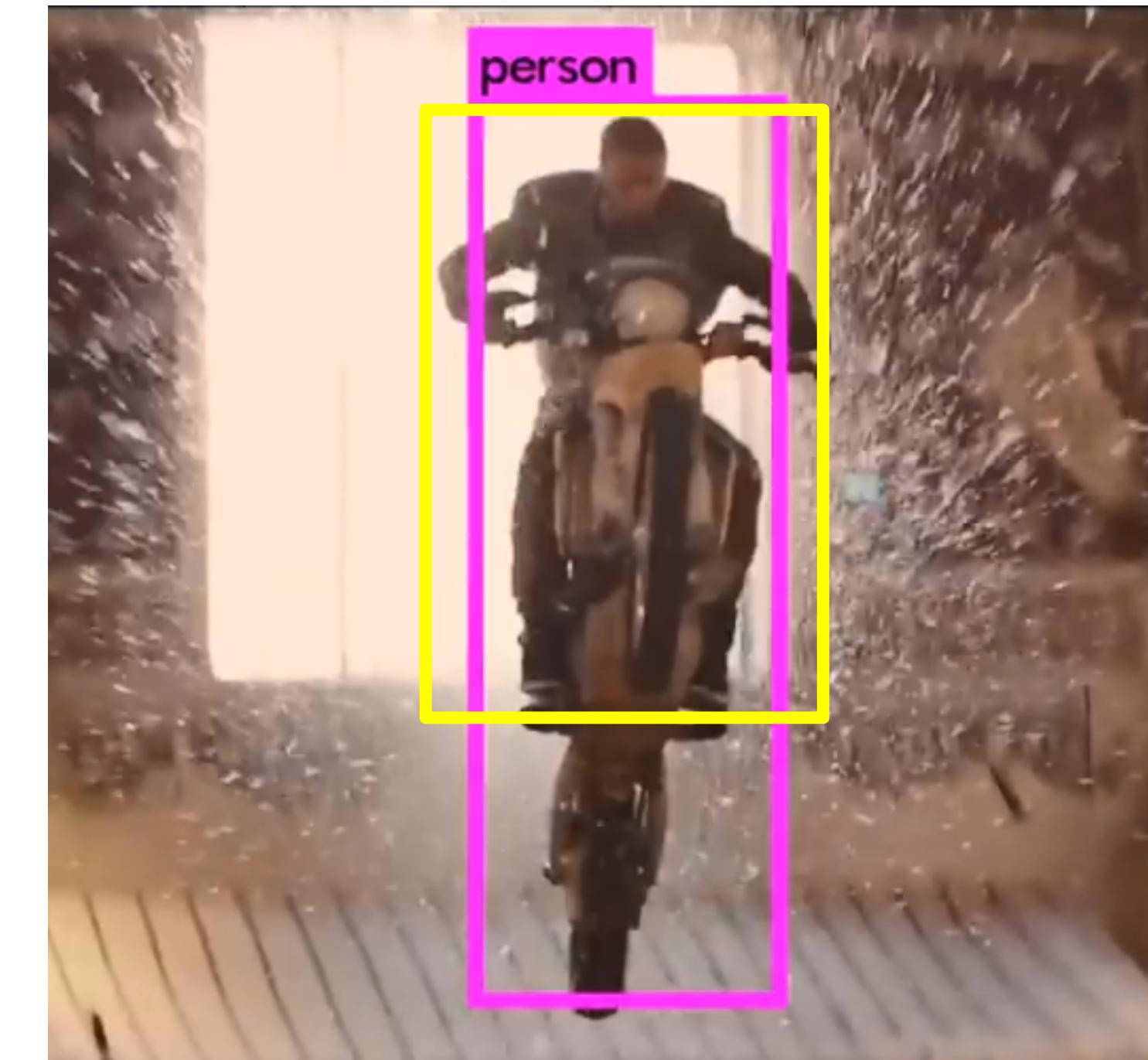


Challenges: How to evaluate object detection?

Images may contain many objects and classes



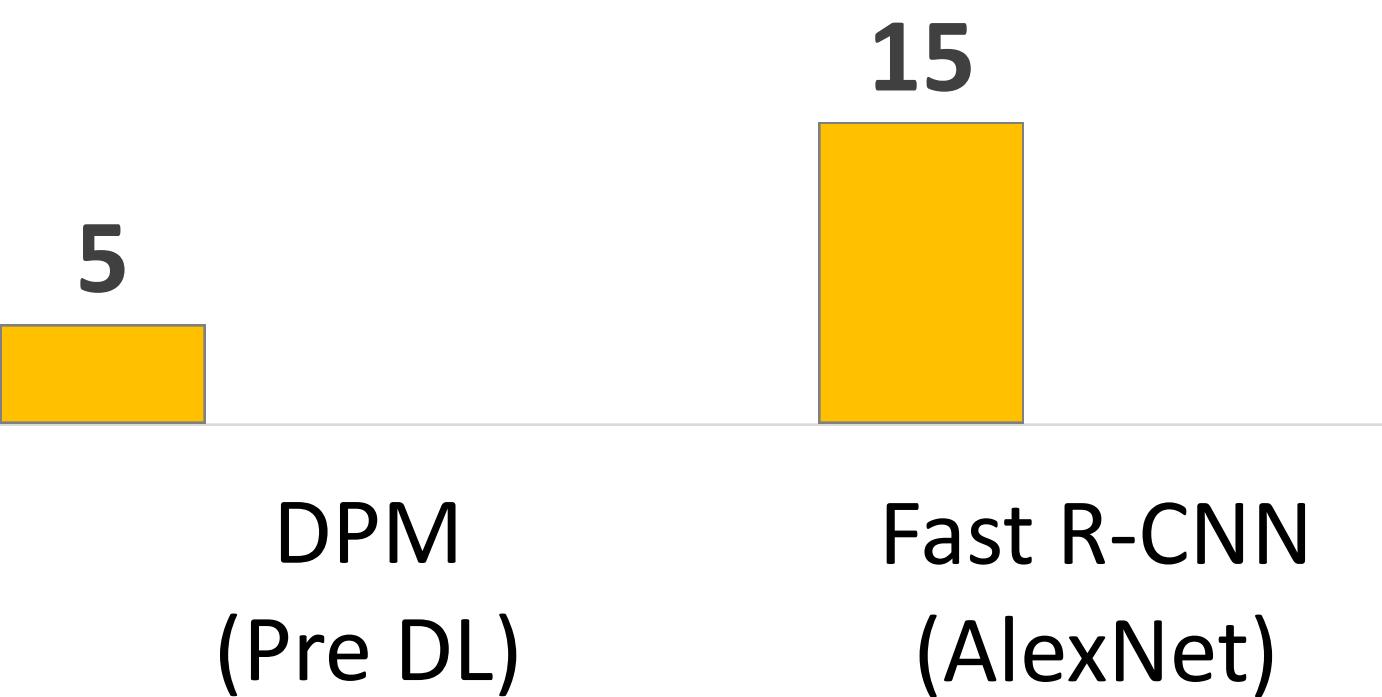
[Image source](#)



- Ground truth
- Detector output

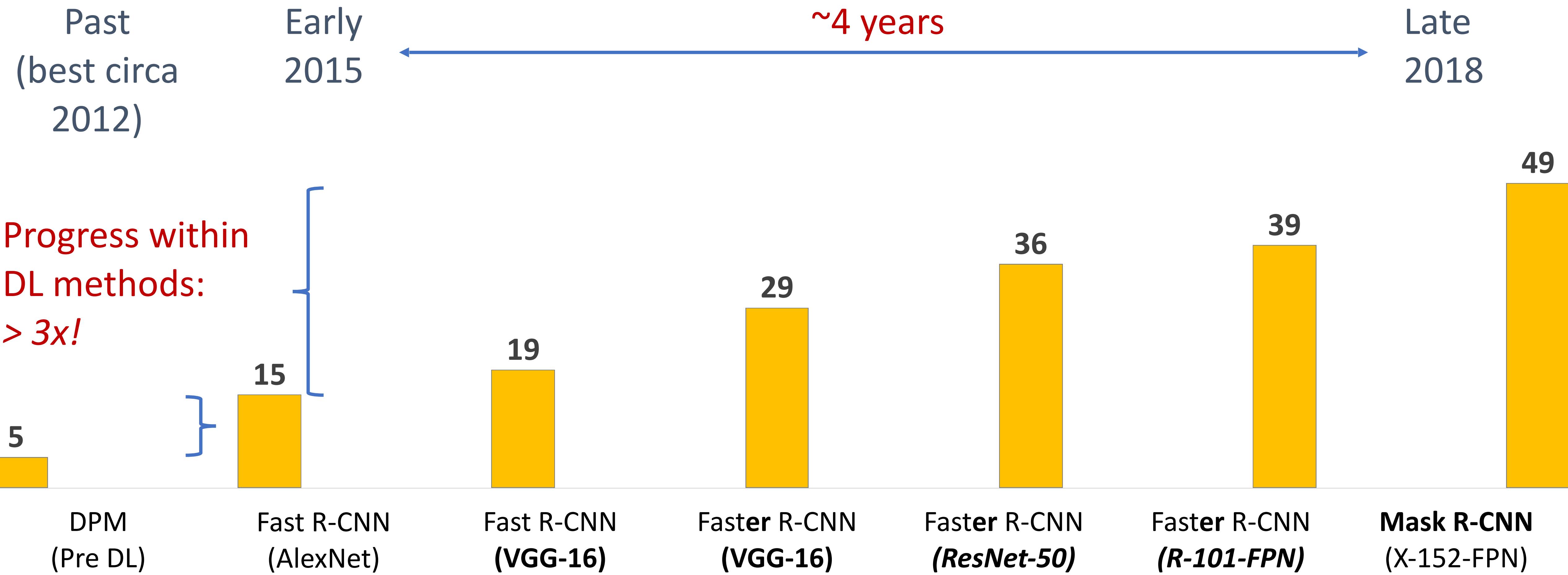
coco Object Detection Average Precision (%)

Past
(best circa
2012)



Movement to
Deep Learning methods:
3x improvement in AP

coco Object Detection Average Precision (%)



Steady Progress on Boxes and Masks

- **R-CNN** [Girshick et al. 2014]
- **SPP-net** [He et al. 2014]
- **Fast R-CNN** [Girshick. 2015]
- **Faster R-CNN** [Ren et al. 2015]
- **R-FCN** [Dai et al. 2016]
- Feature Pyramid Networks + Faster R-CNN [Lin et al. 2017]
- Mask R-CNN [He et al. 2017]
- Training with Large Minibatches (MegDet) [Peng, Xiao, Li, et al. 2017]
- Cascade R-CNN [Cai & Vasconcelos 2018]

DPM ... Fast R-CNN Fast R-CNN Faster R-CNN Faster R-CNN Faster R-CNN Mask R-CNN
(Pre DL) (AlexNet) (VGG-16) (VGG-16) (ResNet-50) (R-101-FPN) (X-152-FPN)

Beyond Boxes and Masks: Human Keypoints



COCO Keypoint Detection Task
[COCO team @ cocodataset.org 2016 - present]

Beyond Boxes and Masks: Human Surfaces



DensePose: Dense Human Pose Estimation In The Wild
[Güler, Neverova, Kokkinos CVPR 2018]

Beyond Boxes and Masks: 3D Shape

Input Image



2D Recognition



3D Meshes



3D Voxels

Mesh R-CNN

[Gkioxari, Malik, Johnson ICCV 2019]

2. Beyond classification

- a. Intro to structured outputs

- b. Object detection (localization)

- c. Segmentation

- d. Human pose estimation

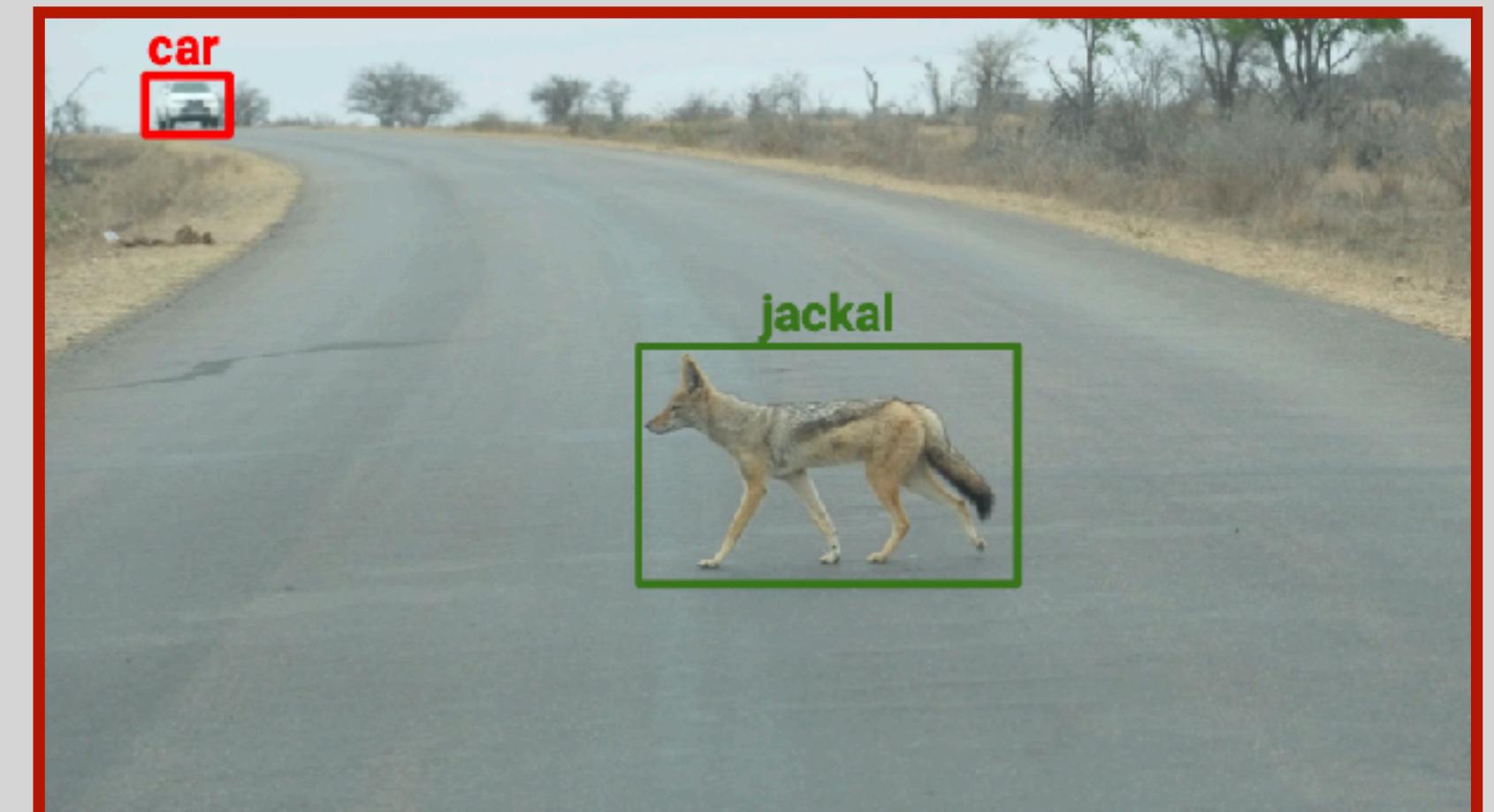
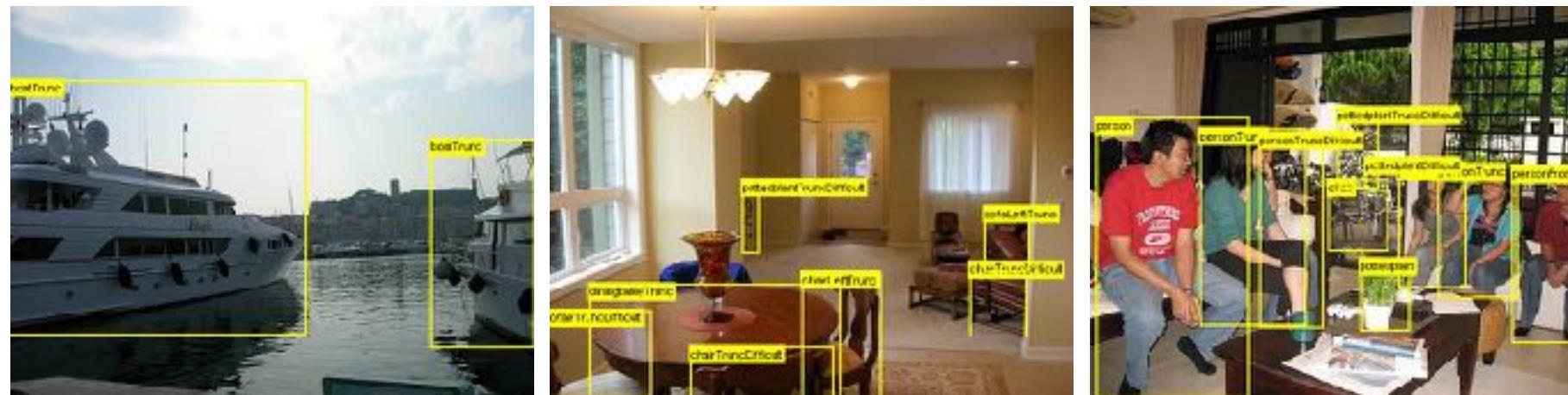


Image credits: Naila Murray

Object detection datasets (benchmarks)

Datasets	Categories	Images	Bounding Boxes
PASCAL-VOC	20	11K	27K
COCO	80 (91 stuff)	328K	2500K
LVIS	1200	164K	2.2M



PASCAL-VOC [2005-2008]



COCO [2014-2015]

Evaluating a detector



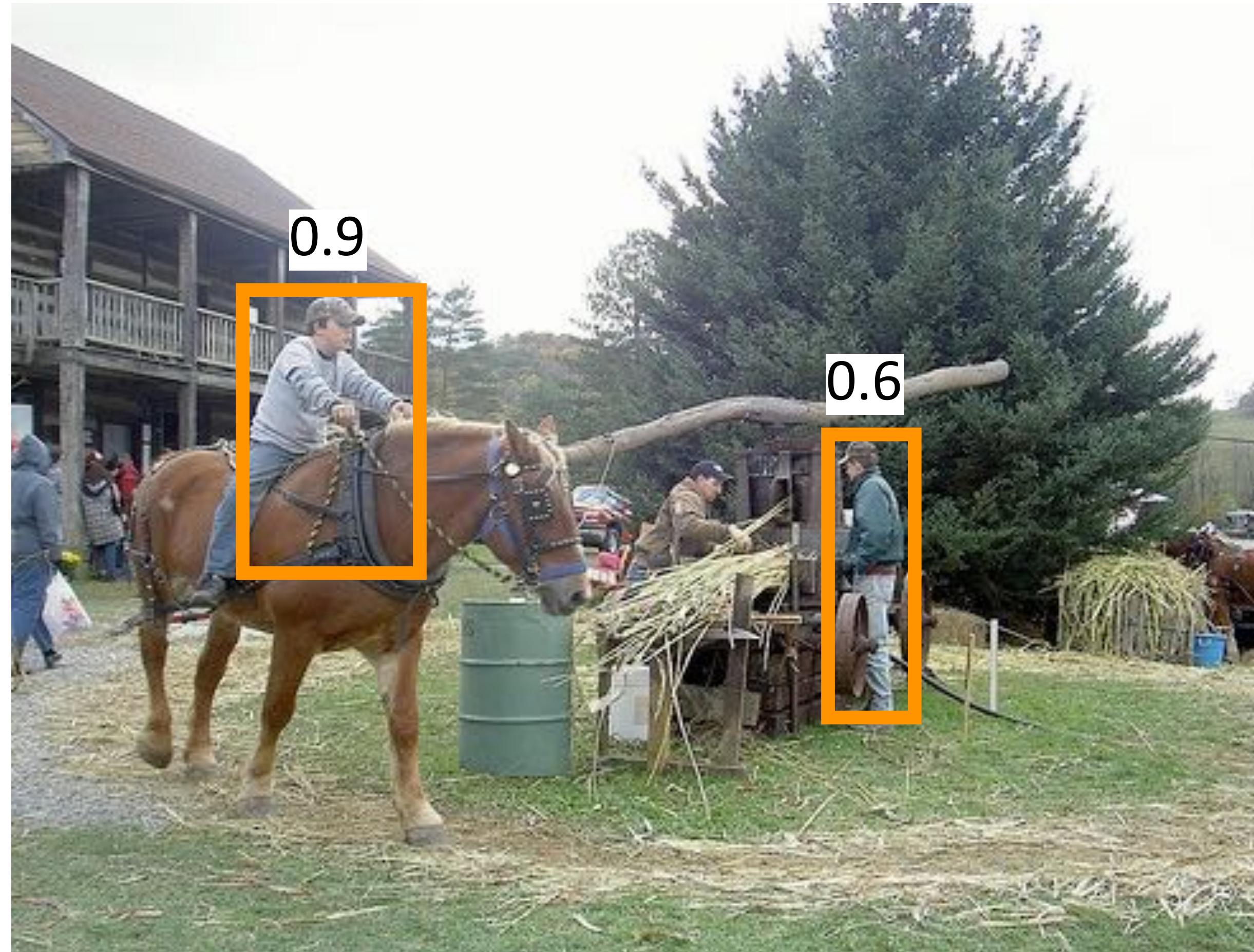
Test image (previously unseen)

First detection ...



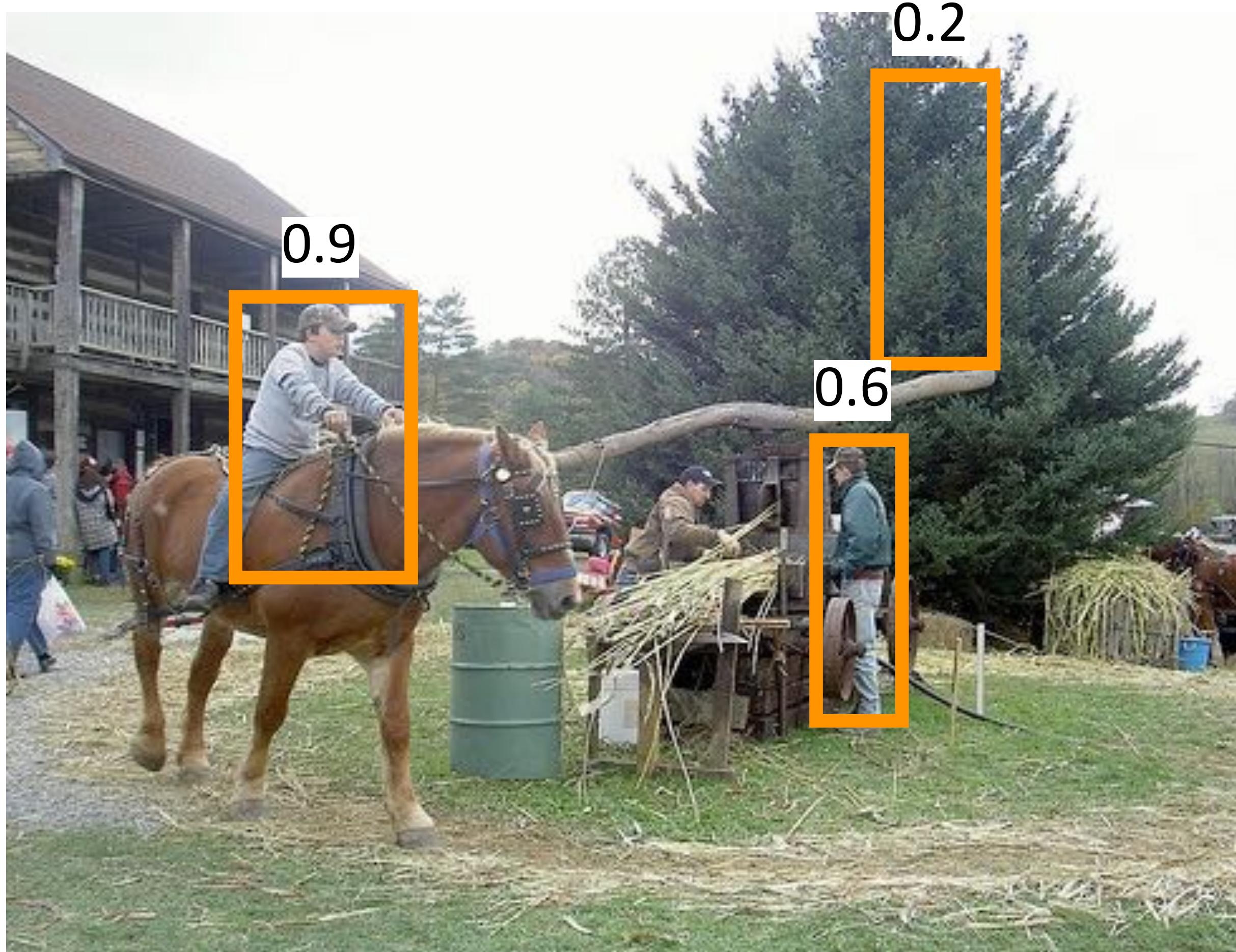
'person' detector predictions

Second detection ...



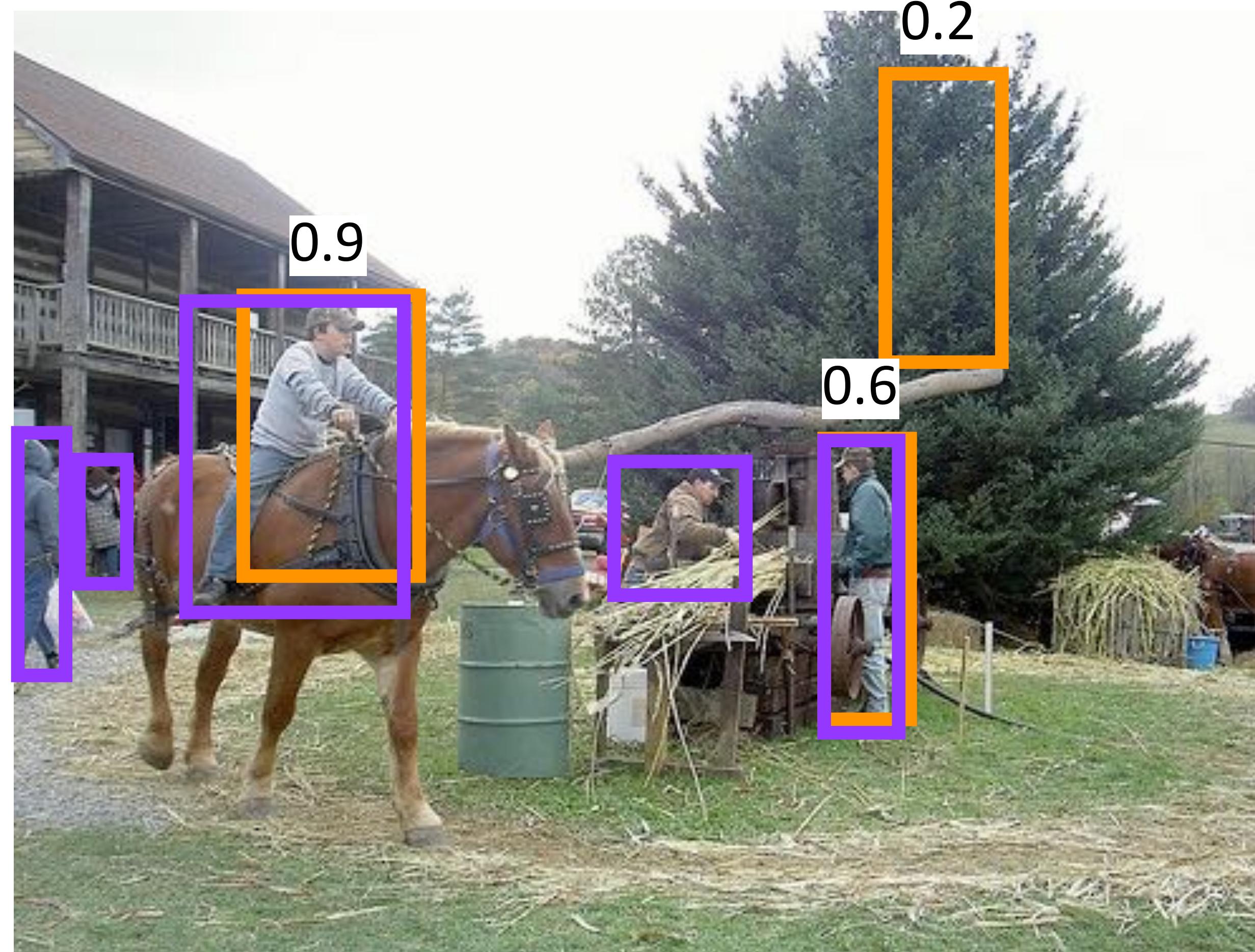
‘person’ detector predictions

Third detection ...



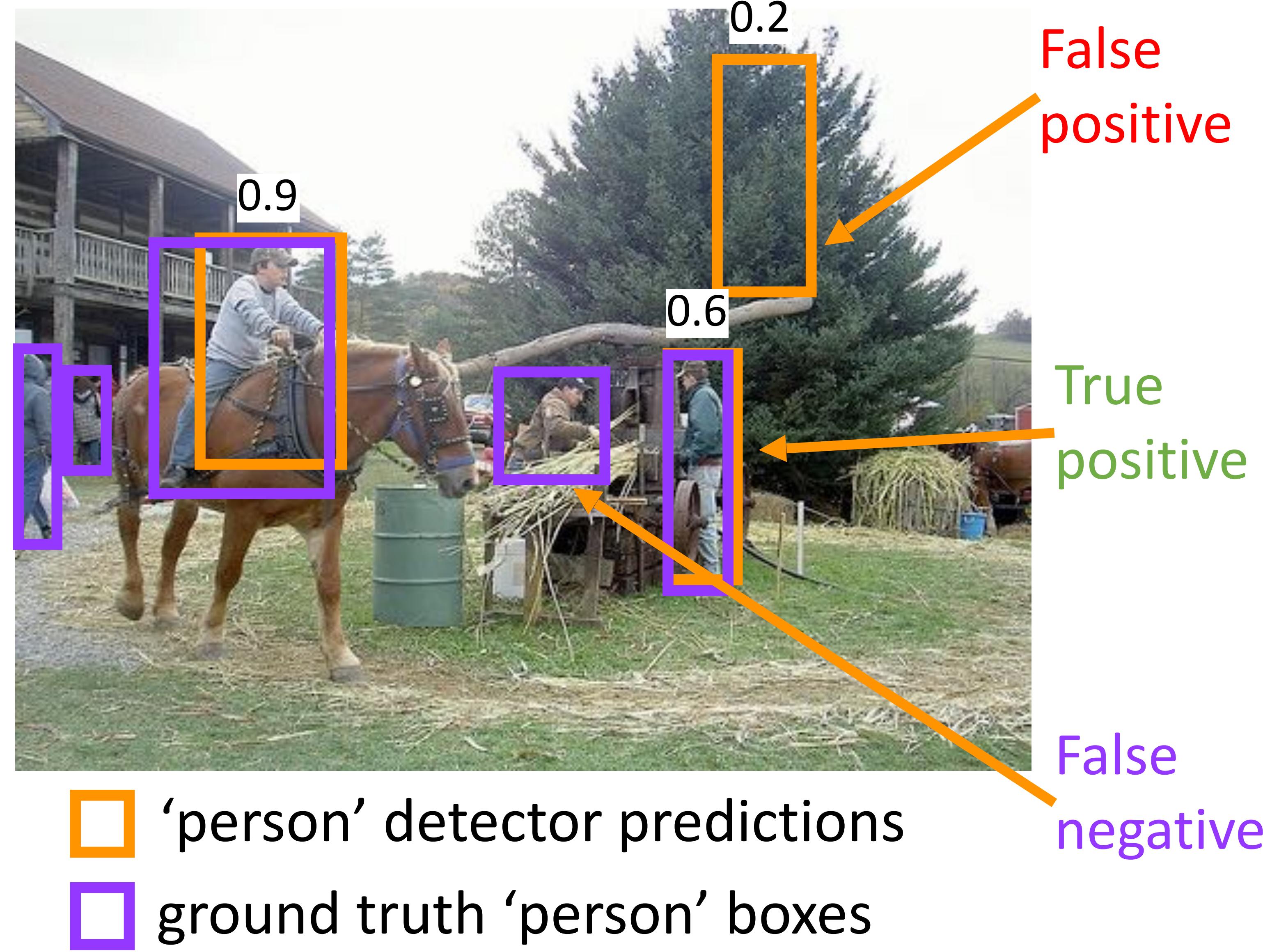
'person' detector predictions

Compare to ground truth

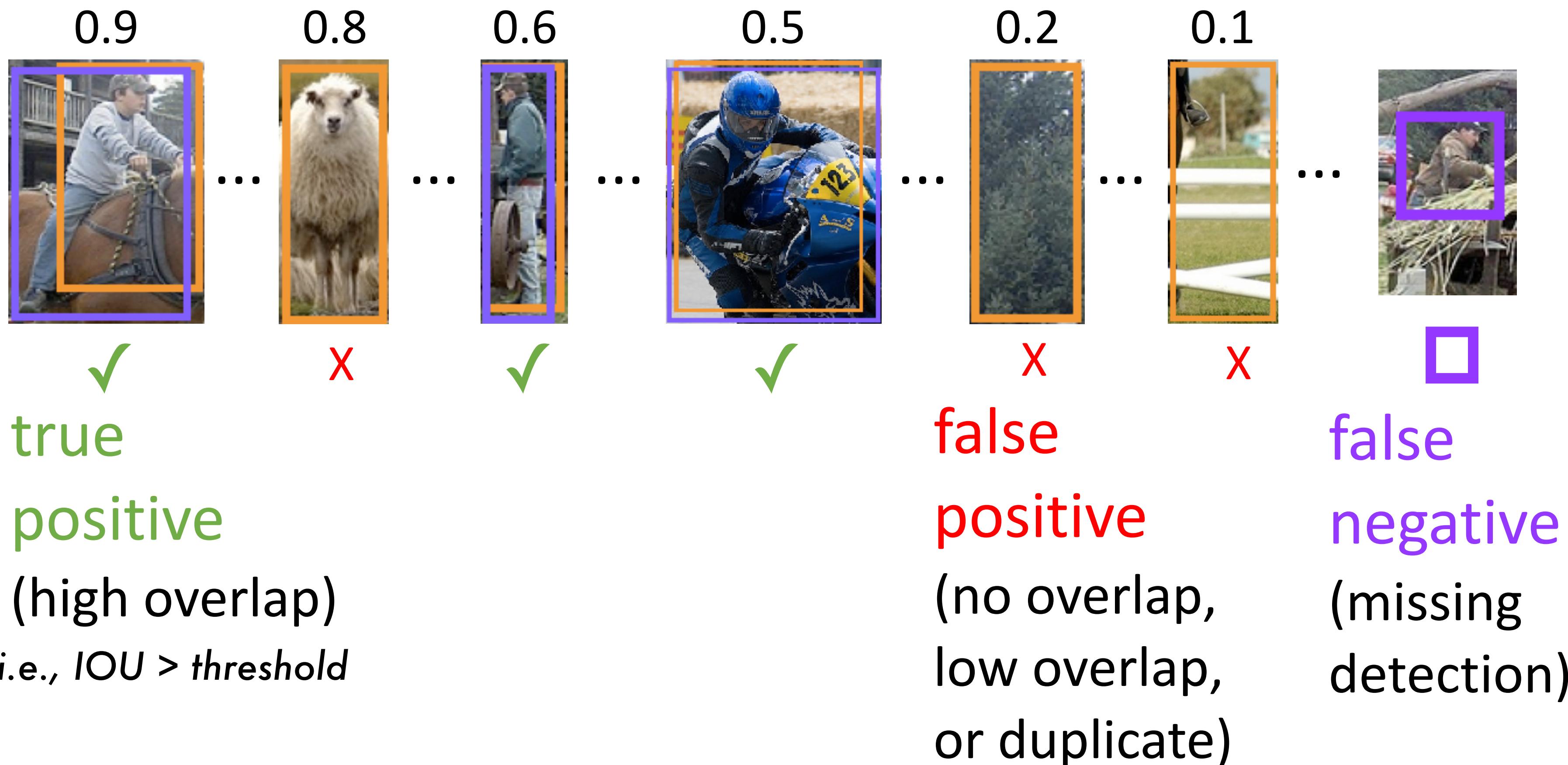


- ‘person’ detector predictions
- ground truth ‘person’ boxes

Compare to ground truth



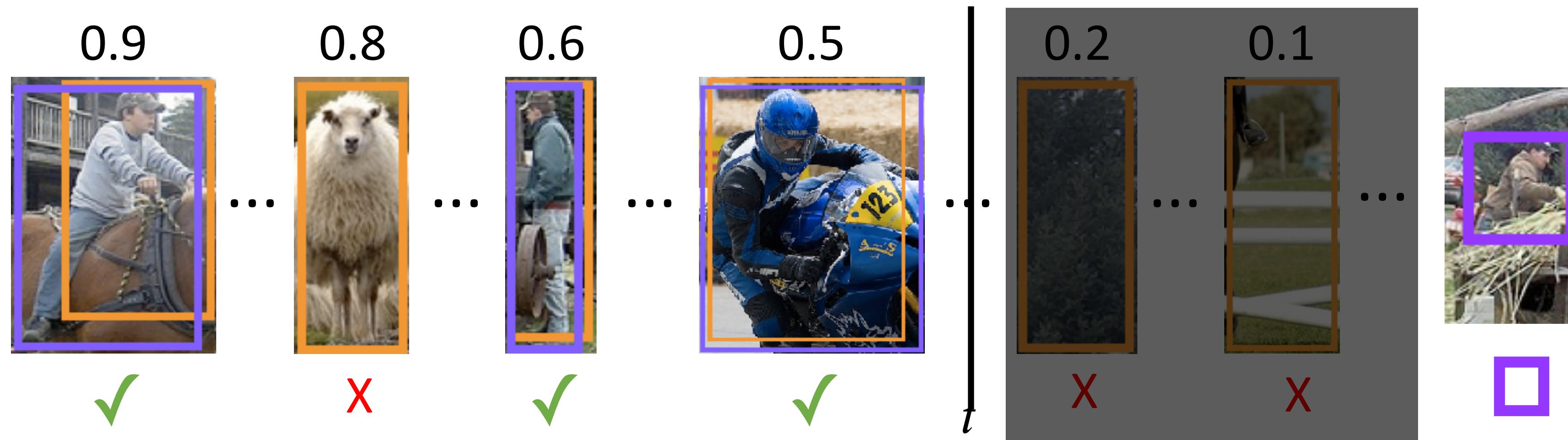
Sort by confidence



$$IOU \text{ (intersection over union)}: \frac{\text{intersection}}{\text{union}}$$

The diagram illustrates the Intersection over Union (IoU) formula. It shows two overlapping rectangles: a blue one and a red one. The intersection is the gray area where they overlap, and the union is the total area covered by either rectangle or both.

Sort by confidence



$$precision@t = \frac{\#true\ positives@t}{\#true\ positives@t + \#false\ positives@t}$$

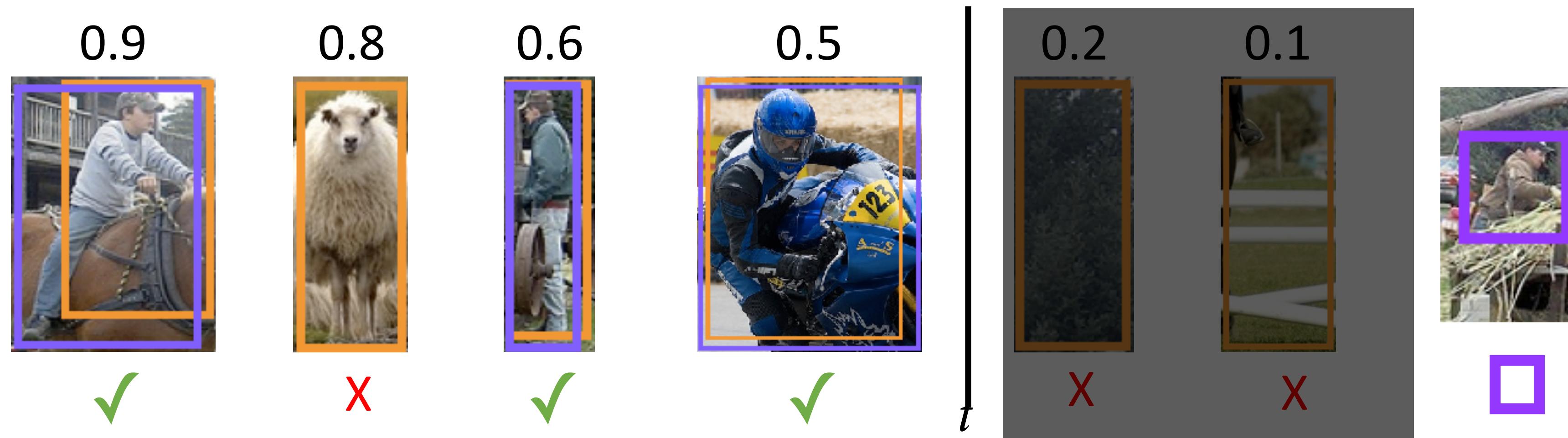
$$\frac{\checkmark}{\checkmark + \times}$$

$$recall@t = \frac{\#true\ positives@t}{\#ground\ truth\ objects}$$

$$\frac{\checkmark}{\square}$$

?

Sort by confidence



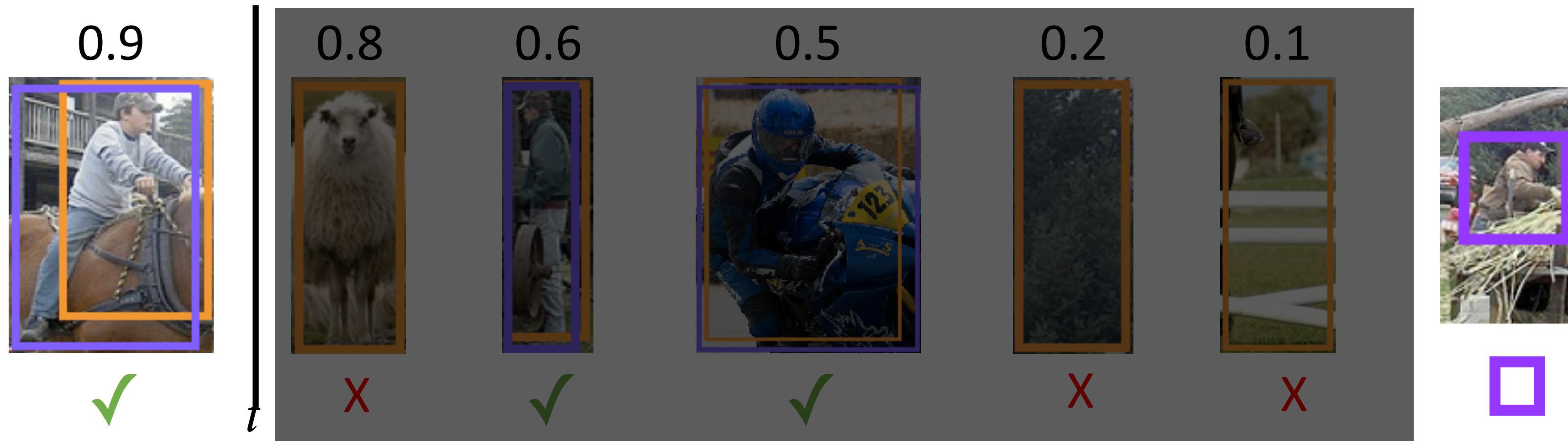
$$precision@t = \frac{\#true\ positives@t}{\#true\ positives@t + \#false\ positives@t}$$

$$\frac{\checkmark}{\checkmark + \times} = 75\%$$

$$recall@t = \frac{\#true\ positives@t}{\#ground\ truth\ objects}$$

$$\frac{\checkmark}{\square} = 75\%$$

Sort by confidence



$$t = 0.9$$

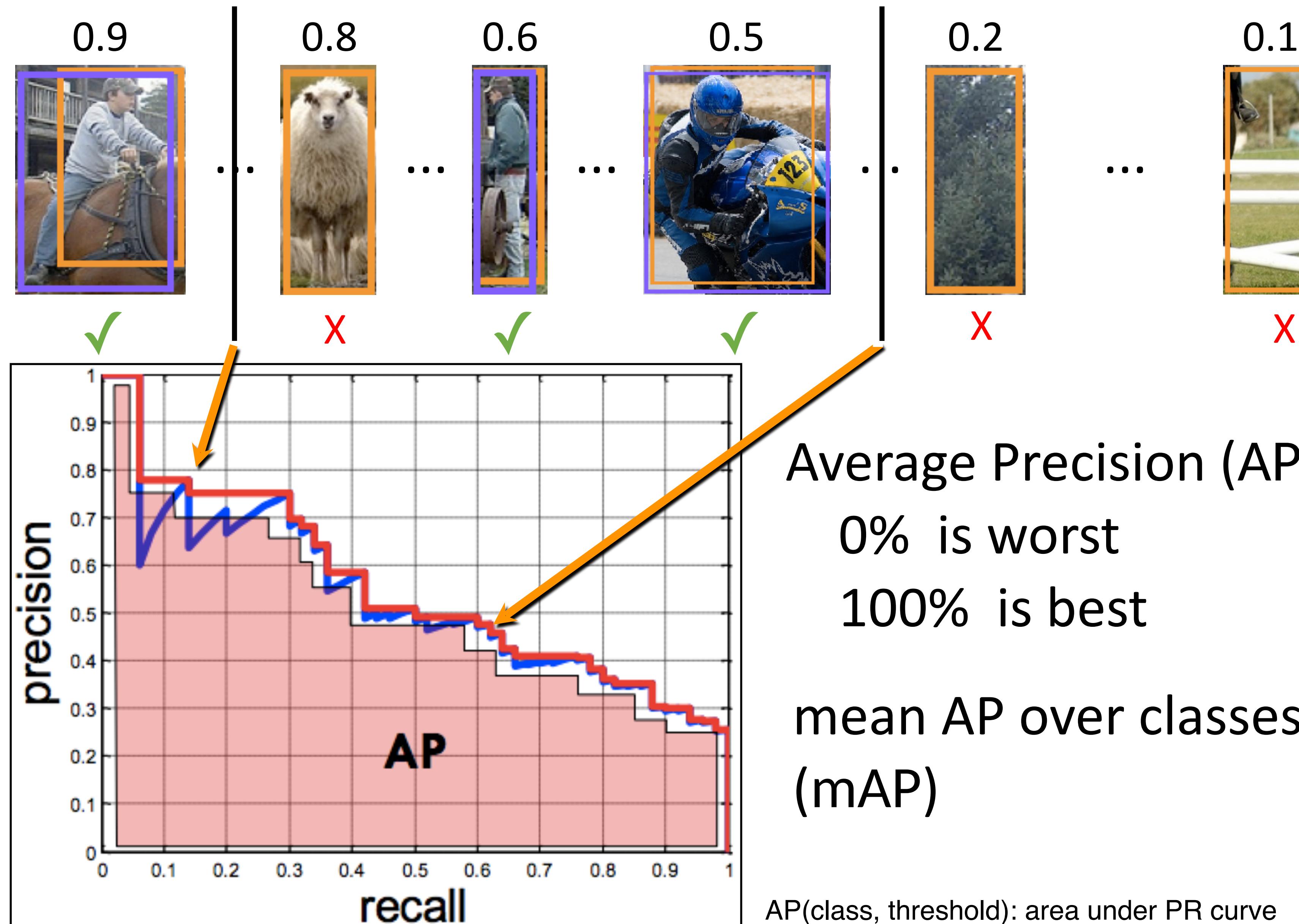
$$precision @ t = \frac{\#true\ positives @ t}{\#true\ positives @ t + \#false\ positives @ t}$$

$$\frac{\checkmark}{\checkmark + \text{X}} = 100\%$$

$$recall @ t = \frac{\#true\ positives @ t}{\#ground\ truth\ objects}$$

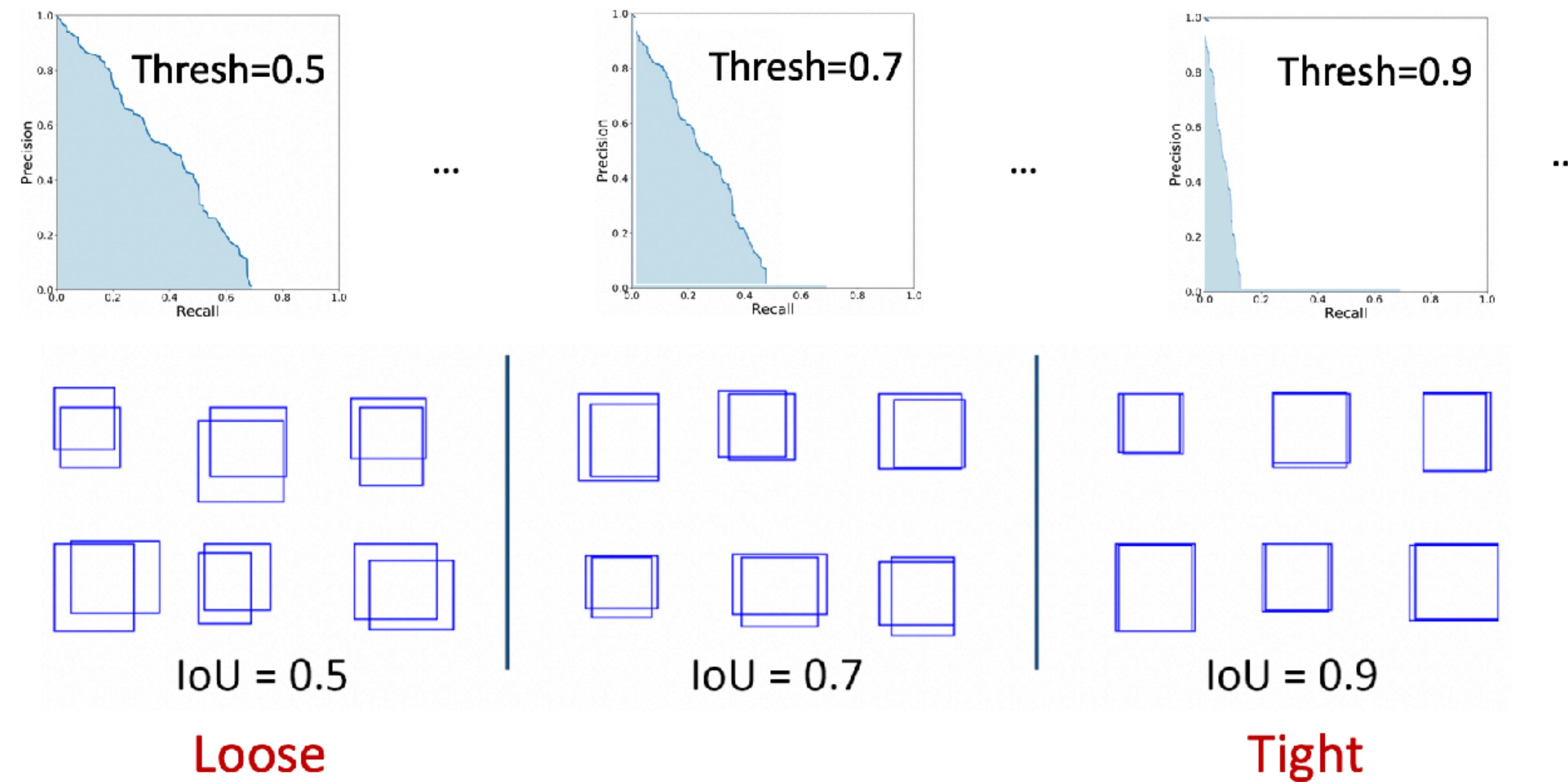
$$\frac{\checkmark}{\square} = 25\%$$

Average Precision for a (class, IOU threshold) pair



Average Precision for a class

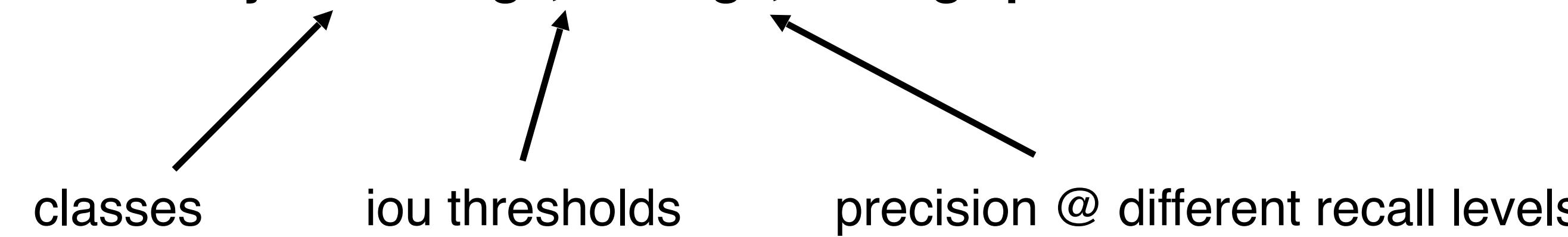
$$AP(class) = \frac{1}{\#thresholds} \sum_{iou \in threshold} AP(class, iou)$$



Overall Average Precision (%)

$$AP = \frac{1}{\#classes} \sum_{class \in classes} AP(class)$$

“AP” is really an average, average, average precision.



Average Precision (AP):

AP

% AP at IoU=.50:.05:.95 (primary challenge metric)

AP^{IoU=.50}

% AP at IoU=.50 (PASCAL VOC metric)

AP^{IoU=.75}

% AP at IoU=.75 (strict metric)

AP Across Scales:

AP^{small}

% AP for small objects: area < 32²

AP^{medium}

% AP for medium objects: 32² < area < 96²

AP^{large}

% AP for large objects: area > 96²

Object detection: naive attempt

This is a chair



Find the chair in this image

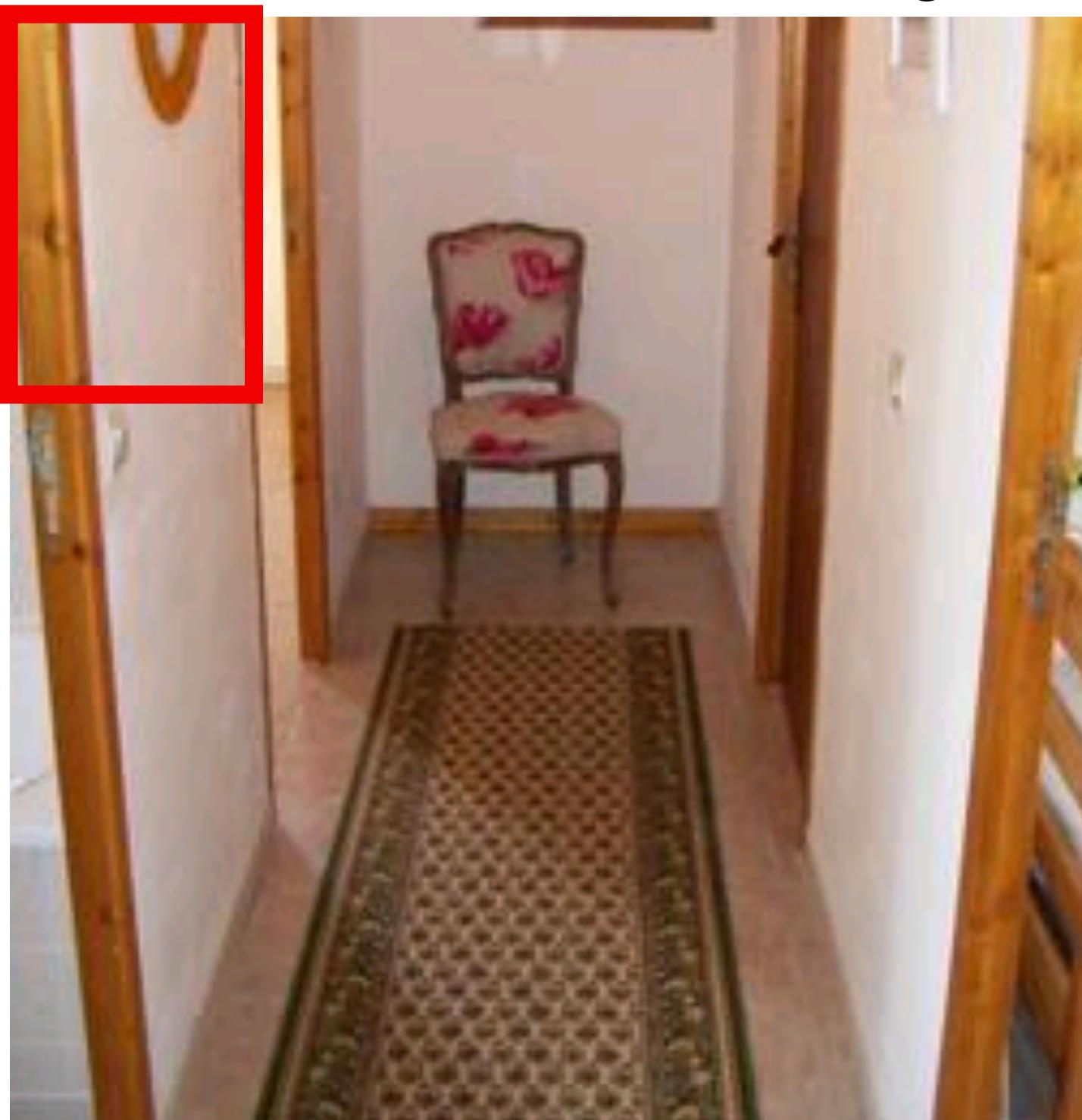


Object detection: naive attempt

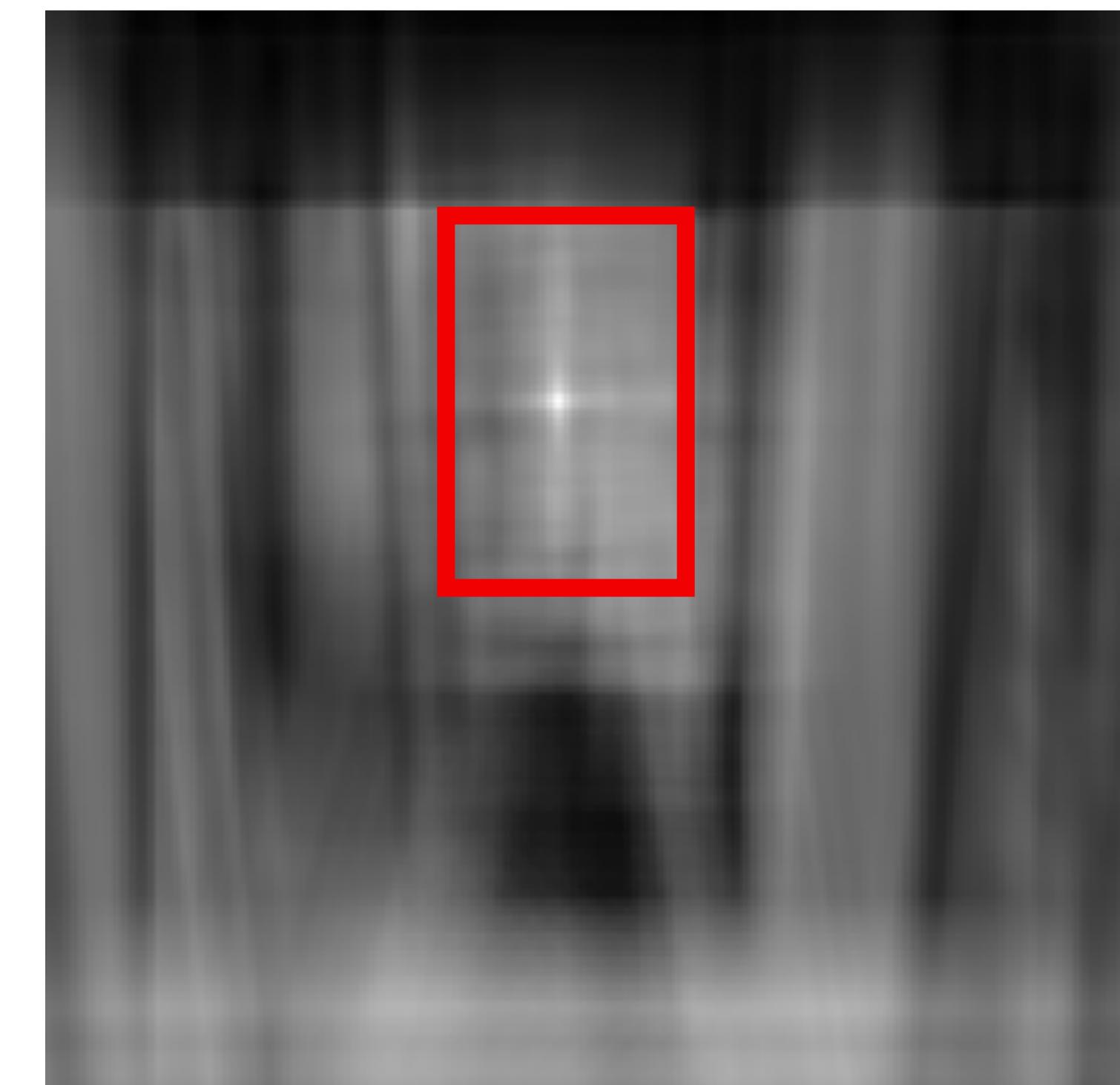
This is a chair



Find the chair in this image



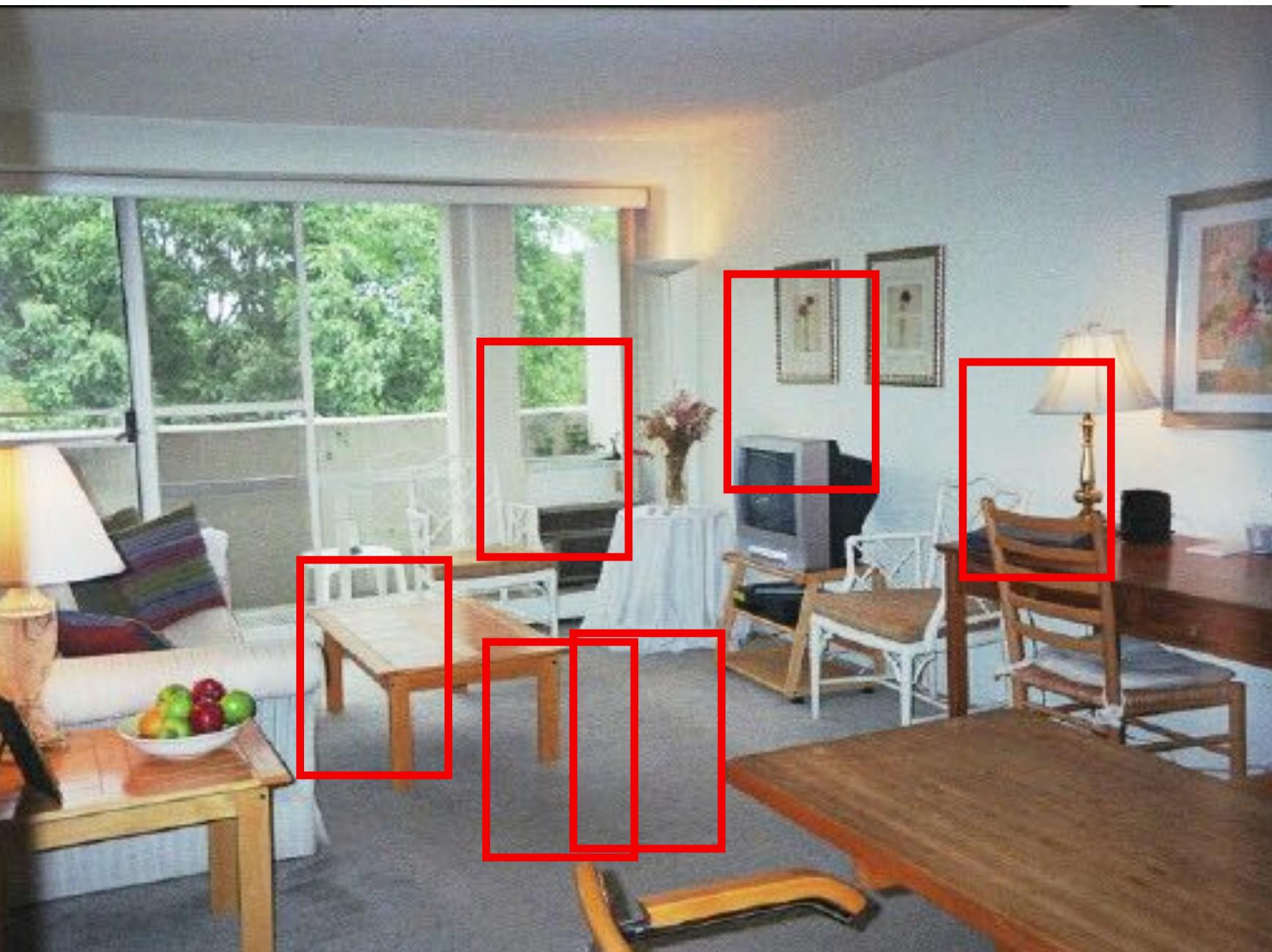
Output of normalized correlation



Object detection: naive attempt



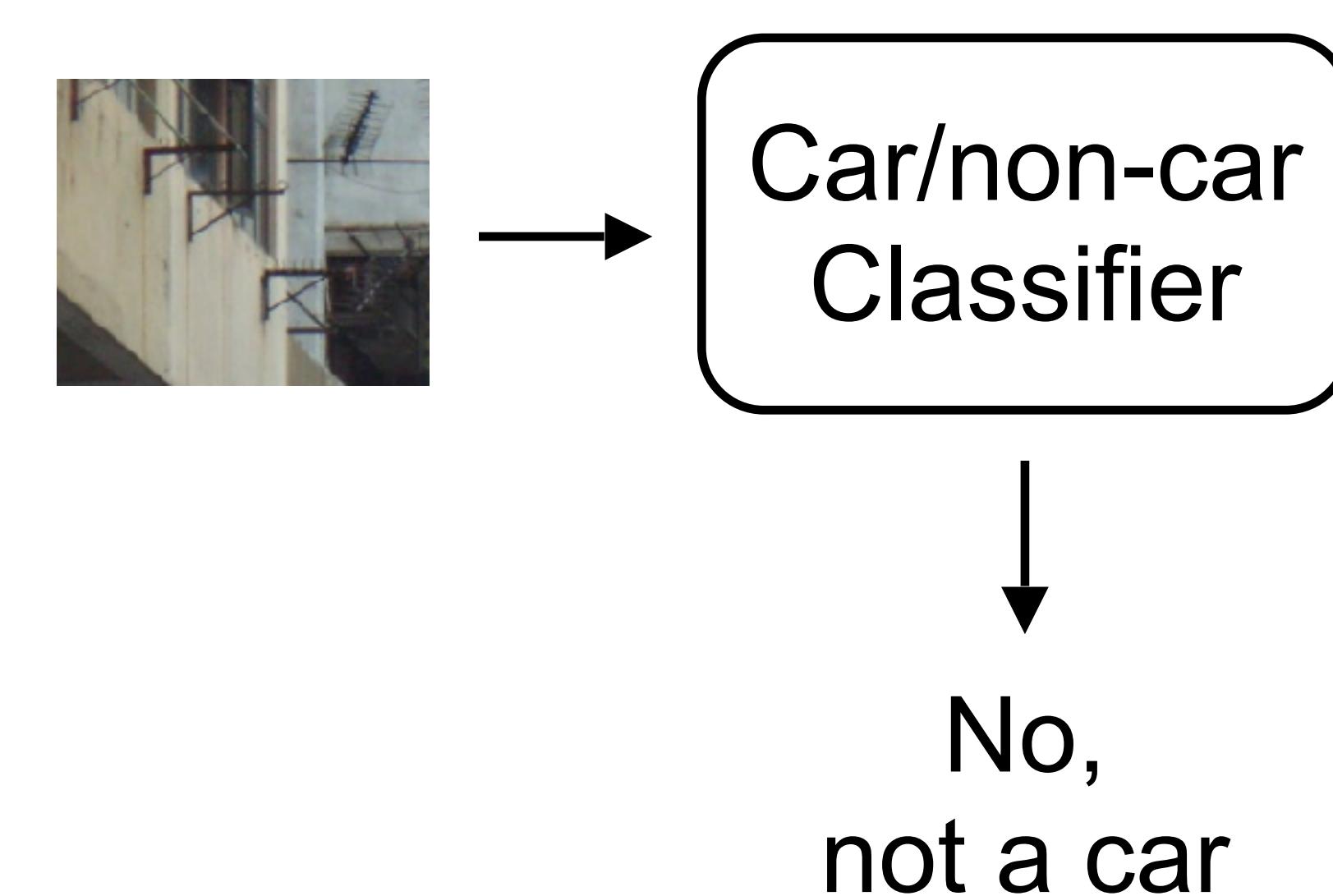
Find the chair in this image



Pretty much garbage
Simple template matching is not going to make it

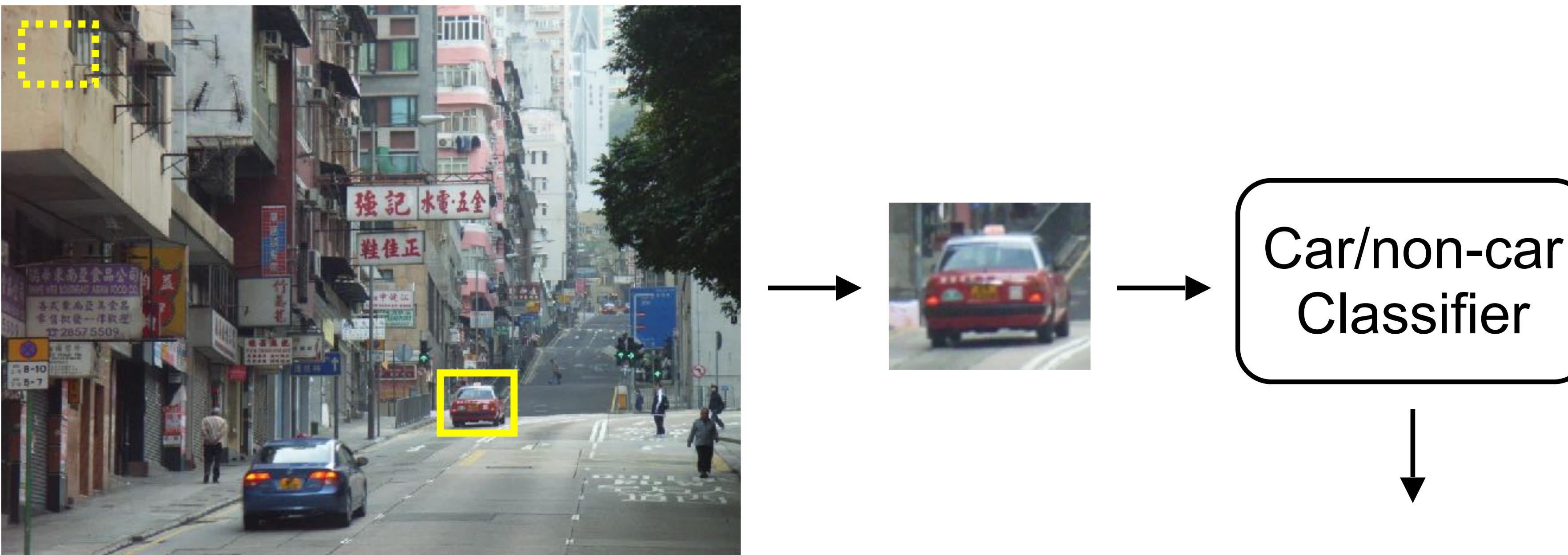
Detection by Classification

- Basic component: binary classifier



Detection by Classification

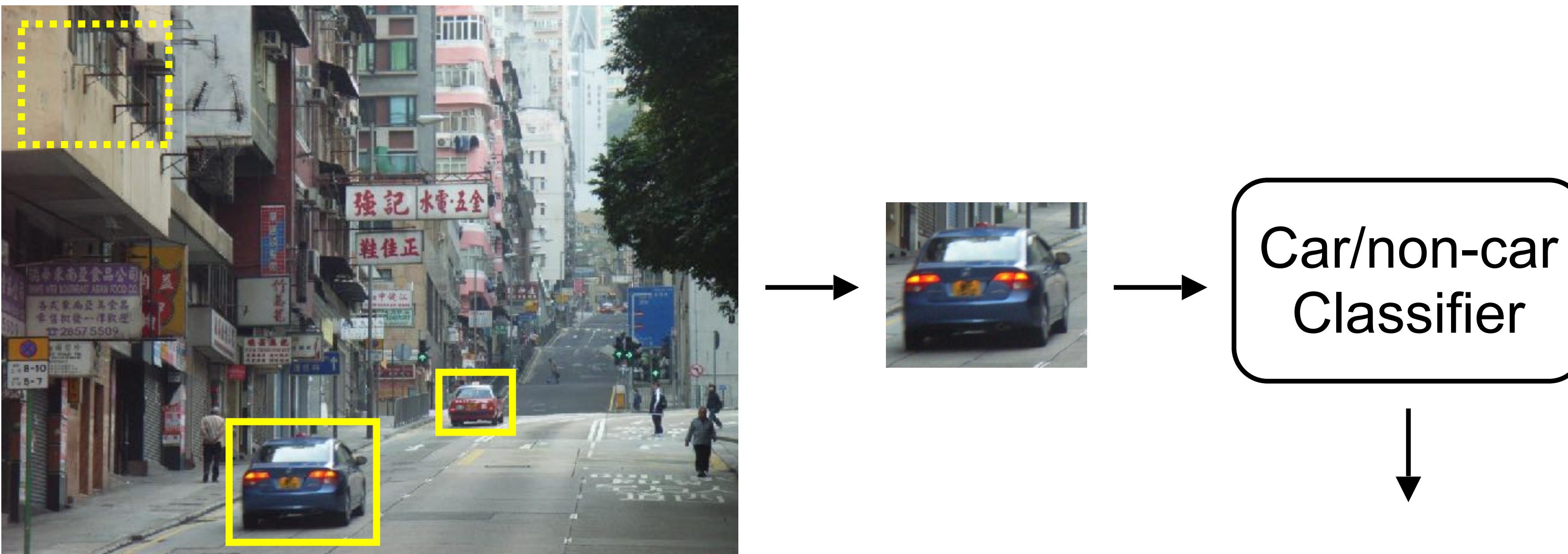
- Detect objects in clutter by search



- **Sliding window:** exhaustive search over position and scale

Detection by Classification

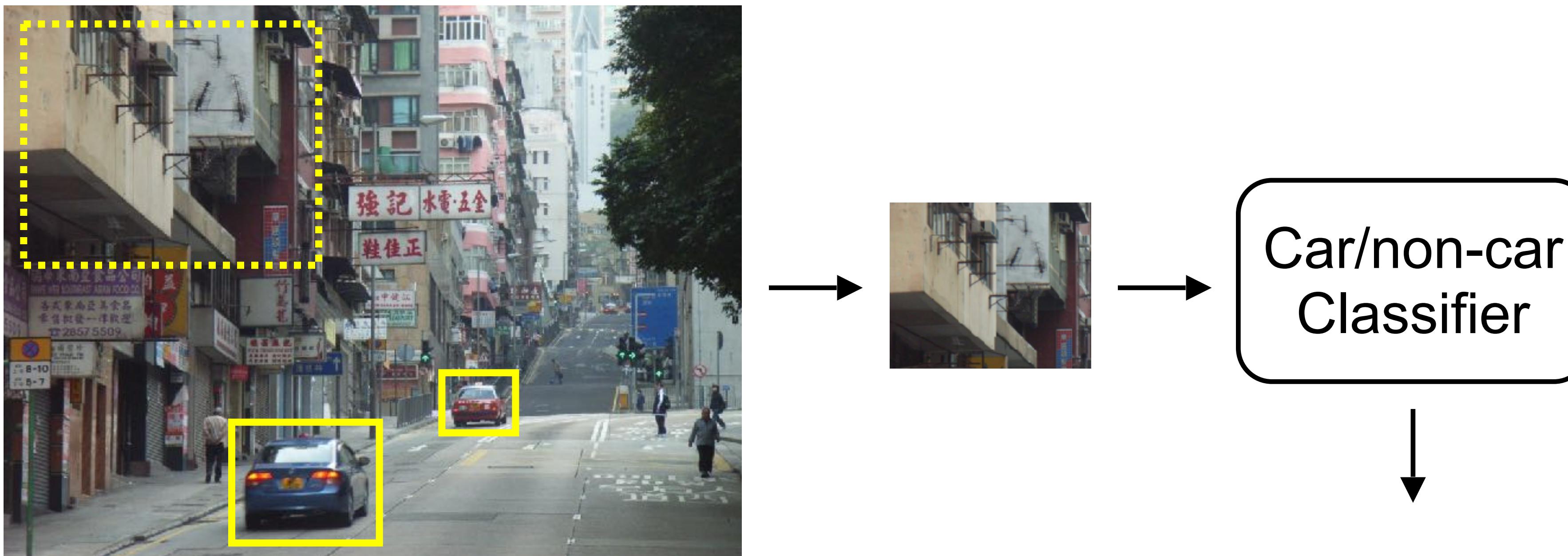
- Detect objects in clutter by search



- **Sliding window:** exhaustive search over position and scale

Detection by Classification

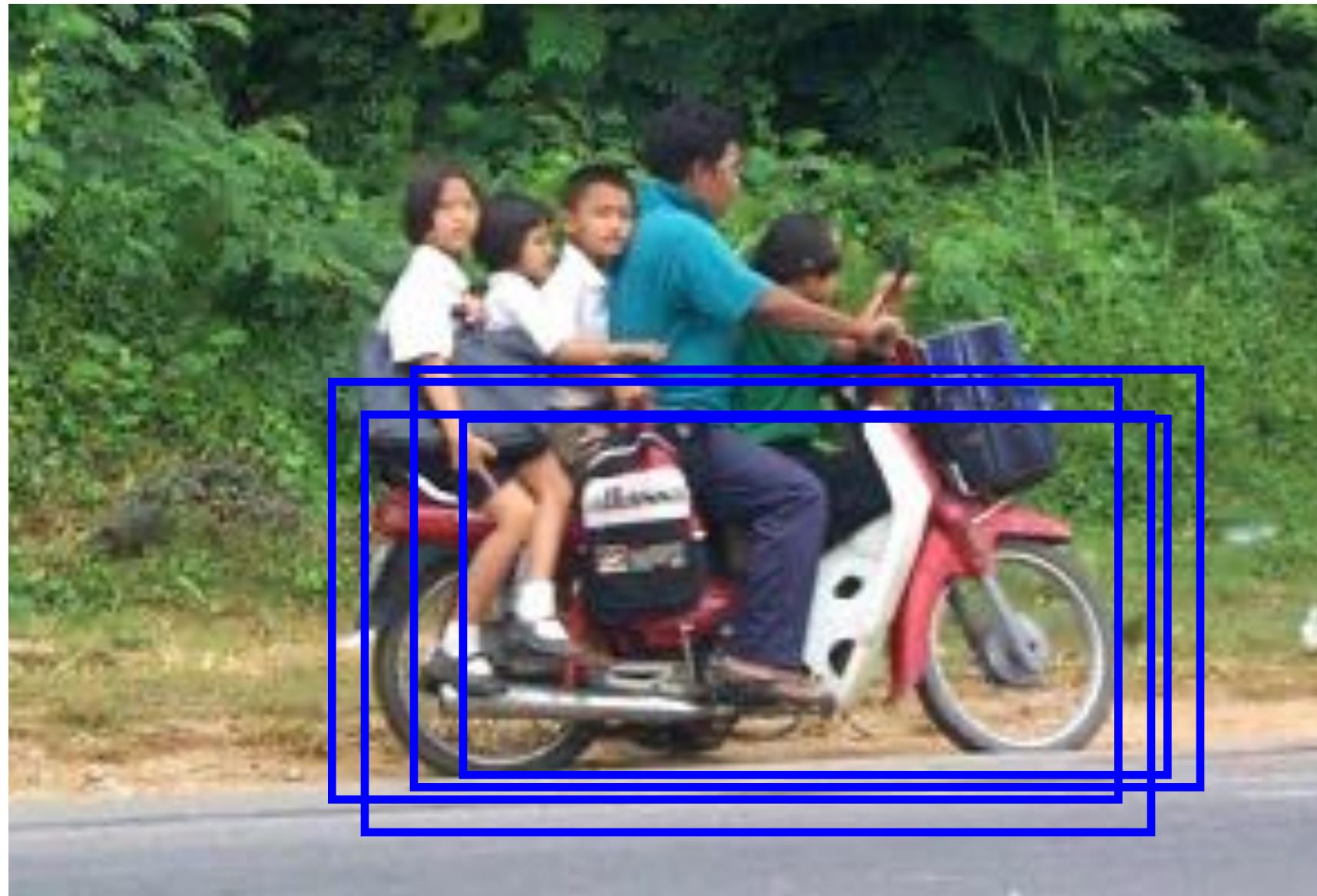
- Detect objects in clutter by search



- **Sliding window:** exhaustive search over position and scale
(can use same size window over a spatial pyramid of images)

Test: Non-maximum suppression (NMS)

- Scanning-window detectors typically result in multiple responses for the same object



- To remove multiple responses, a simple greedy procedure called “Non-maximum suppression” is applied:

NMS:

- Sort all detections by detector confidence
- Choose most confident detection d_i ; remove all d_j s.t. $\text{overlap}(d_i, d_j) > T$
- Repeat Step 2. until convergence

Test: Non-maximum suppression (NMS)

- Scanning-window detectors typically result in multiple responses for the same object



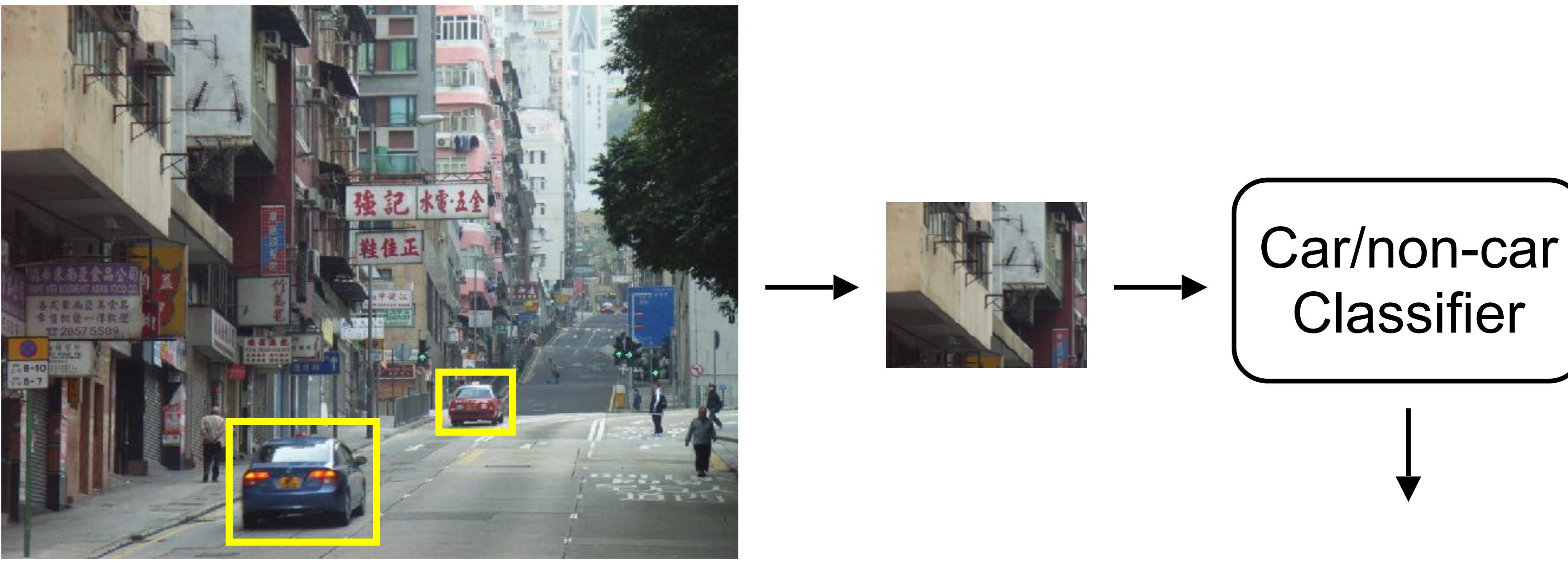
- To remove multiple responses, a simple greedy procedure called “Non-maximum suppression” is applied:

NMS:

- Sort all detections by detector confidence
- Choose most confident detection d_i ; remove all d_j s.t. $\text{overlap}(d_i, d_j) > T$
- Repeat Step 2. until convergence

Detection by Classification

- Detect objects in clutter by search

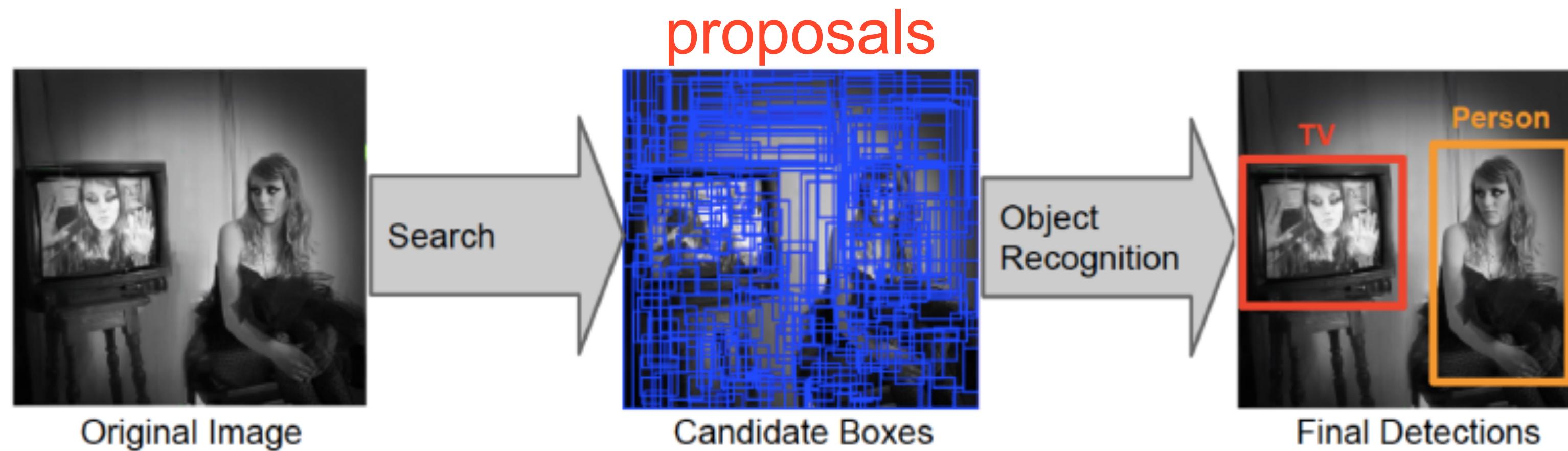


Problem: too many windows to run a classifier

- **Sliding window:** exhaustive search over position and scale
(can use same size window over a spatial pyramid of images)

Object proposals

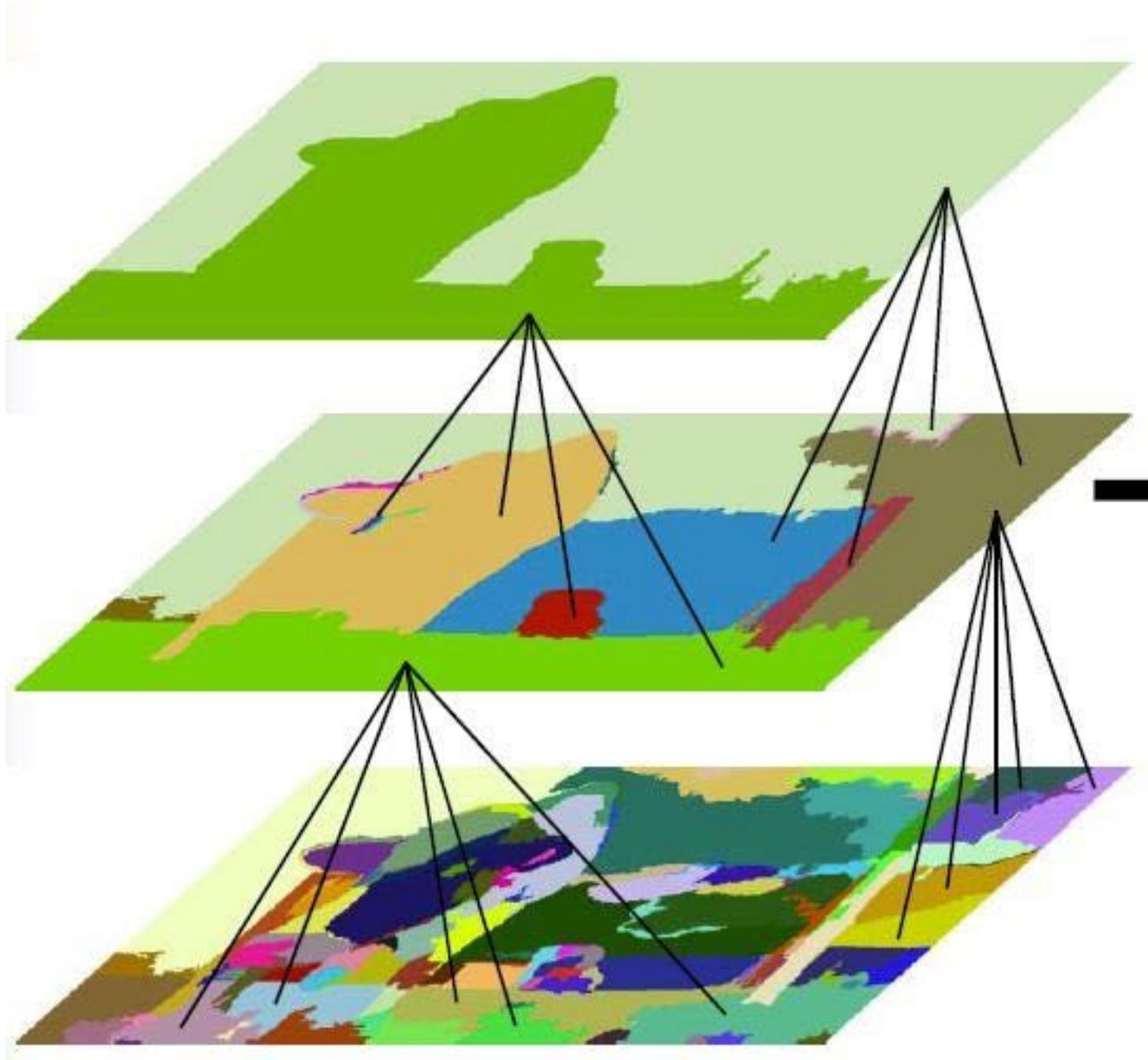
Generate and evaluate a few hundred region proposals.



- Proposal mechanism can:
 - take advantage of low-level perceptual organization cues,
 - be category-specific or category-independent, handcrafted or trained.
- Classifier can be slower but more powerful.

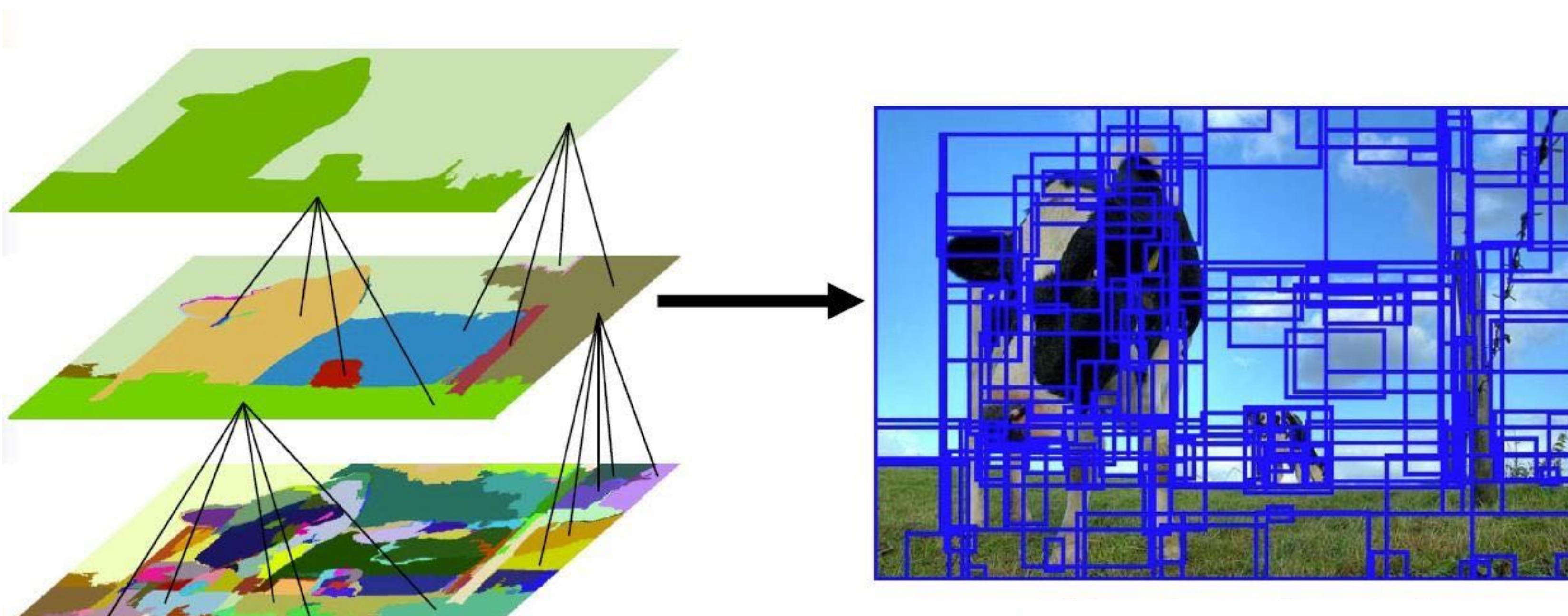
Region proposals: Selective search

1. Merge two most similar regions based on similarity.
2. Update similarities between the new region and its neighbors.
3. Go back to step 1. until the whole image is a single region.

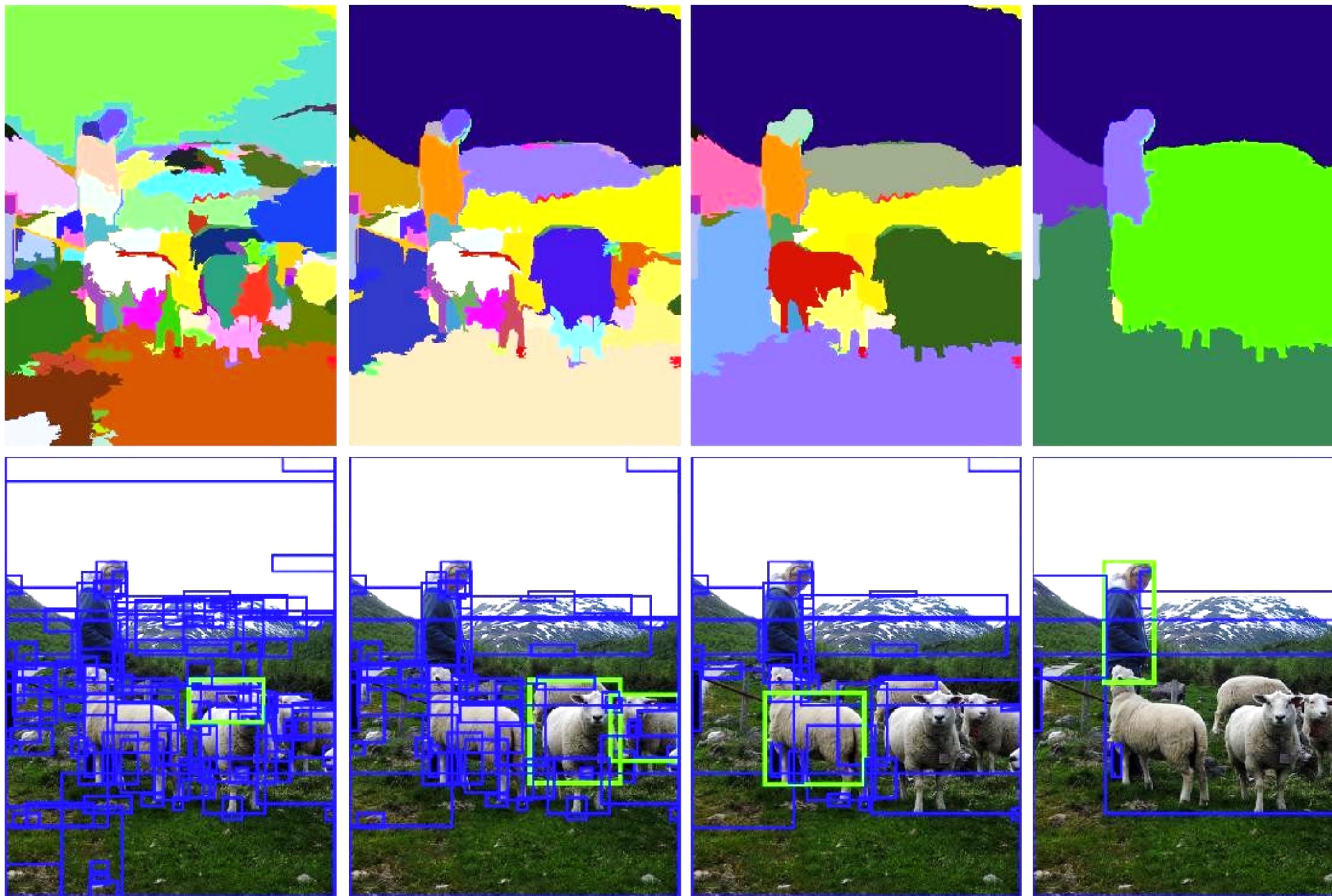


Region proposals: Selective search

Take bounding boxes of all generated regions
and treat them as possible object locations.



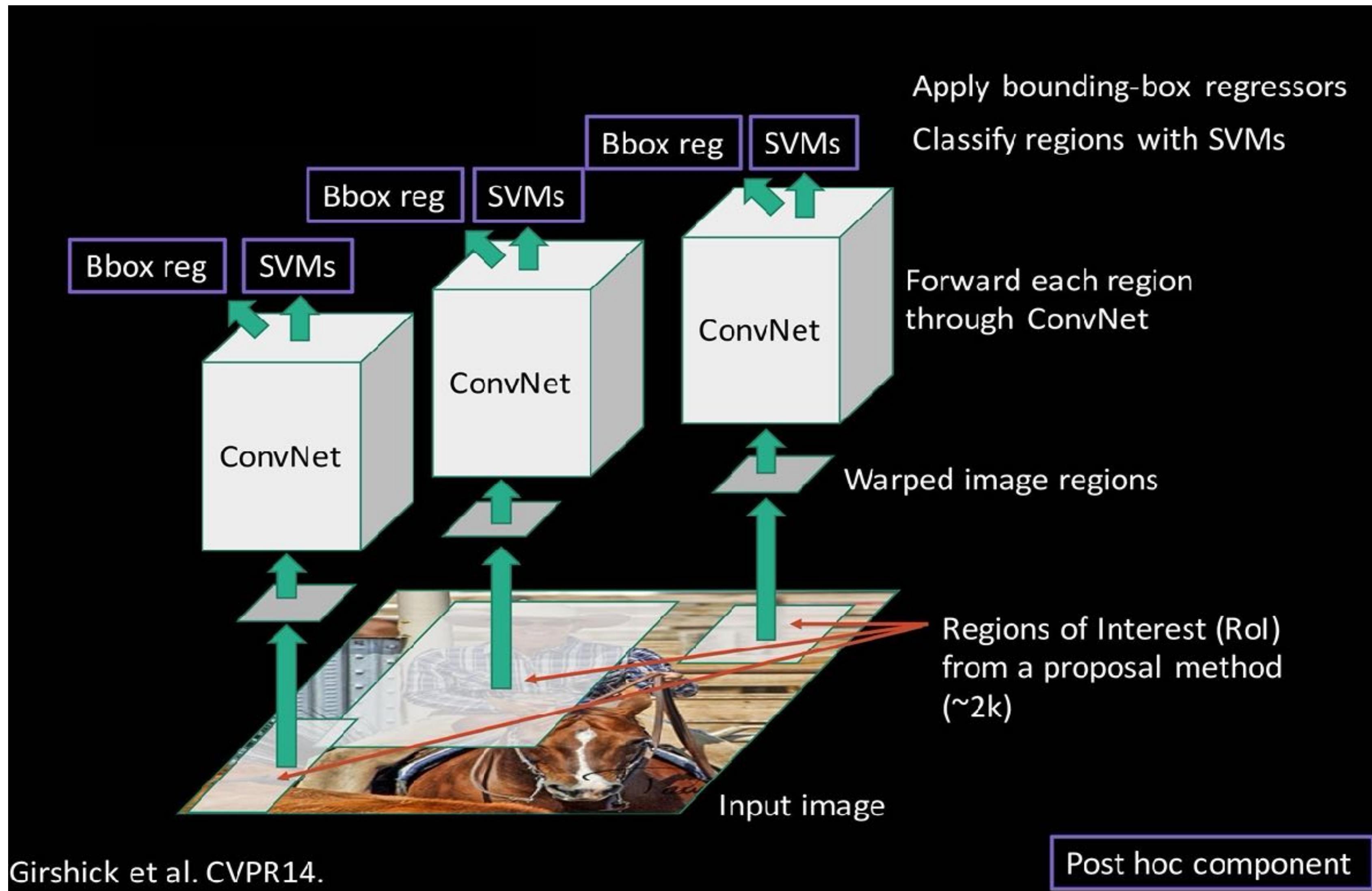
Region proposals: Selective search



[K. van de Sande, J. Uijlings, T. Gevers, and A. Smeulders, ICCV 2011]

Object detection: CNN-based methods

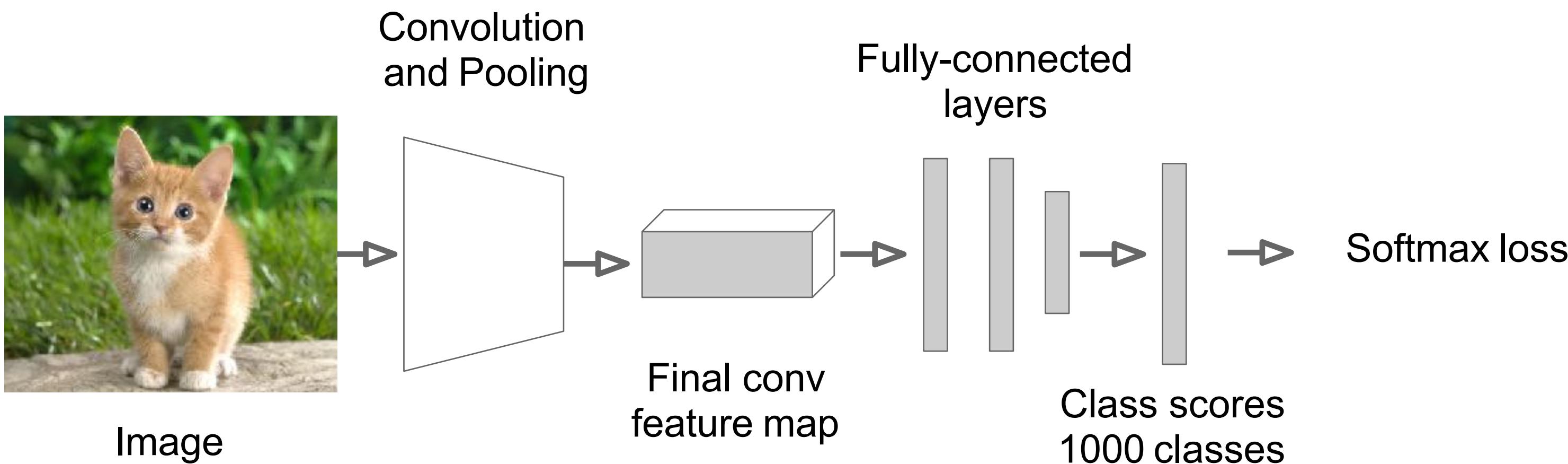
R-CNN: Region-based CNN



Girshick et al. CVPR14.

R-CNN Training

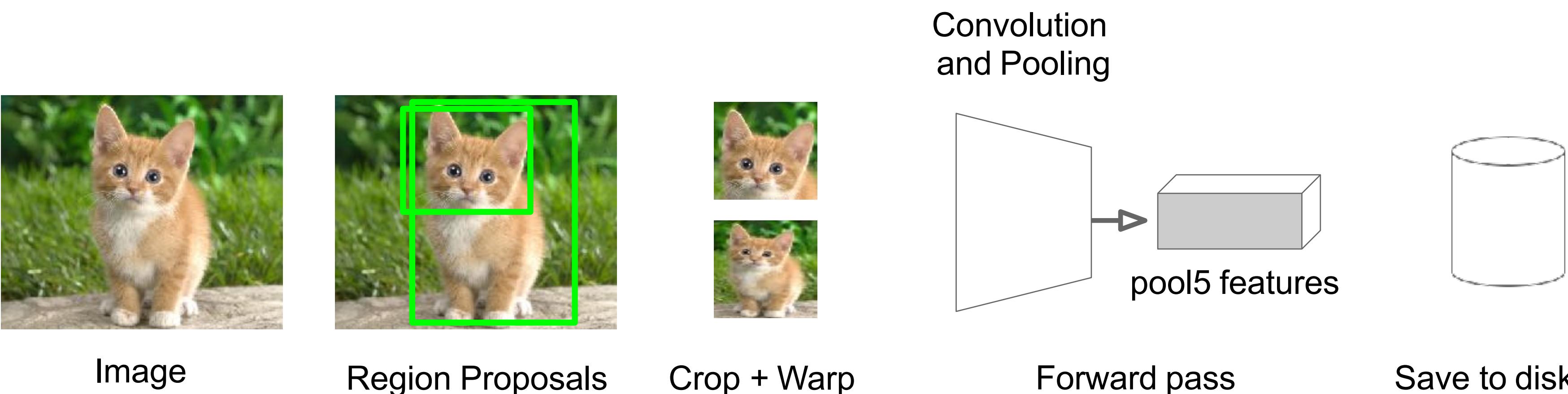
Step 1: Train (or download) a classification model for ImageNet (AlexNet)



R-CNN Training

Step 2: Extract features

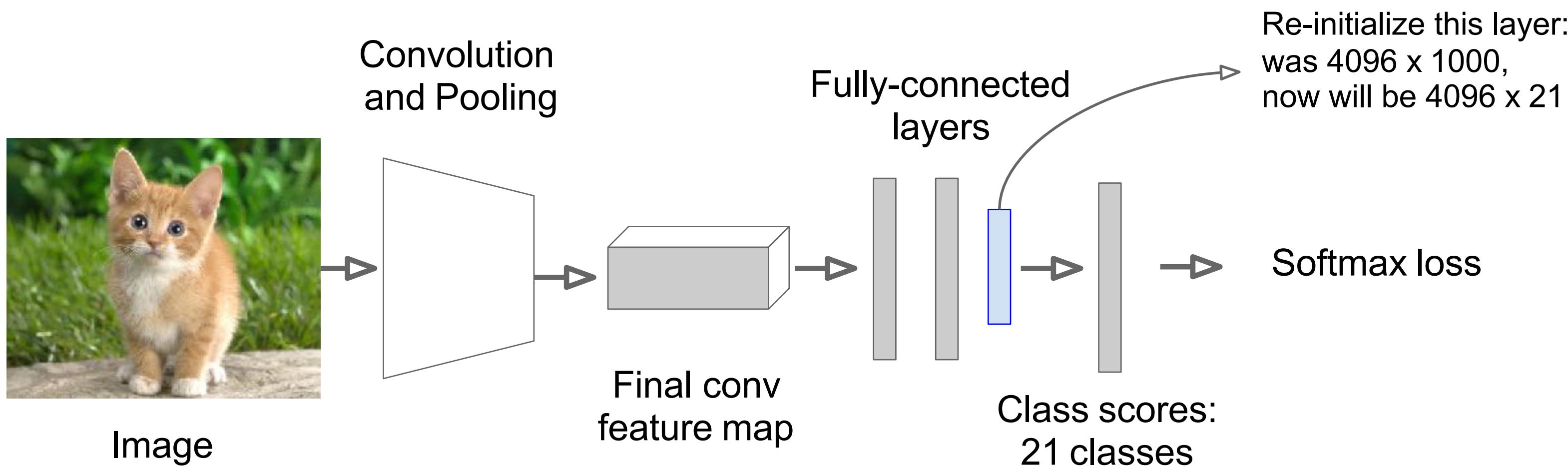
- Extract region proposals for all images
- For each region: warp to CNN input size, run forward through CNN, save pool5 features to disk
- Have a big hard drive: features are ~200GB for PASCAL dataset!



R-CNN Training

Step 3: Fine-tune model for detection

- Instead of 1000 ImageNet classes, want 20 object classes + background
- Throw away final fully-connected layer, reinitialize from scratch
- Keep training model using positive / negative regions from detection images

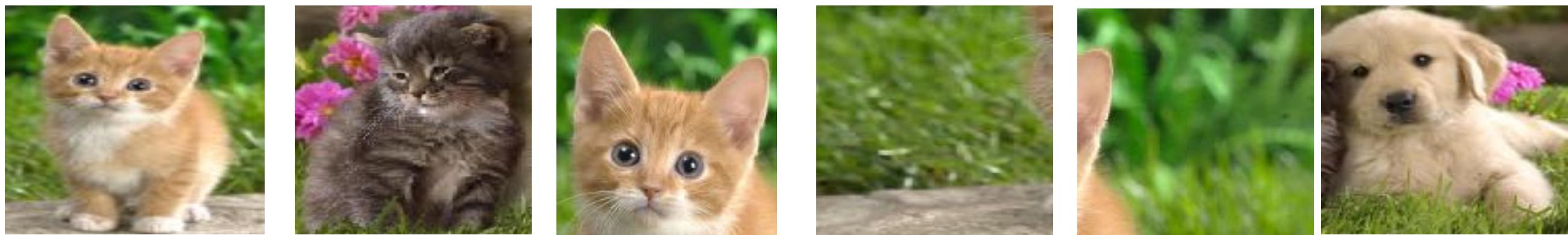


R-CNN Training

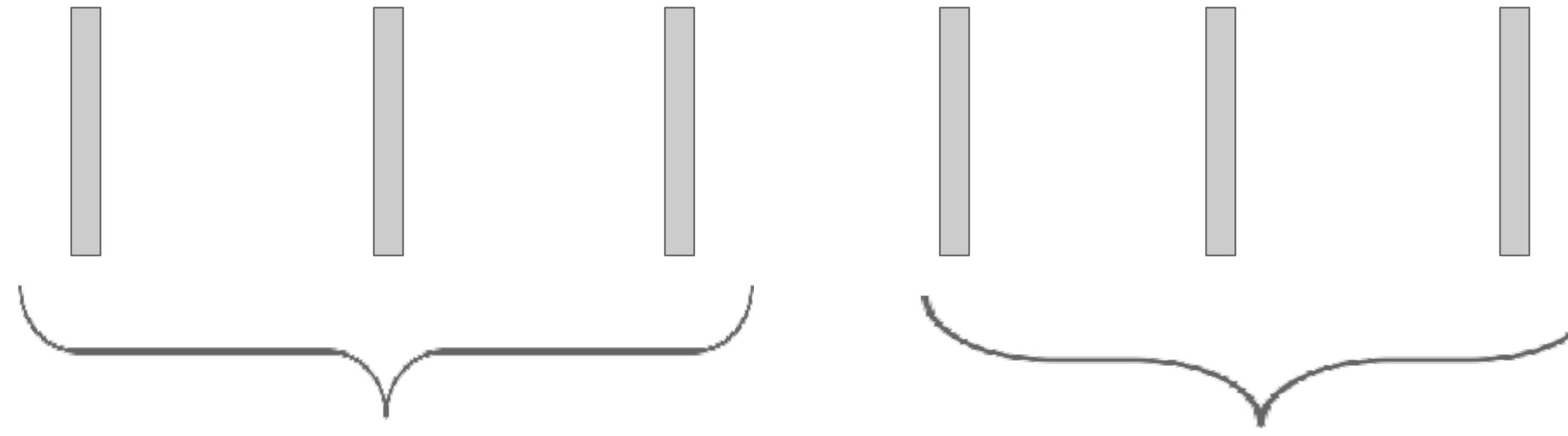
(Go back to Step 2: Re-extract features with the fine-tuned CNN)

Step 4: Train one binary SVM per class to classify region features

Training image regions



Cached region features



+ Positive samples for **cat** SVM

- Negative samples for **cat** SVM

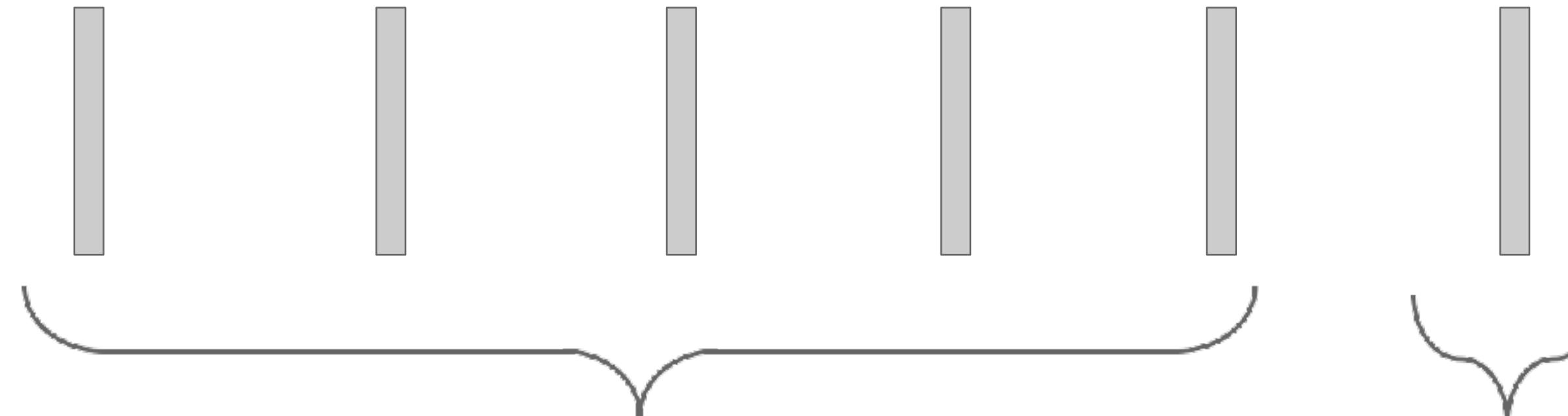
R-CNN Training

Step 4: Train one binary SVM per class to classify region features

Training image regions



Cached region features



- Negative samples for **dog** SVM

+ Positive samples for **dog** SVM

R-CNN Training

Step 5 (bbox regression): For each class, train a linear regression model to map from cached features to offsets to GT boxes to make up for “slightly wrong” proposals

Training image regions



Cached region features



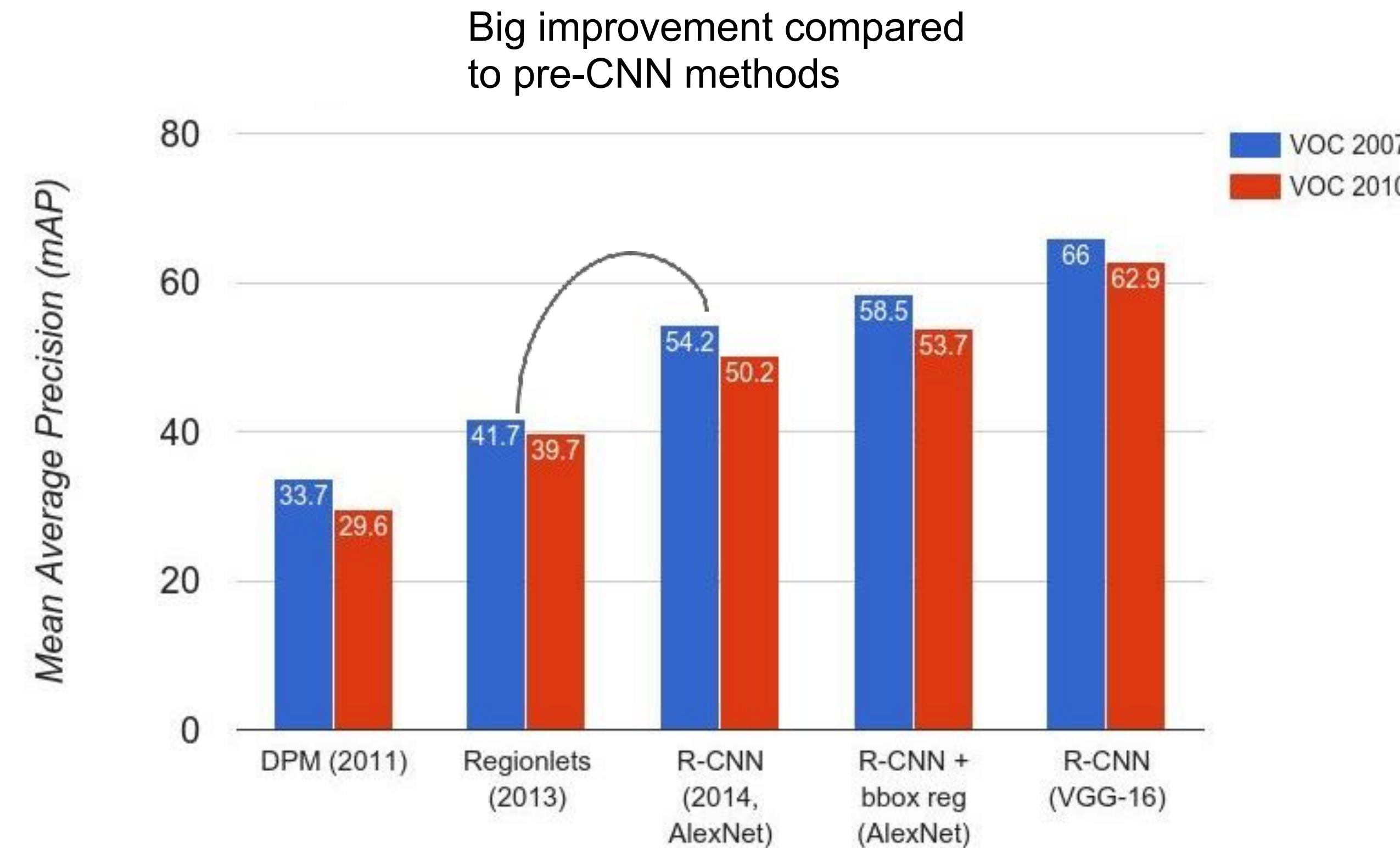
Regression targets
(dx , dy , dw , dh)
Normalized coordinates

(0, 0, 0, 0)
Proposal is good

(.25, 0, 0, 0)
Proposal too far to left

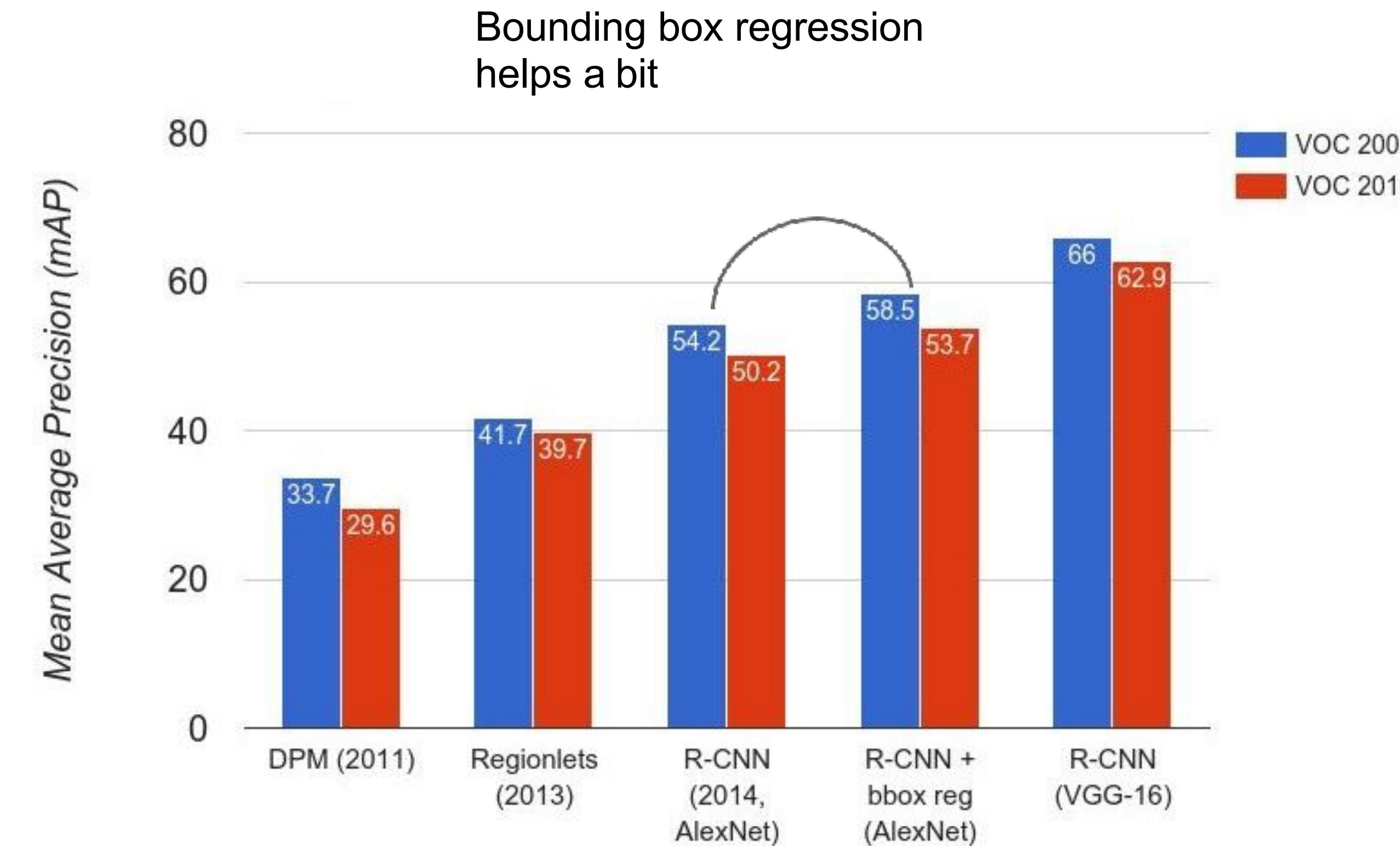
(0, 0, -0.125, 0)
Proposal too wide

R-CNN Results

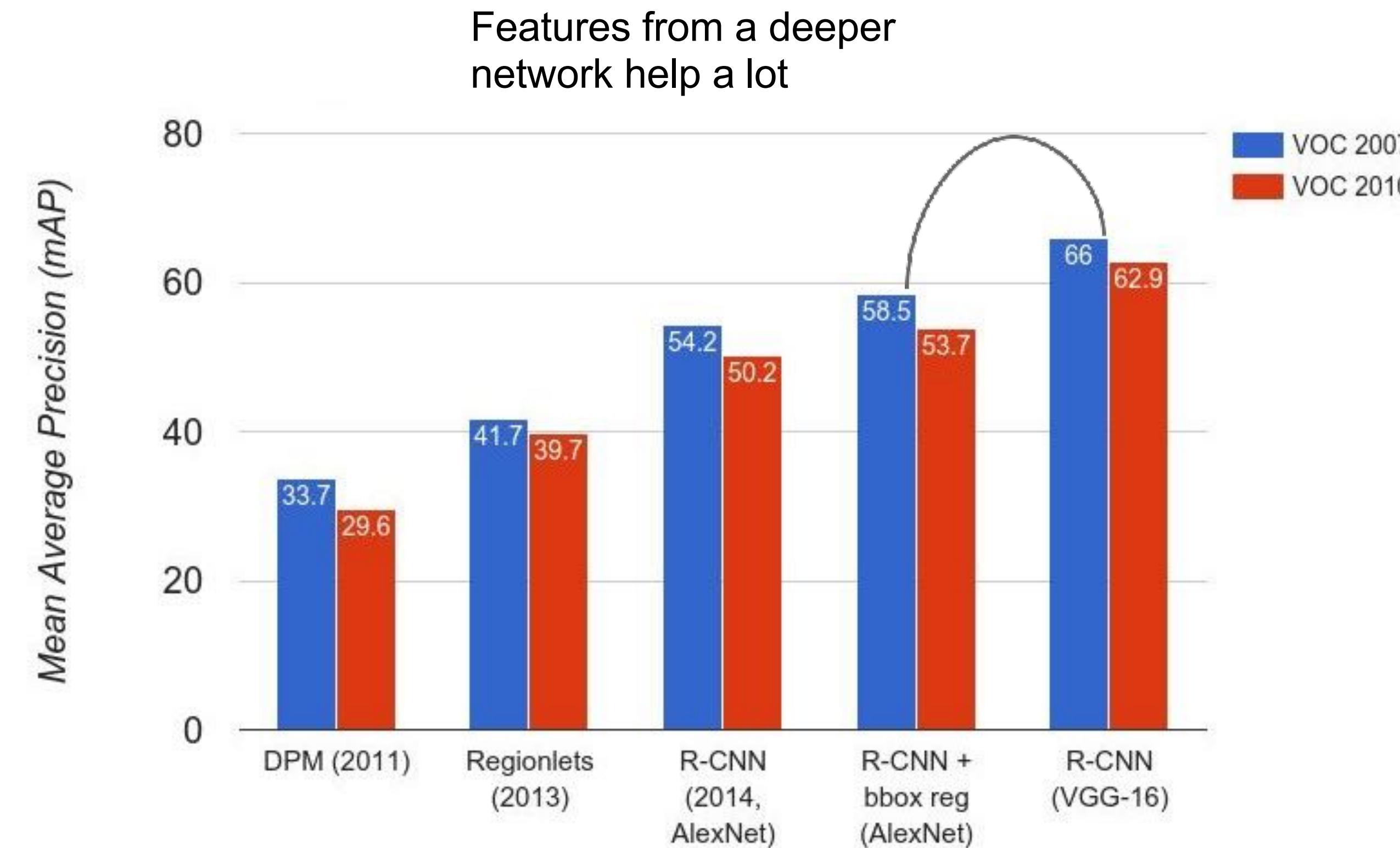


Wang et al, "Regionlets for Generic Object Detection", ICCV 2013

R-CNN Results



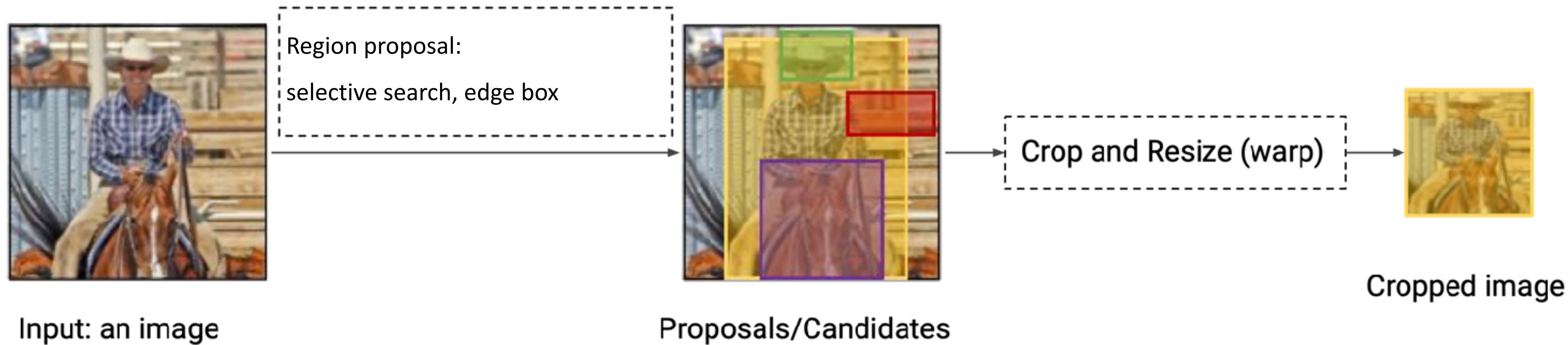
R-CNN Results



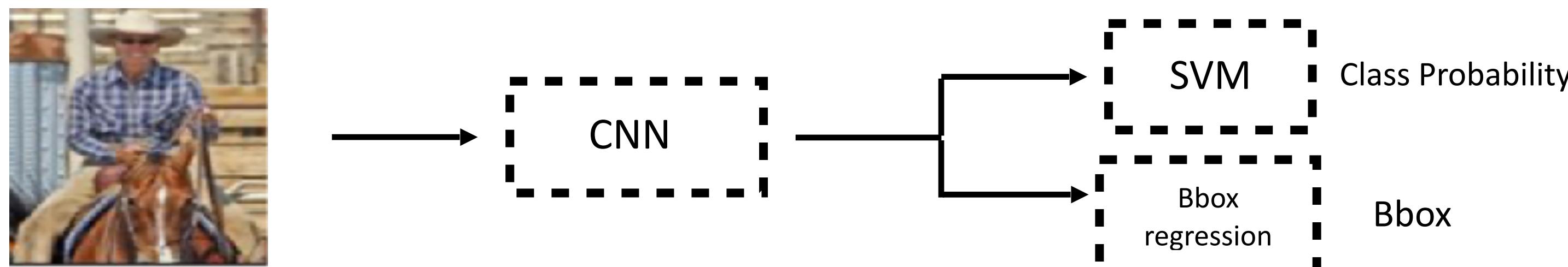
R-CNN [CVPR 2014] Summary

Two-stage detector

- Propose large number of regions potentially with objects



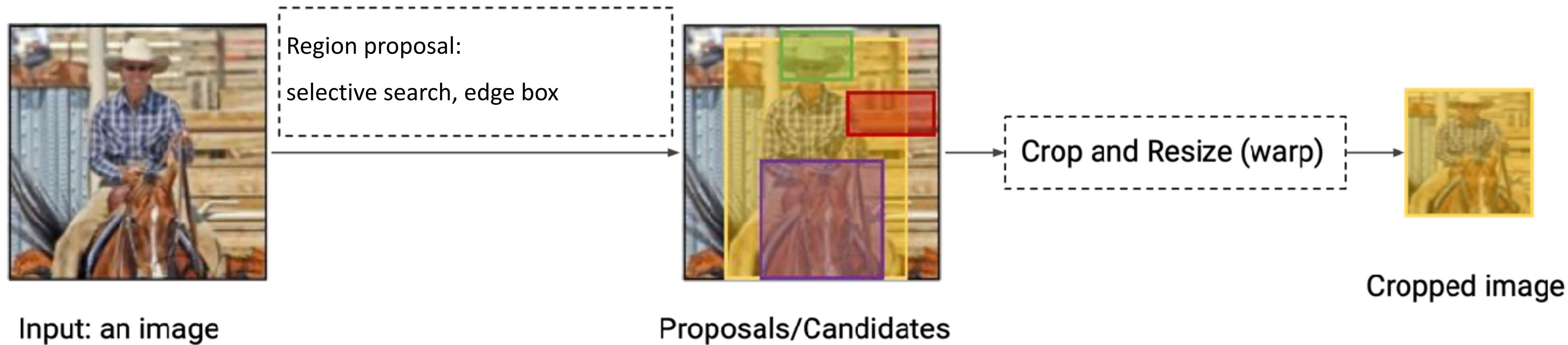
- Classify each proposed region



R-CNN [CVPR 2014] Limitations

Two-stage detector

- Propose large number of regions potentially with objects

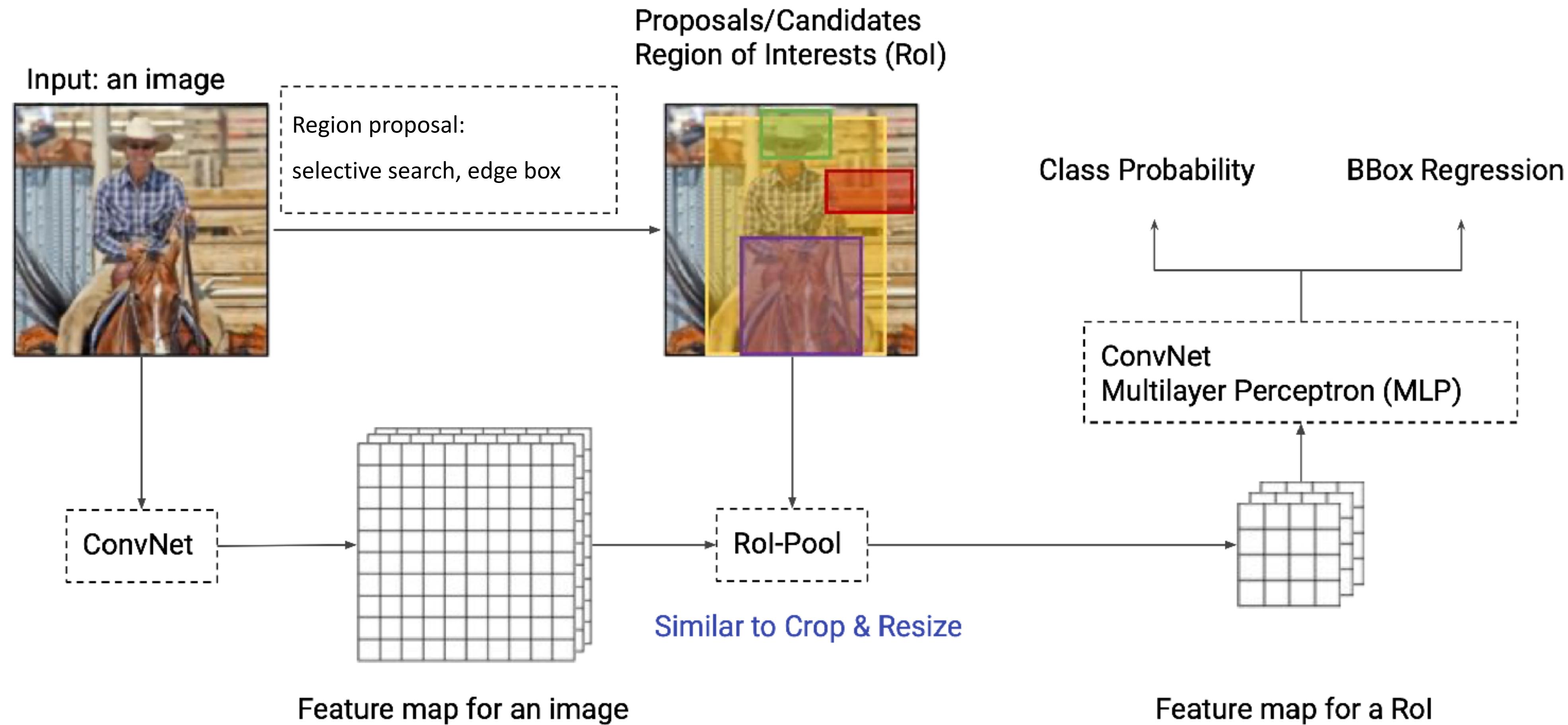


- Slow at test-time:** need to run full forward pass of CNN for each region proposal
- Not end-to-end:** SVMs and regressors are post-hoc, CNN features not updated in response to SVMs and regressors
- Complex multistage training pipeline**



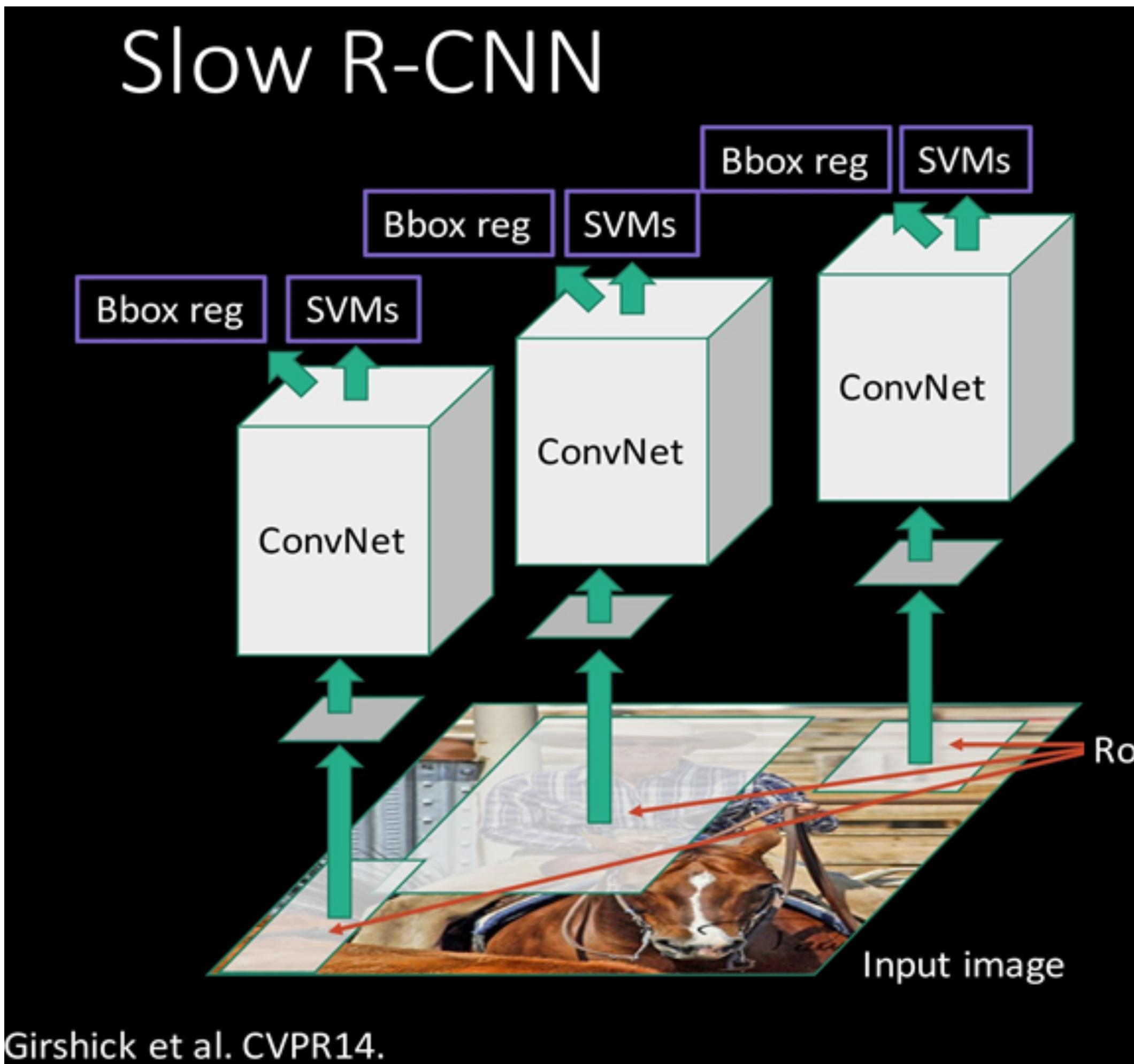
Fast R-CNN [ICCV 2015]

- Small accuracy improvement
- Timing excluding region proposal
 - ~10x faster for training
 - ~100x faster for testing (< 1 sec / image)



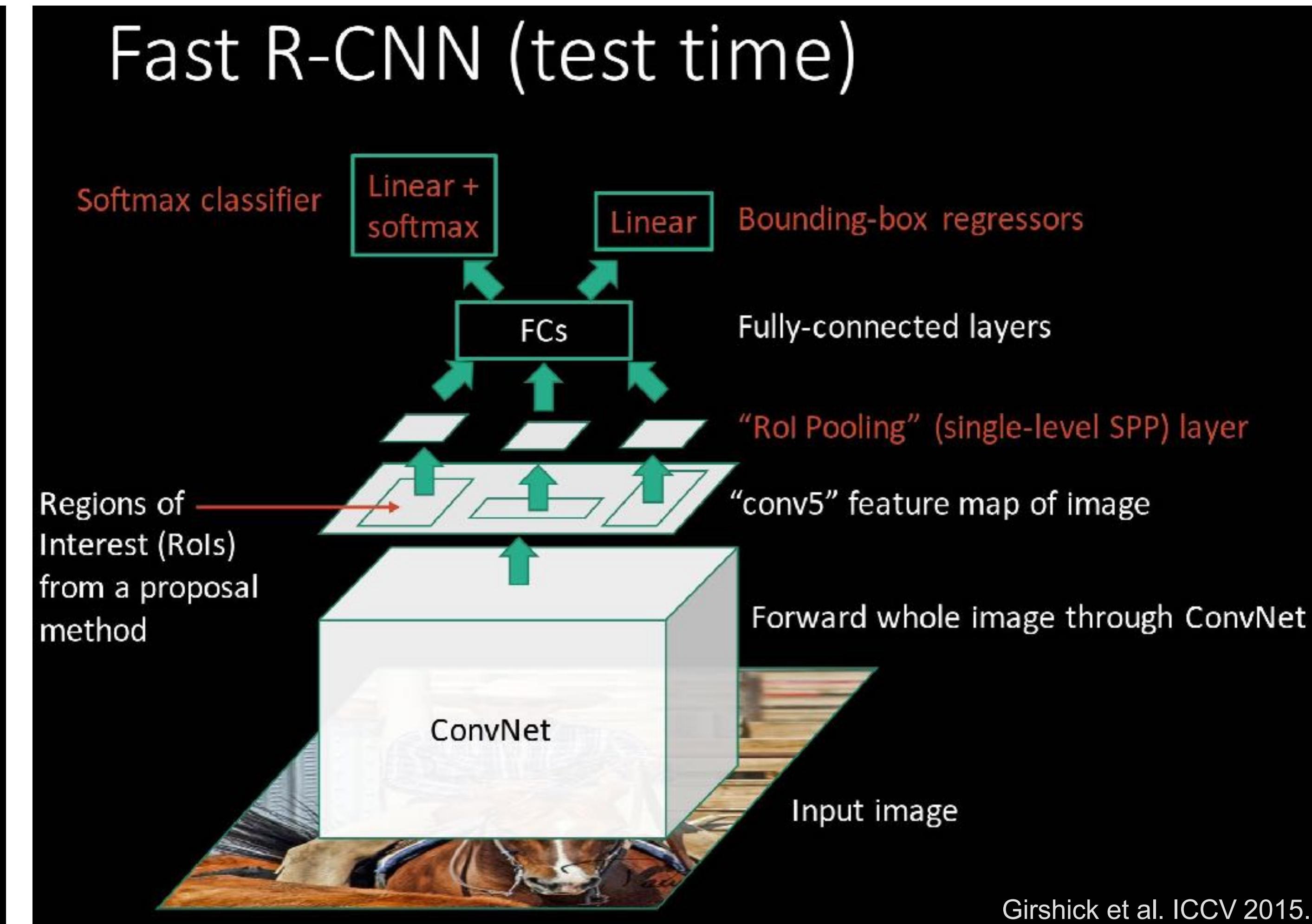
R-CNN Problems

Problem #1: Slow at test-time due to **independent forward passes** of the CNN



Fast R-CNN Solutions

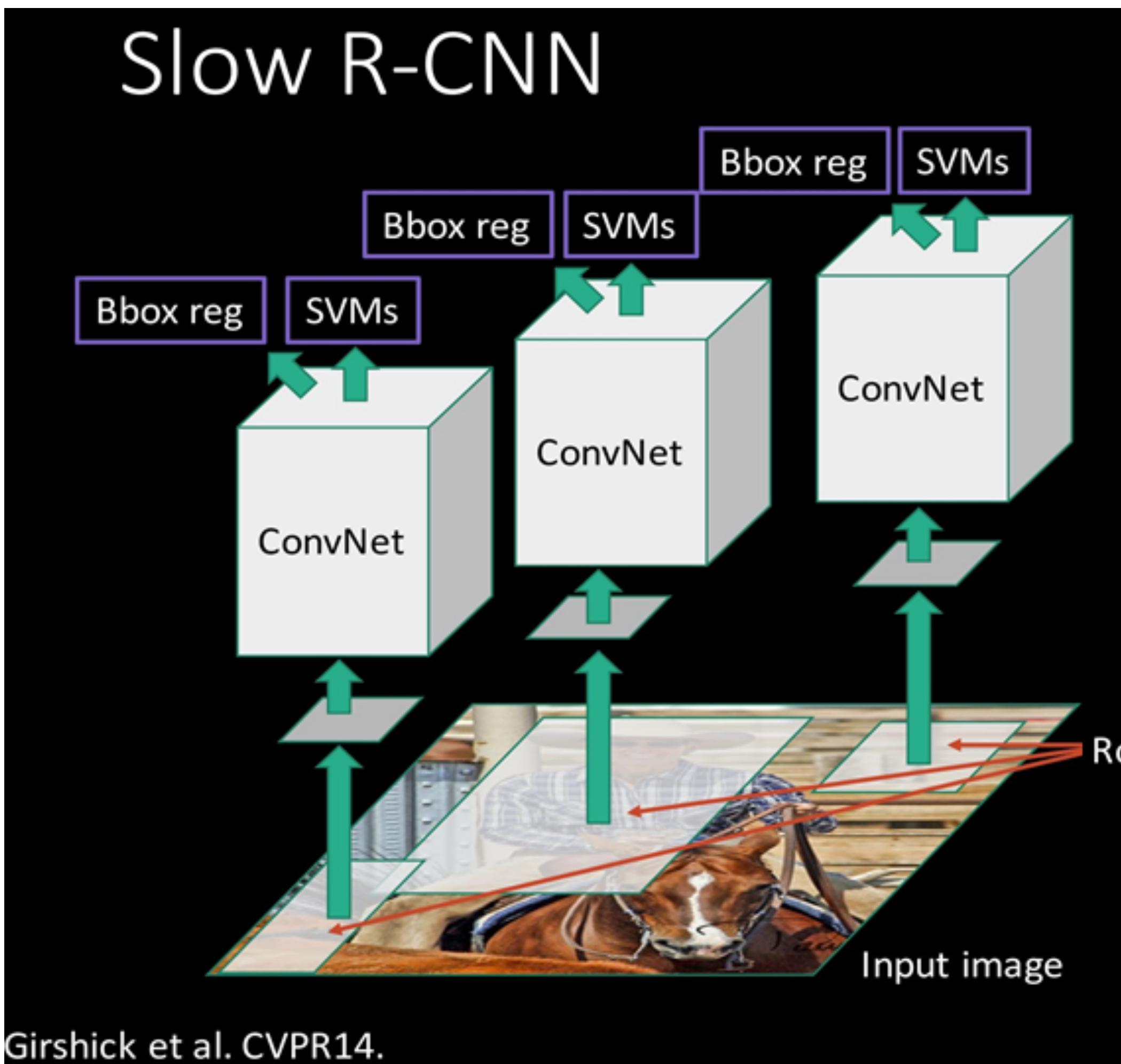
Solution: Share computation of convolutional layers between proposals for an image



R-CNN Problems

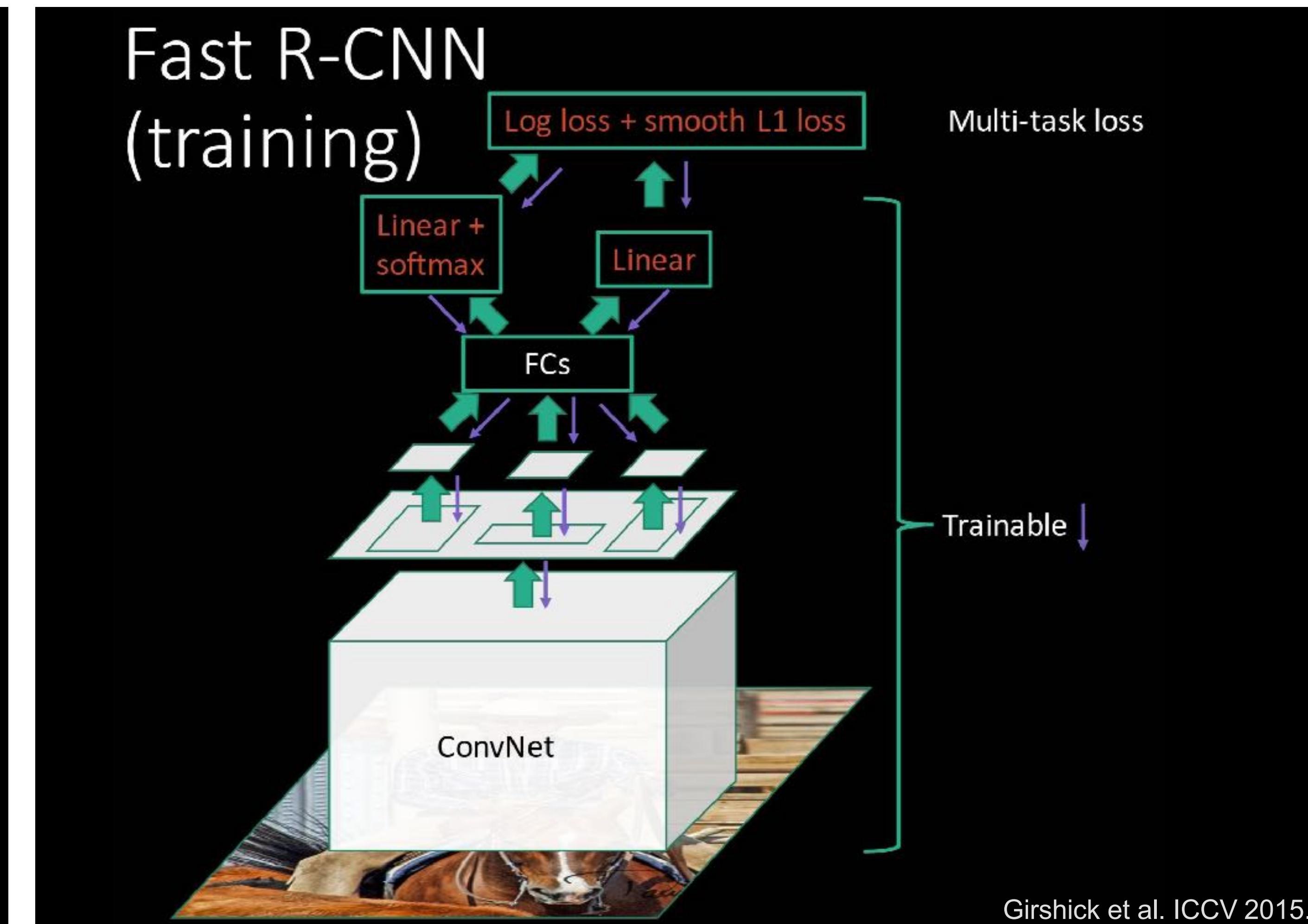
Problem #2: Post-hoc training: CNN not updated in response to final classifiers and regressors.

Problem #3: Complex training pipeline.



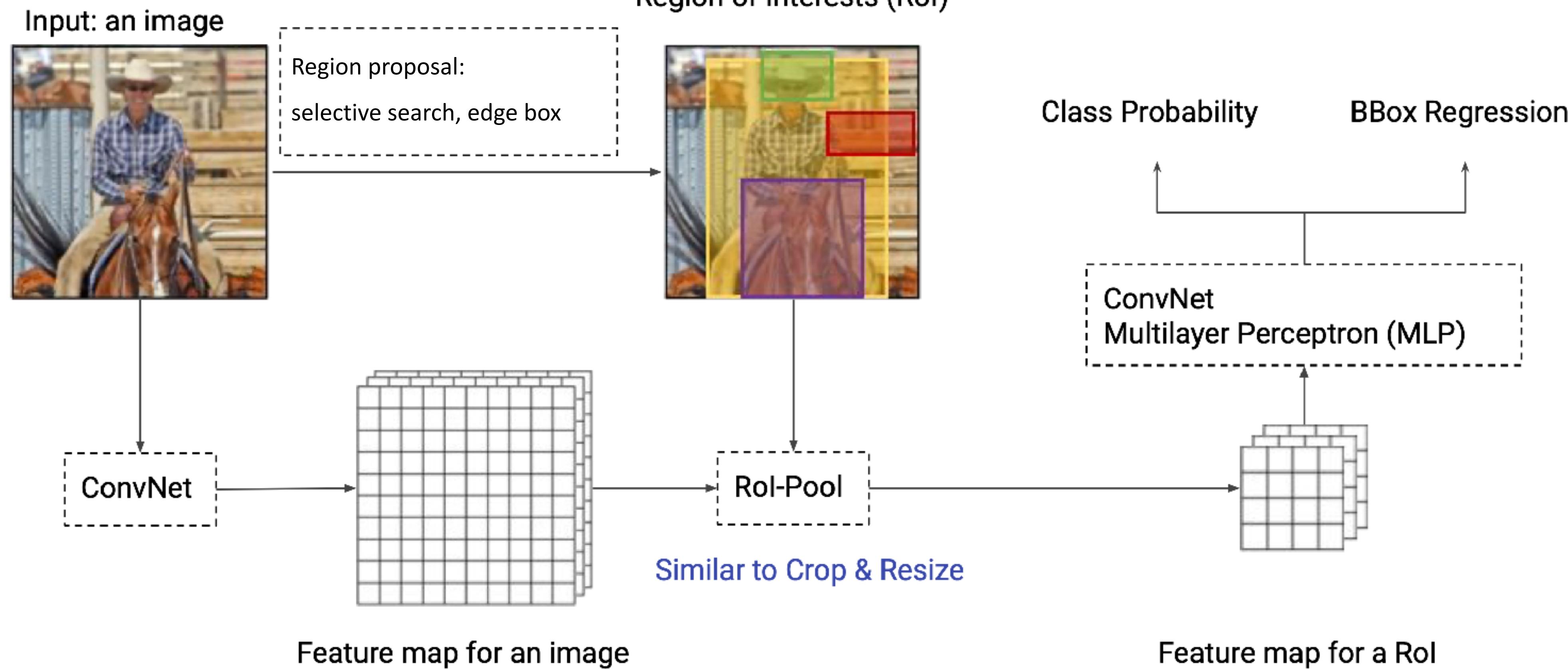
Fast R-CNN Solutions

Solution: Just train the whole system end-to-end all at once!

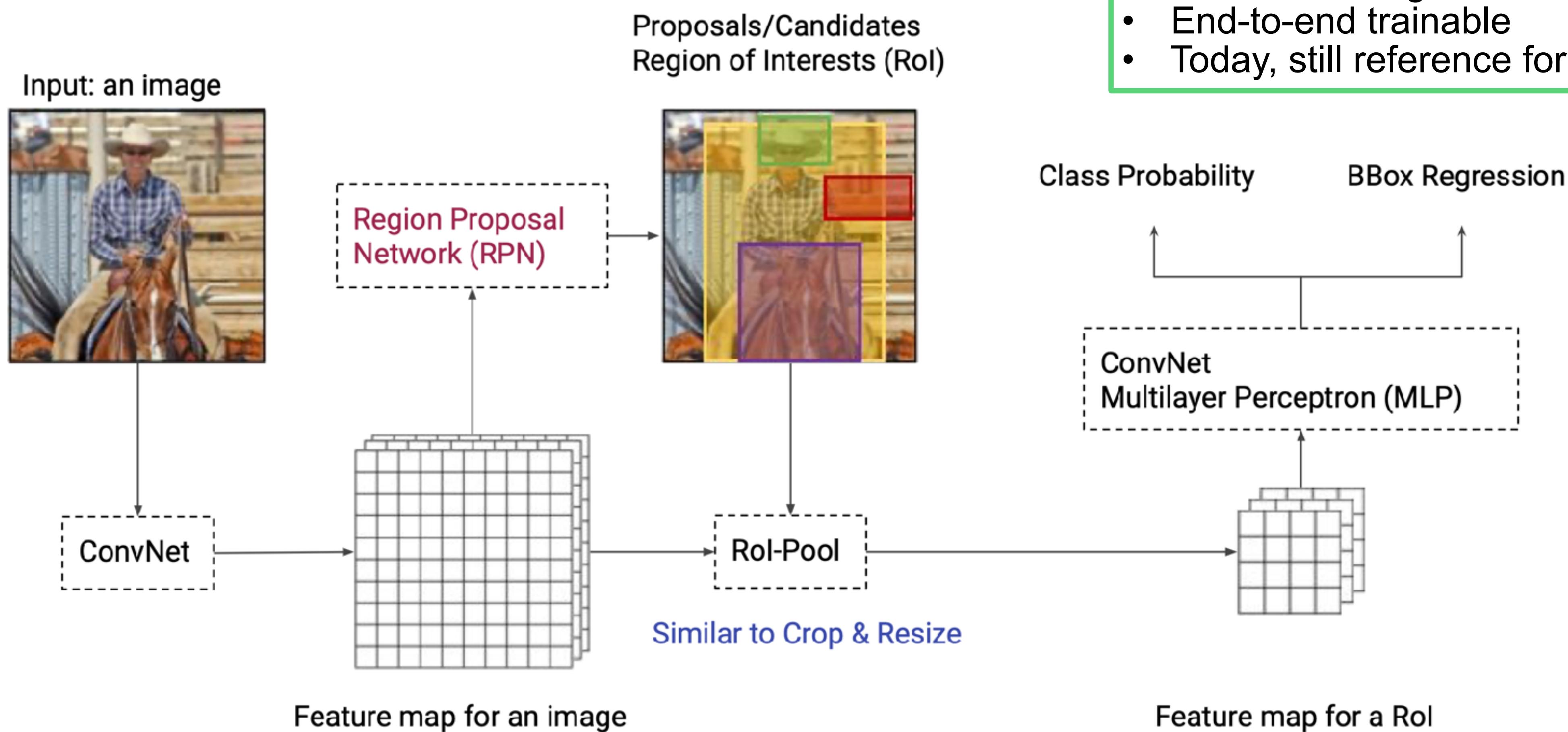


Fast R-CNN [ICCV 2015]

Region proposal is still independent
and can be slow (1-2 sec)



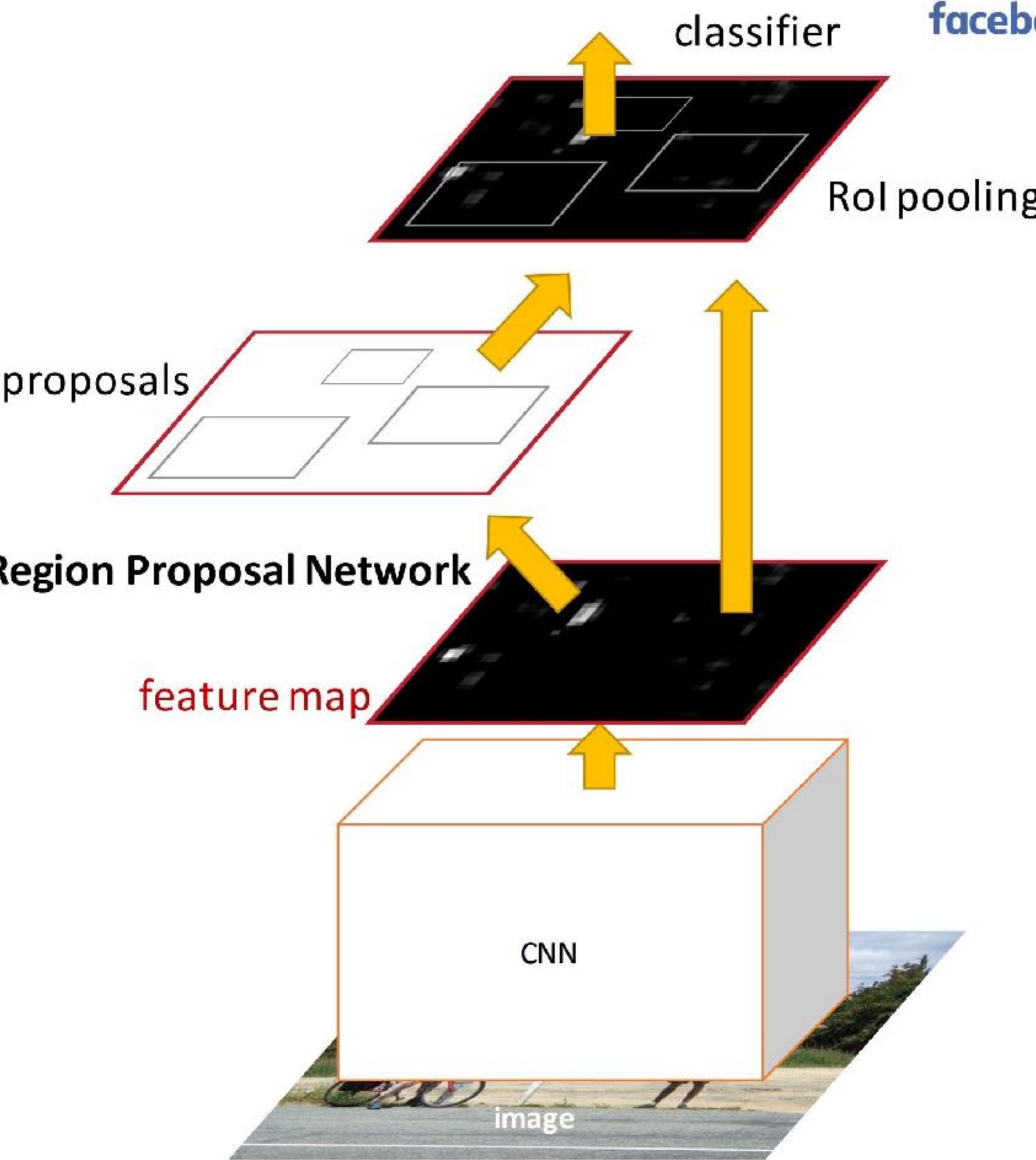
Faster R-CNN [NeurIPS 2015]



Faster R-CNN

facebook

- Insert a **Region Proposal Network (RPN)** after the last convolutional layer.
- RPN trained to produce region proposals directly; no need for external region proposals!
- After RPN, use “RoI Pooling” and an upstream classifier and bbox regressor just like Fast R-CNN.



RPN: Region Proposal Network

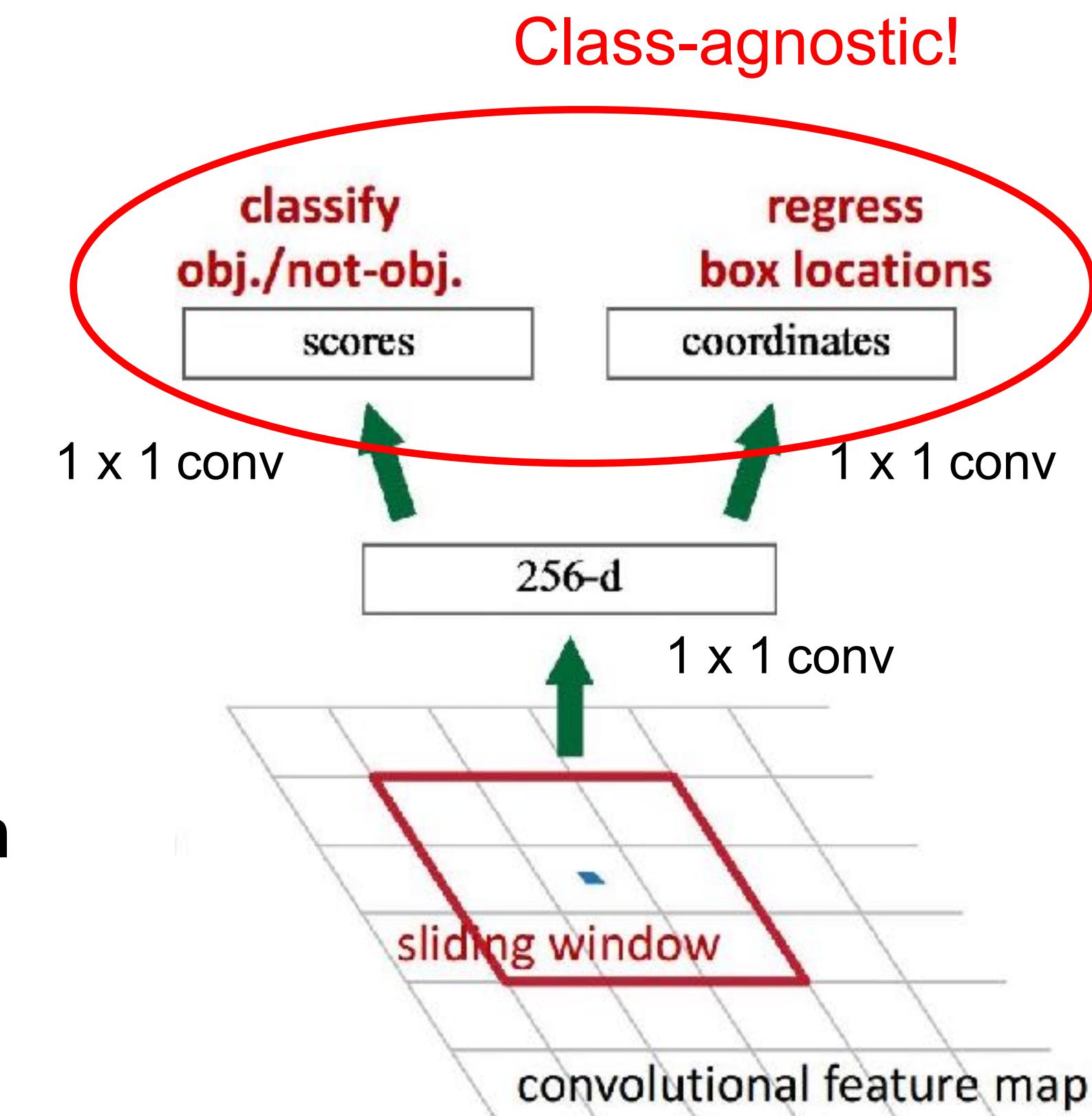
Slide a small window on the feature map

Build a small network for:

- classifying object or not-object, and
- regressing bbox locations

Position of the sliding window provides localization information with reference to the image

Box regression provides finer localization information with reference to this sliding window



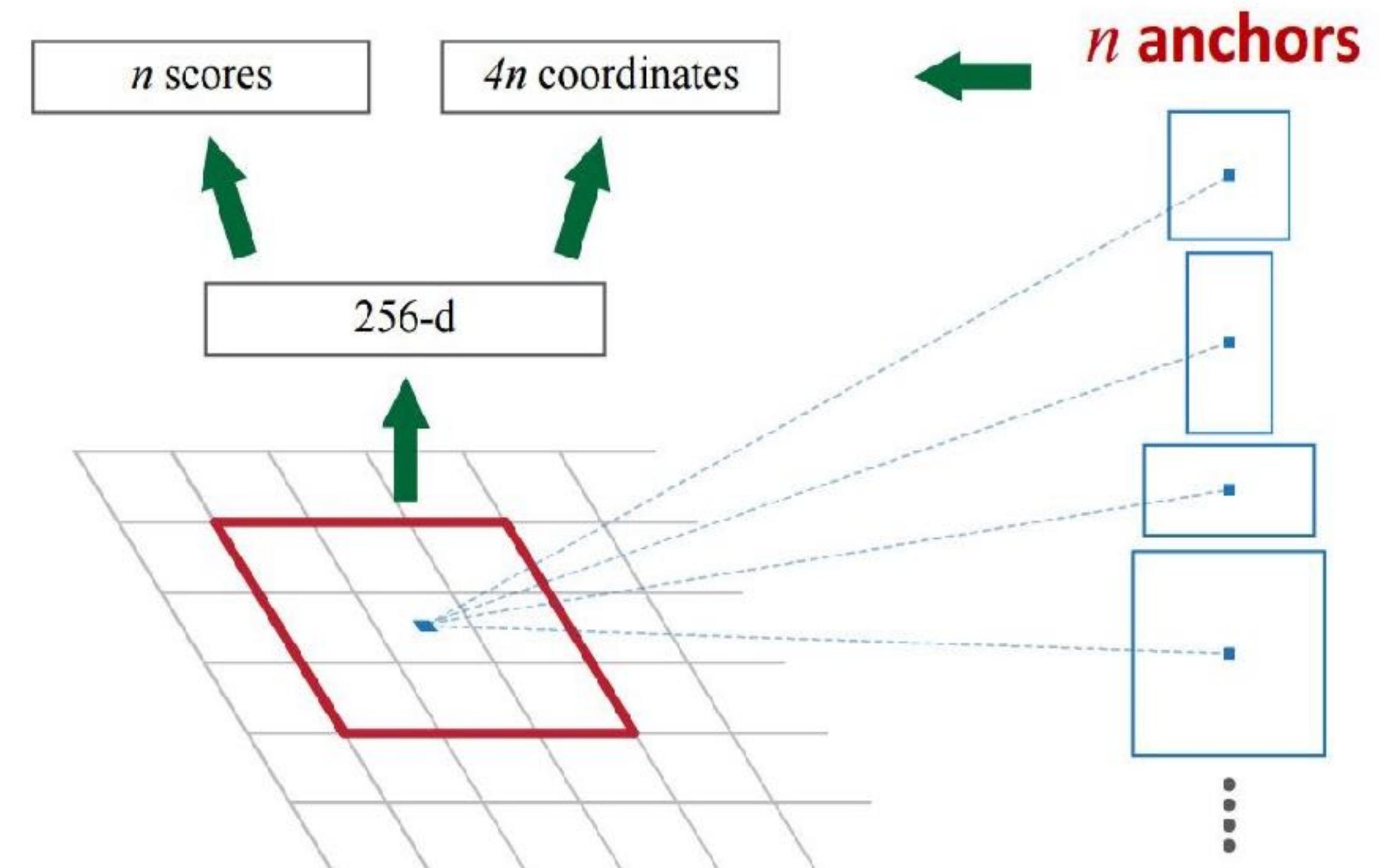
RPN: Region Proposal Network

Use **N anchor boxes** at each location.

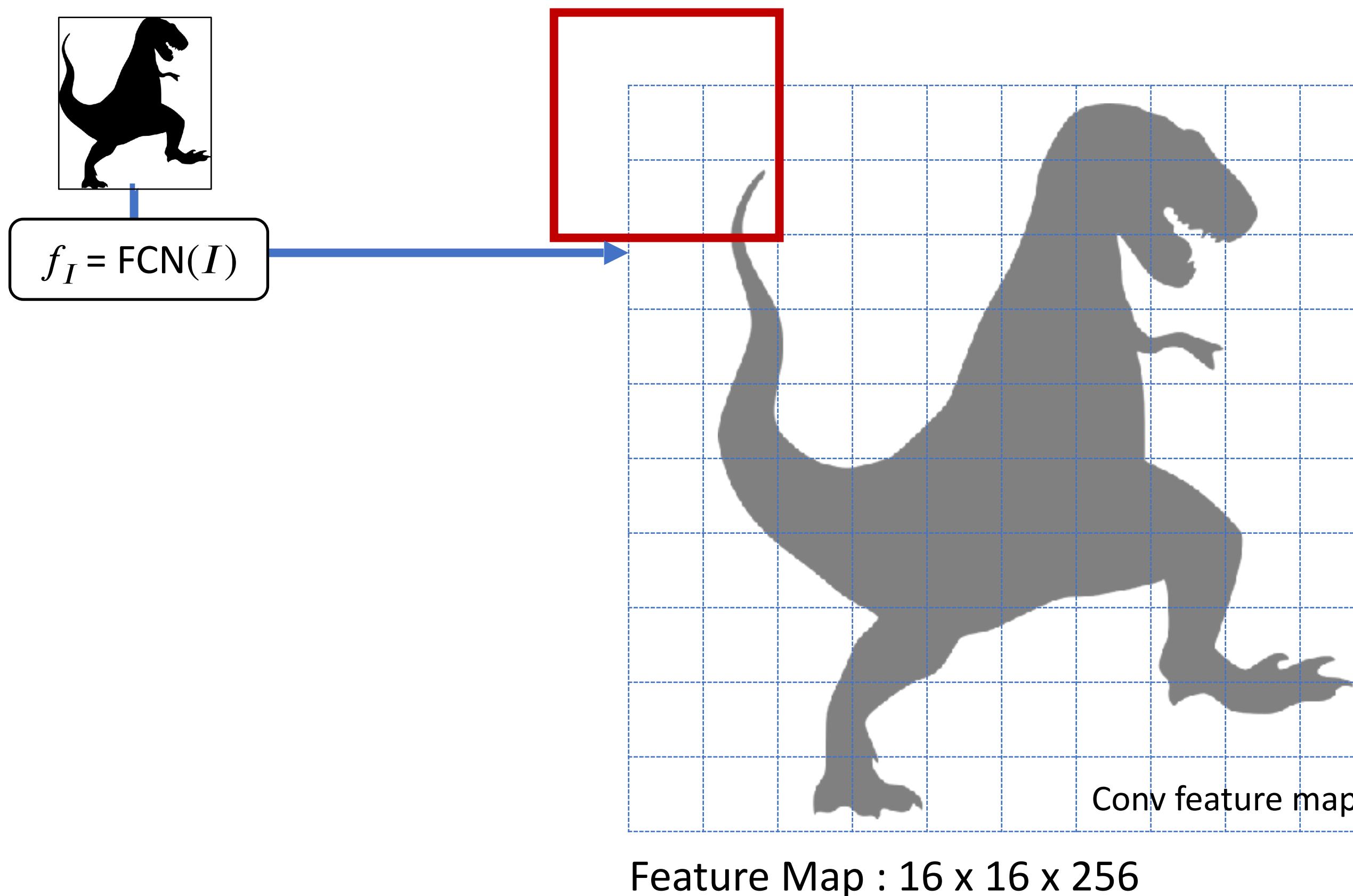
Anchors are **translation invariant**: use the same ones at every location.

Regression gives offsets from anchor boxes.

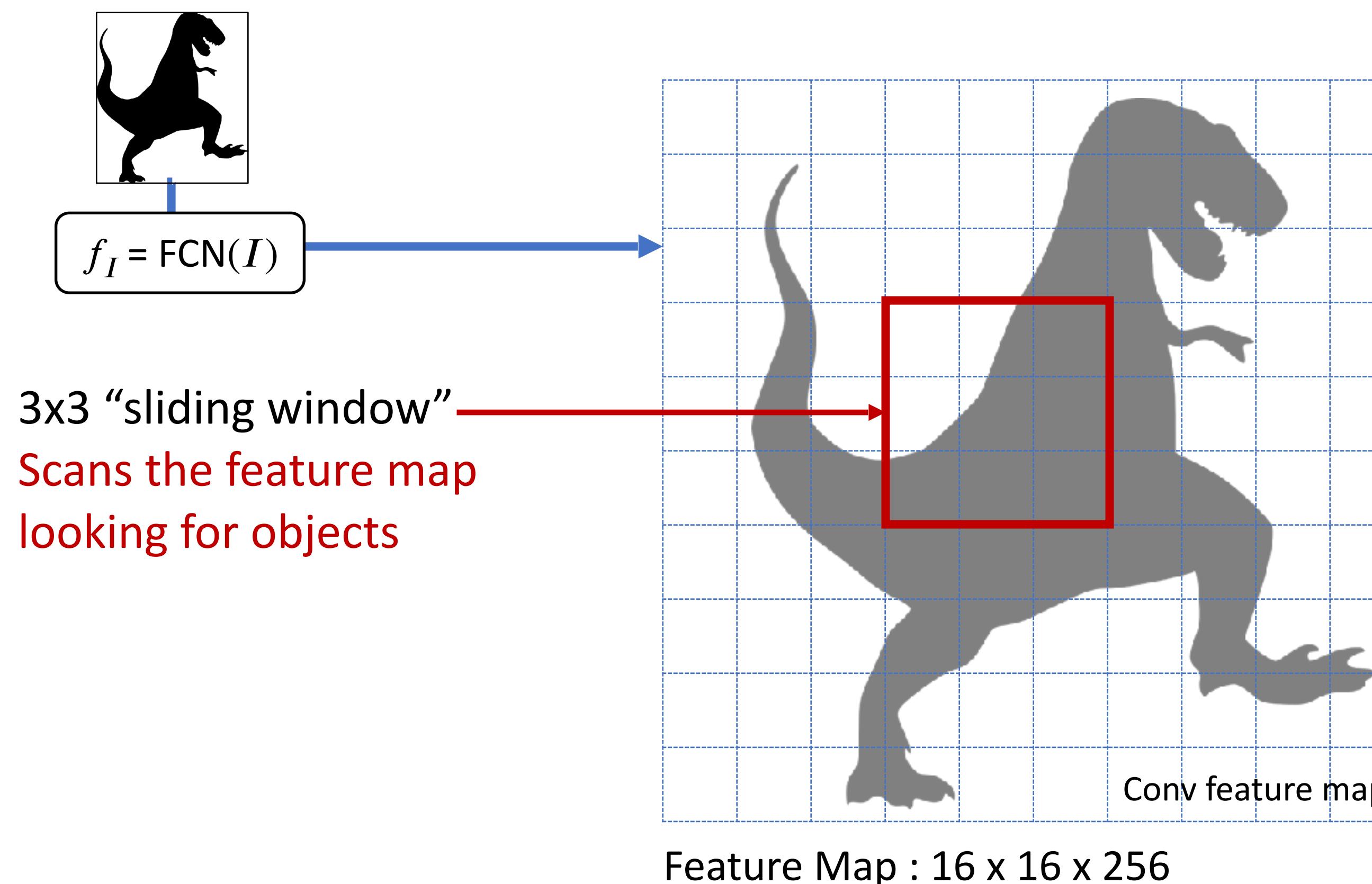
Classification gives the probability that each (regressed) anchor shows an object.



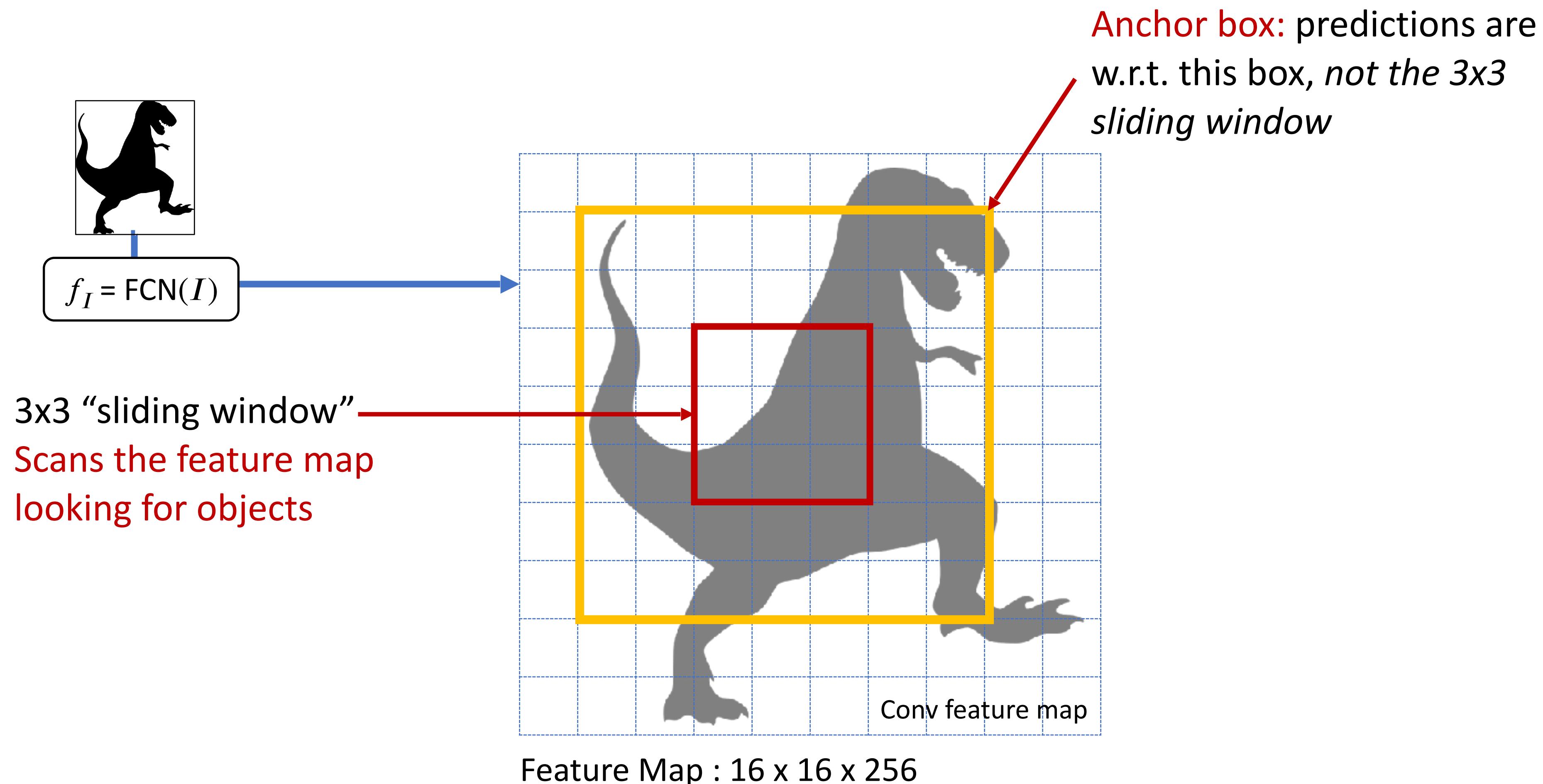
RPN: Region Proposal Network



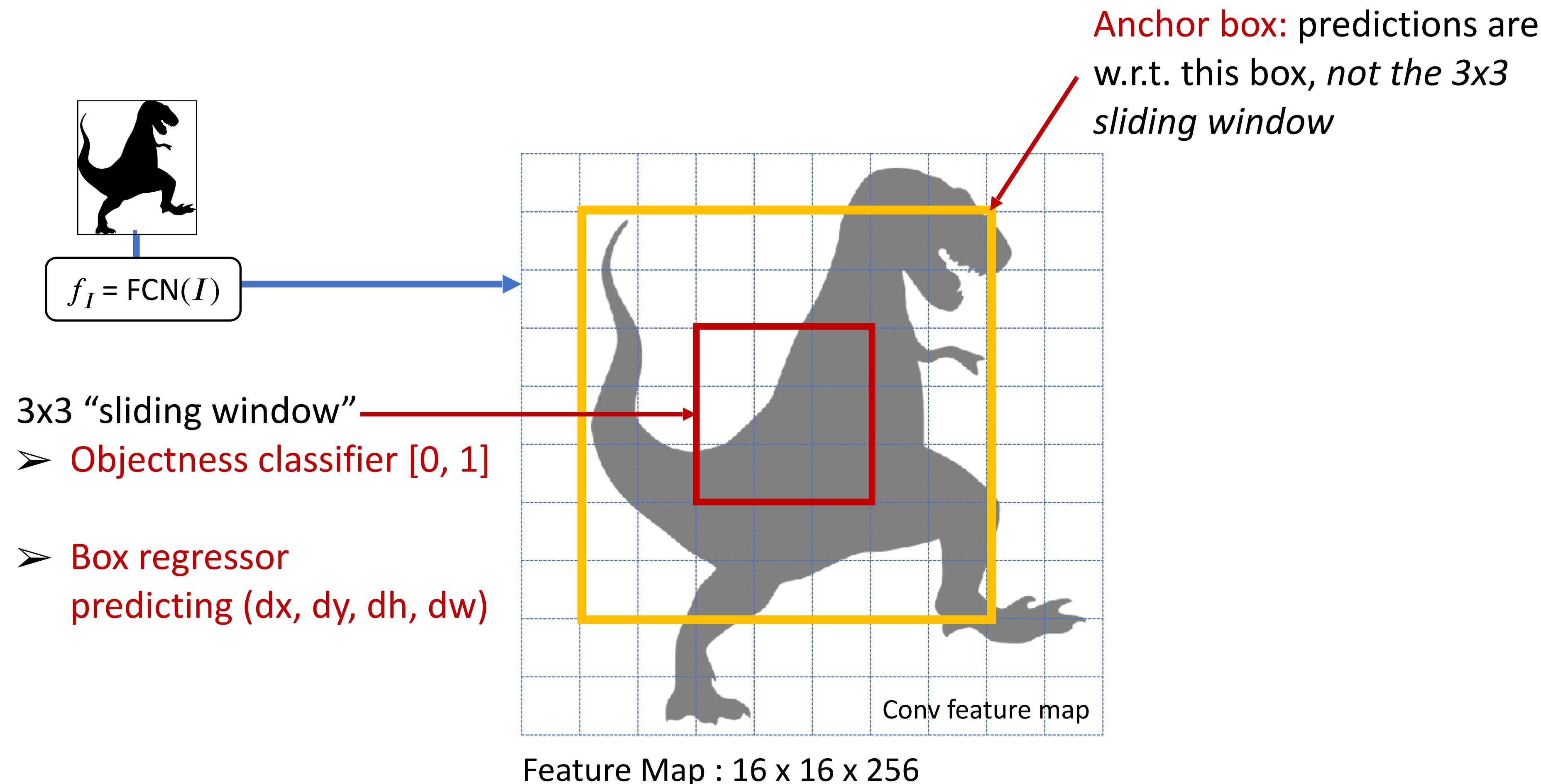
RPN: Region Proposal Network



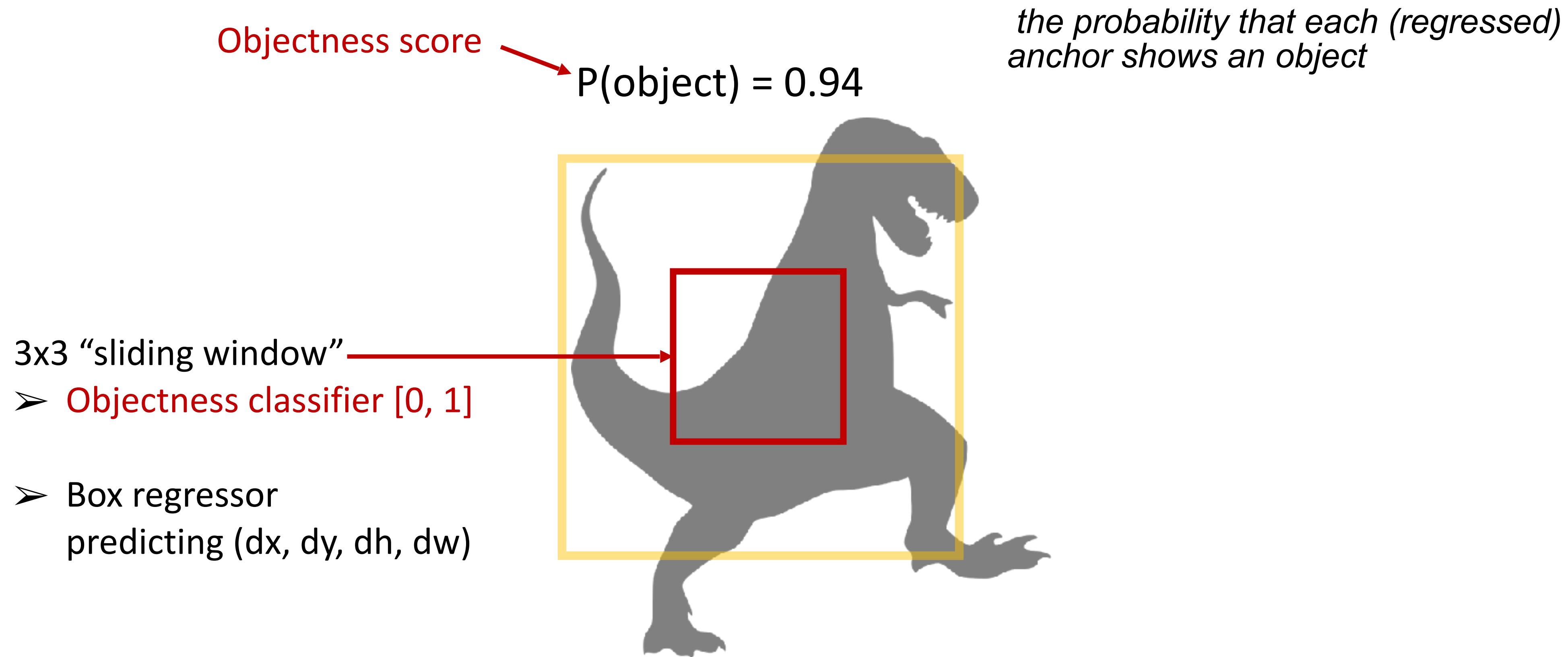
RPN: Anchor Box



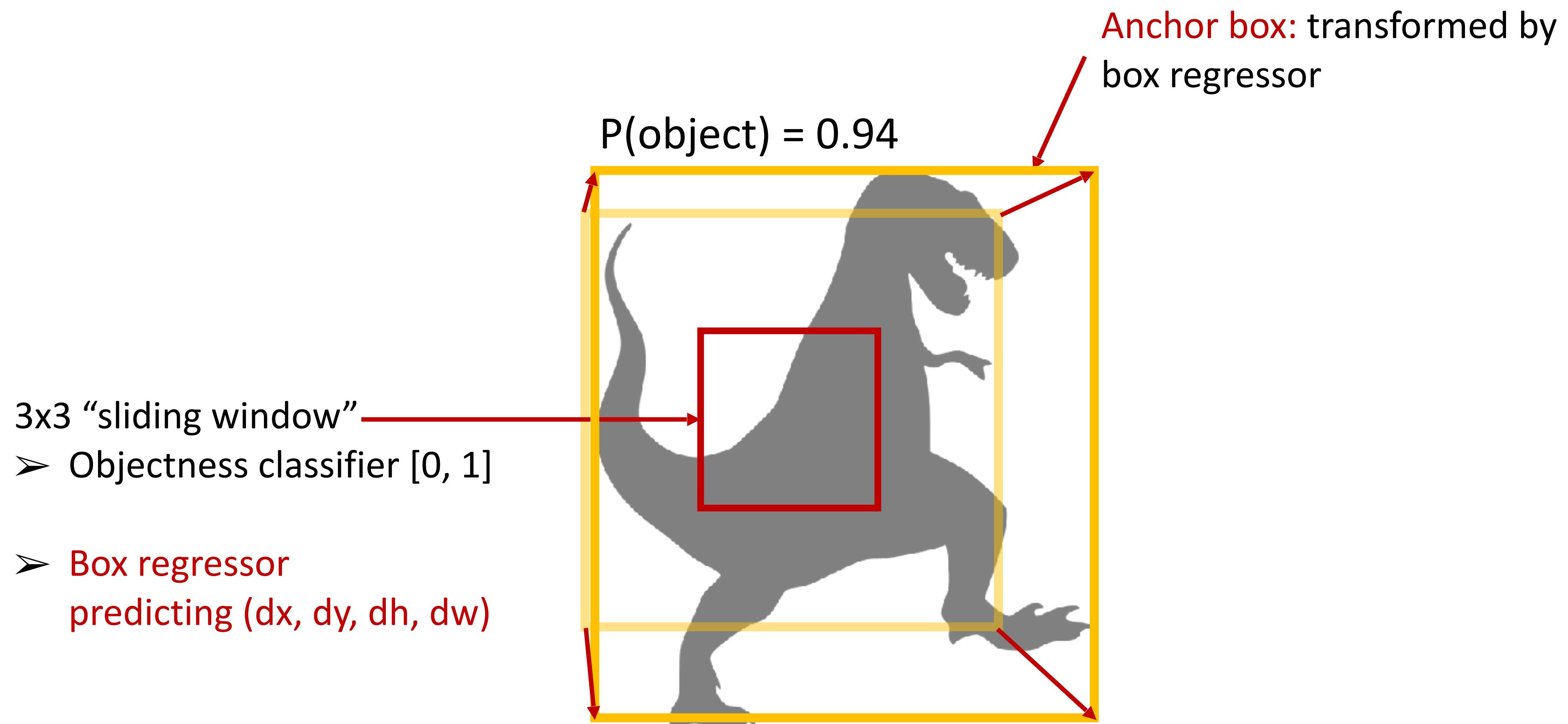
RPN: Anchor Box



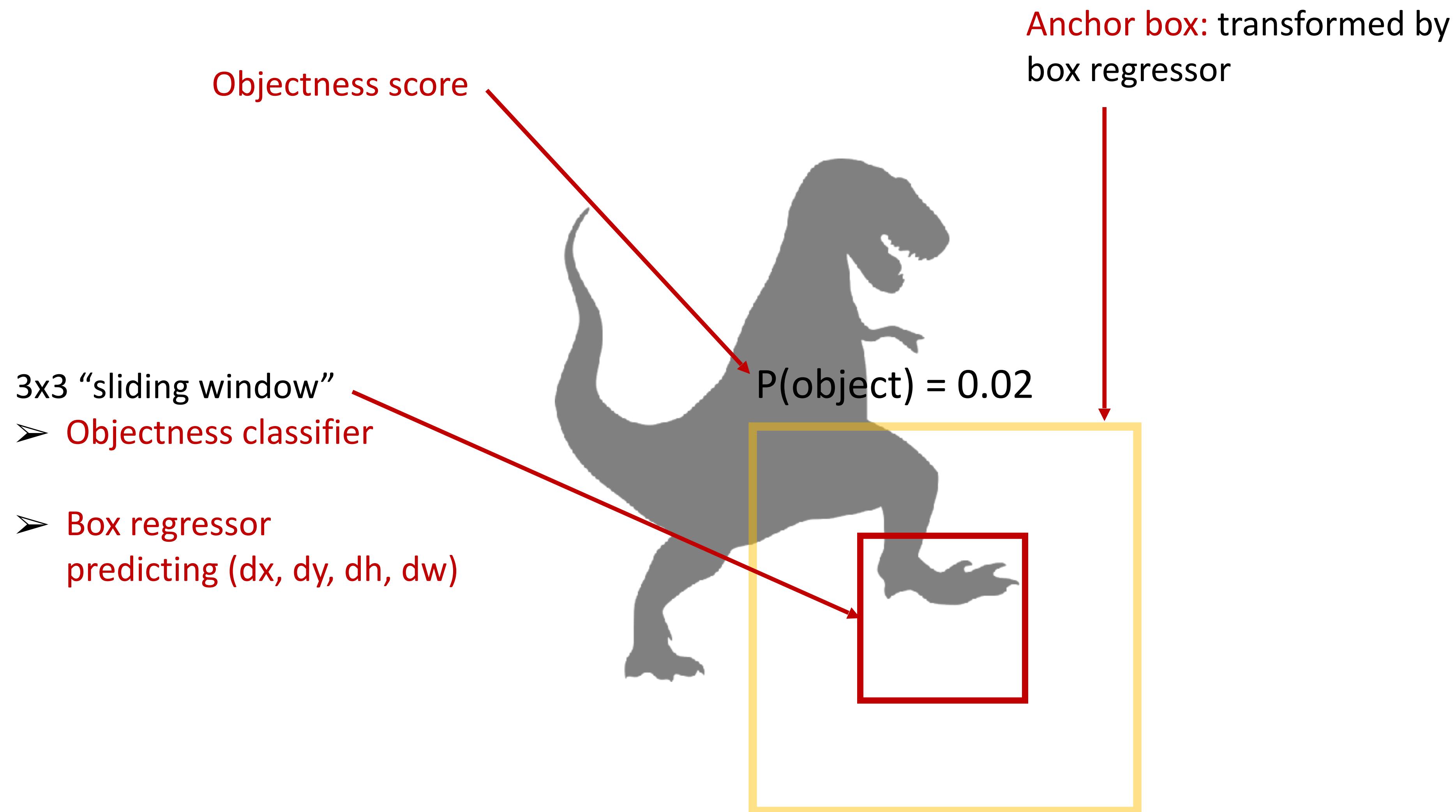
RPN: Prediction (on object)



RPN: Prediction (on object)



RPN: Prediction (off object)

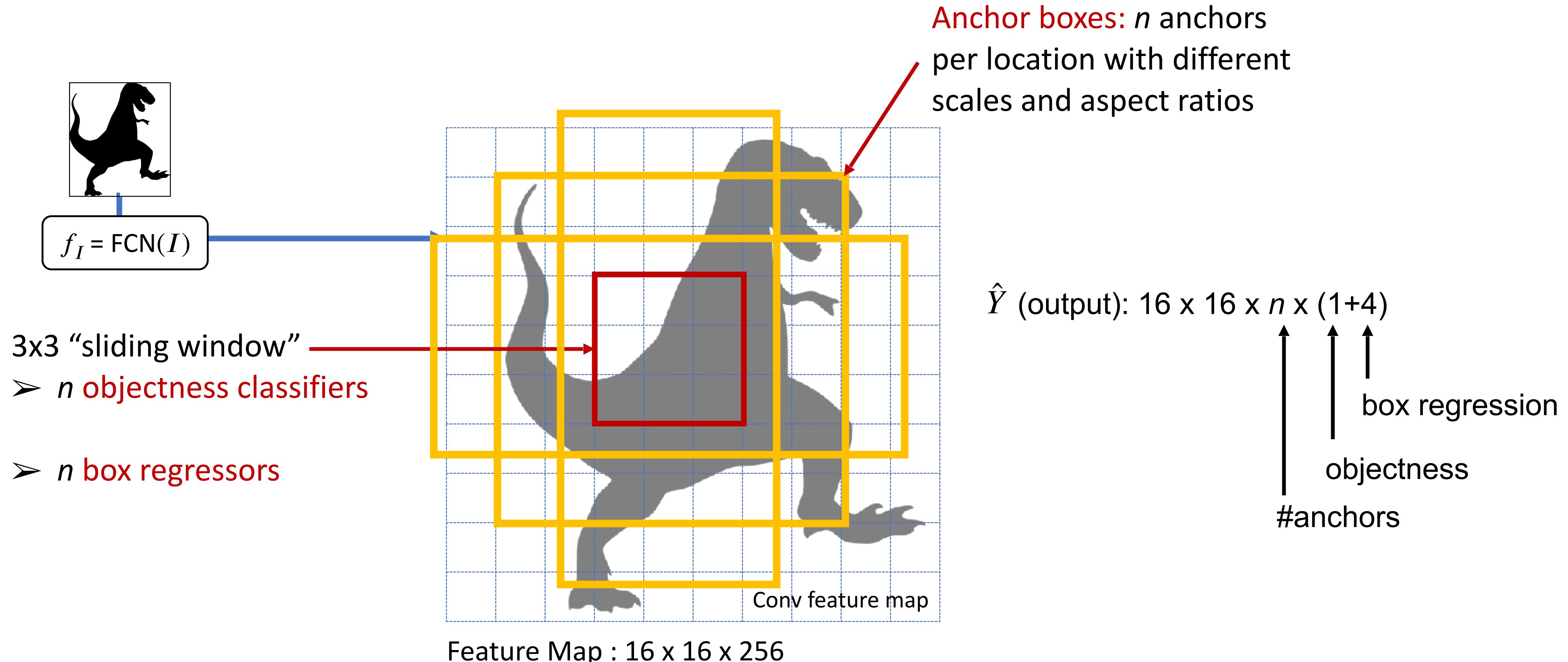


RPN: Multiple Anchors

Y (ground truth) : [1 1 1 1 1 1 0 0 0 0 1 0 0 0 0]

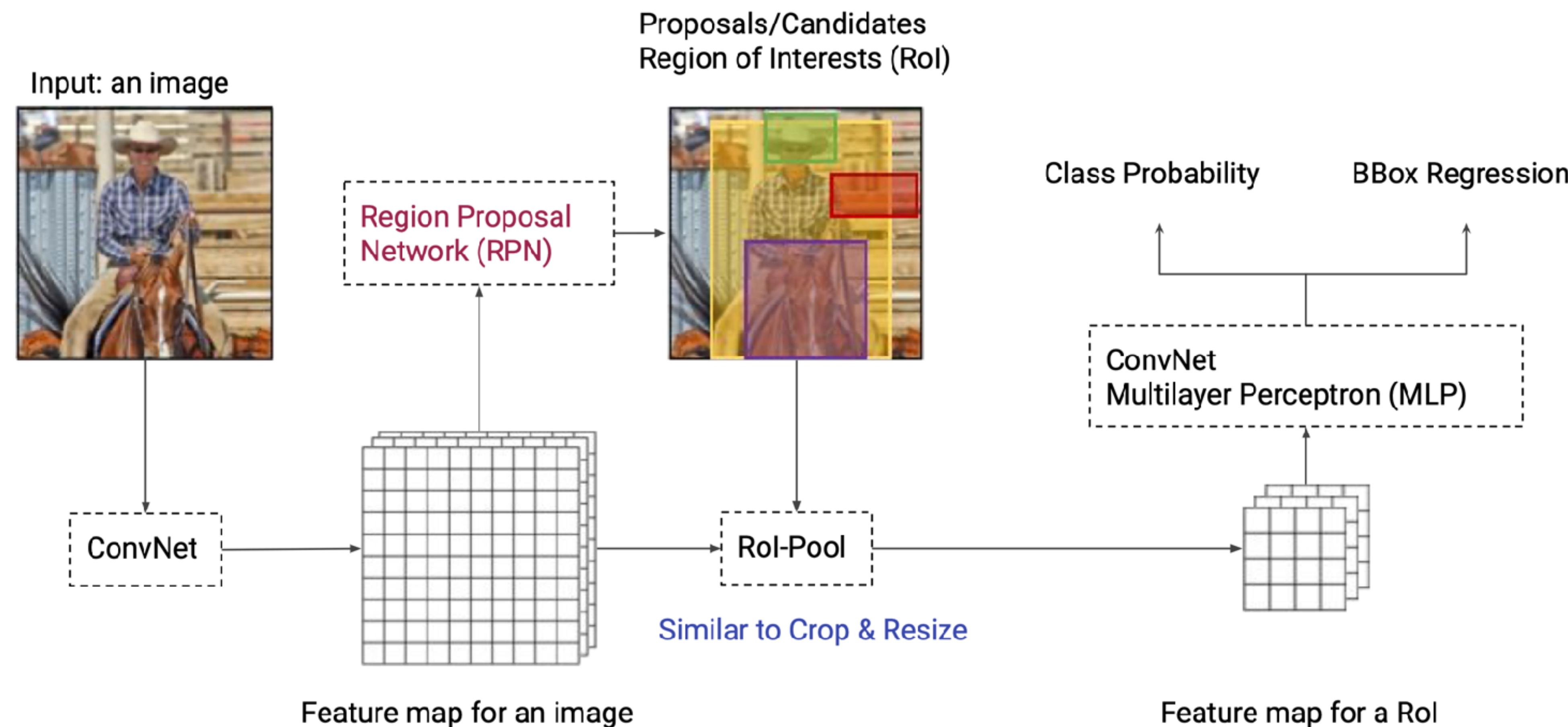
M (mask) : [1 dx dy dh dw 0 - - - 0 - - - -]

$$Loss = \sum_i M_i \cdot \mathcal{L}(\hat{Y}_i, Y_i)$$



Faster R-CNN

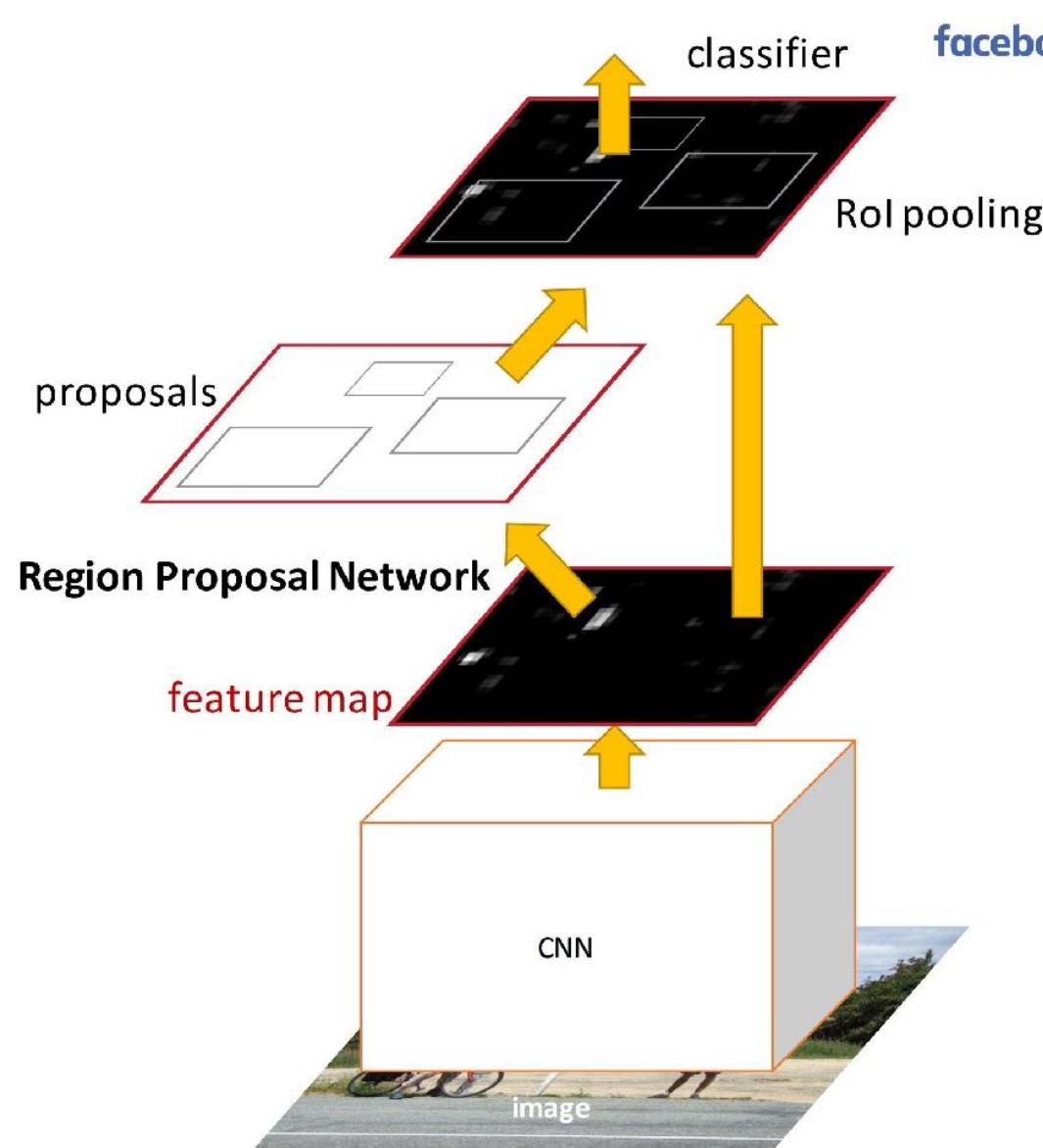
- Still two-stage



Object detection: 2-stage vs 1-stage

CNN Detectors

2-stage detectors

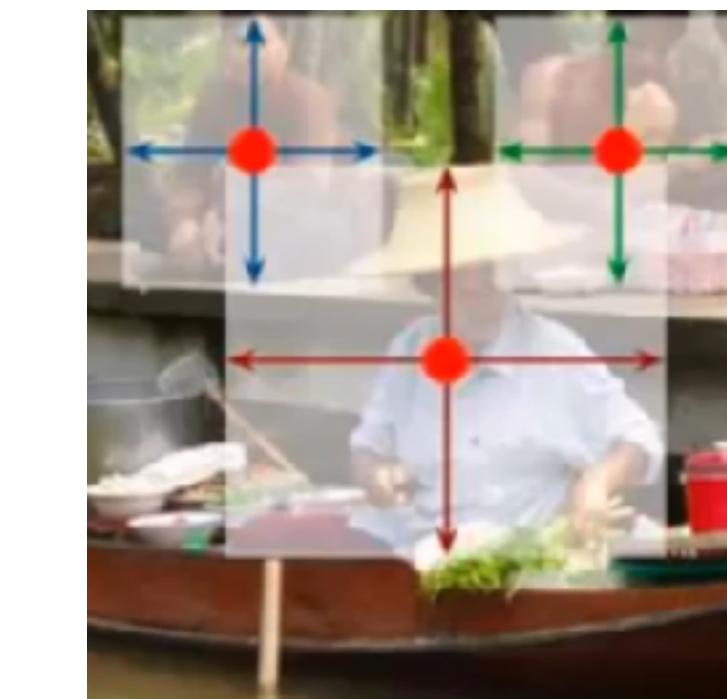


Fast(er) RCNN, Mask-RCNN,
SNIPPER, PANet, TridentNet

1-stage detectors



Anchor-free
detectors



- No object proposals
- Use **anchors**
- Faster but less accurate

YOLO, SSD, RetinaNet,
EfficientDet

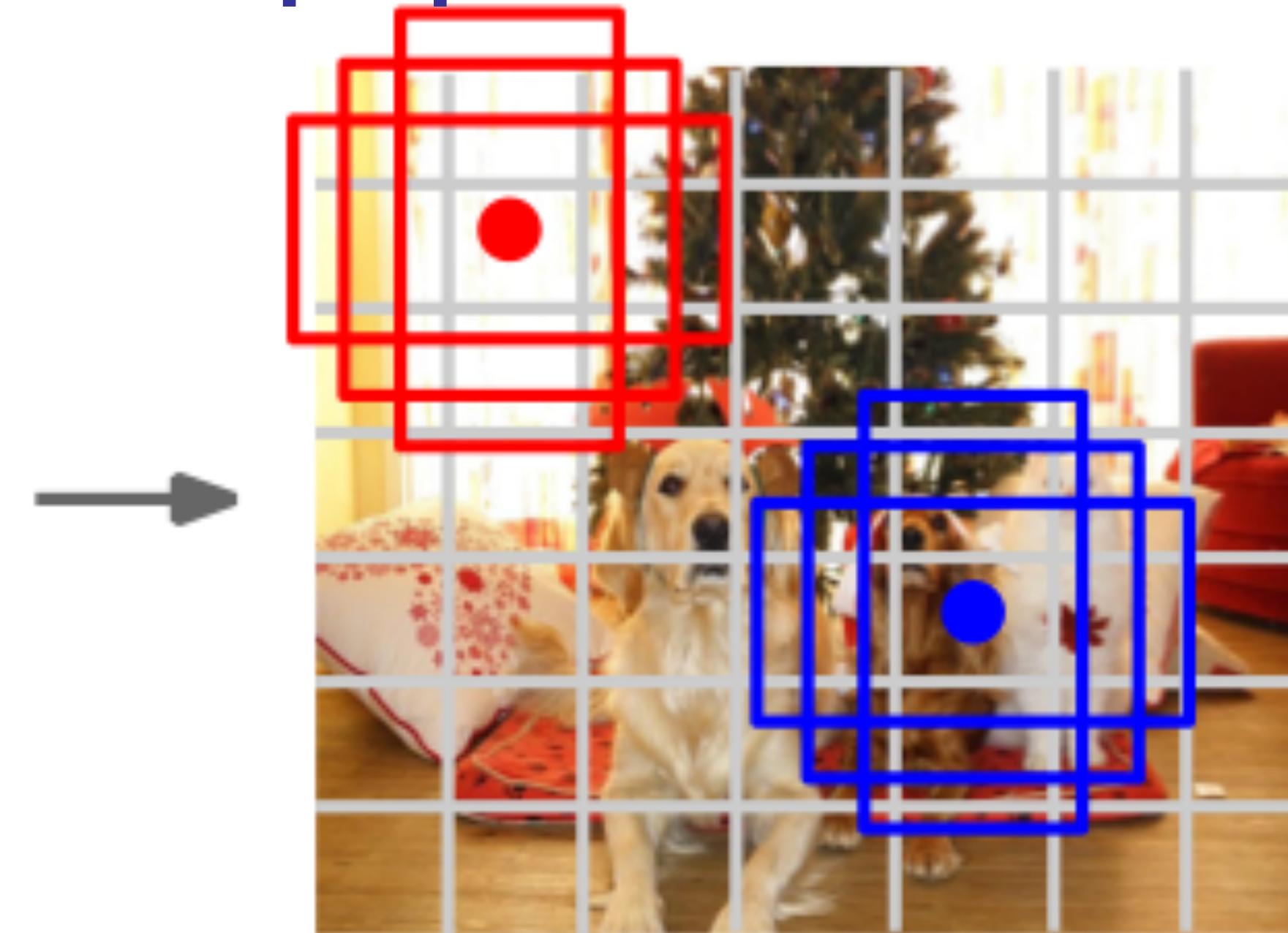
CornerNet, CenterNet,
FCOS, ExtremeNet

1-stage object detection: YOLO/SSD

Detection without proposals



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

Image a set of **base boxes**
centered at each grid cell
Here $B = 3$

Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", CVPR 2016
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

1-stage object detection: YOLO/SSD

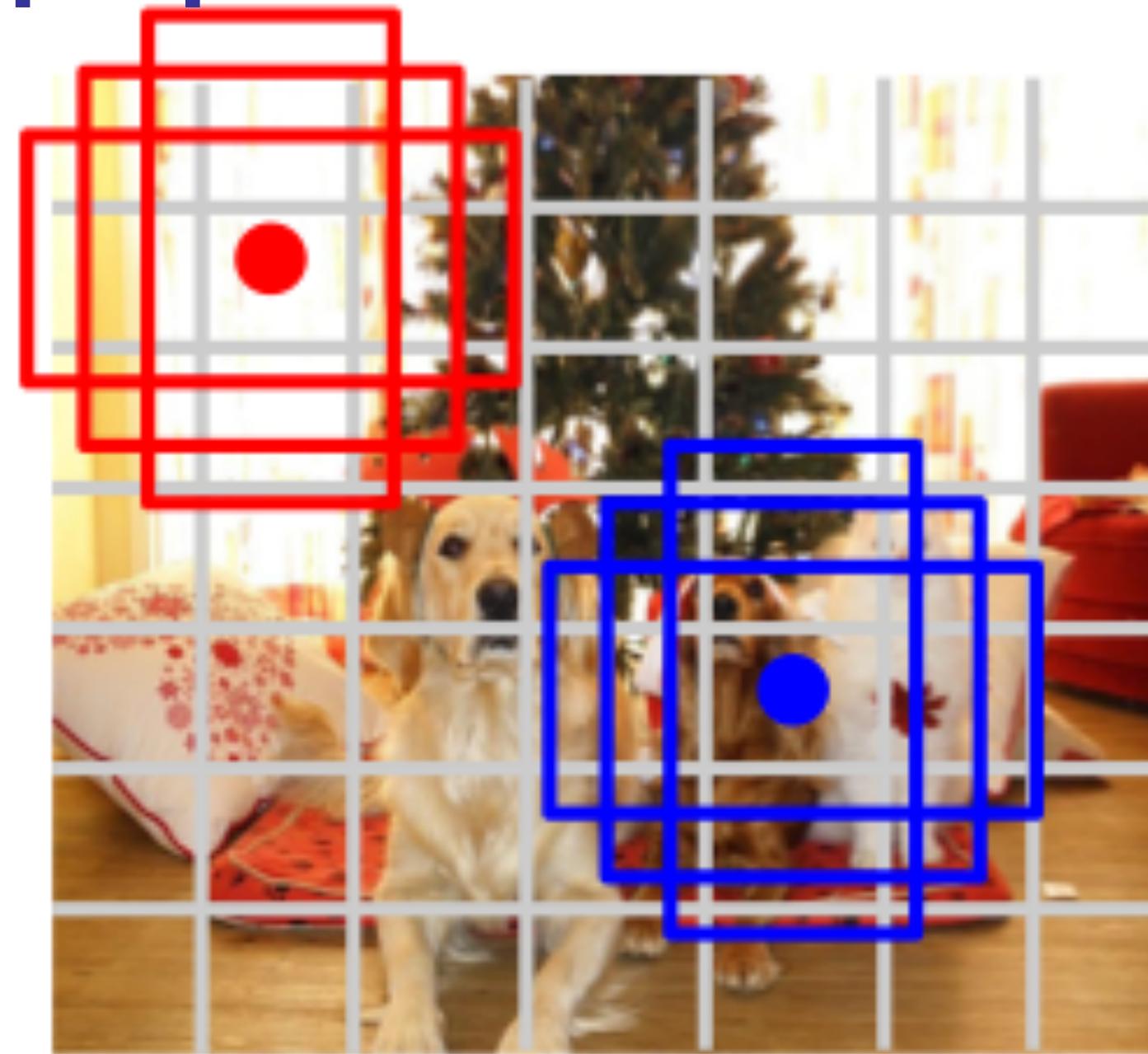
Detection without proposals

Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers:
 $(dx, dy, dh, dw, \text{confidence})$
- Predict scores for each of C classes (including background as a class)

Output:

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



Divide image into grid
 7×7

Image a set of **base boxes**
centered at each grid cell
Here $B = 3$

From input image to scores with a single network. **Faster but not as accurate as RCNN.**

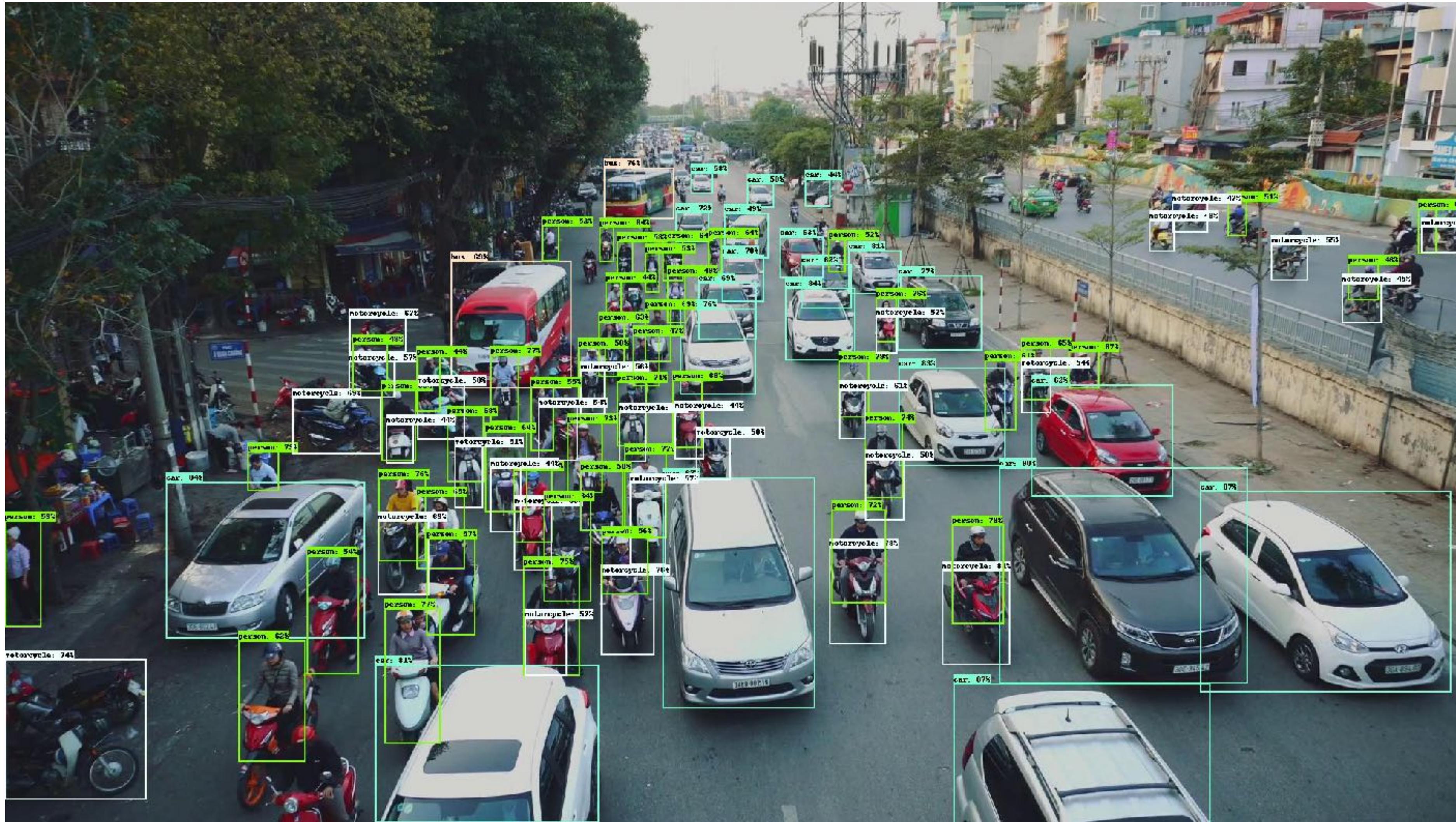
See also: Lin et al., Focal loss for dense object detection, ICCV 2017.

Yolo v2 Demo video



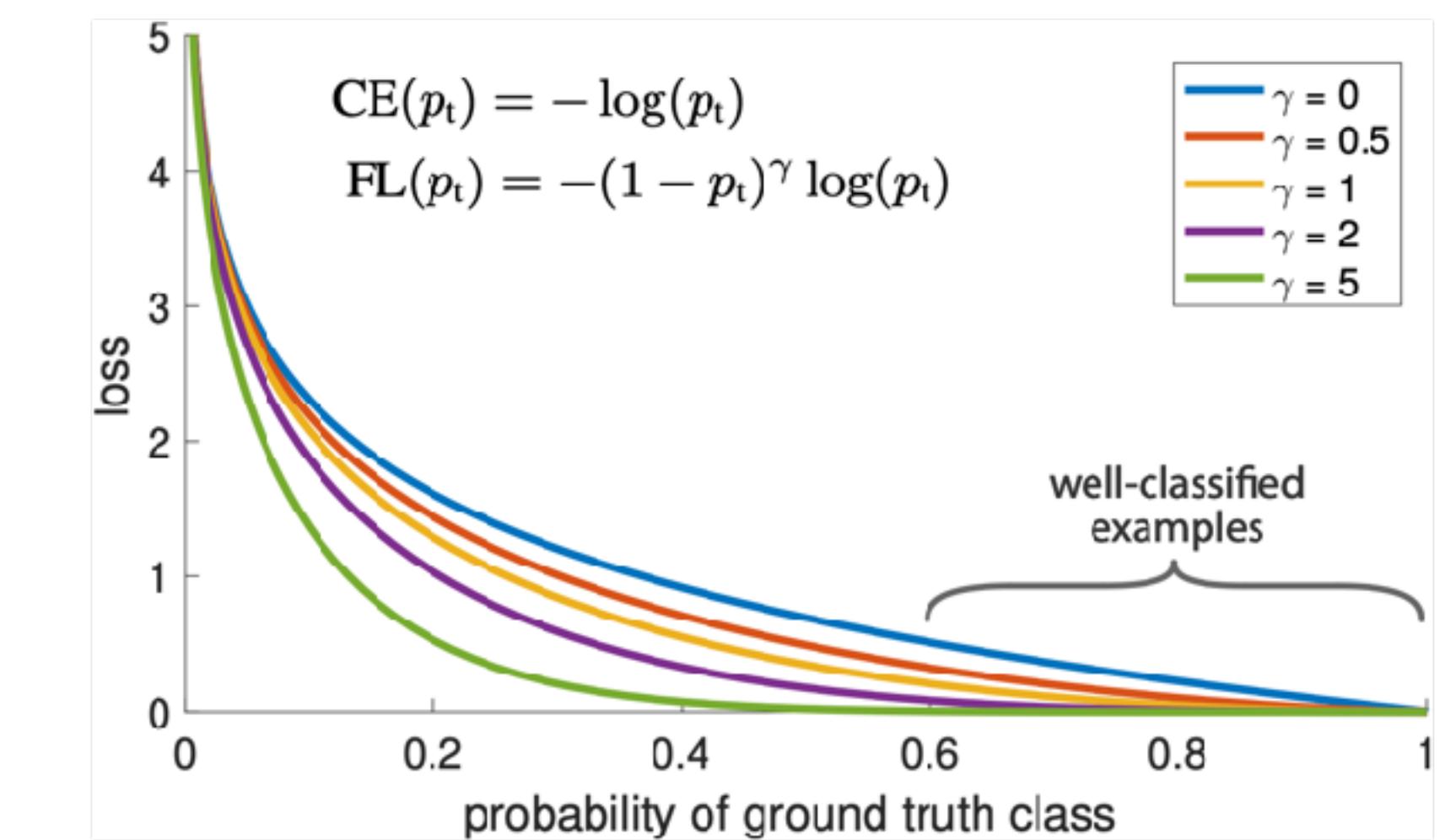
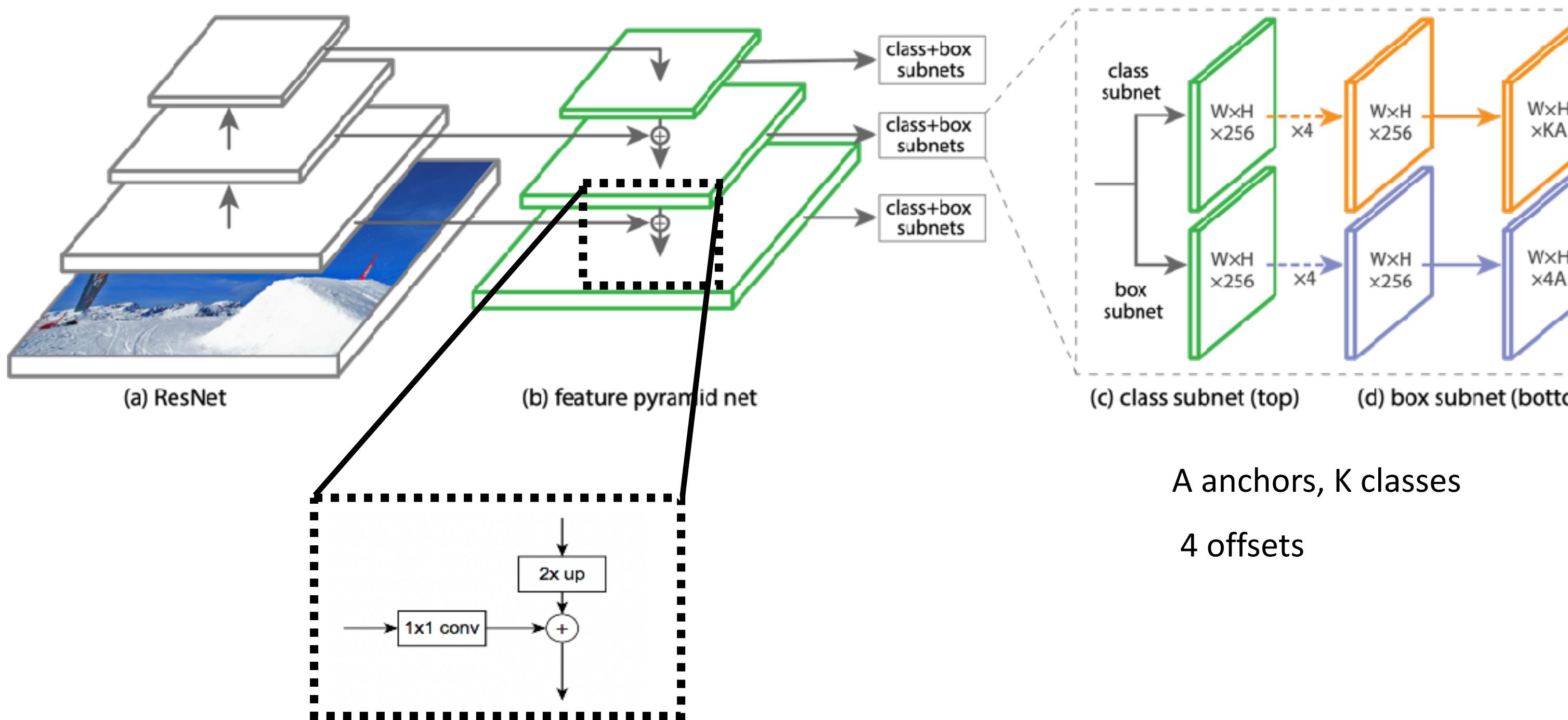
Scale in object detection

Problem with YOLO: Single cell can corresponds to multiple objects, even with multiple anchors, still **too coarse**.



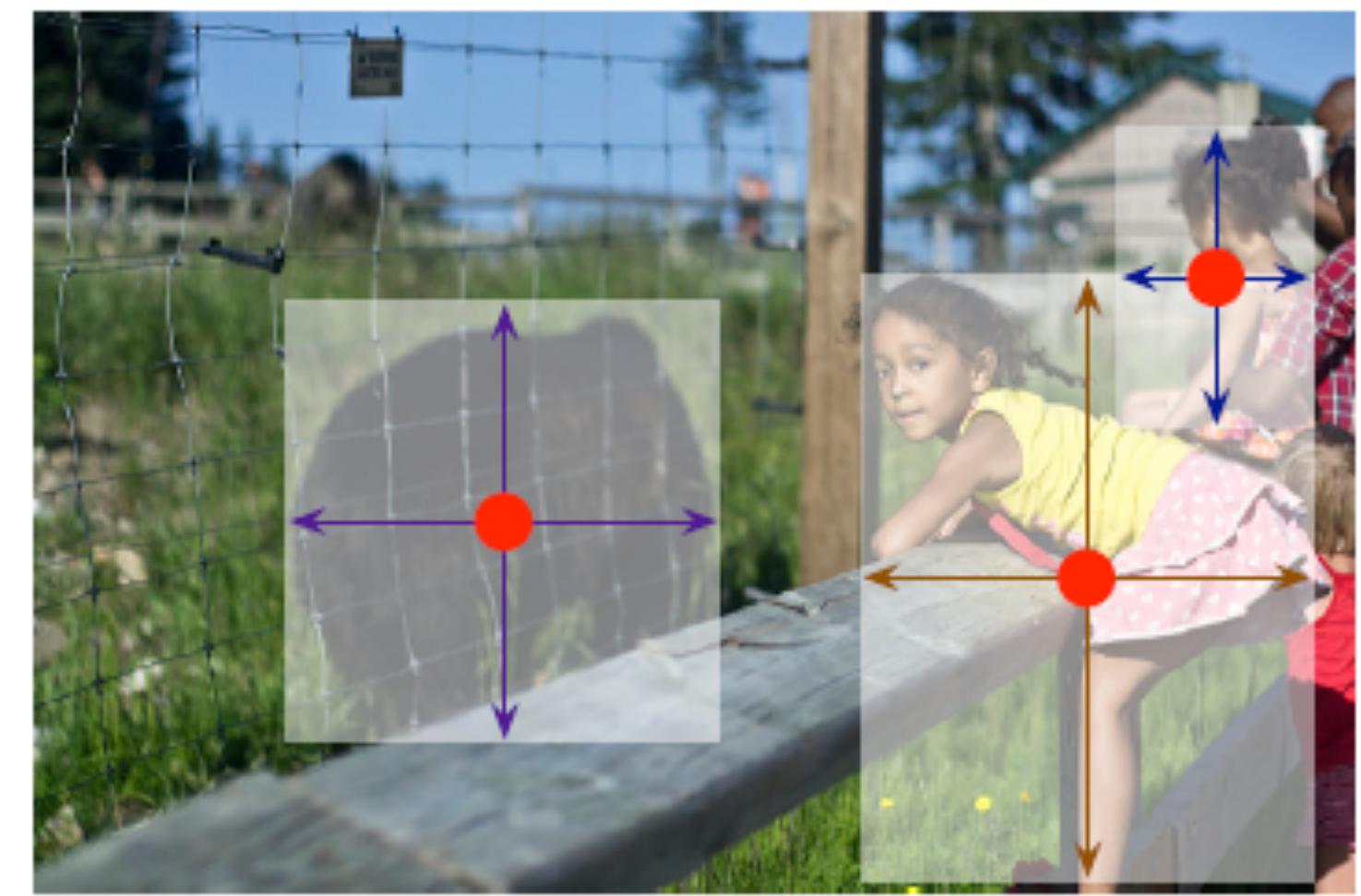
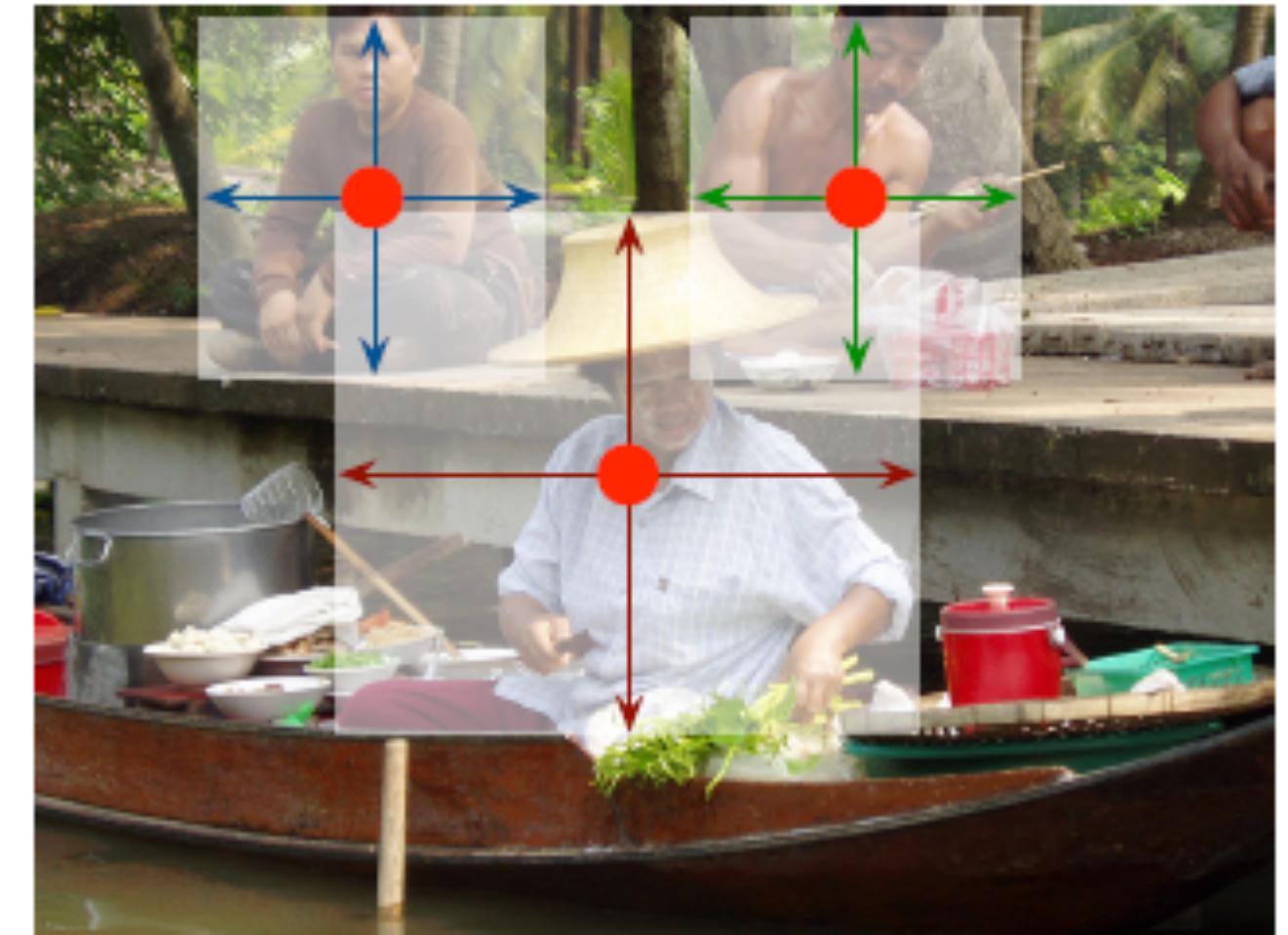
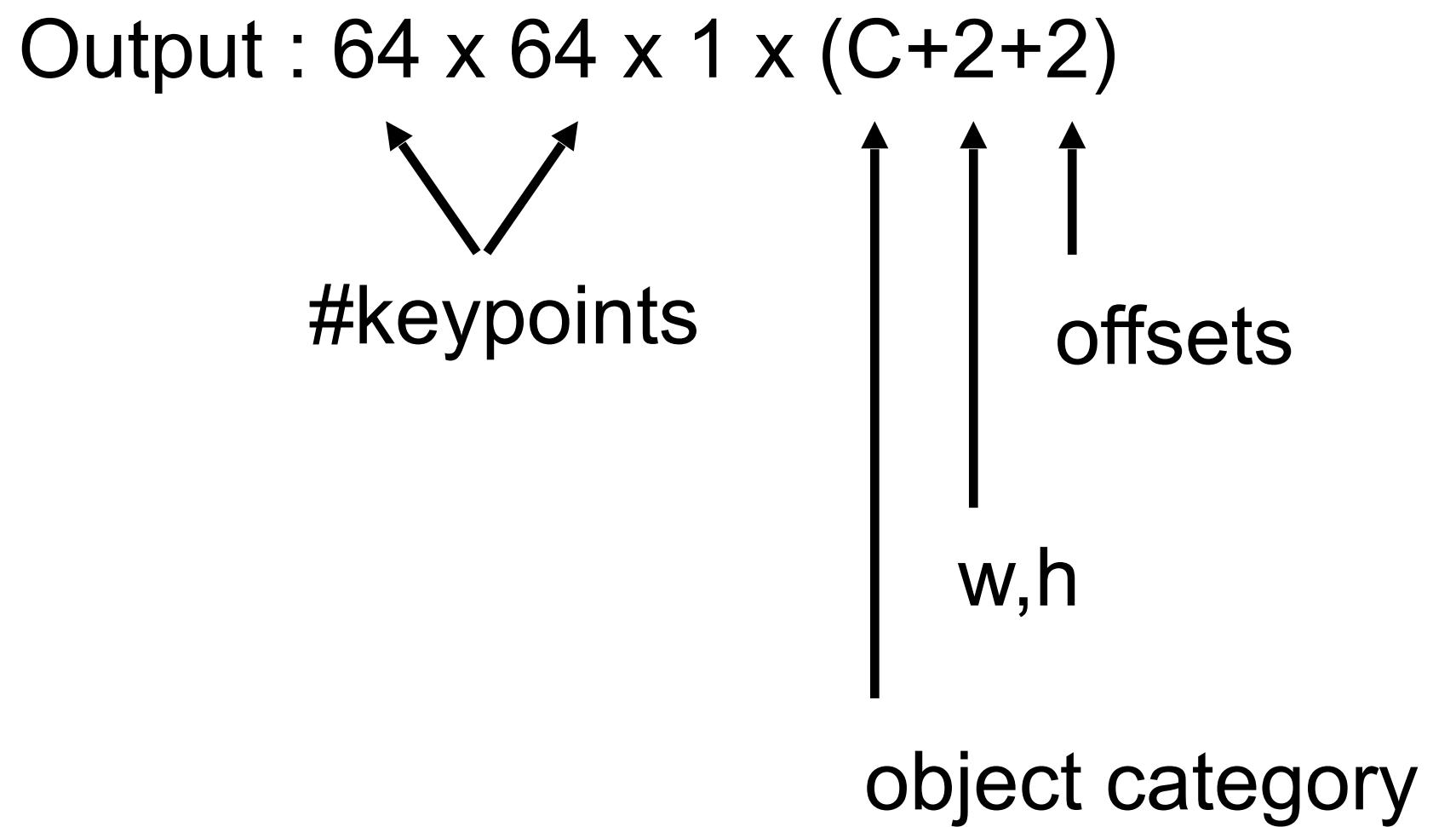
1-stage object detection: RetinaNet

- Pre-define anchor boxes on **multiple scales**, e.g., Feature Pyramid Networks (FPNs).
- 6 anchors per location, 100 - 200k anchor boxes to classify per image (dense detection).
- Focal loss for soft-version hard sample mining.



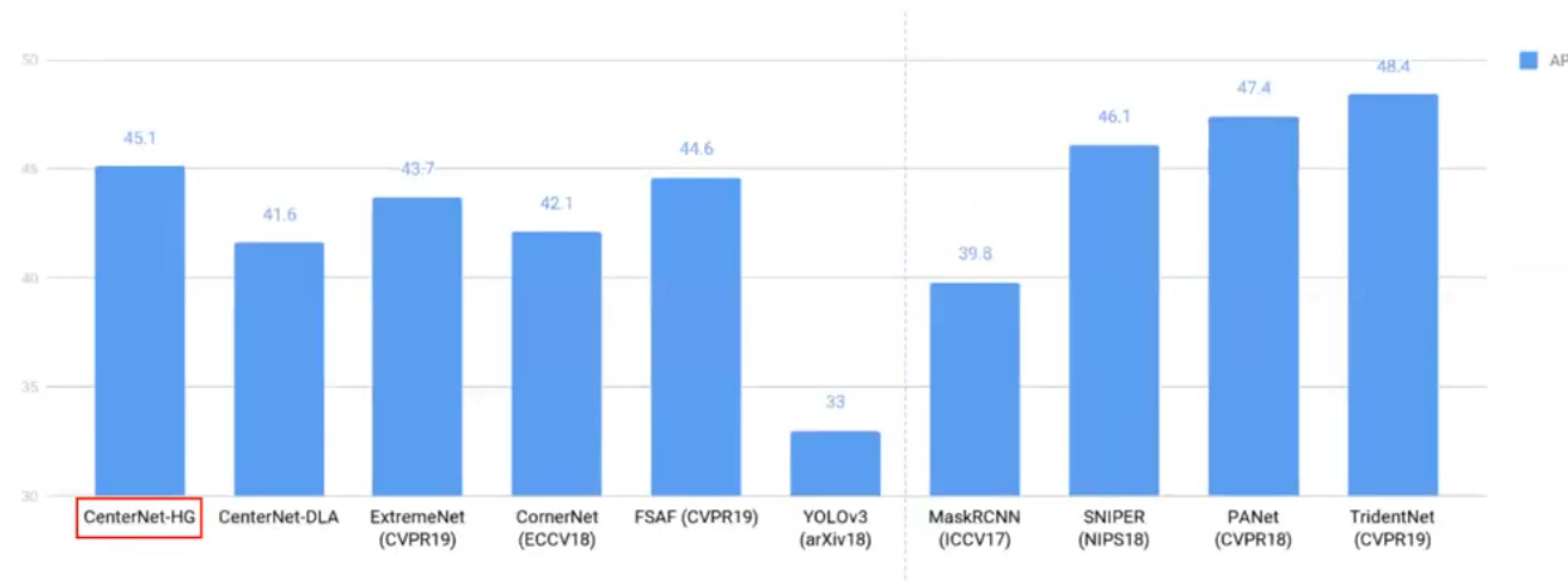
1-stage object detection: CenterNet (anchor-free)

- Represent objects by a single point + (width, height)
- Regress other parameters such as
 - Bounding box
 - 3D box
 - human pose
 - ...



State-of-the-art comparison: MS COCO

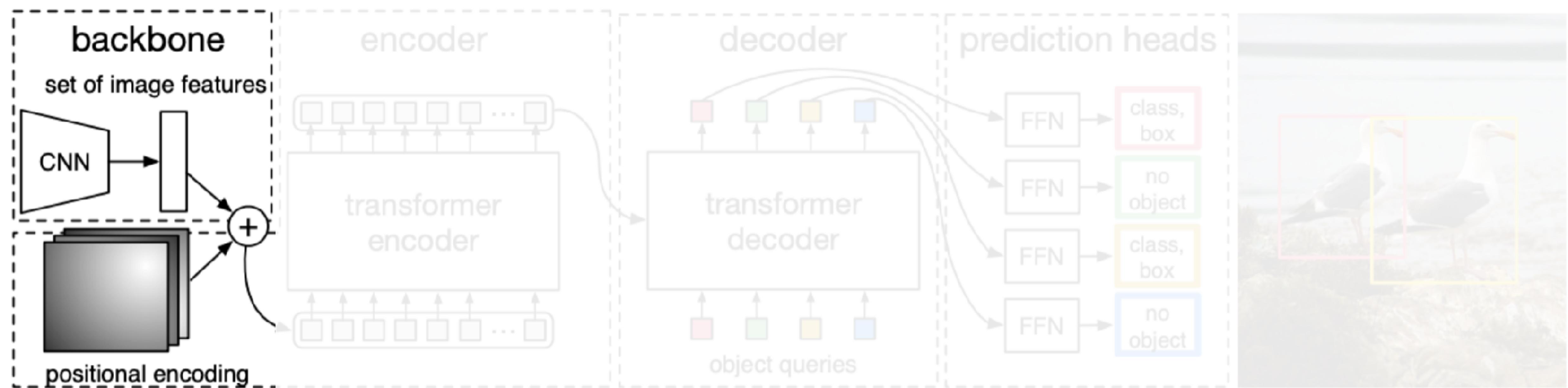
1-stage detectors



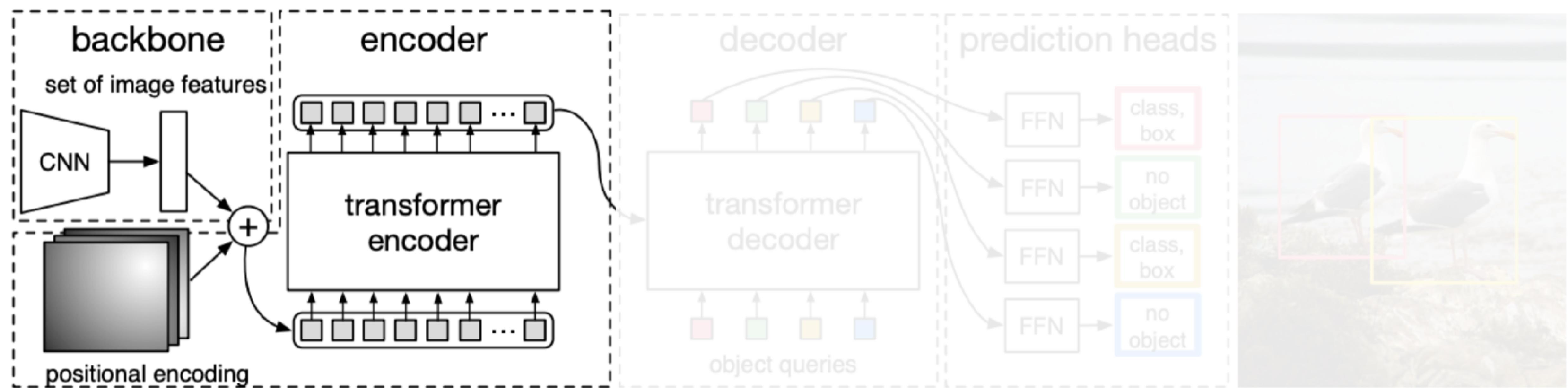
2-stage detectors

Object detection: Transformer-based methods

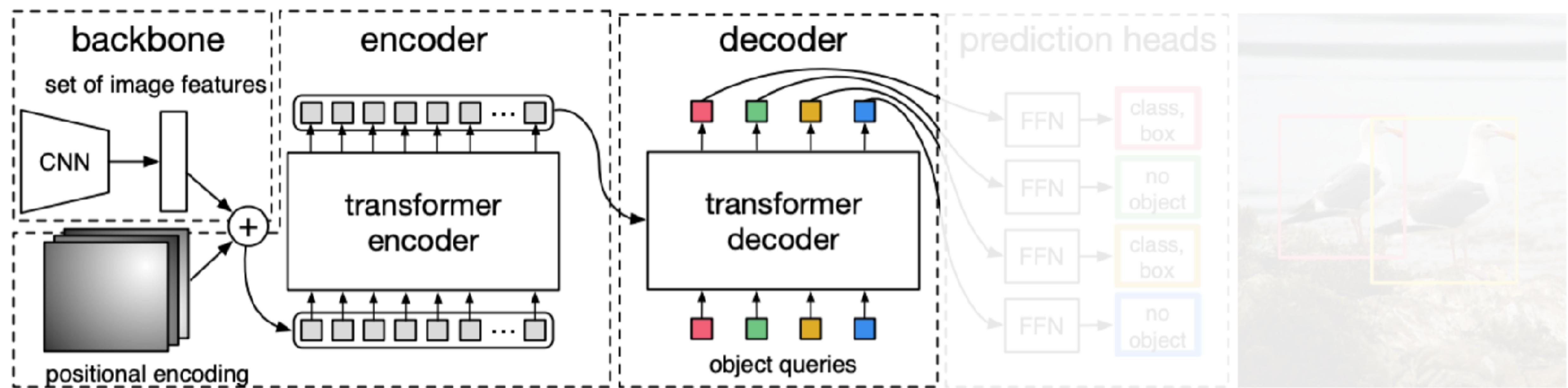
DETR: Object detection with transformers



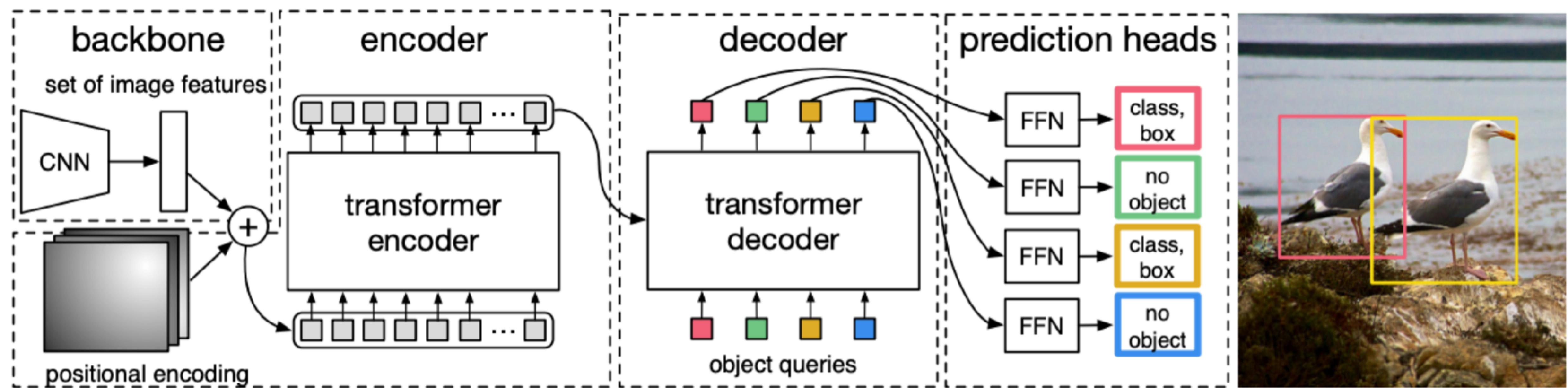
DETR: Object detection with transformers



DETR: Object detection with transformers

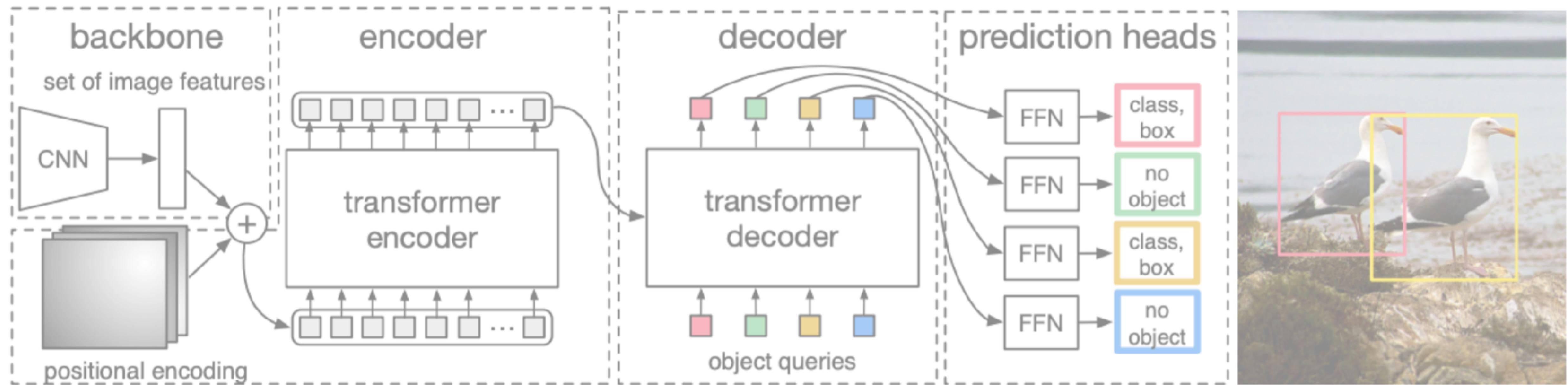
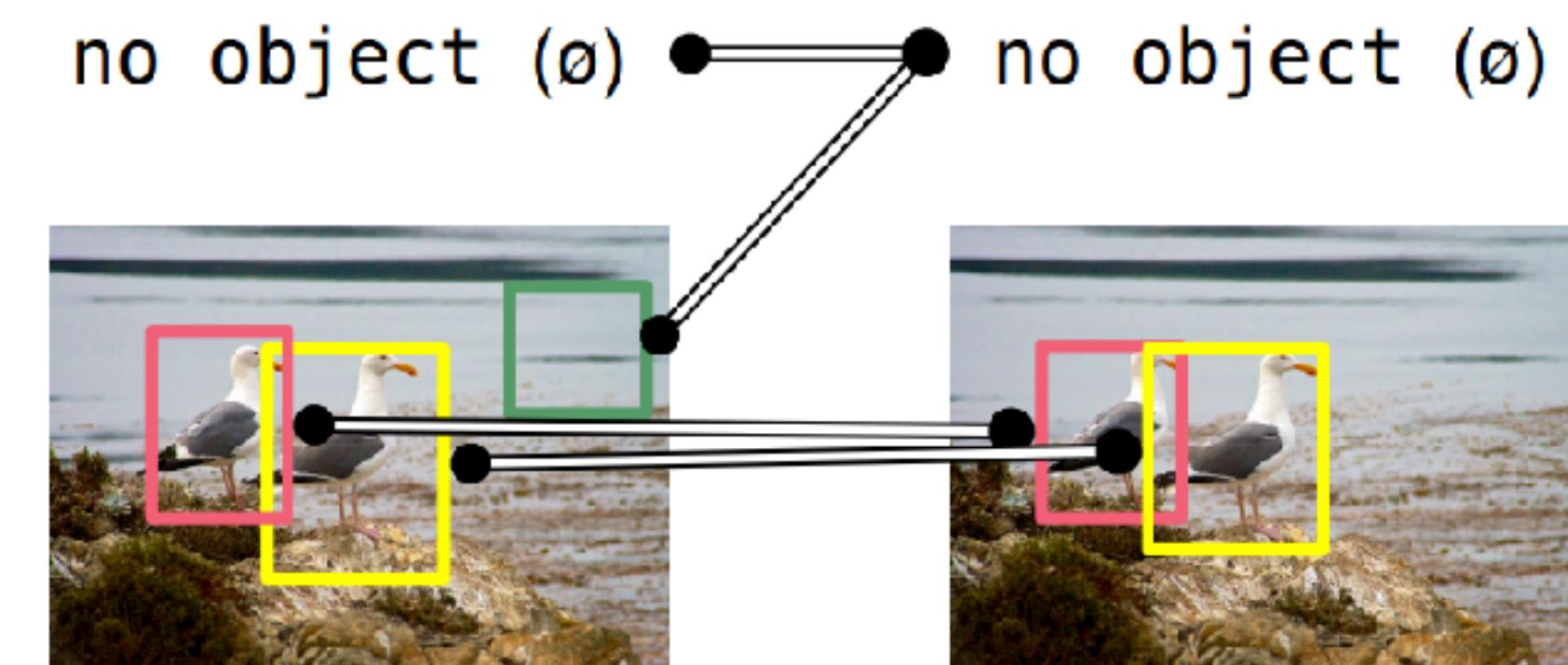


DETR: Object detection with transformers



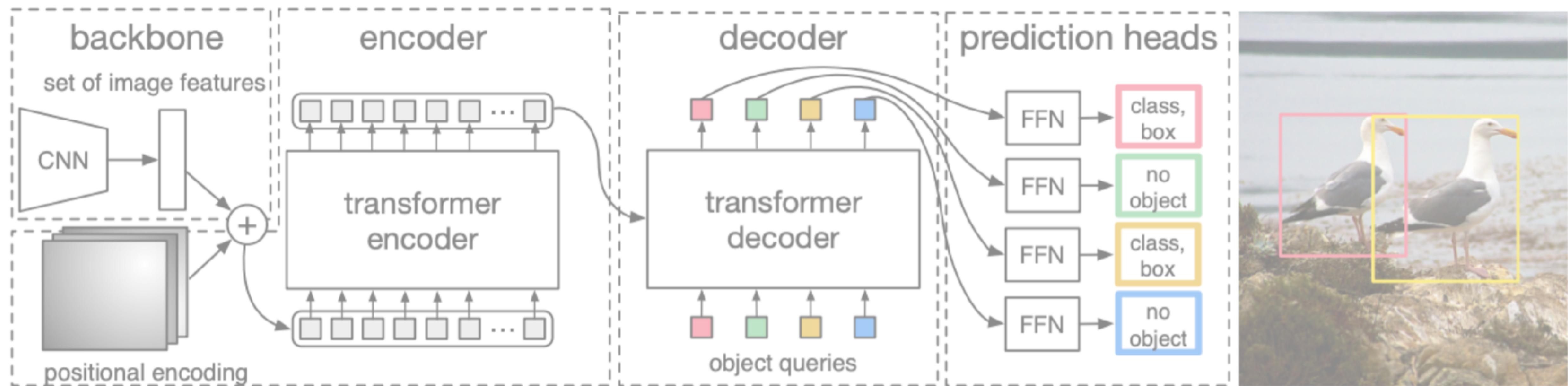
DETR: Training

- Match each box proposal to ground truth
- Use Hungarian algorithm to find permutation minimizing matching loss



DETR: Results COCO Val

Model	Epochs	mAP	mAP (small)	mAP (medium)	mAP (large)
Faster RCNN-FPN	109	42.0	26.6	45.4	53.4
DETR	500	42.0	20.5	45.8	61.1



2. Beyond classification

- a. Intro to structured outputs
- b. Object detection (localization)
- c. Segmentation
- d. Human pose estimation

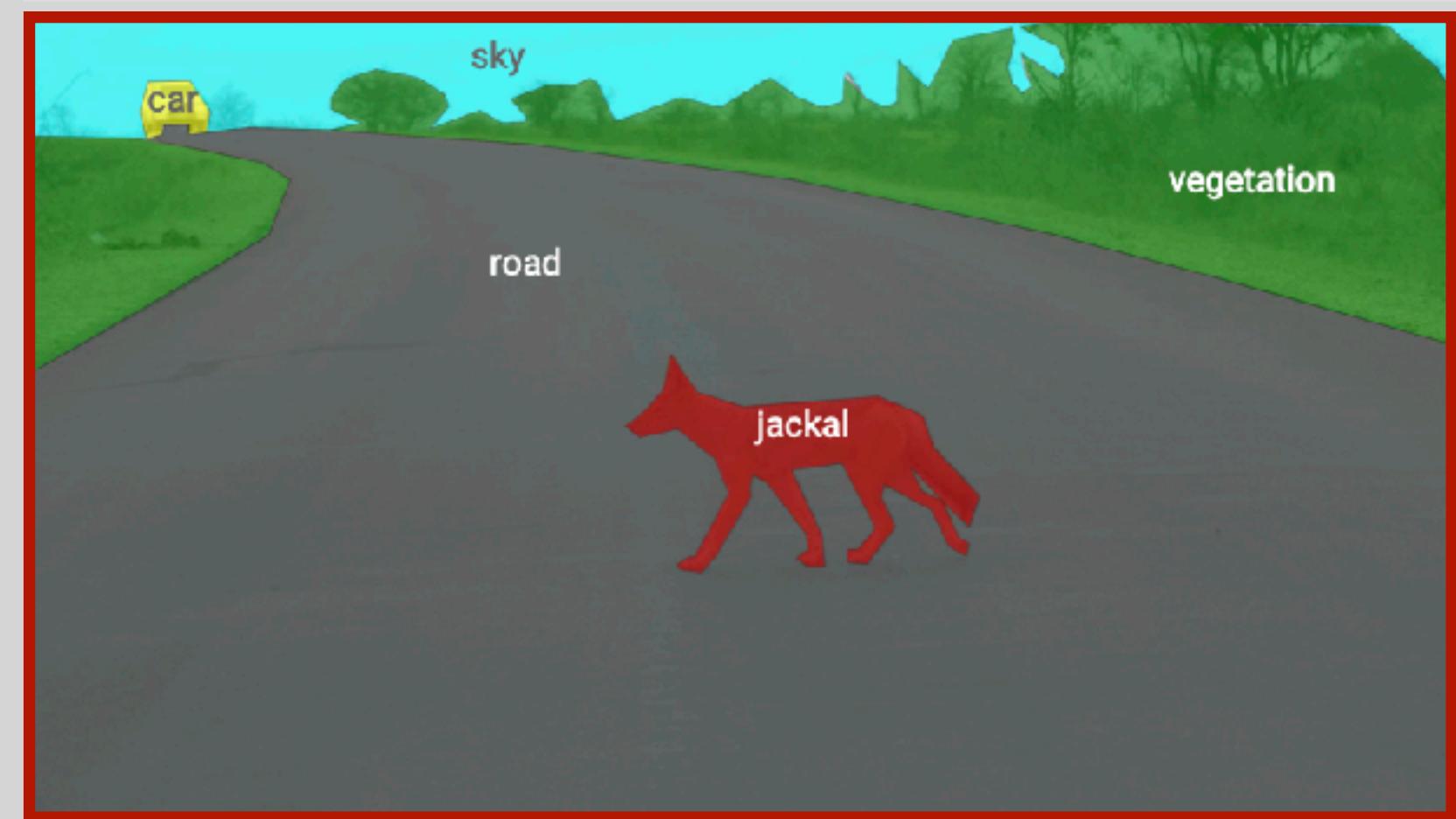
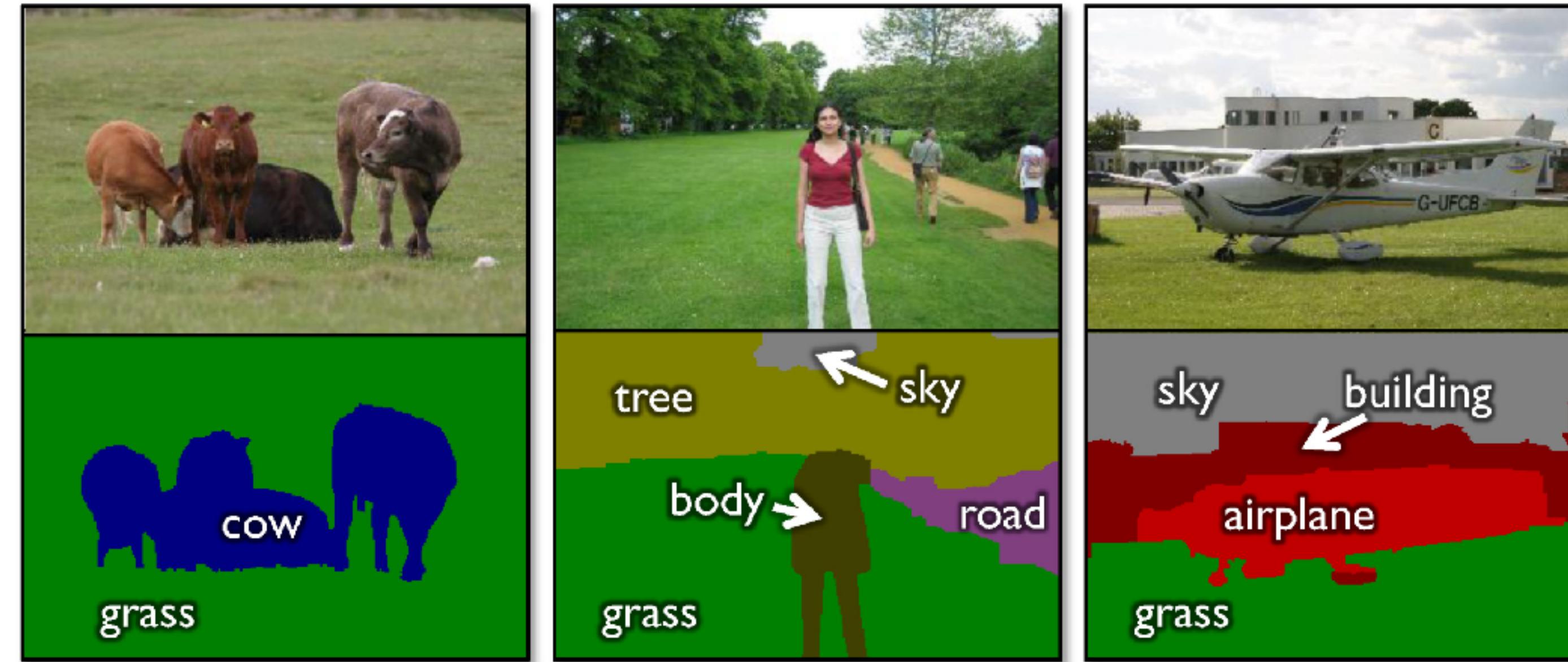


Image credits: Naila Murray

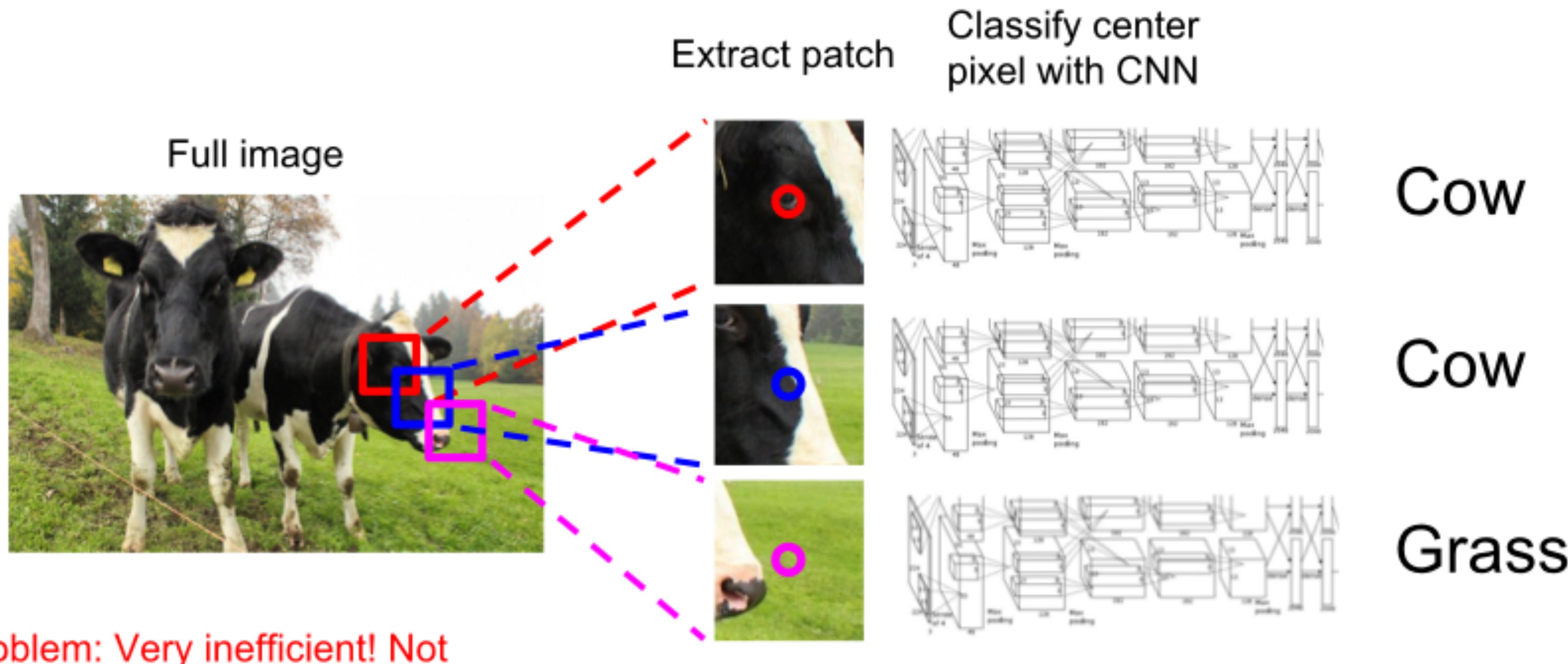
Semantic segmentation

- Label each pixel in the image with a category label
- Don't differentiate instances, only care about pixels



object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Semantic segmentation: sliding window



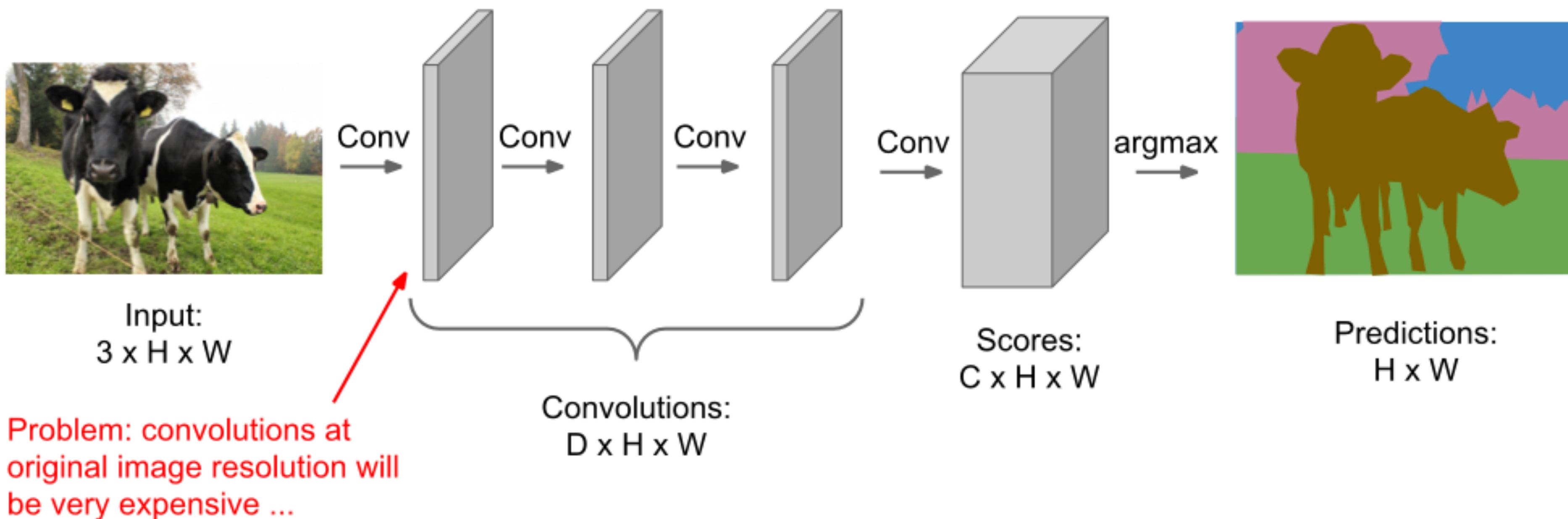
Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

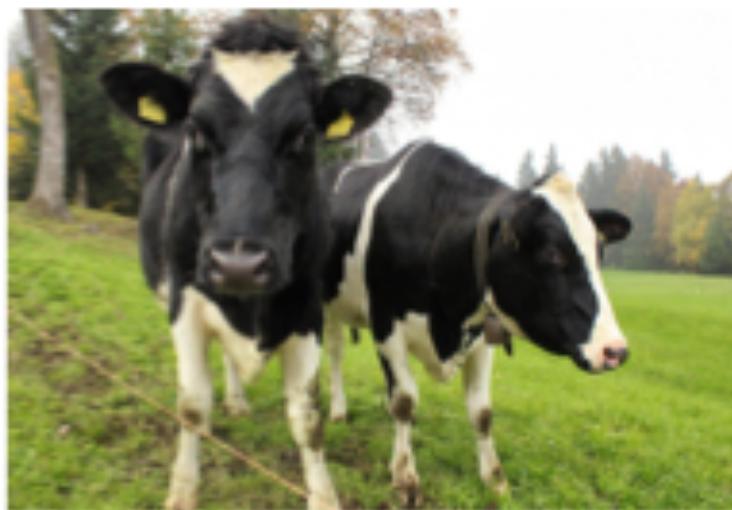
Semantic segmentation: fully convolutional

Design a network as a bunch of convolutional layers
to make predictions for pixels all at once!



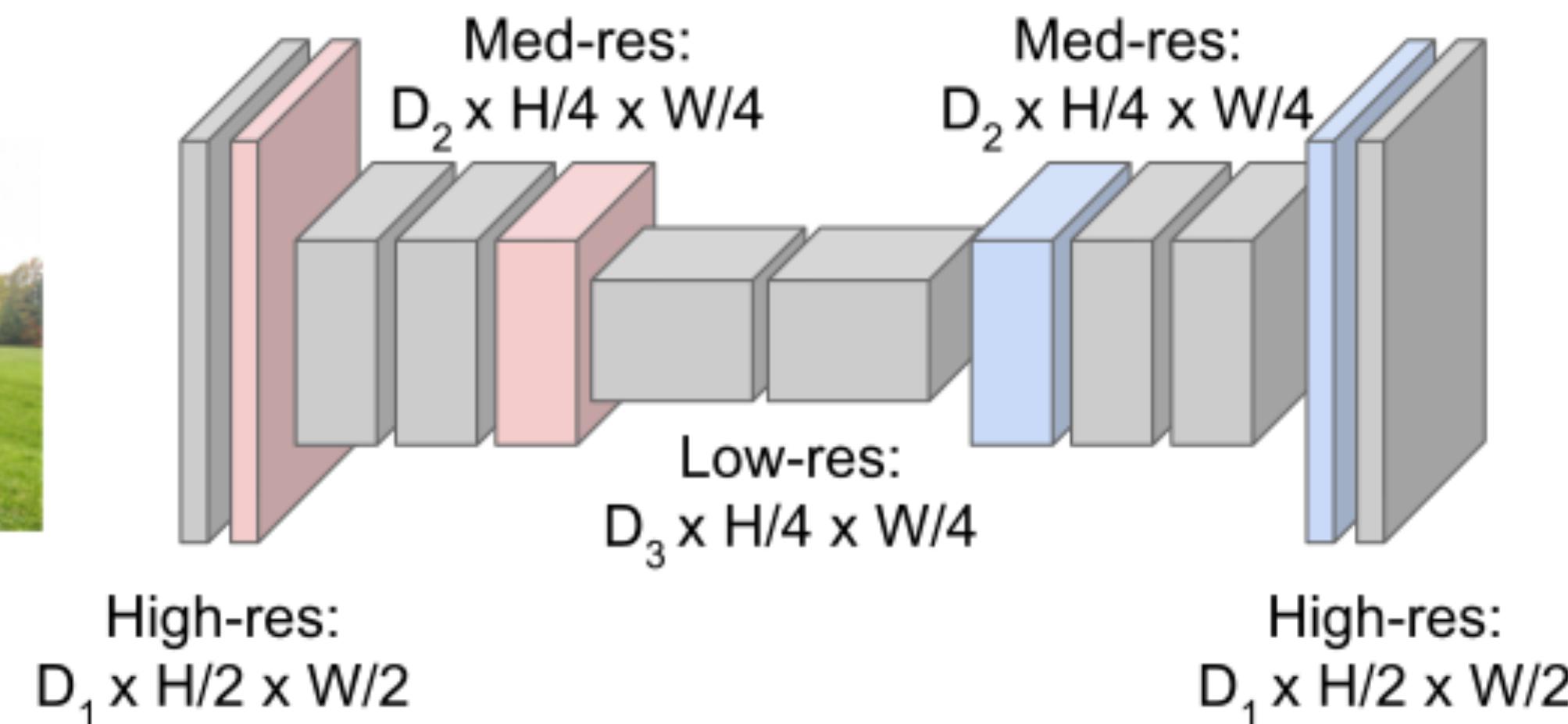
Semantic segmentation: fully convolutional

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



Upsampling:
???



Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al. "Learning Deconvolution Network for Semantic Segmentation". ICCV 2015

In-network upsampling: “Unpooling”

Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

Rest of the network

Max Unpooling

Use positions from pooling layer

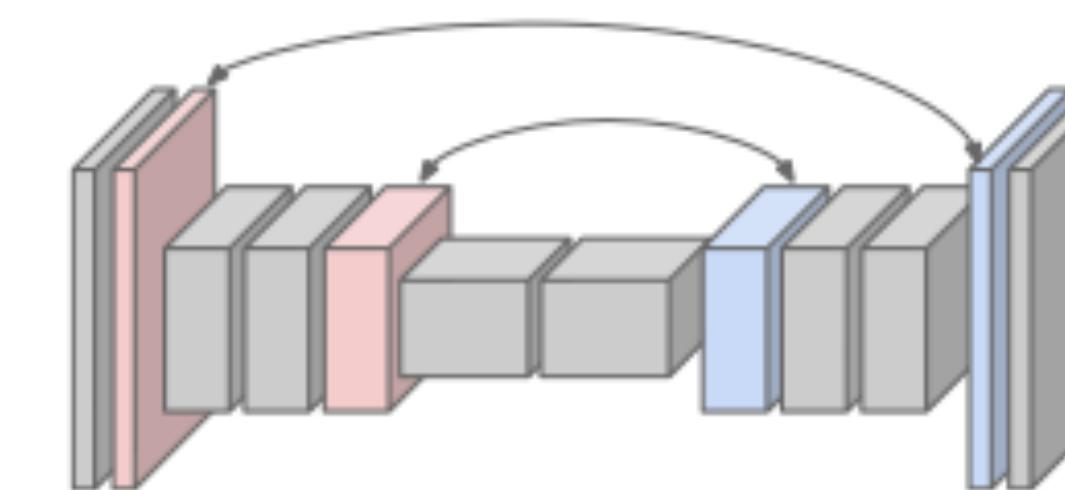
1	2
3	4

Input: 2 x 2

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

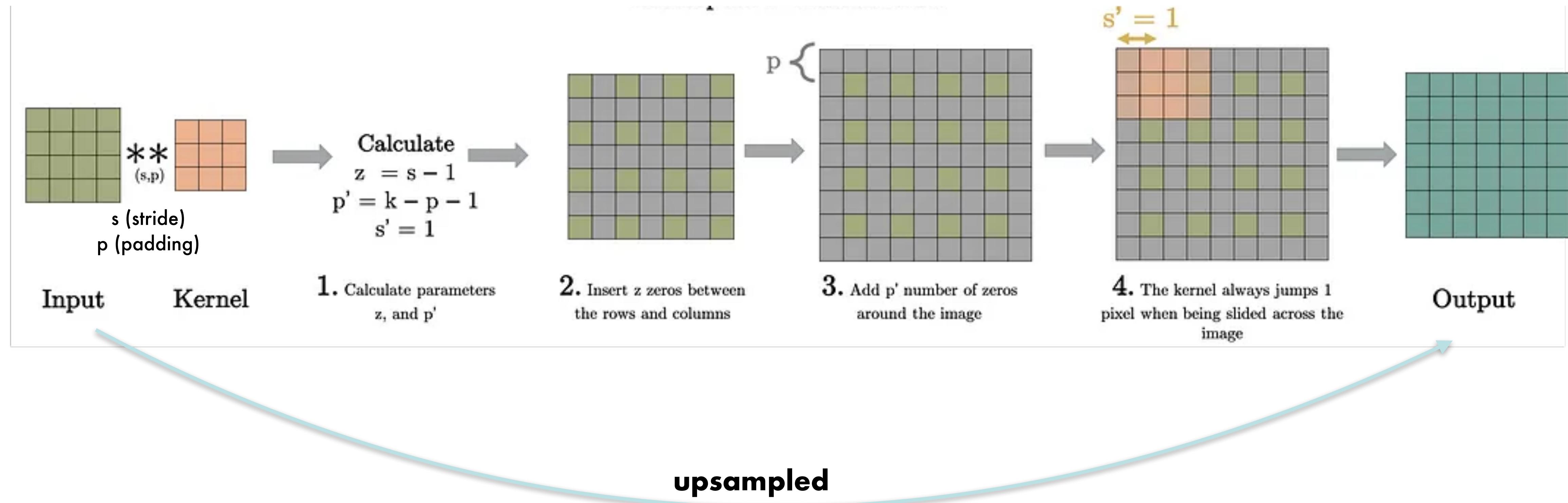
Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers



Learnable upsampling

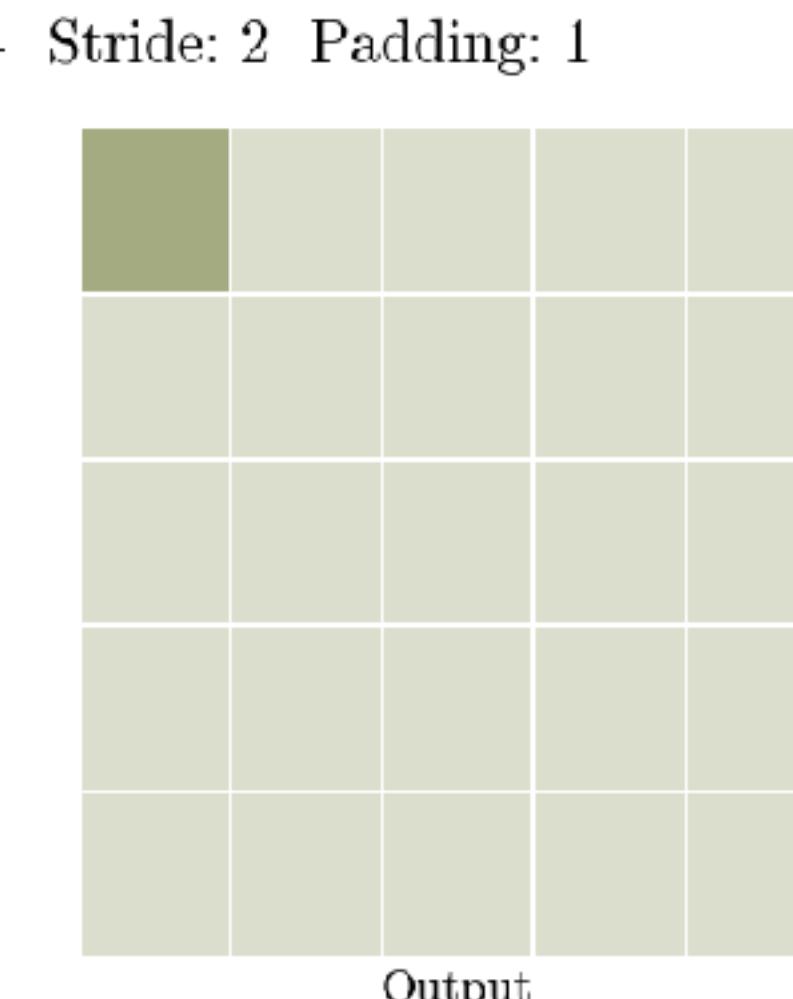
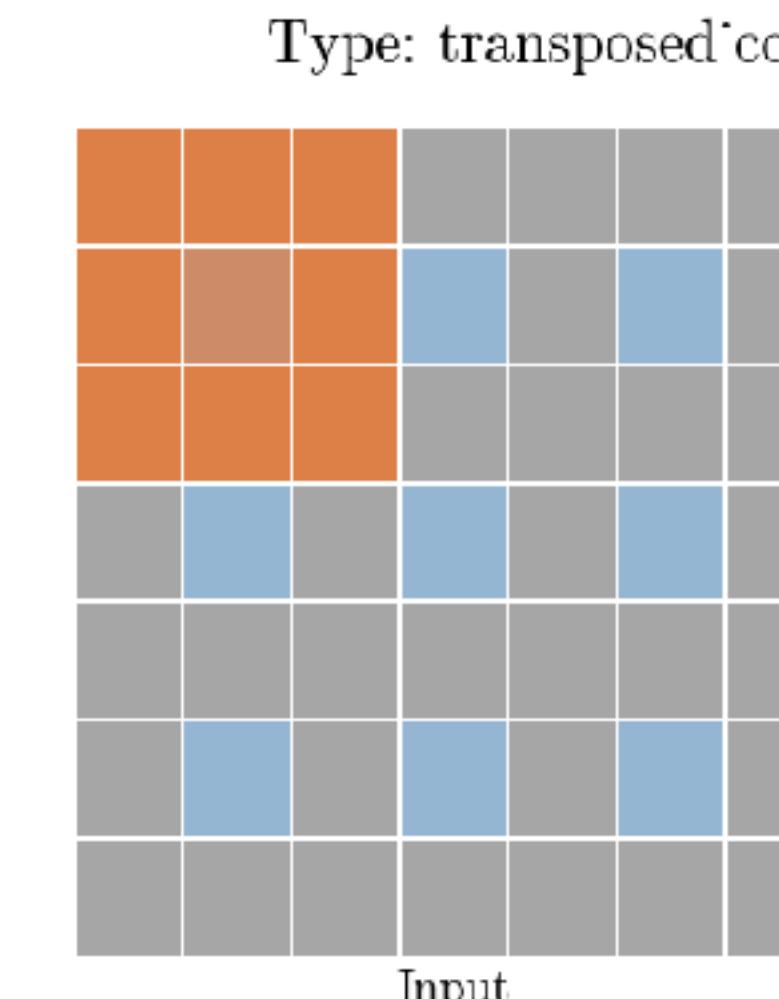
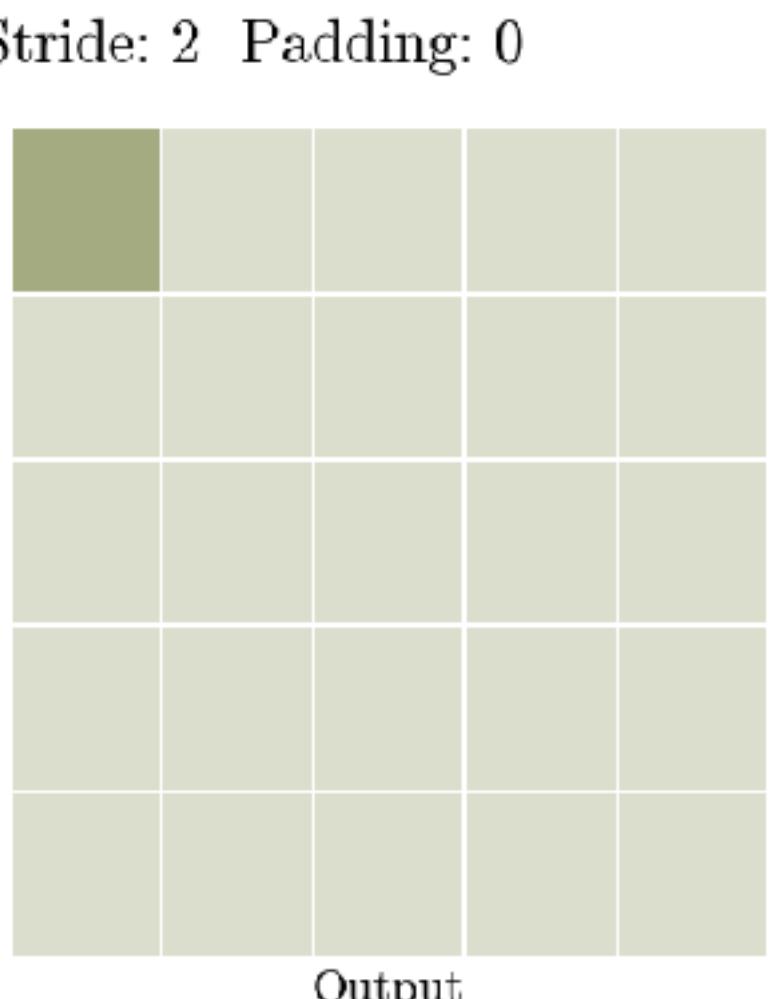
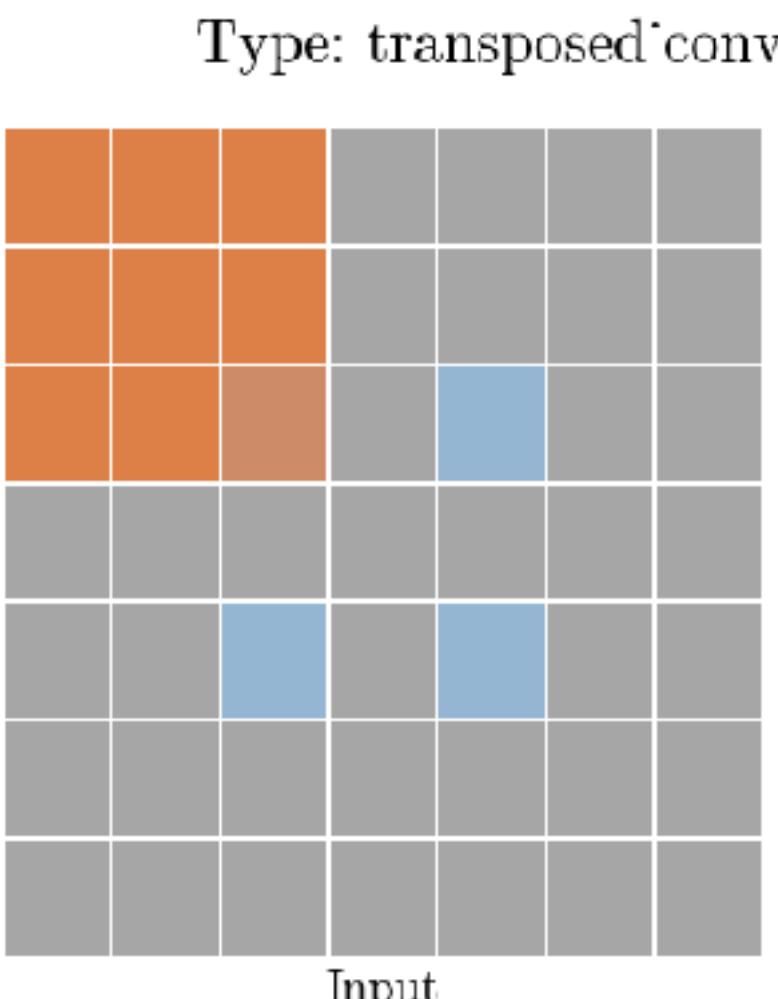
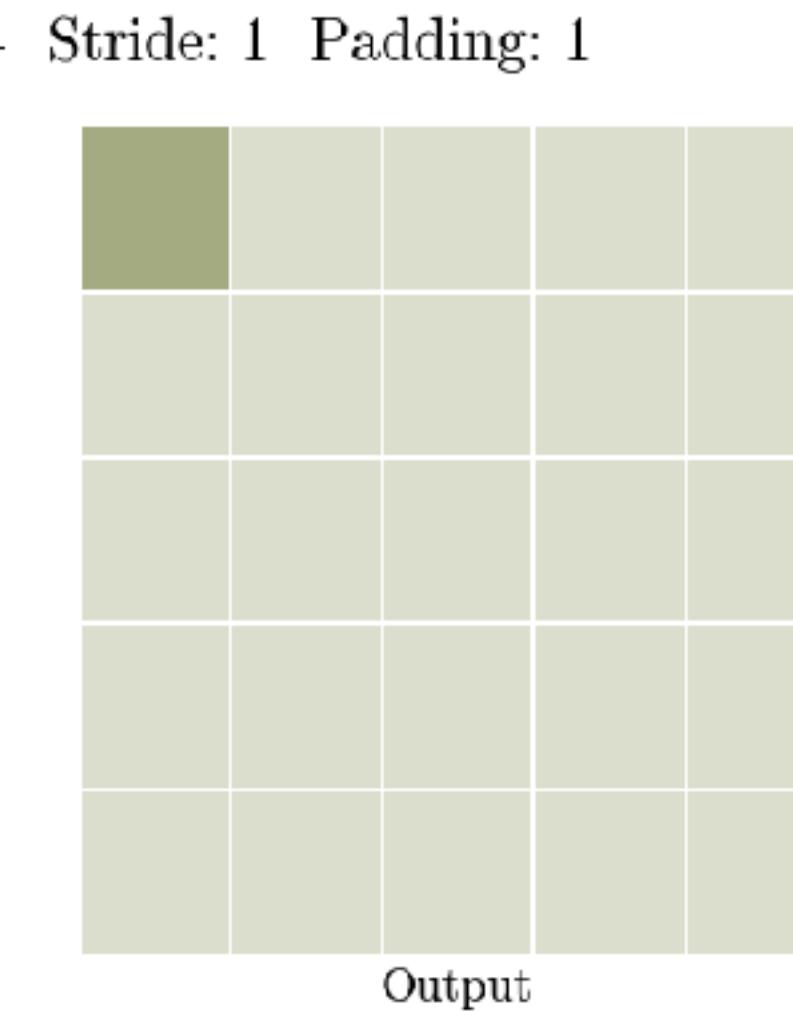
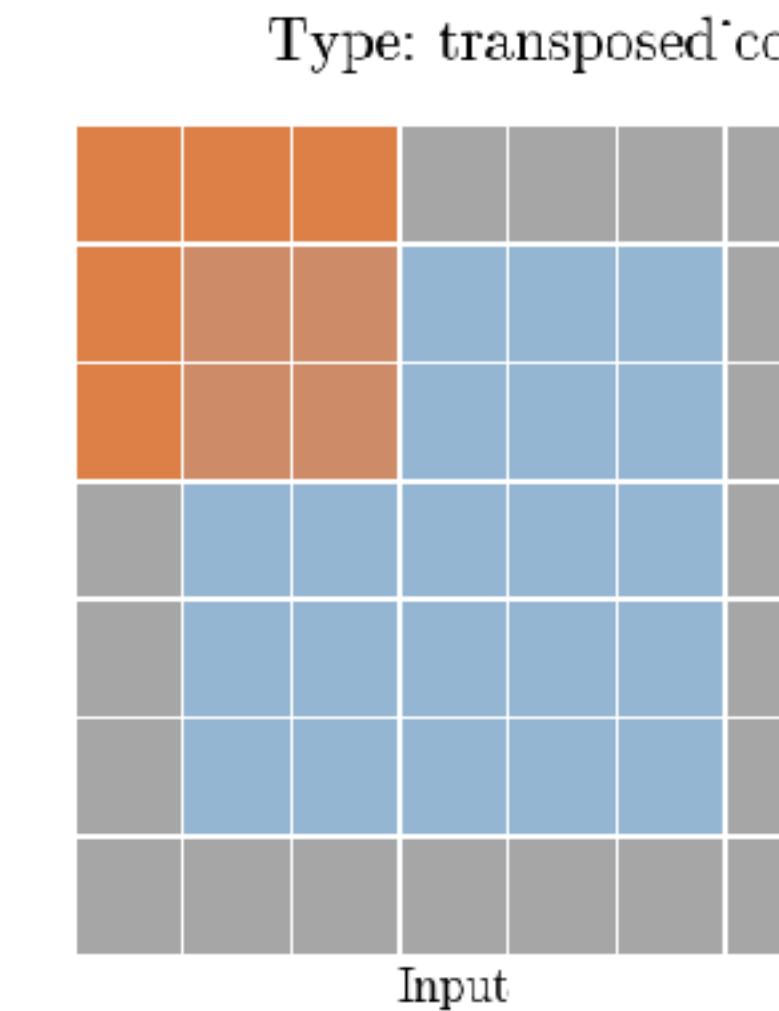
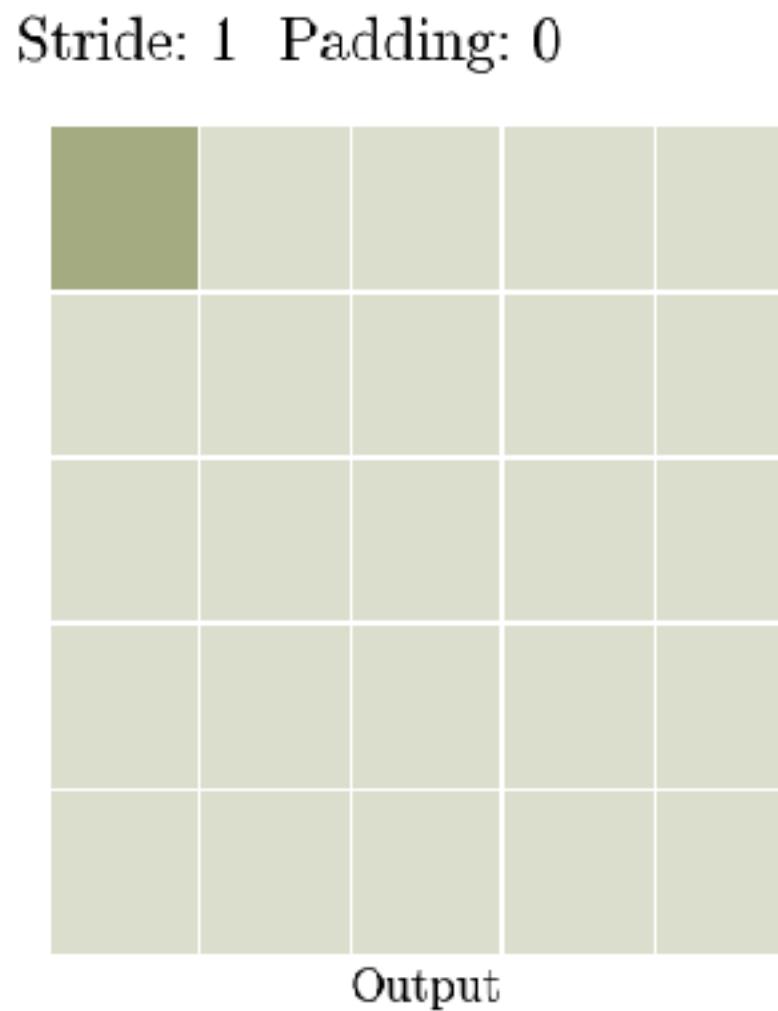
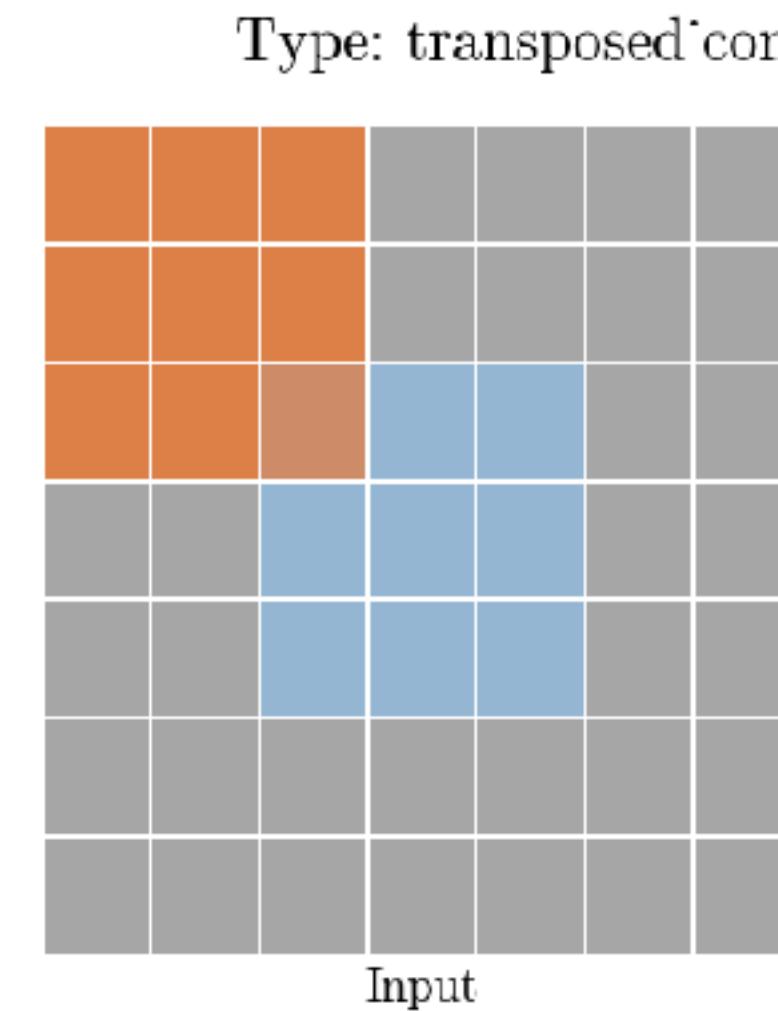
Other names: Deconvolution (bad); Upconvolution; Fractionally strided convolution; Backward strided convolution; Transpose convolution



Does exactly what a standard convolutional layer does but on a modified input feature map

Learnable upsampling

Other names: Deconvolution (bad); Upconvolution; Fractionally strided convolution; Backward strided convolution; Transpose convolution



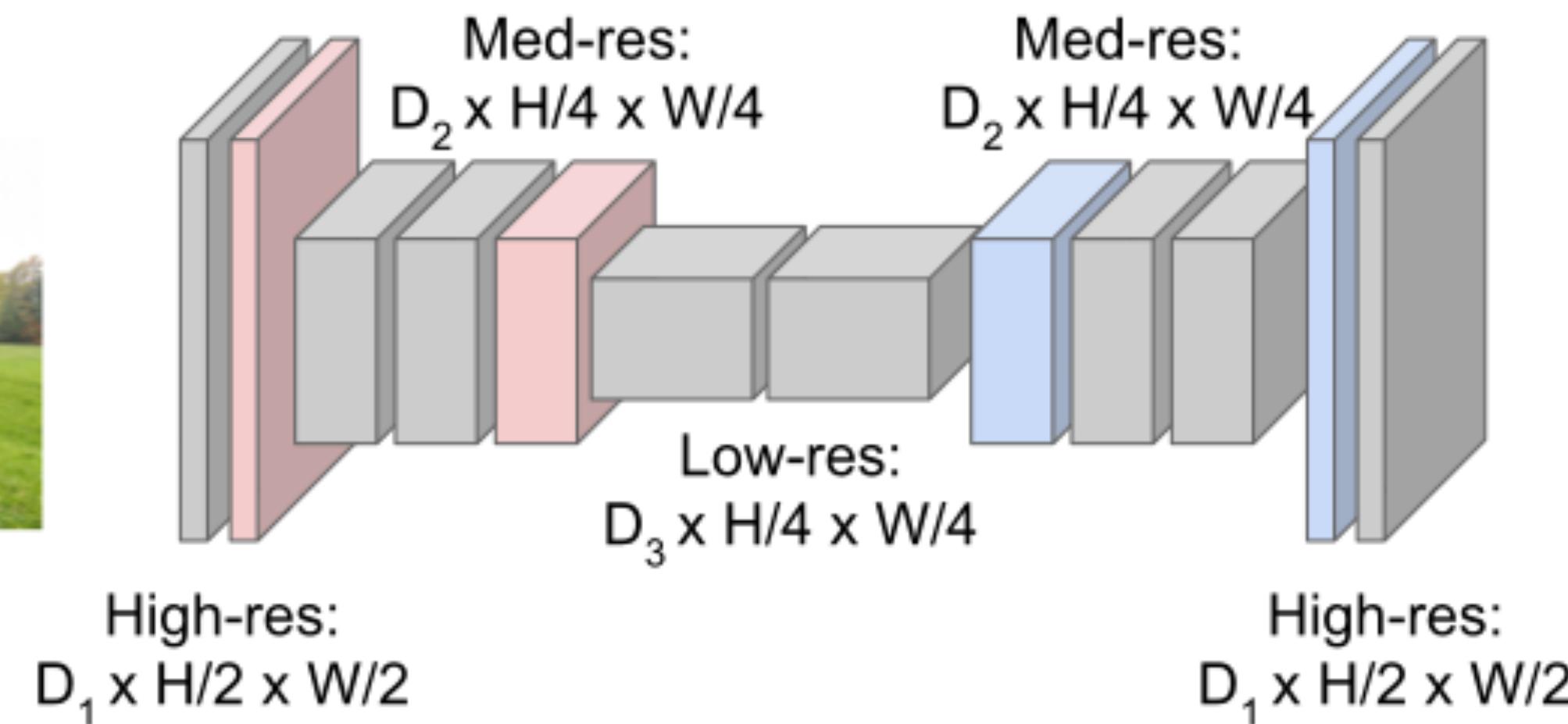
Semantic segmentation: fully convolutional

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



Upsampling:
Unpooling or strided
transpose convolution

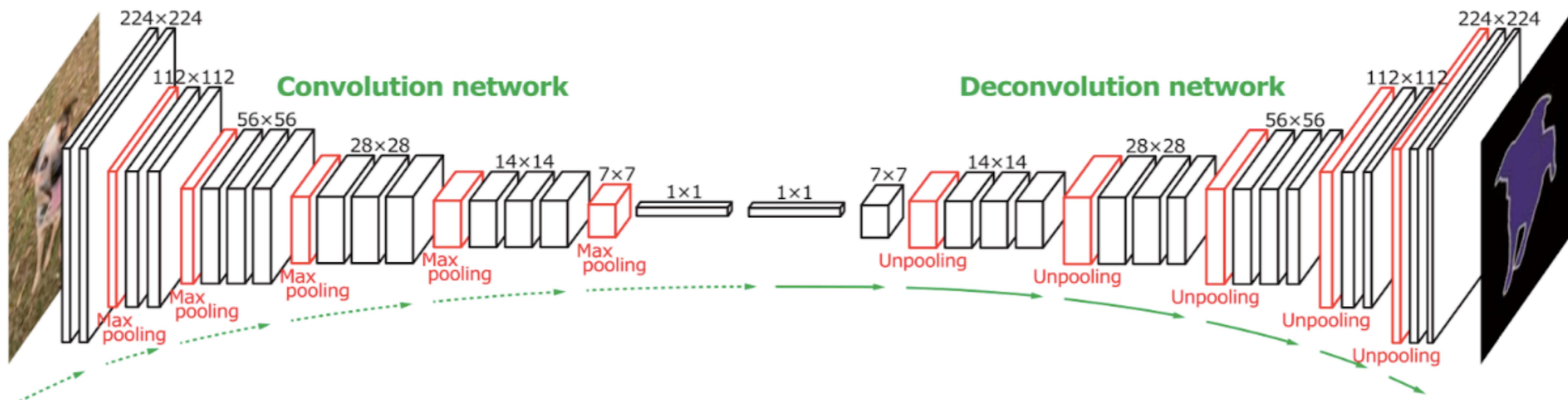


Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al. "Learning Deconvolution Network for Semantic Segmentation". ICCV 2015

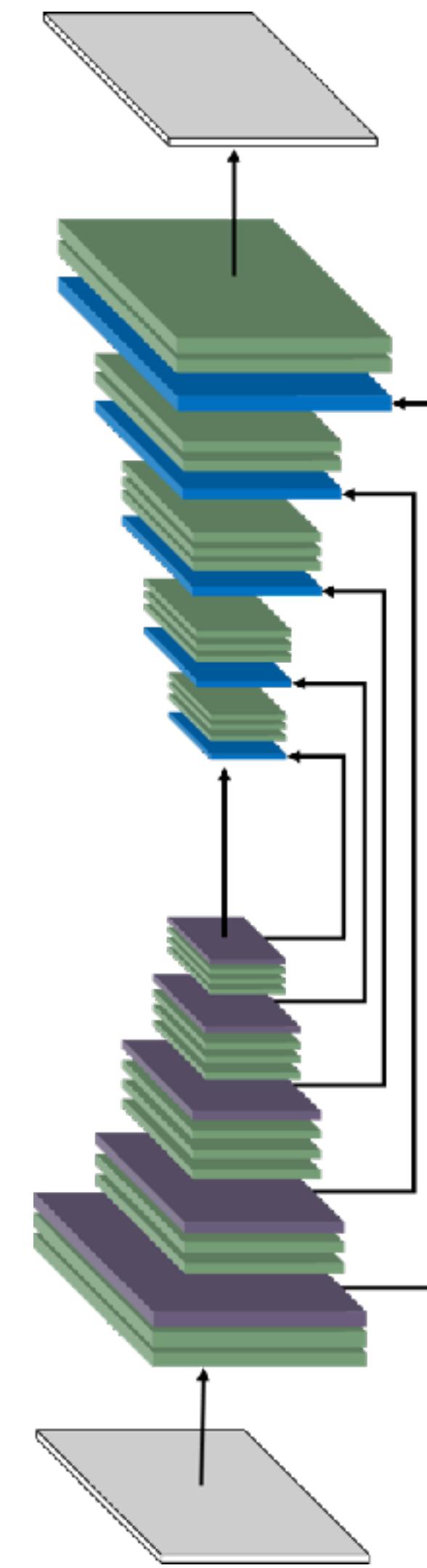
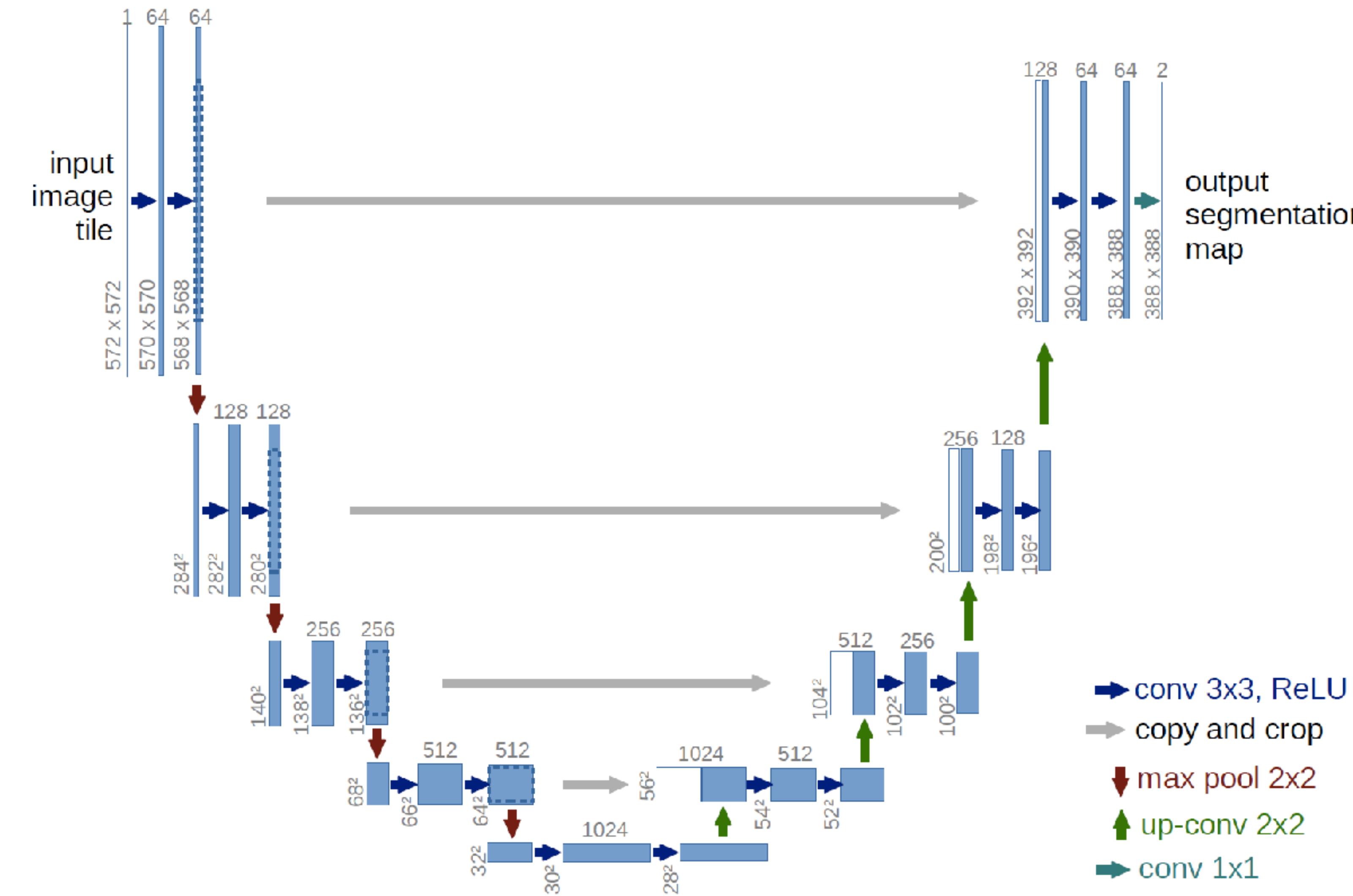
Semantic segmentation: Auto-encoder



Why is this a bad idea for segmentation?

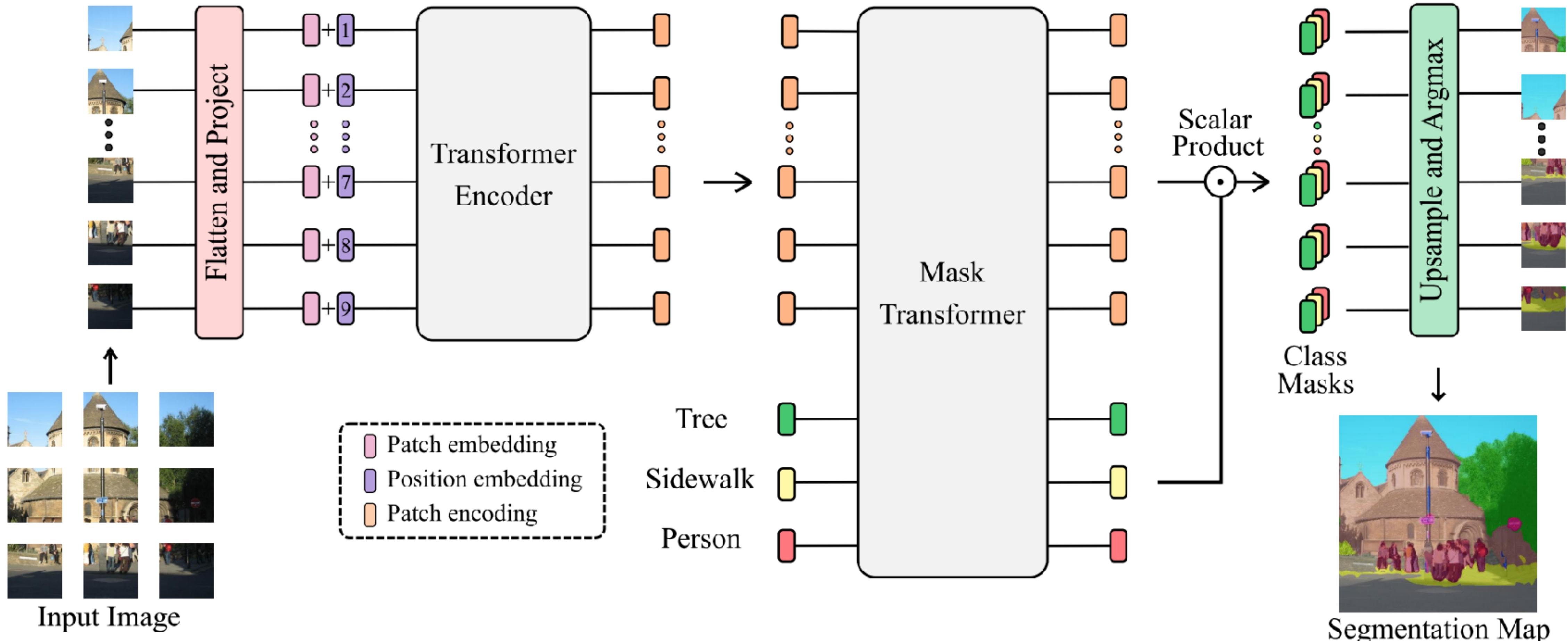
Semantic segmentation: U-Net

or “Hourglass”



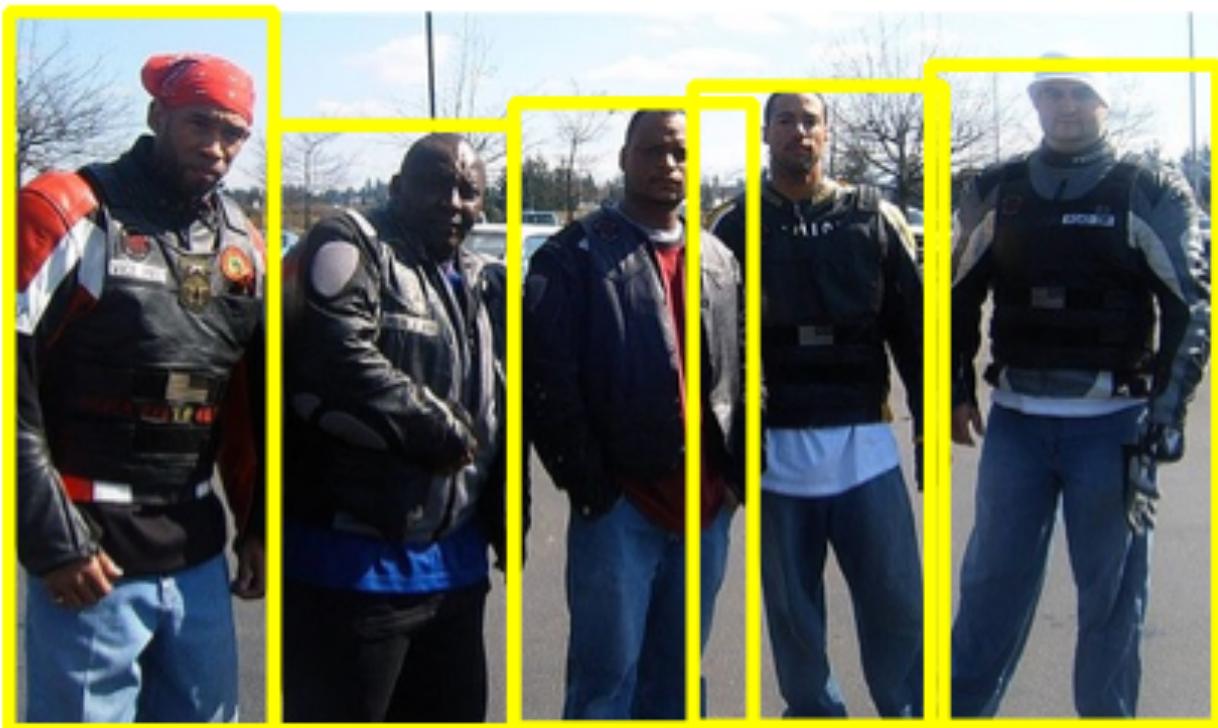
Semantic segmentation: Segmenter

Transformer architecture for image segmentation



(Object) Instance segmentation

- Differentiate instances
- Object detection + segmentation



Object Detection



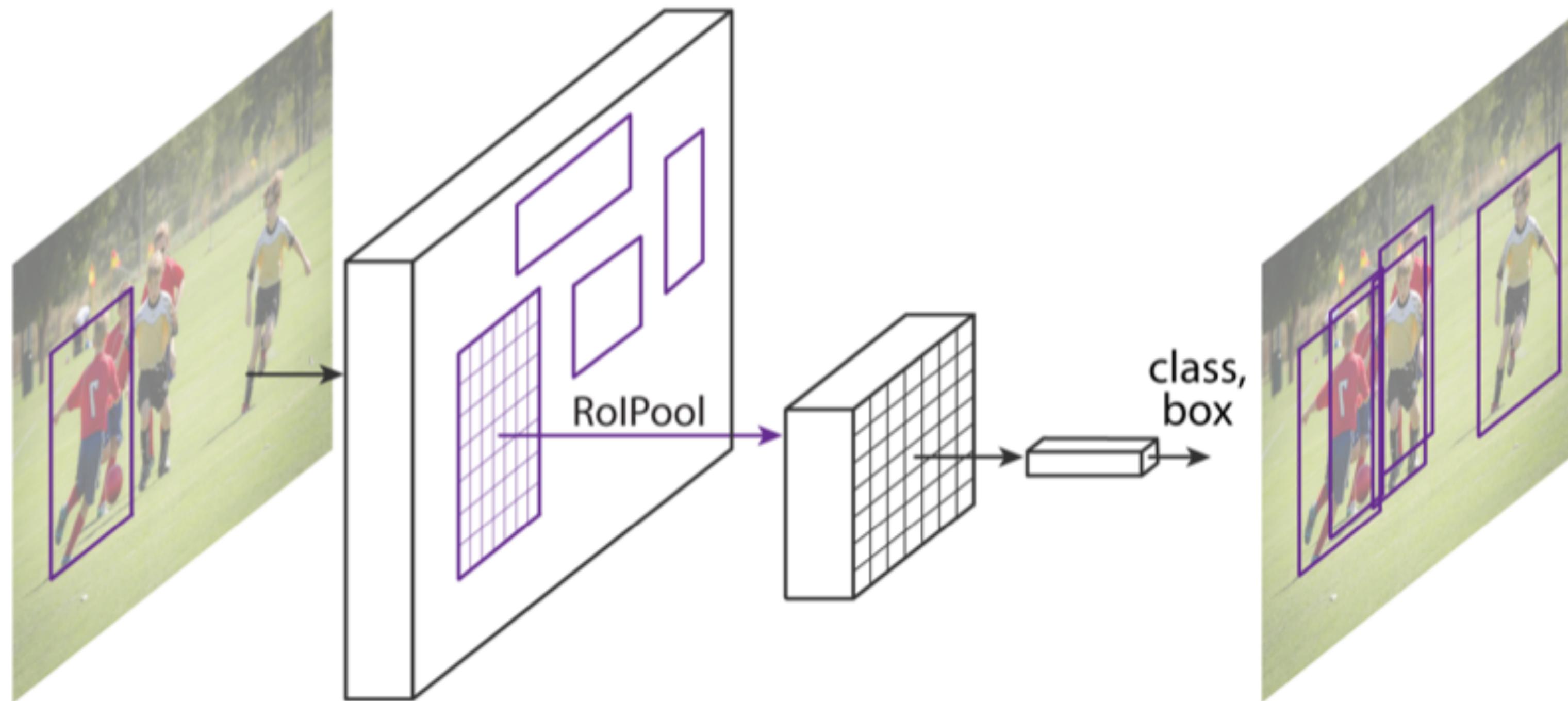
Semantic Segmentation



Instance Segmentation

Remember:

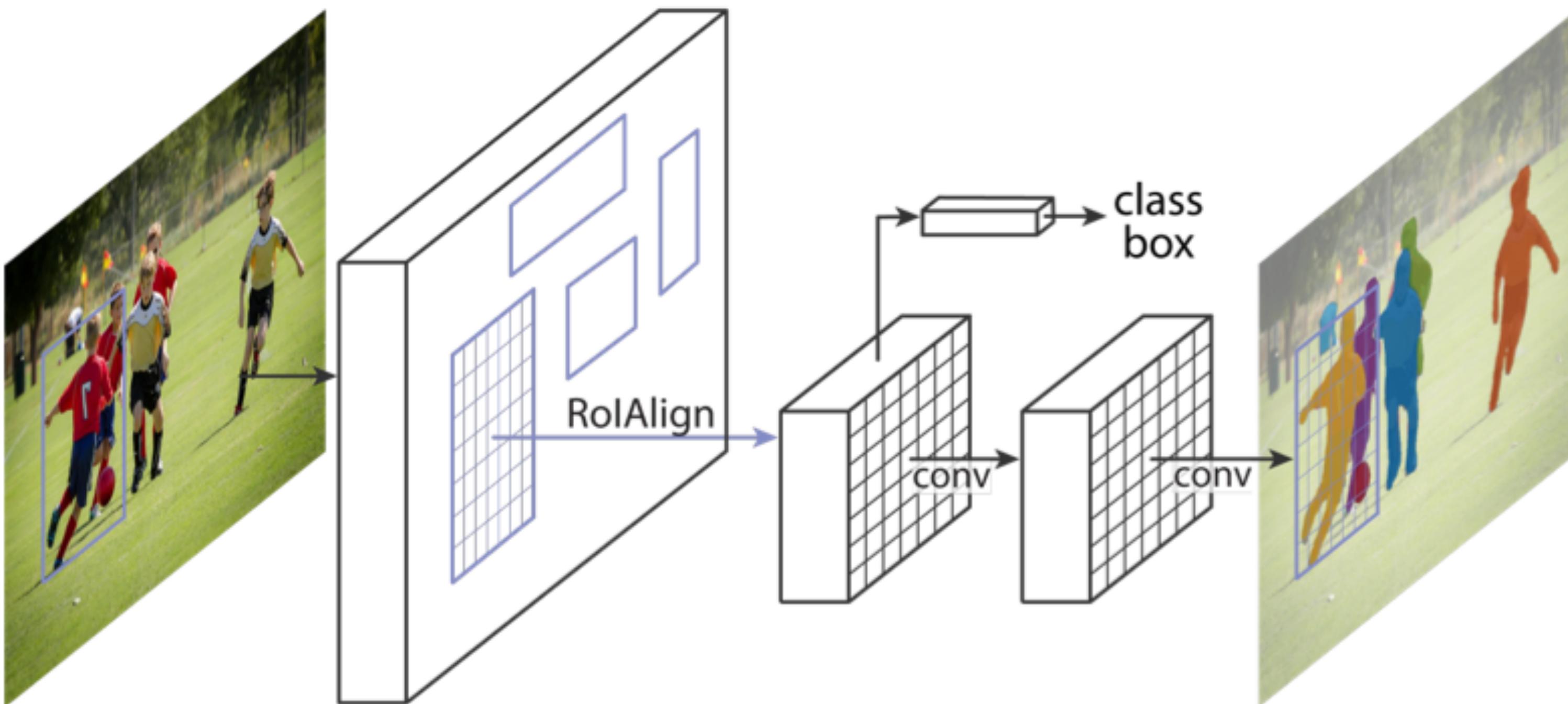
Object detection using Fast(er) R-CNN



Ross Girshick. “Fast R-CNN”. ICCV 2015. Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun.
“Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks”. NIPS 2015.

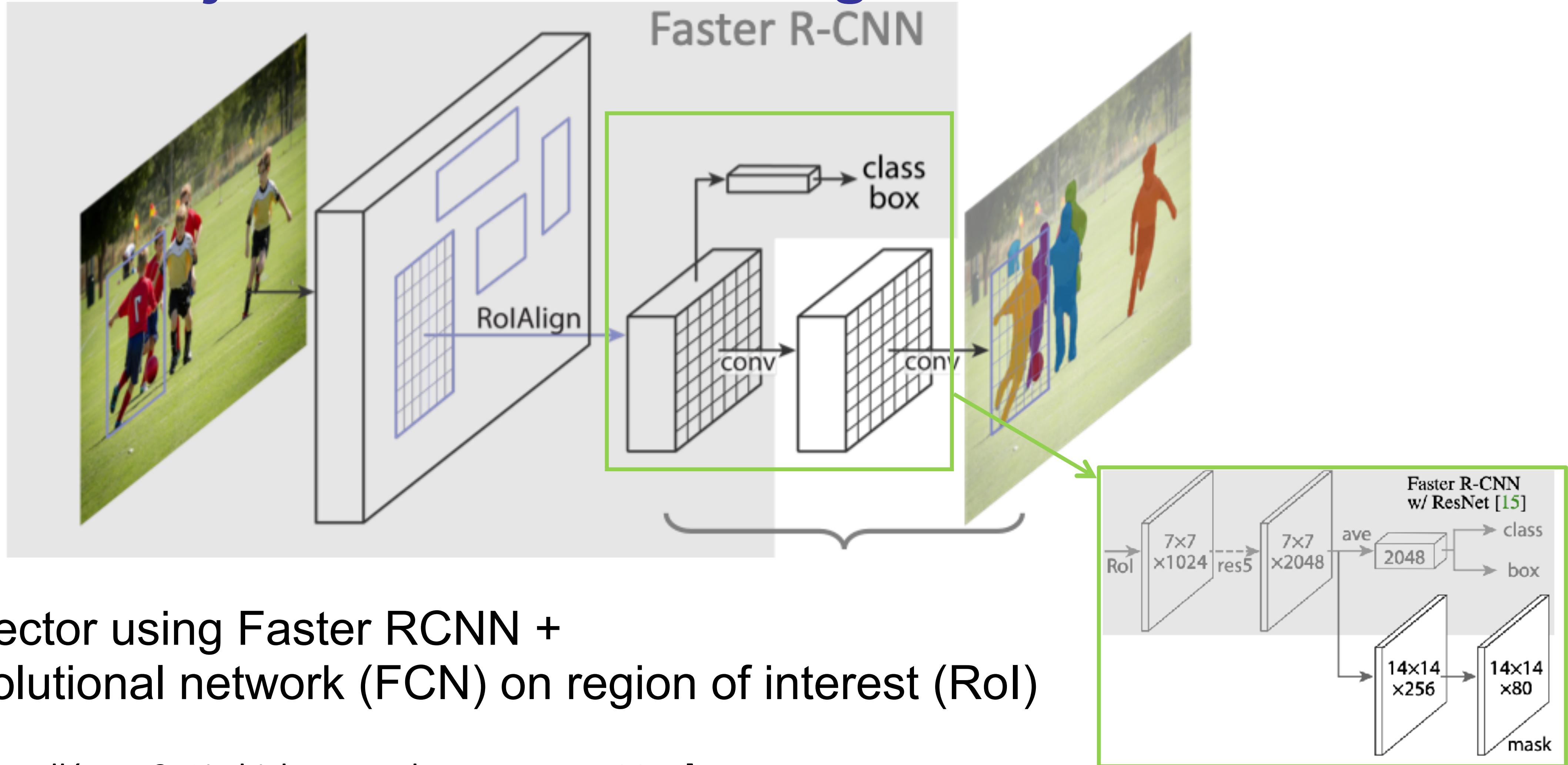
Mask R-CNN

Object detection and segmentation



Mask R-CNN

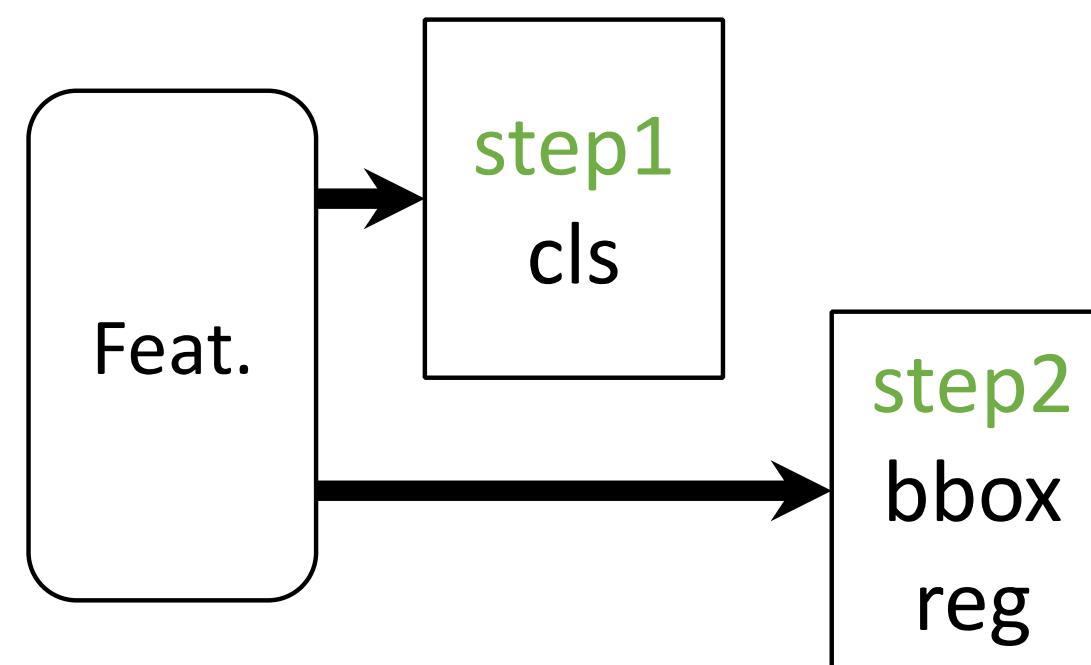
Object detection and segmentation



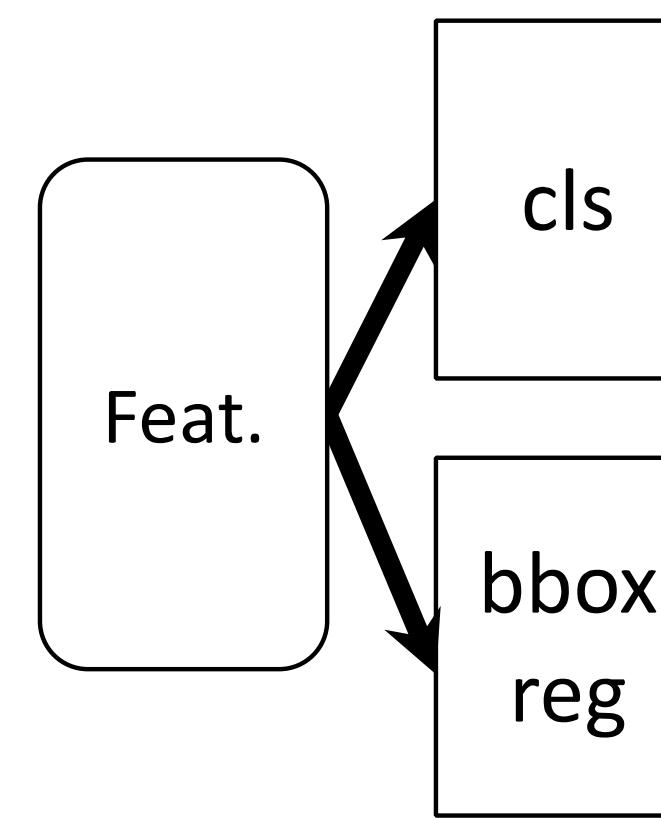
Mask R-CNN

Combining loss functions

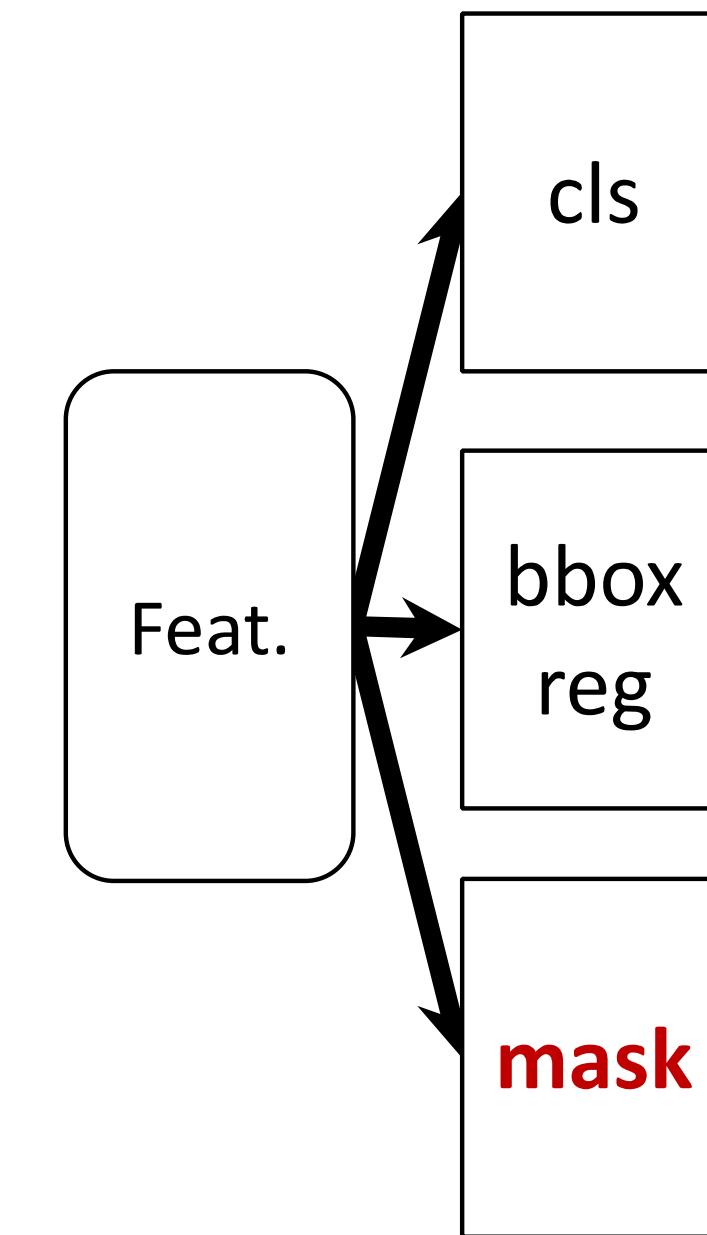
- Easy, fast to implement and train



(slow) R-CNN



Fast/er R-CNN



Mask R-CNN

Mask R-CNN

Example results

object
surrounded by
same-category
objects



Mask R-CNN results on COCO

Mask R-CNN

Example results

disconnected
object

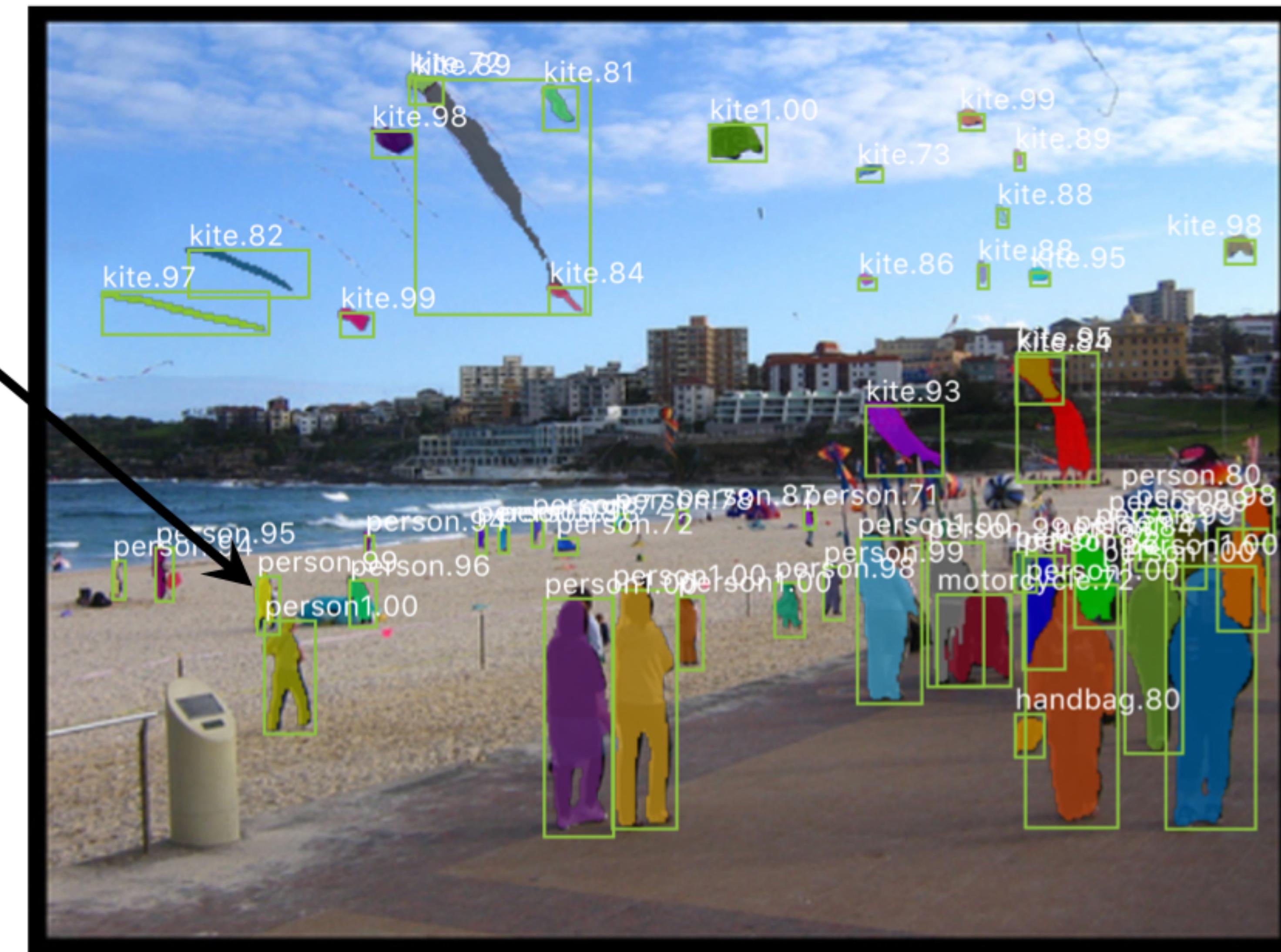


Mask R-CNN results on COCO

Mask R-CNN

Example results

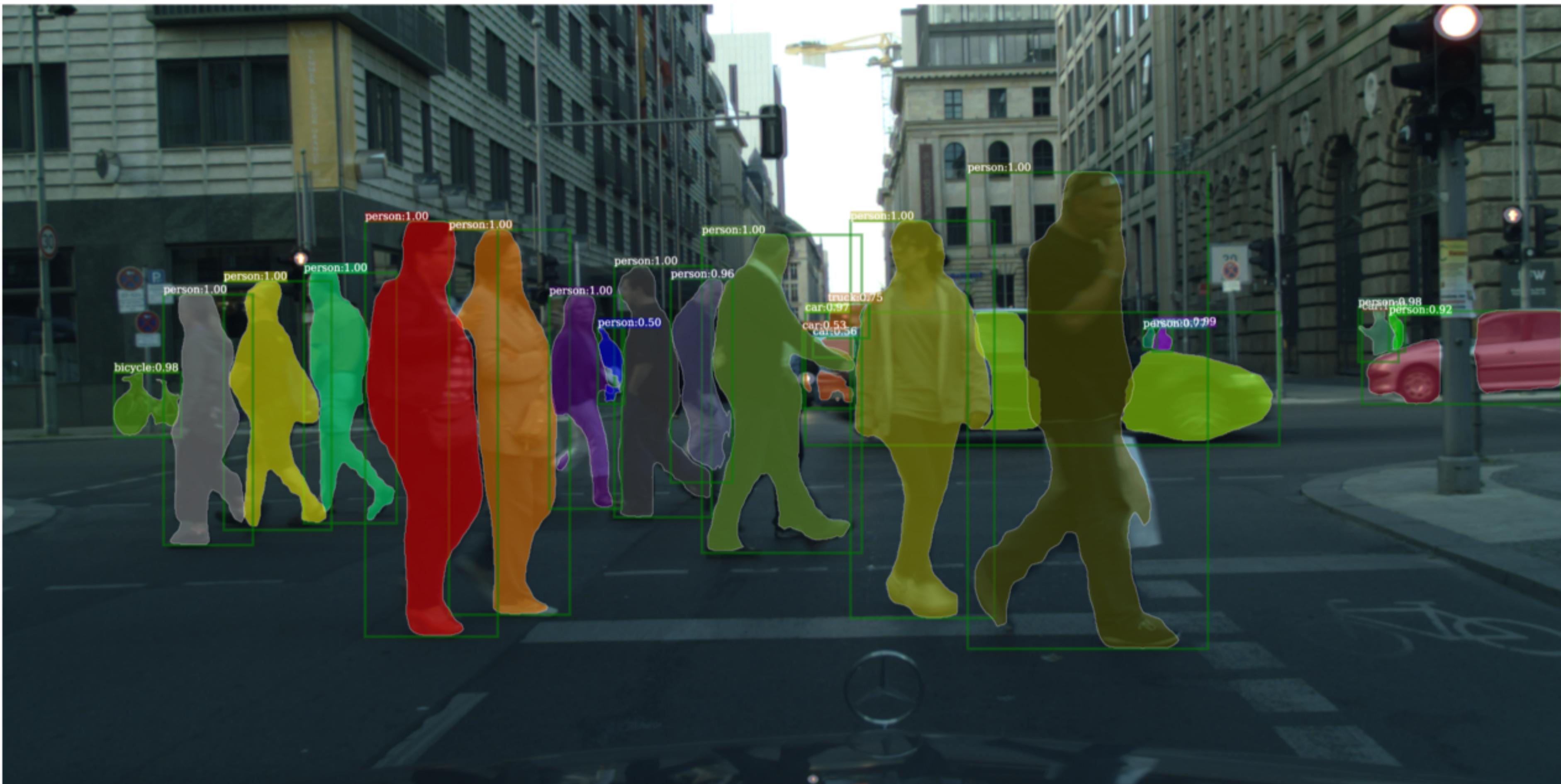
small
objects



Mask R-CNN results on COCO

Mask R-CNN

Example results

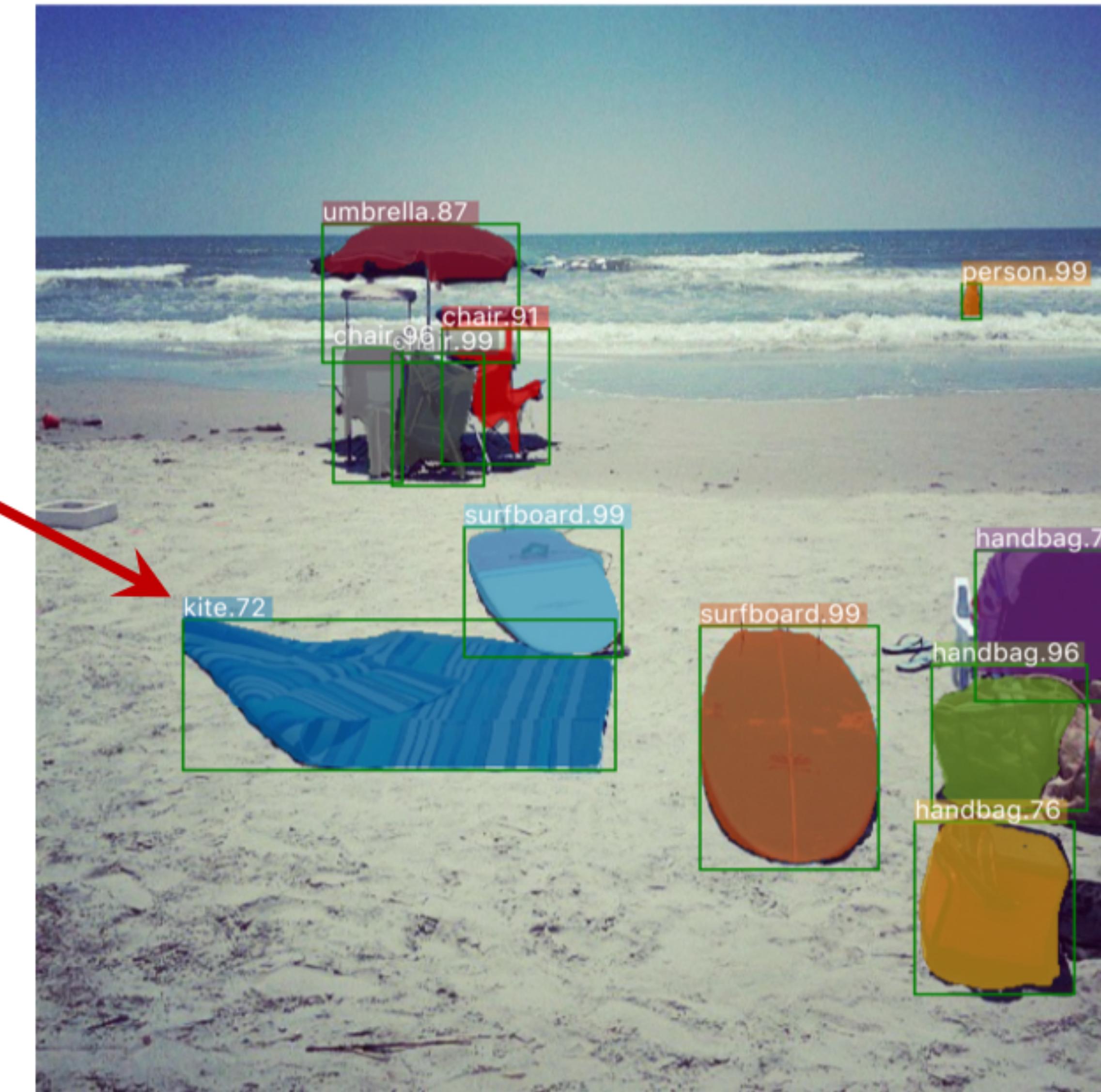


Mask R-CNN results on CityScapes

Mask R-CNN

Example failures: recognition

not a kite



Mask R-CNN results on COCO

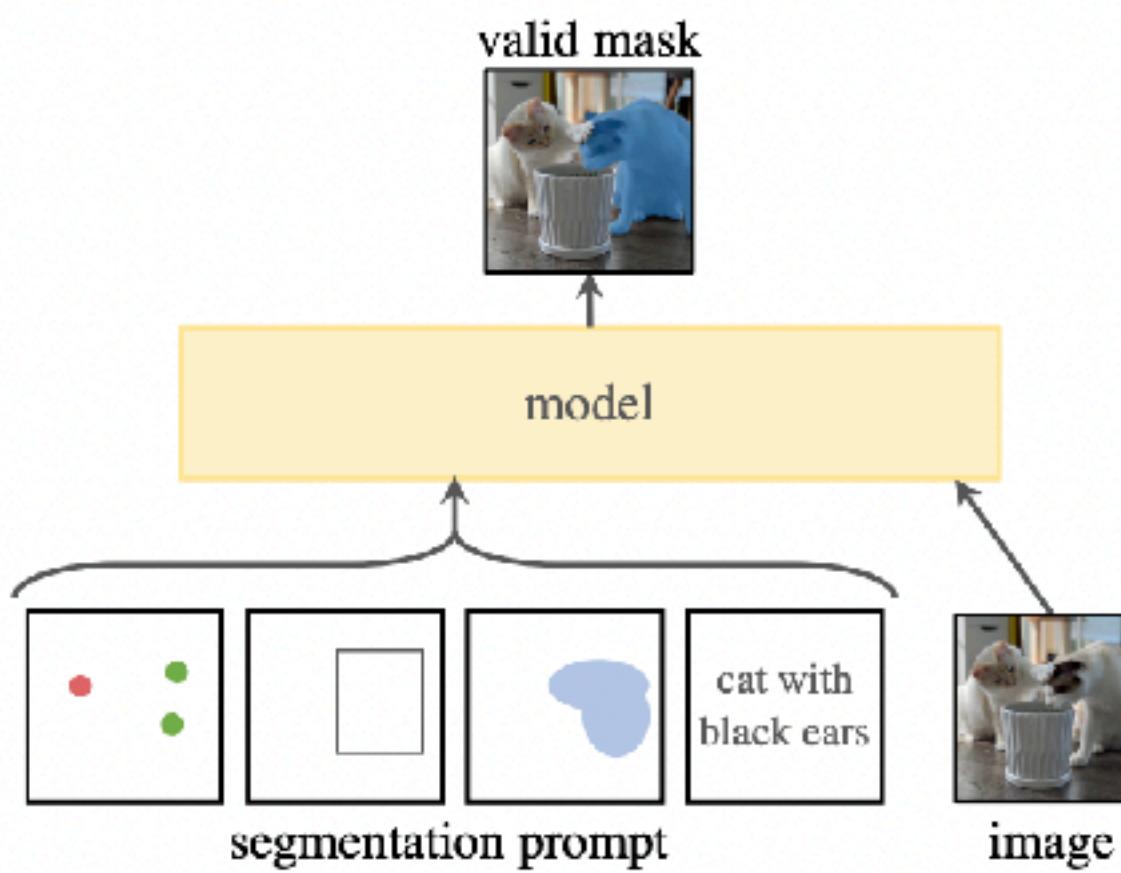
Promptable segmentation

Segment Anything

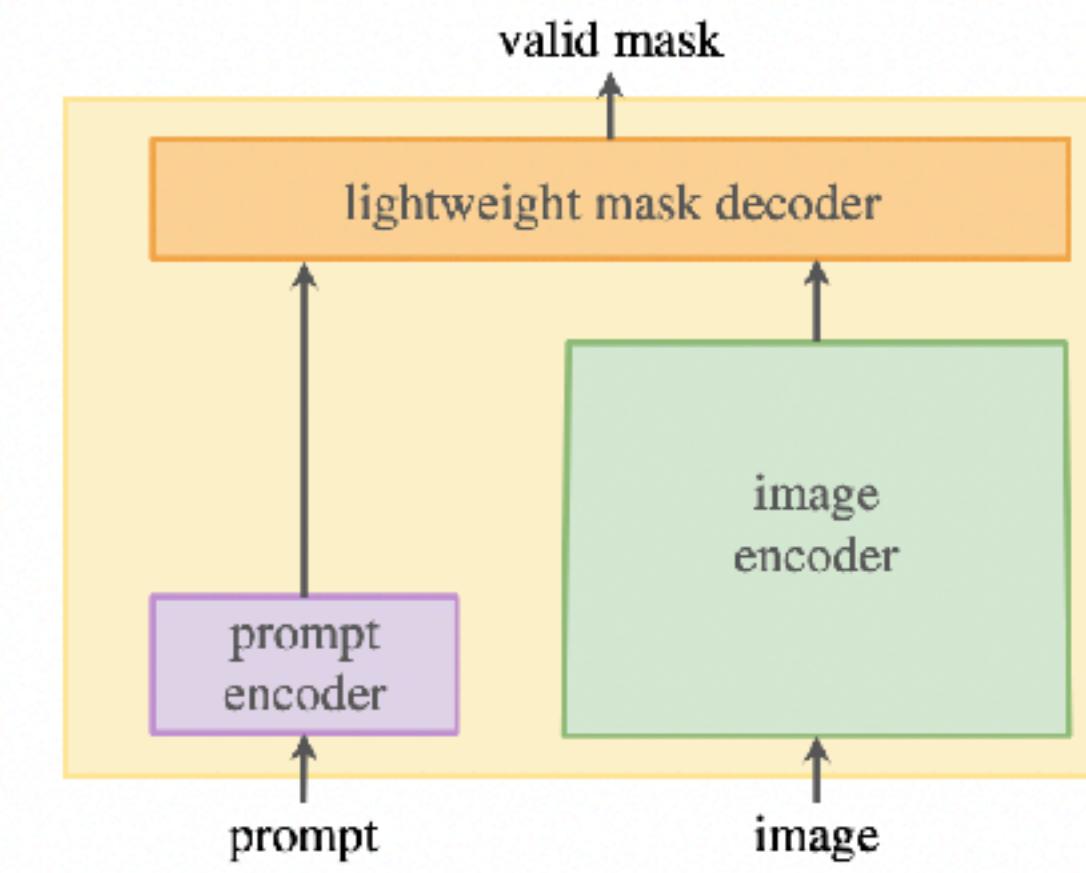
Alexander Kirillov^{1,2,4} Eric Mintun² Nikhila Ravi^{1,2} Hanzi Mao² Chloe Rolland³ Laura Gustafson³
Tete Xiao³ Spencer Whitehead Alexander C. Berg Wan-Yen Lo Piotr Dollár⁴ Ross Girshick⁴
¹project lead ²joint first author ³equal contribution ⁴directional lead

Meta AI Research, FAIR

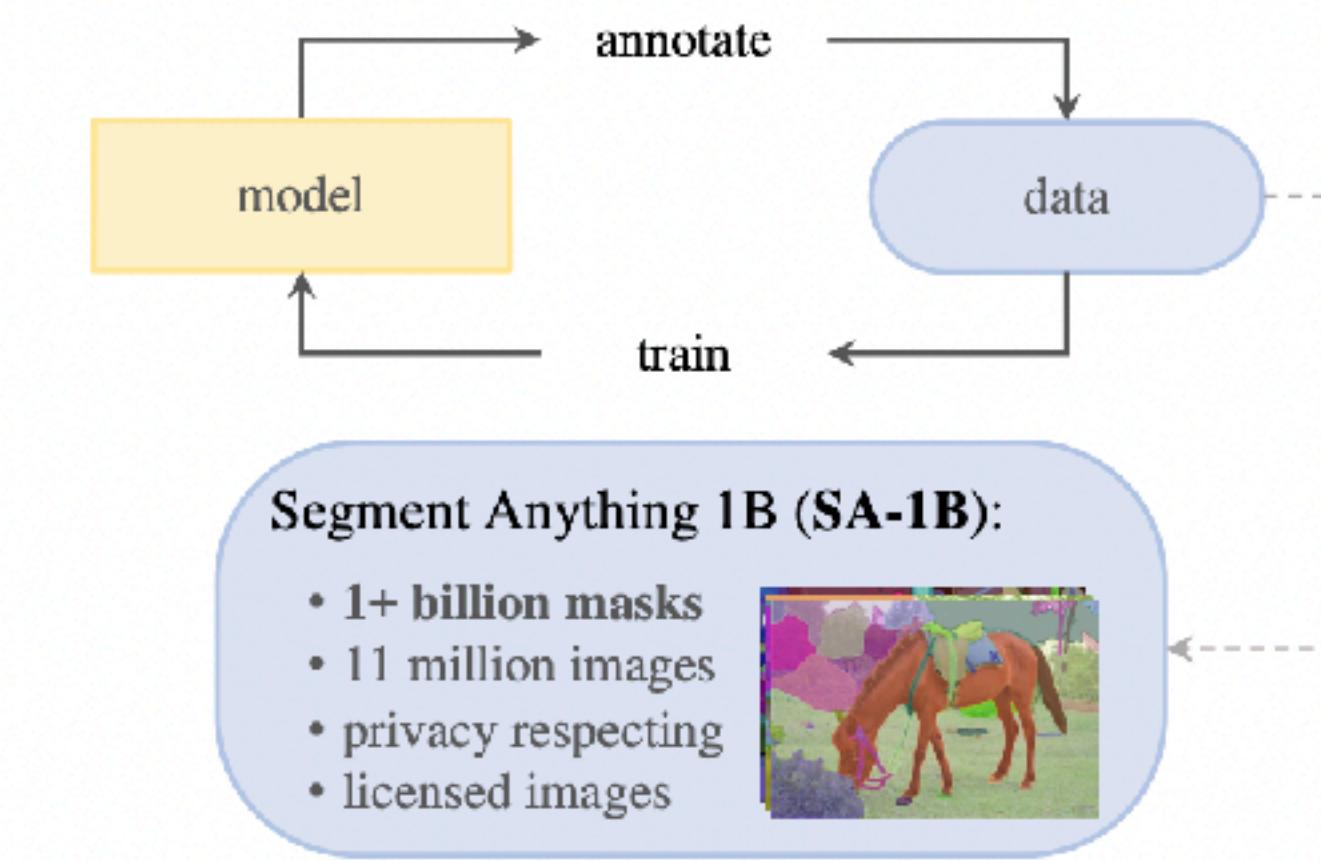
CS.CV] 5 Apr 2023



(a) Task: promptable segmentation



(b) Model: Segment Anything Model (SAM)

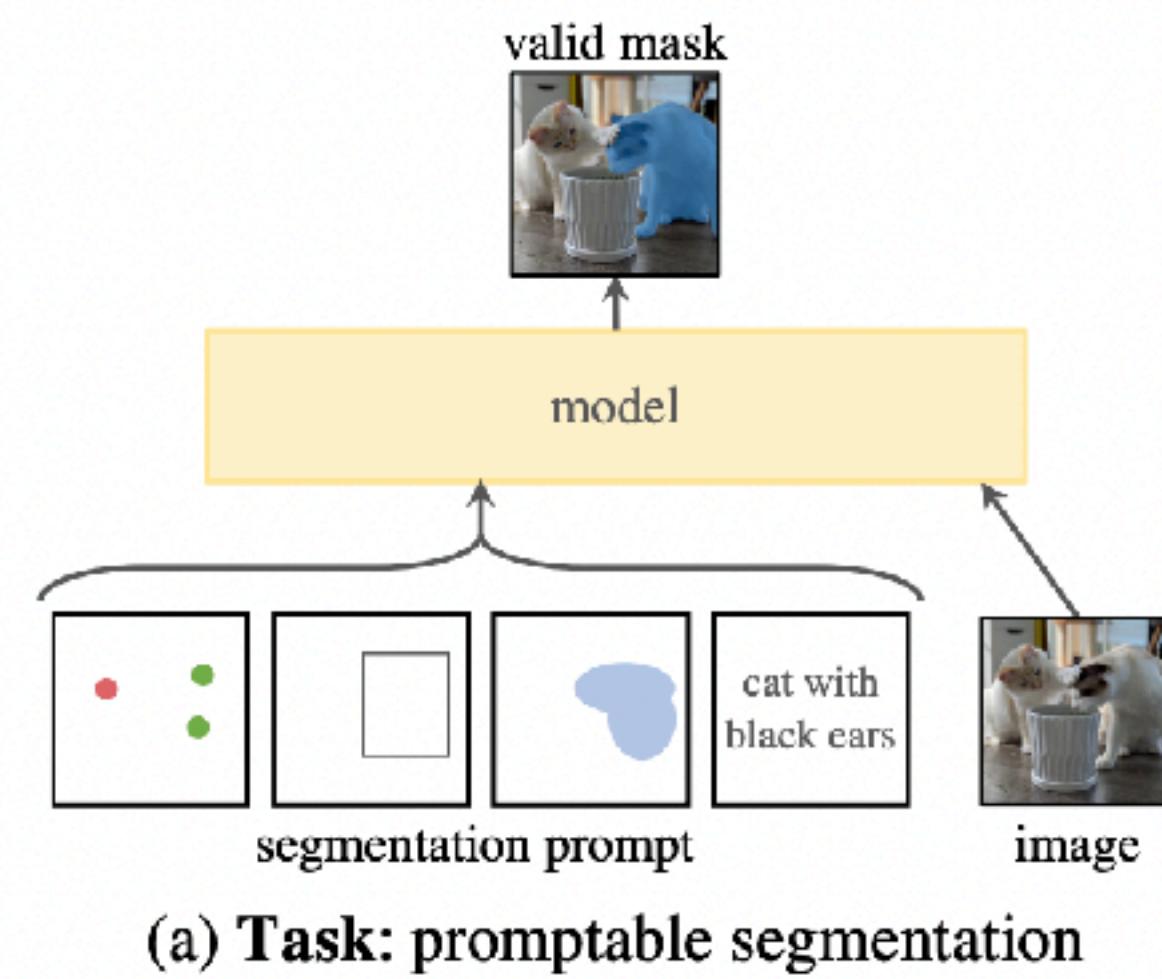


(c) Data: data engine (top) & dataset (bottom)

Figure 1: We aim to build a foundation model for segmentation by introducing three interconnected components: a promptable segmentation *task*, a segmentation *model* (SAM) that powers data annotation and enables zero-shot transfer to a range of tasks via prompt engineering, and a *data* engine for collecting SA-1B, our dataset of over 1 billion masks.

Segment Anything Model (SAM)

<https://segment-anything.com/demo>



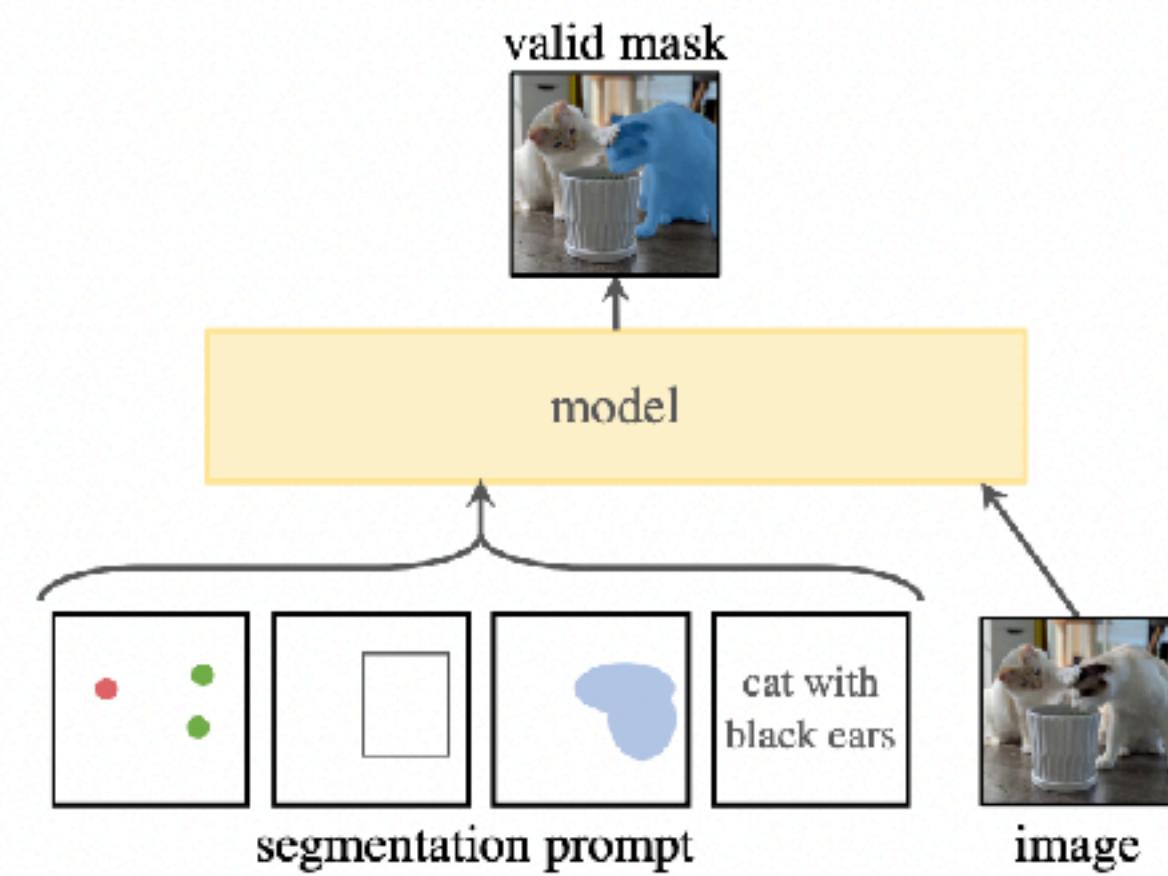
Prompting with a point



Prompting with a dense grid of points

Segment Anything Model (SAM)

Not semantic segmentation (no category)



(a) Task: promptable segmentation

Could be used for **instance segmentation** by integrating an object detector



Prompting with detected boxes

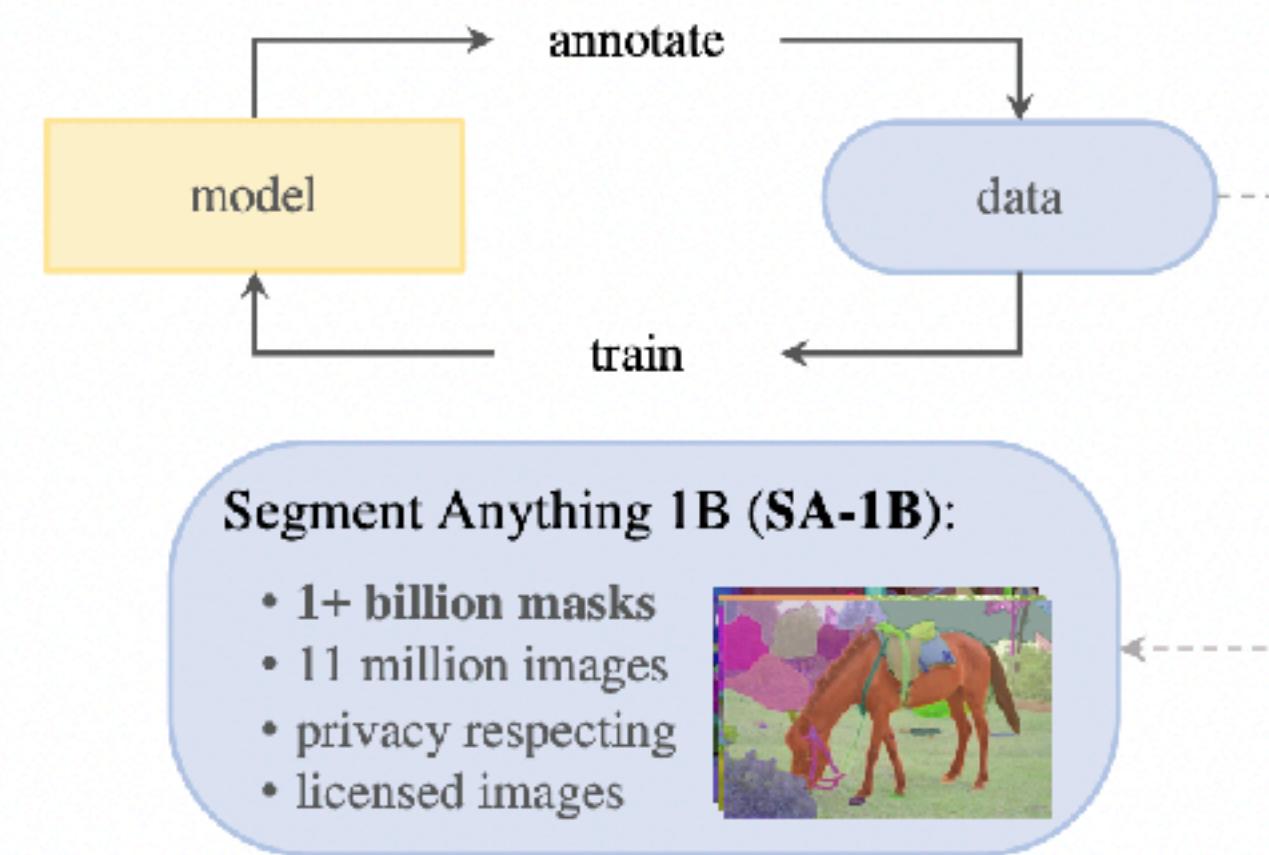
Segment Anything Dataset (SA-1B)

- 11M images
- 1B+ masks (99.1% of masks fully automatic)
- Collected through interactive interface



3-stage annotation:

- Assisted-manual stage (+30sec/image to annotate, reduced to 14sec after 6 x retraining, 4.3M masks from 12K images)
- Semi-automatic stage (bbox for less prominent objects, up to 34sec. 5 x retraining, 5.9M masks in 180K images)
- Fully-automatic stage.

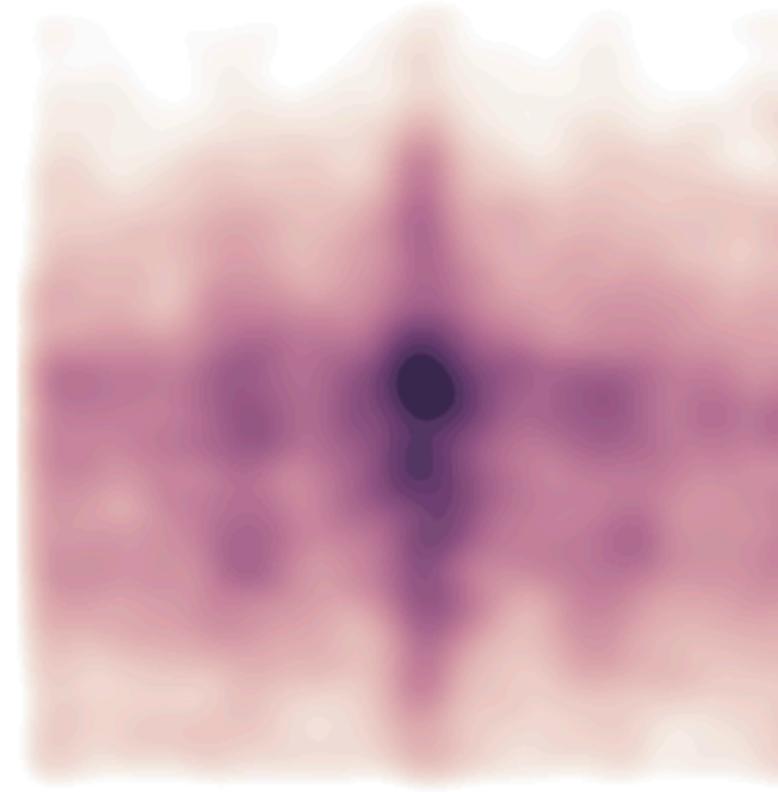


(c) Data: data engine (top) & dataset (bottom)

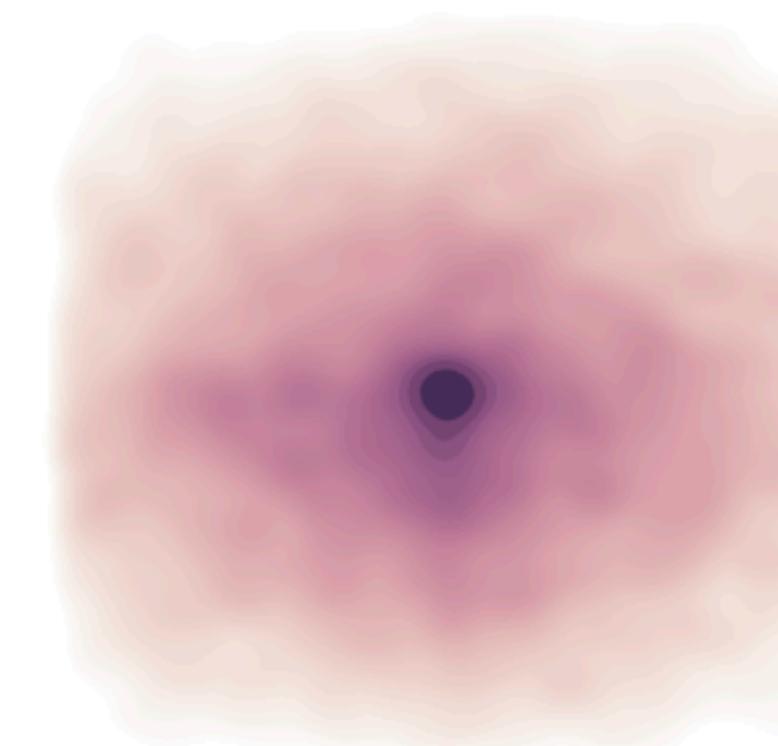
Segment Anything

- Spatial distribution of object centers
- Common photographer bias
- Greater coverage of image corners in SA-1B

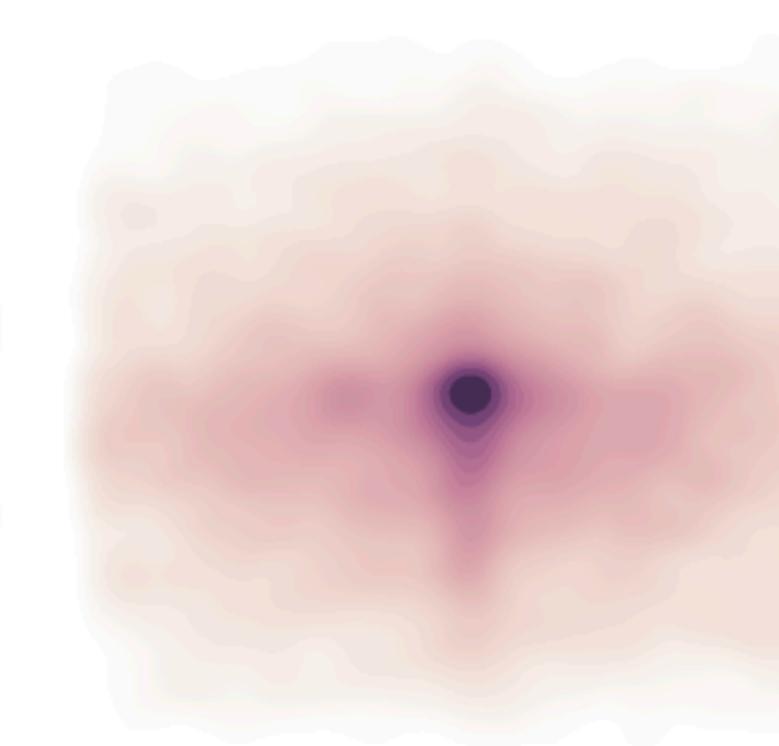
SA-1B



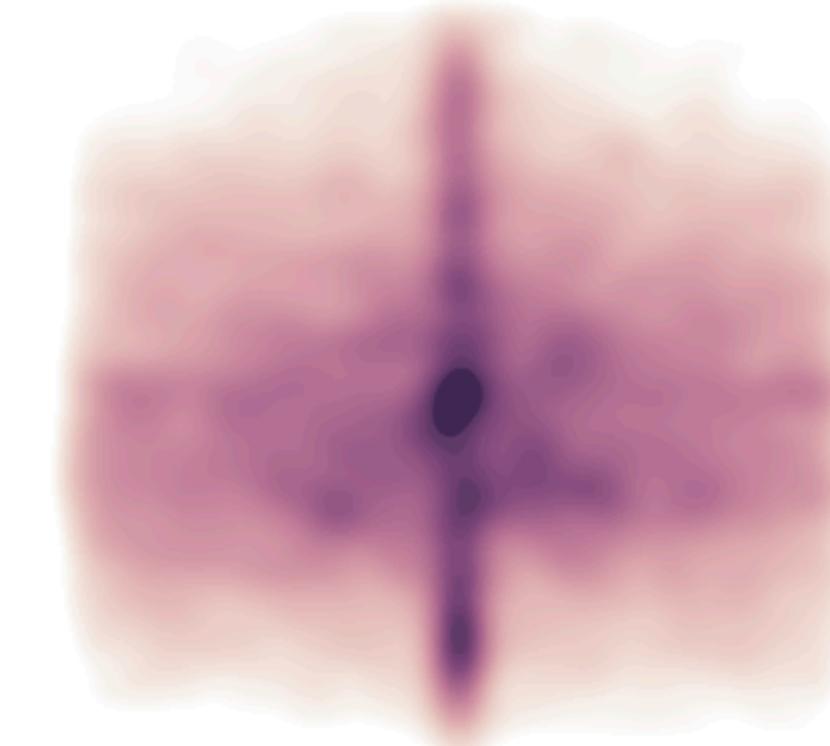
LVIS v1



COCO



ADE20K



Open Images

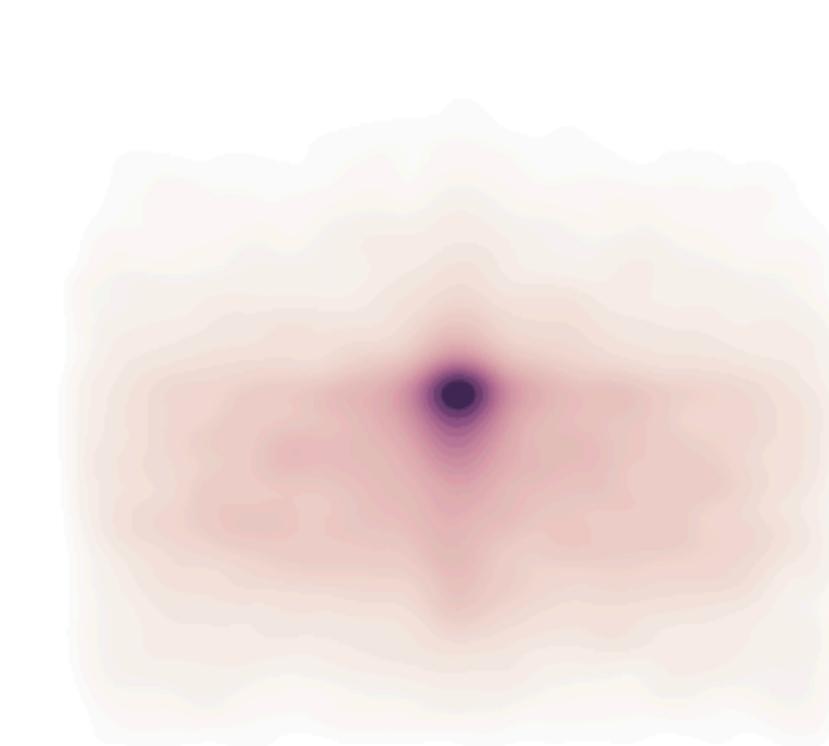


Image-size normalized mask center distributions

2. Beyond classification

- a. Intro to structured outputs
- b. Object detection (localization)
- c. Segmentation
- d. Human pose estimation



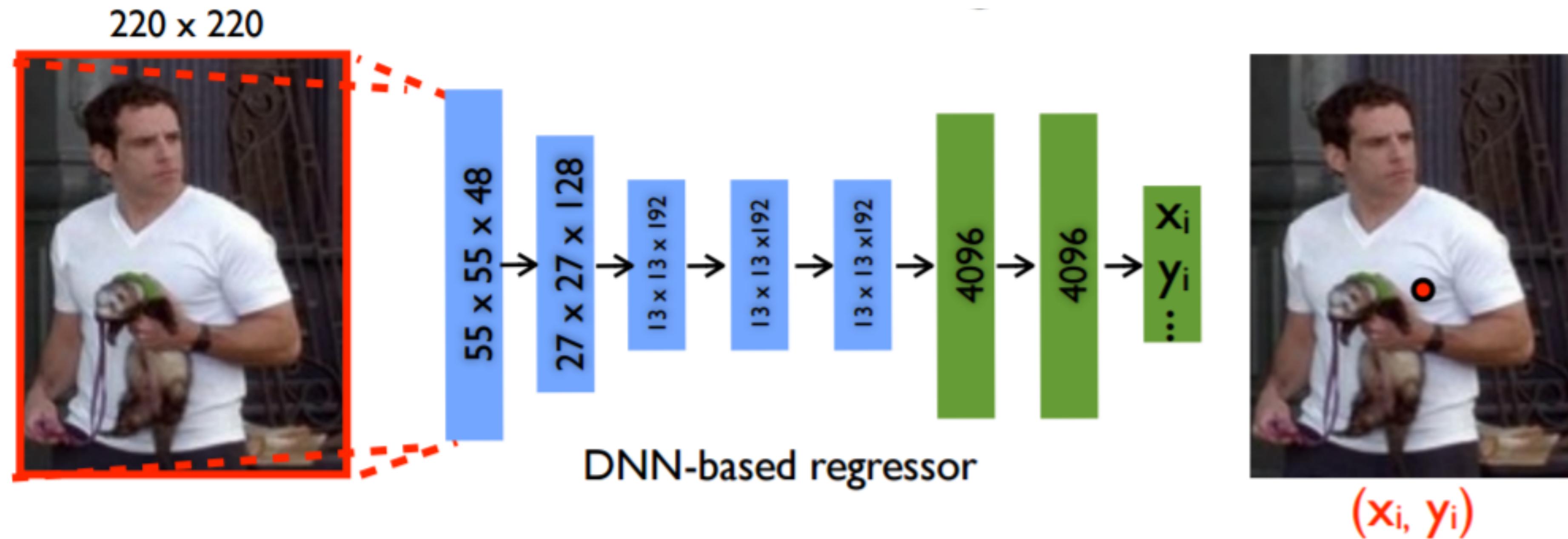
2D Human pose estimation



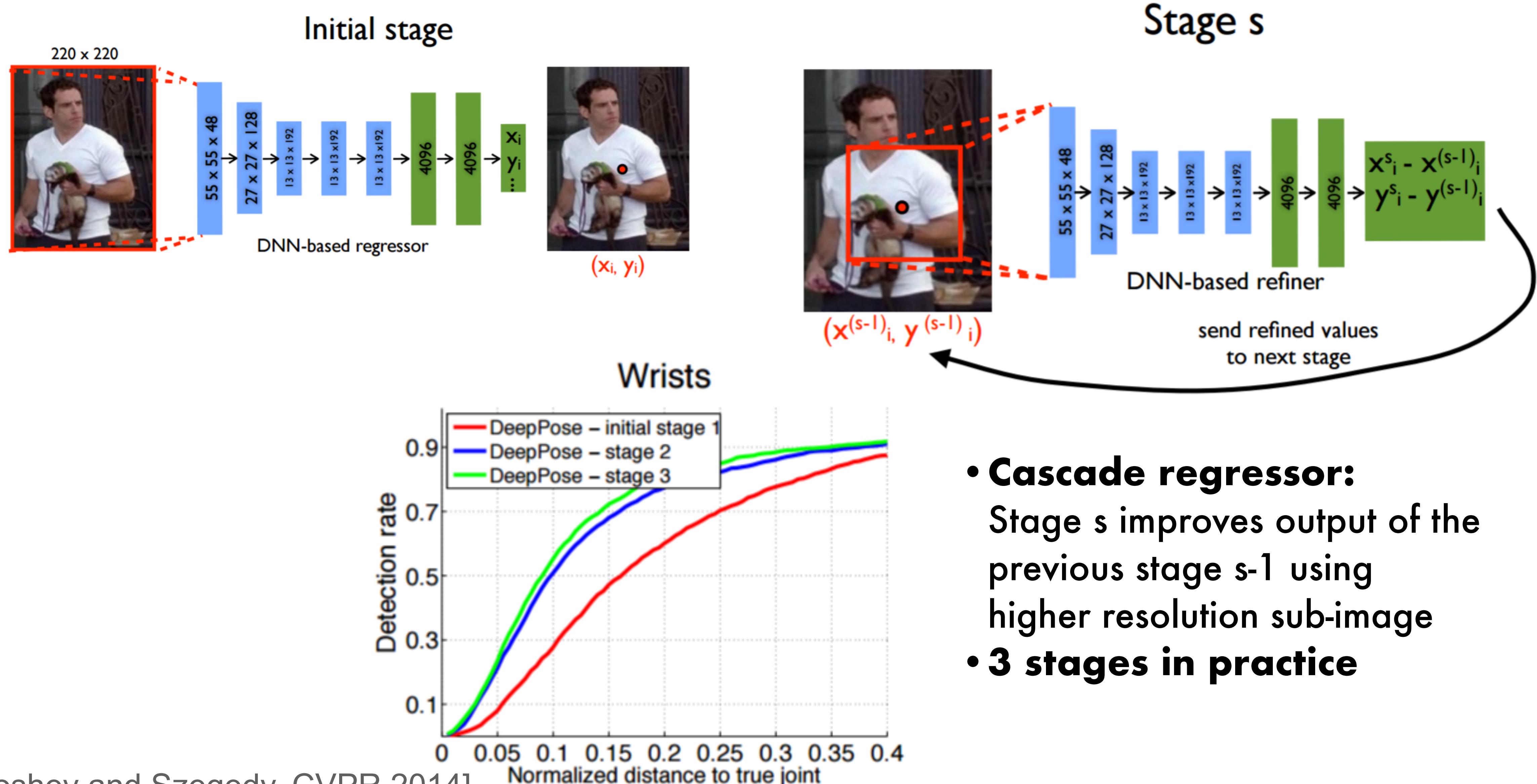
Source: <https://www.youtube.com/watch?v=2DiQUX11YaY>

DeepPose: Human Pose Estimation via Deep Neural Networks

Trains CNN to regress locations (x_i, y_i) for each joint i



DeepPose: Human Pose Estimation via Deep Neural Networks

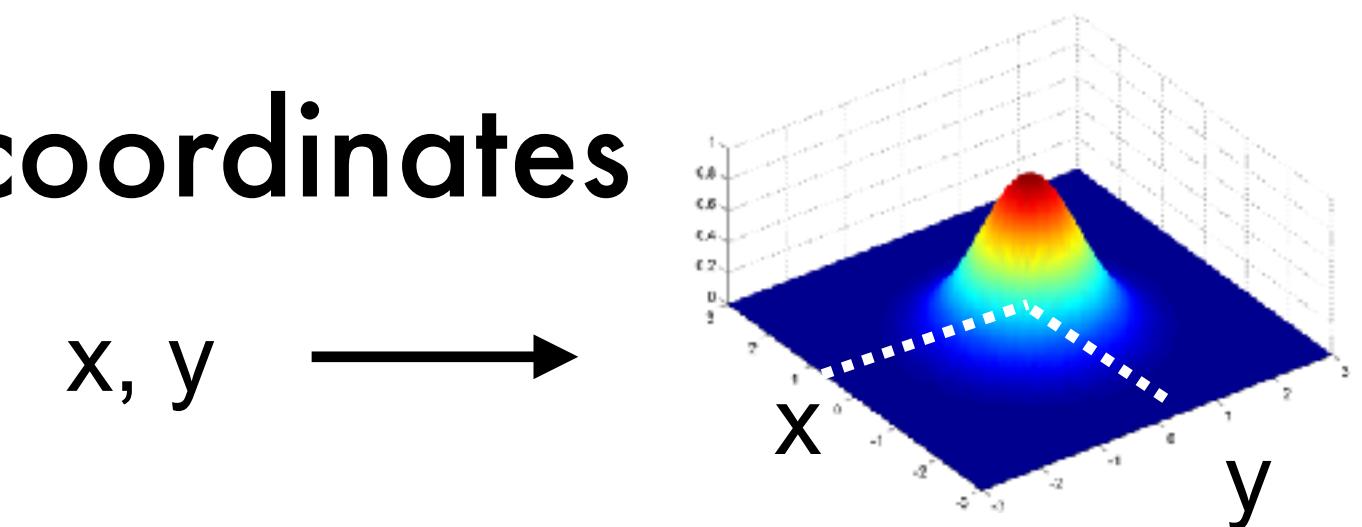


DeepPose: Human Pose Estimation via Deep Neural Networks



Convolutional Pose Machines

- Regression to joint “heatmaps”: 2D gaussians around joint coordinates
- Heatmaps enable to handle spatial ambiguity



Input Image

CNN



Heatmap
for right elbow

Convolutional Pose Machines

- Regression to joint “heatmaps”: 2D gaussians around joint coordinates
- Heatmaps enable to handle spatial ambiguity
- Multi-stage refinement



Input Image

(a) Stage 1

(b) Stage 2

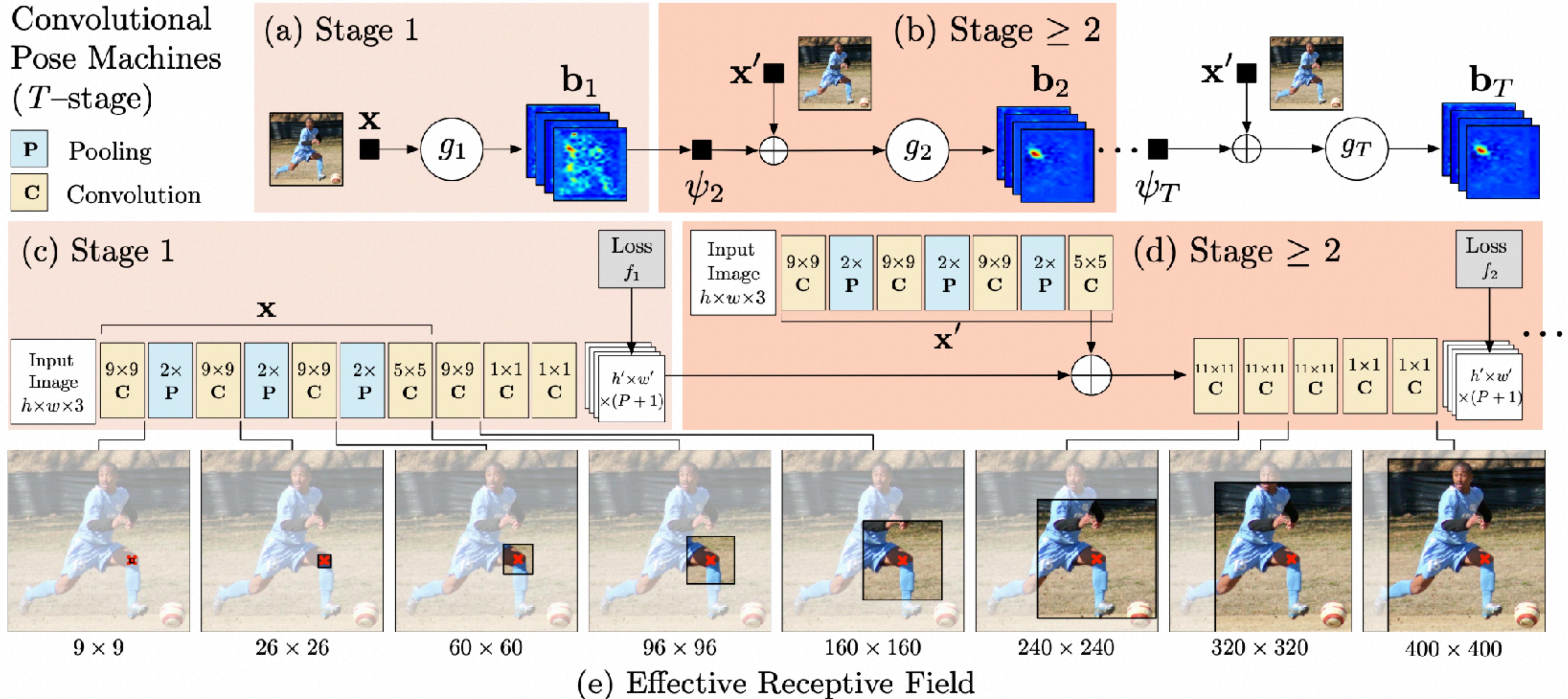
(c) Stage 3

Convolutional Pose Machines

- Intermediate supervision at every stage; Increasing context

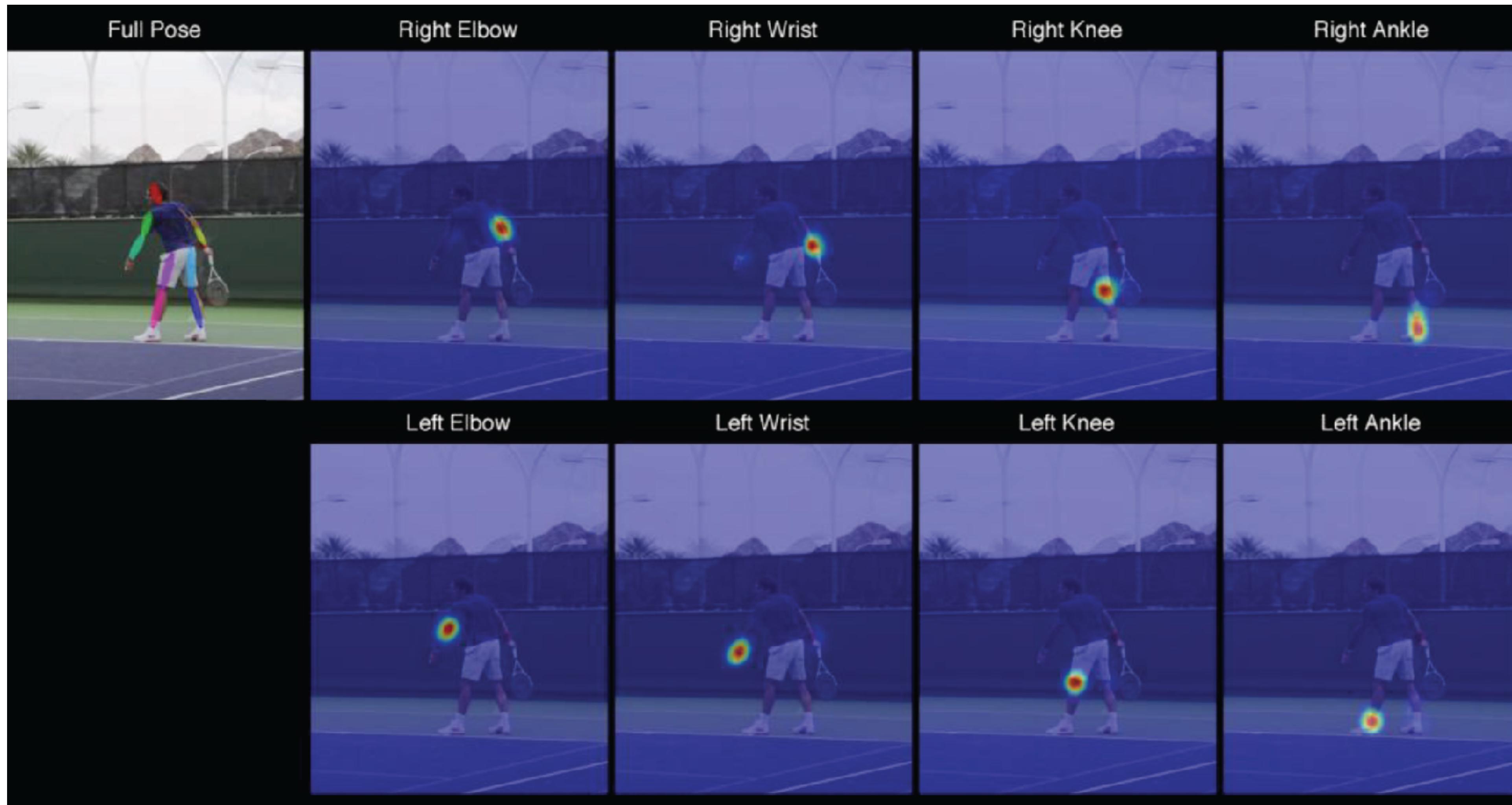
Convolutional
Pose Machines
(T -stage)

P Pooling
C Convolution



Convolutional Pose Machines

Qualitative results



[Wei, Ramakrishna, Kanade and Sheikh, CVPR 2016]

Convolutional Pose Machines

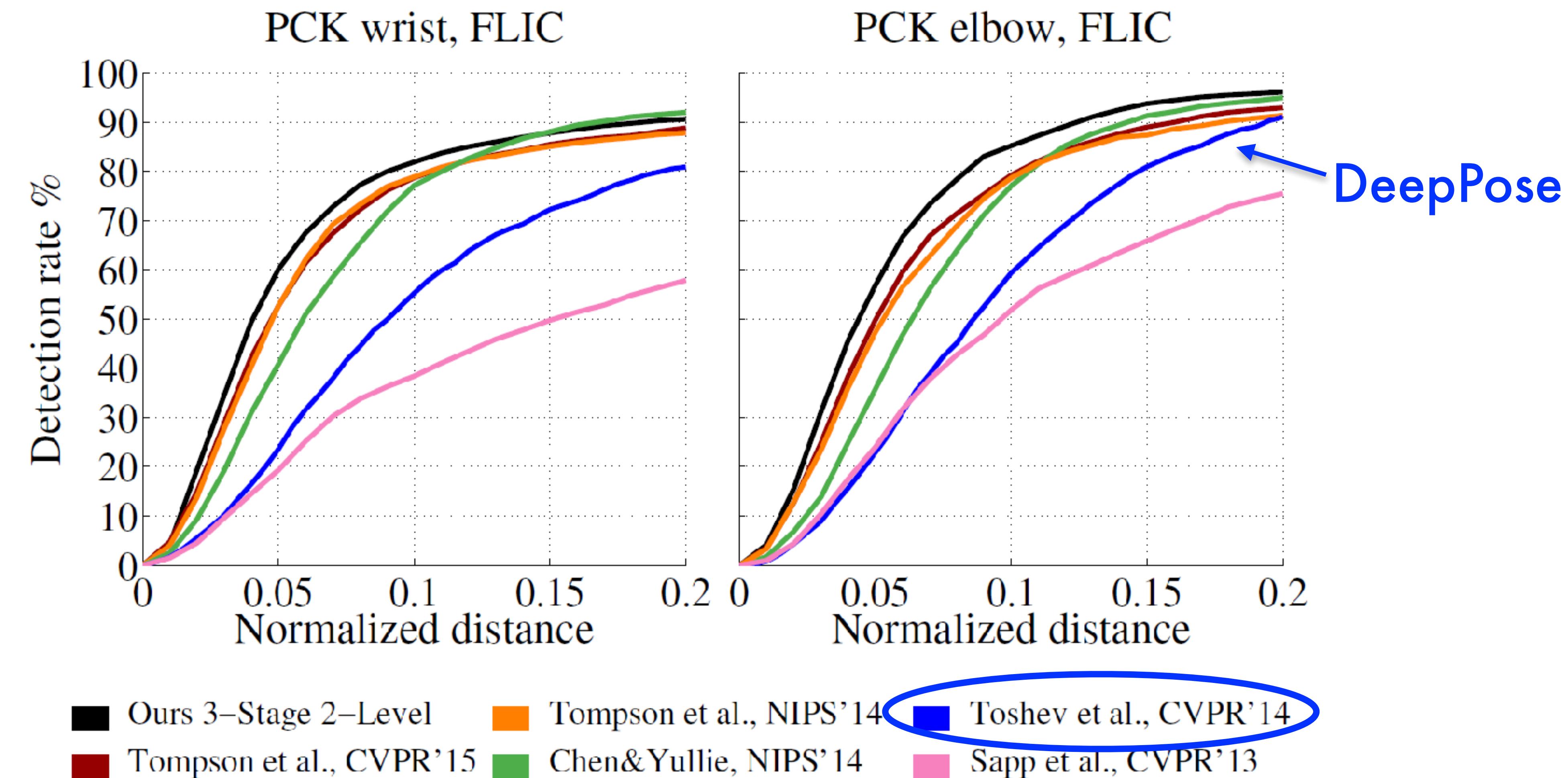
Qualitative results



[Wei, Ramakrishna, Kanade and Sheikh, CVPR 2016]

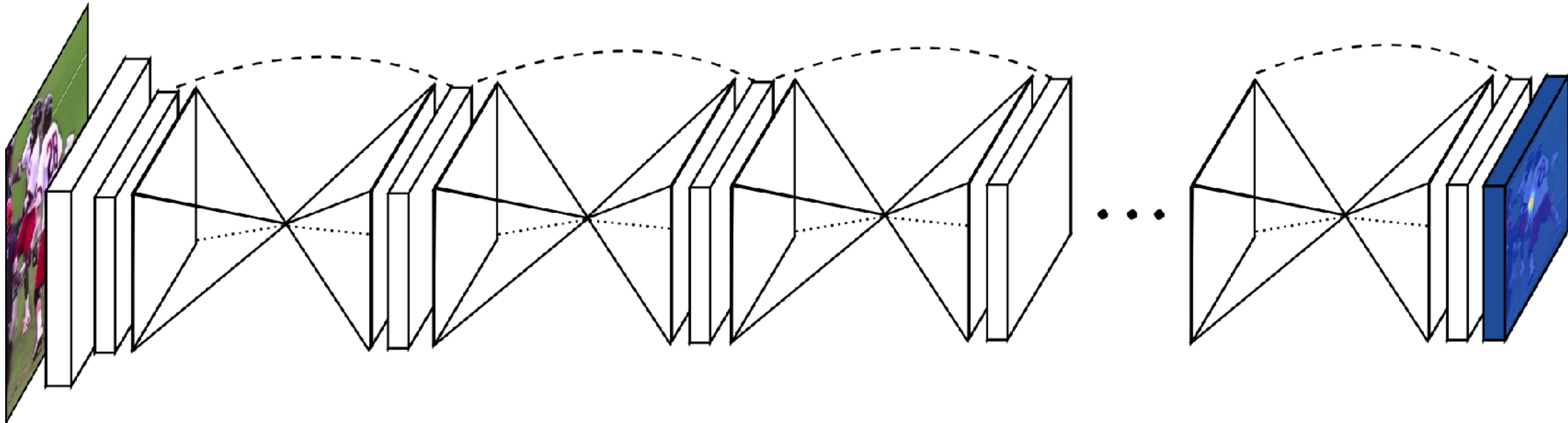
Convolutional Pose Machines

Quantitative comparison



Stacked Hourglass Networks

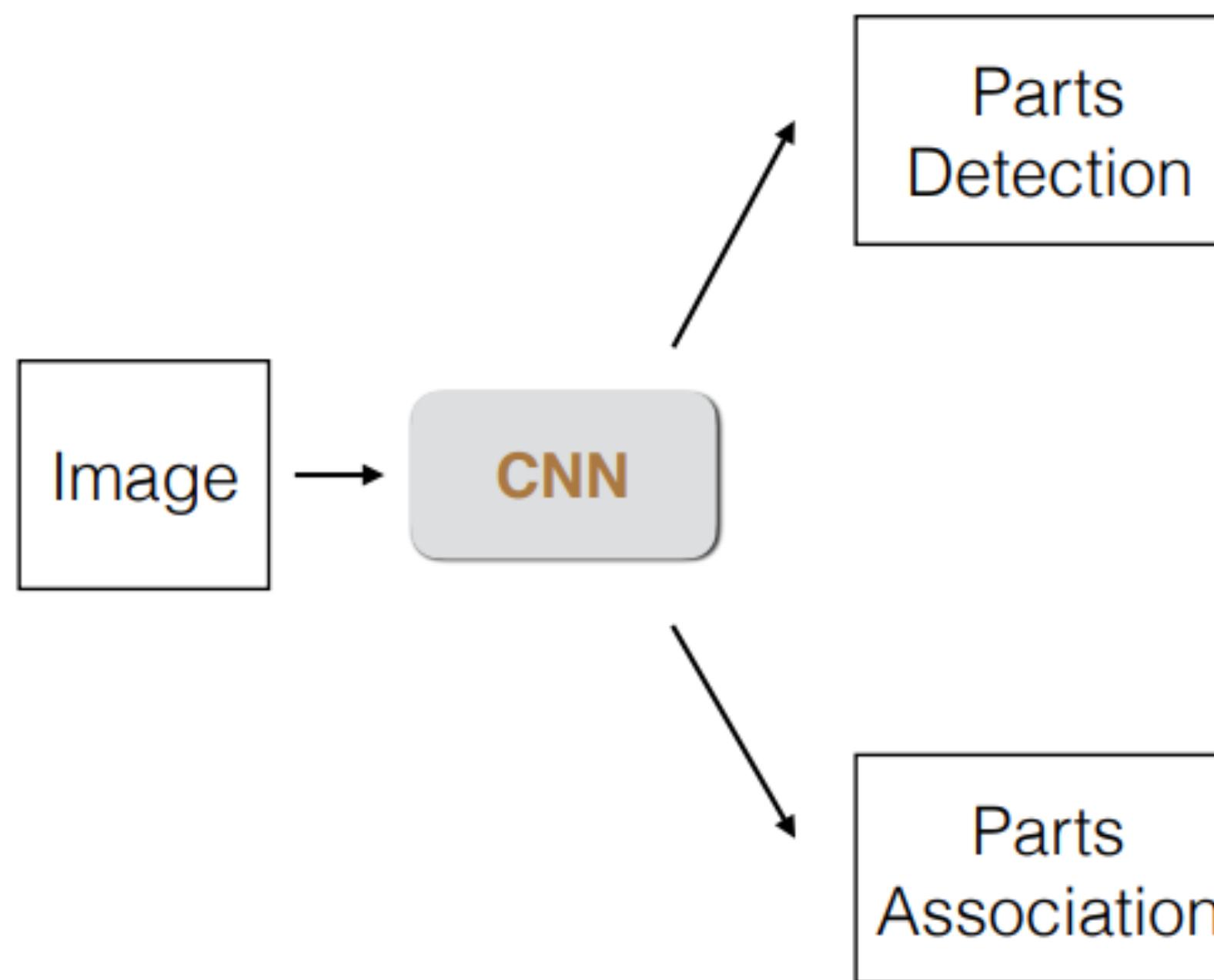
Remember U-Net



- Also heatmap regression
- Also multi-stage refinement - but full context (receptive field = entire image)

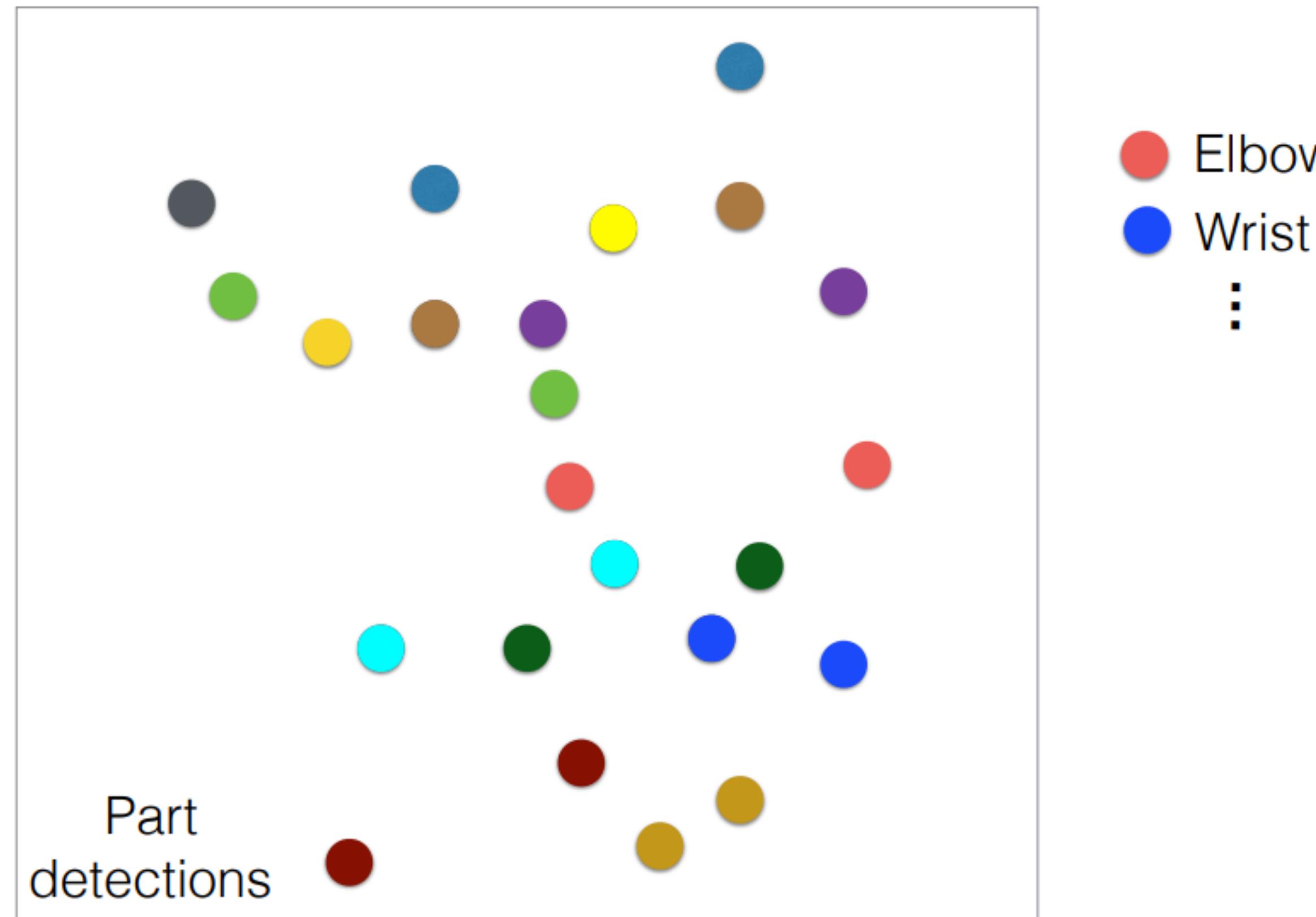
OpenPose: Multi-person pose estimation

Novelty: Jointly Learning Parts Detection and Parts Association



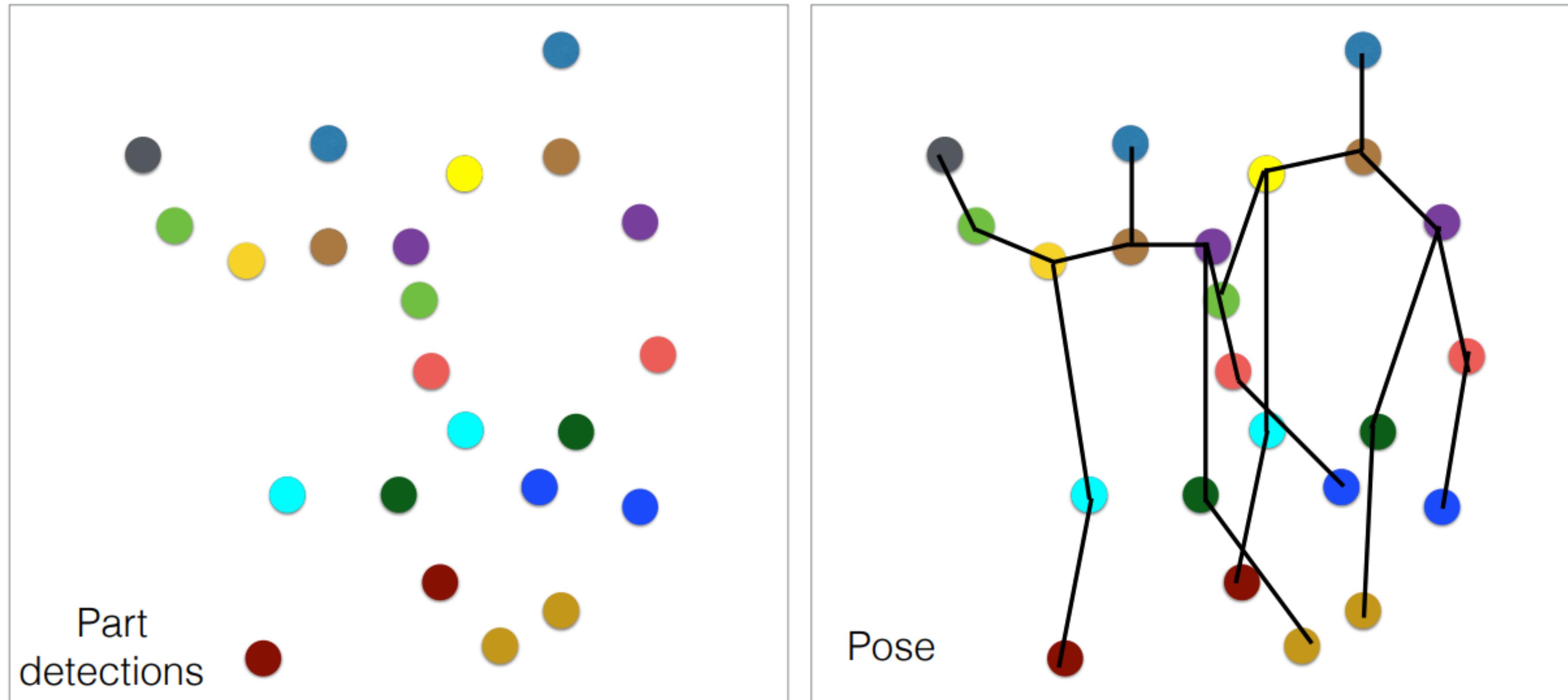
OpenPose

Part-Person Association for Multi-Person Pose Estimation



OpenPose

Part-Person Association for Multi-Person Pose Estimation



OpenPose

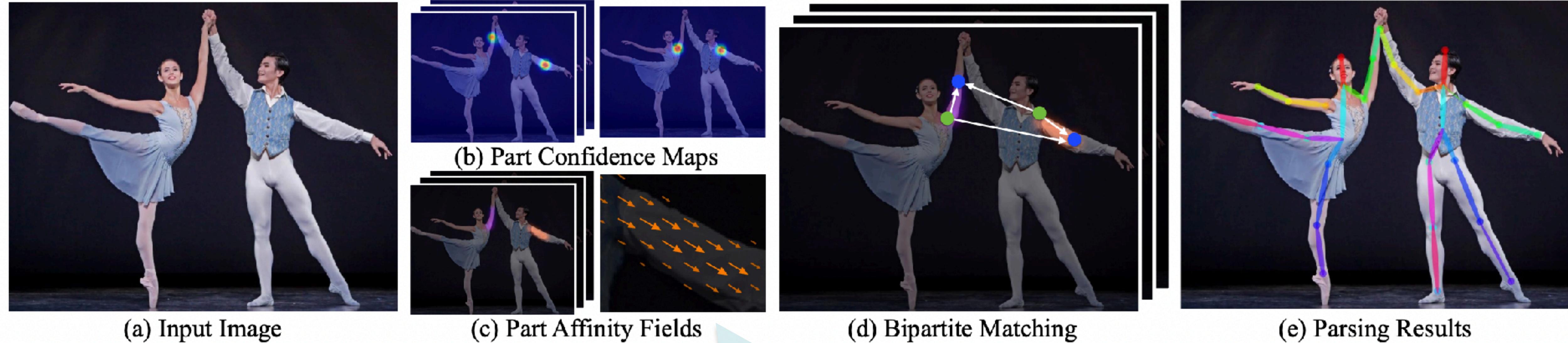
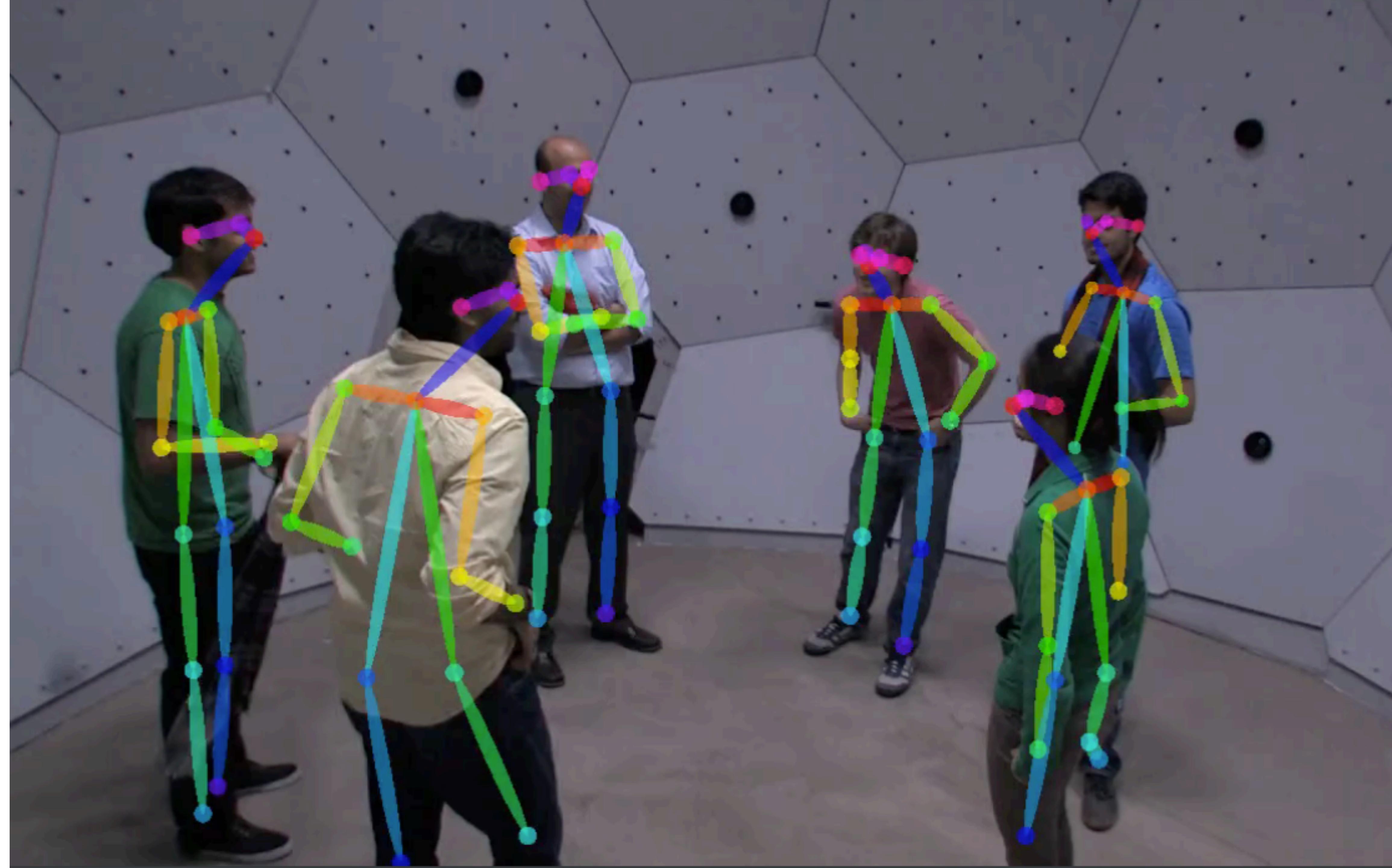


Figure 2. Overall pipeline. Our method takes the entire image as the input for a two-branch CNN to jointly predict confidence maps for body part detection, shown in (b), and part affinity fields for parts association, shown in (c). The parsing step performs a set of bipartite matchings to associate body parts candidates (d). We finally assemble them into full body poses for all people in the image (e).

Key Idea: Encode the Part Affinity Score on the Image Plane
=> Part Affinity Fields encode direction and position



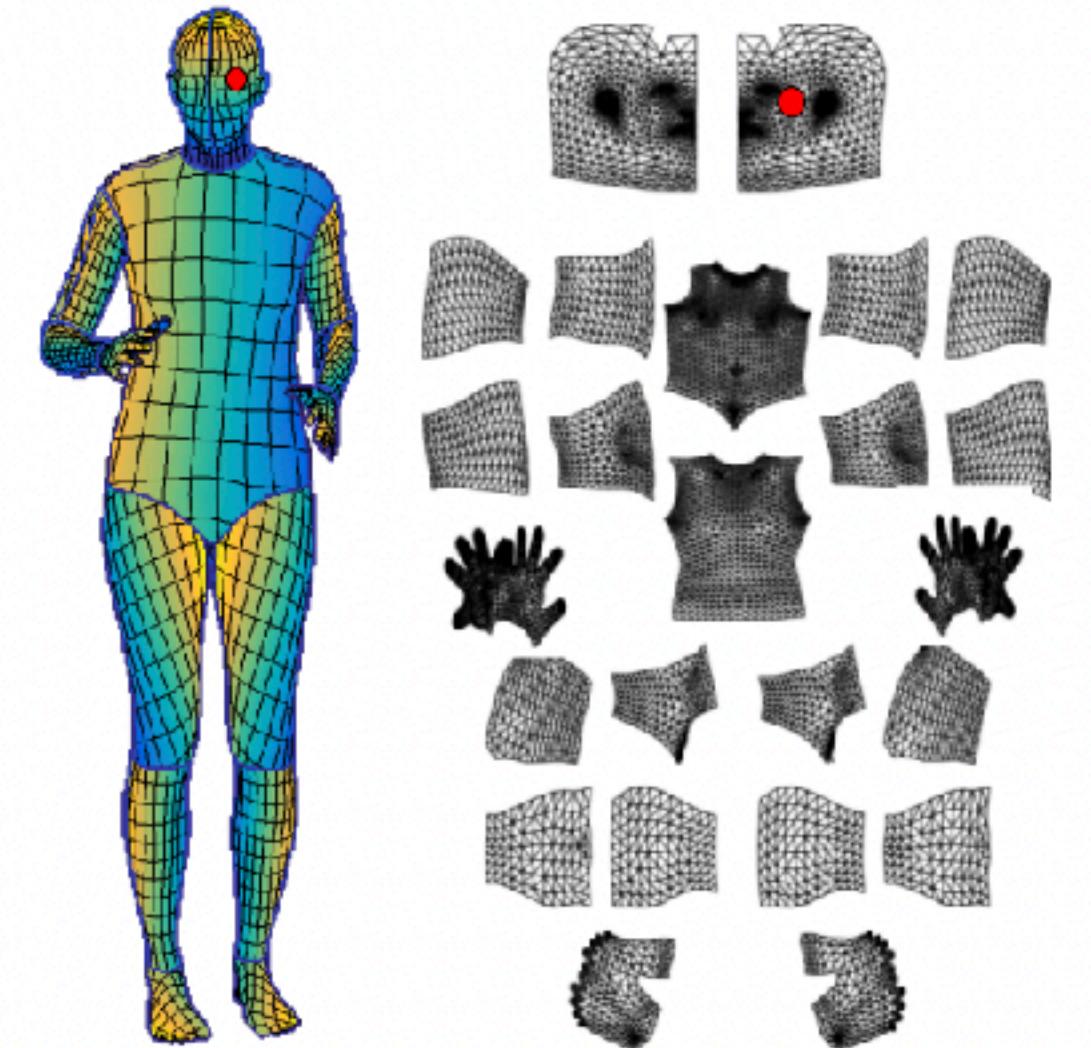
DensePose: Dense Human Pose Estimation



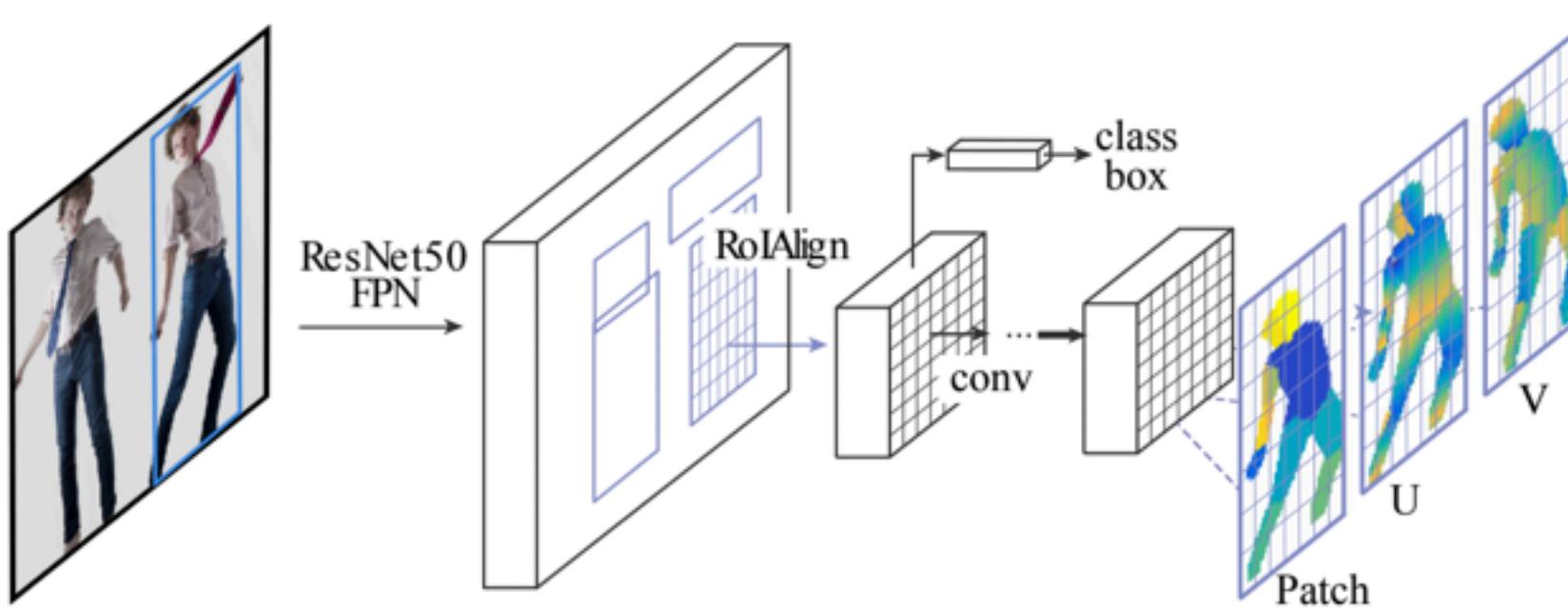
DensePose-RCNN Results



DensePose COCO Dataset



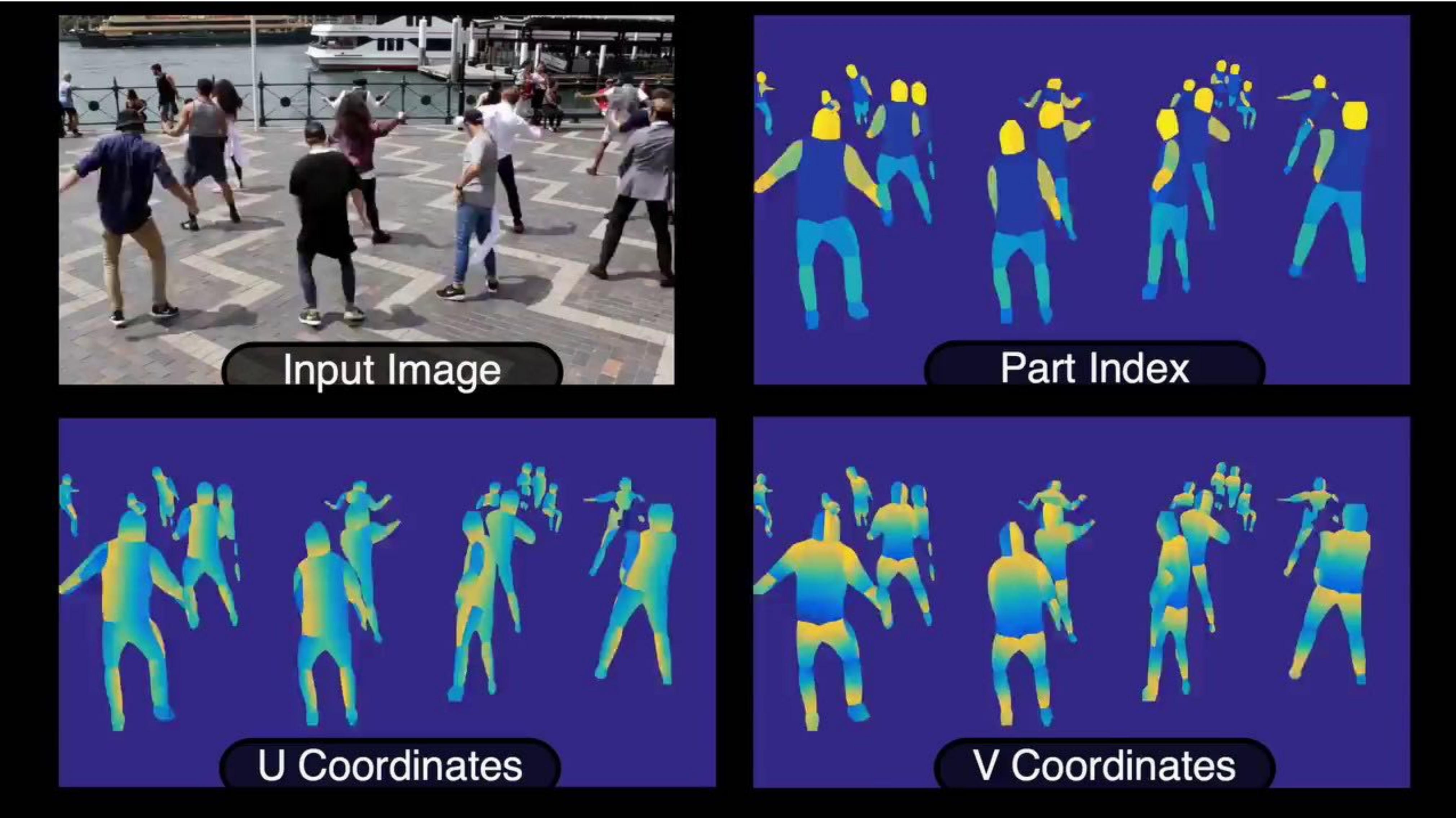
Dense pose estimation aims at mapping all human pixels of an RGB image to the 3D surface of the human body.



regresses to continuous surface coordinates

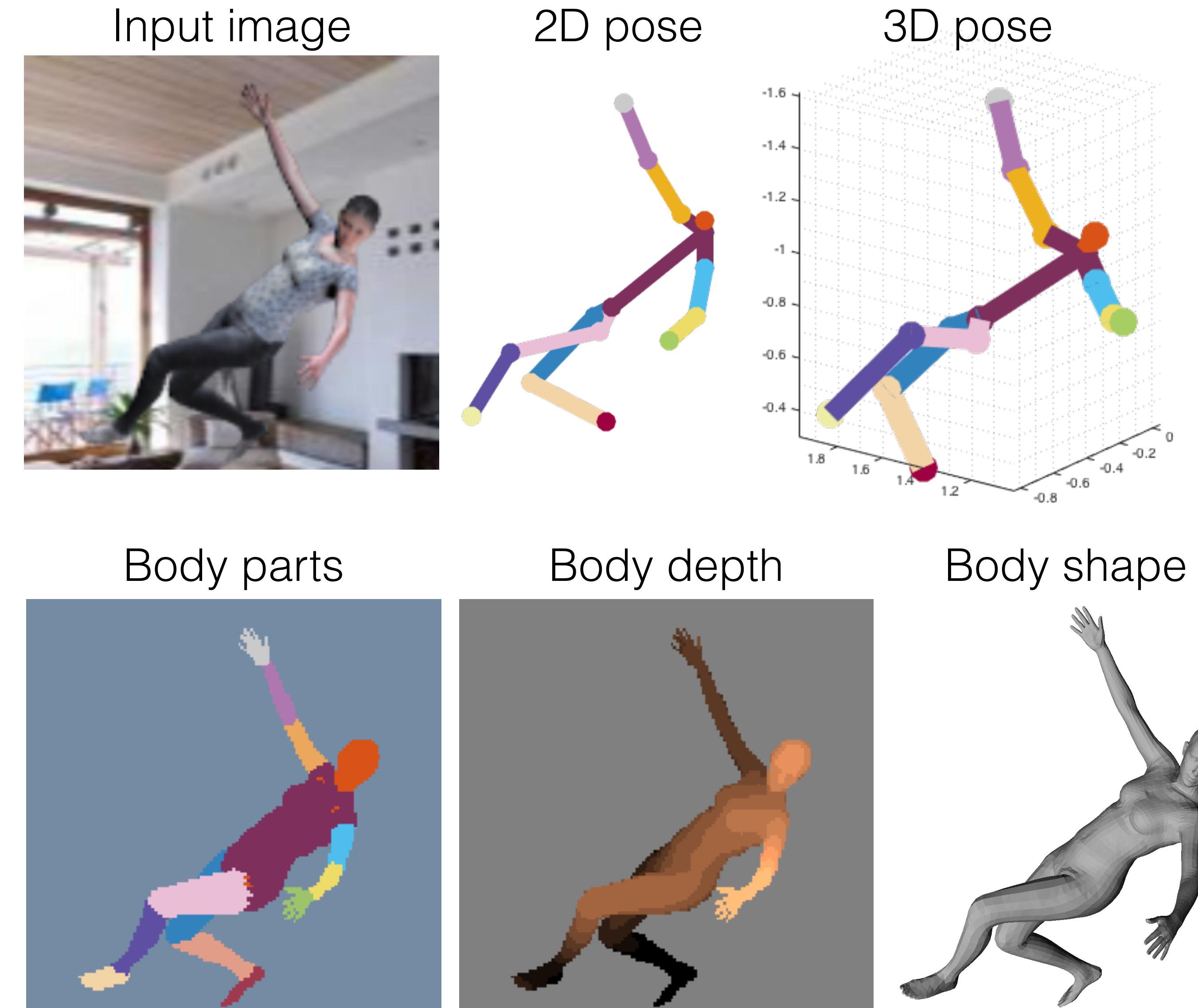
Guler et al. DensePose, CVPR 2018

DensePose



Human pose estimation beyond 2D keypoints

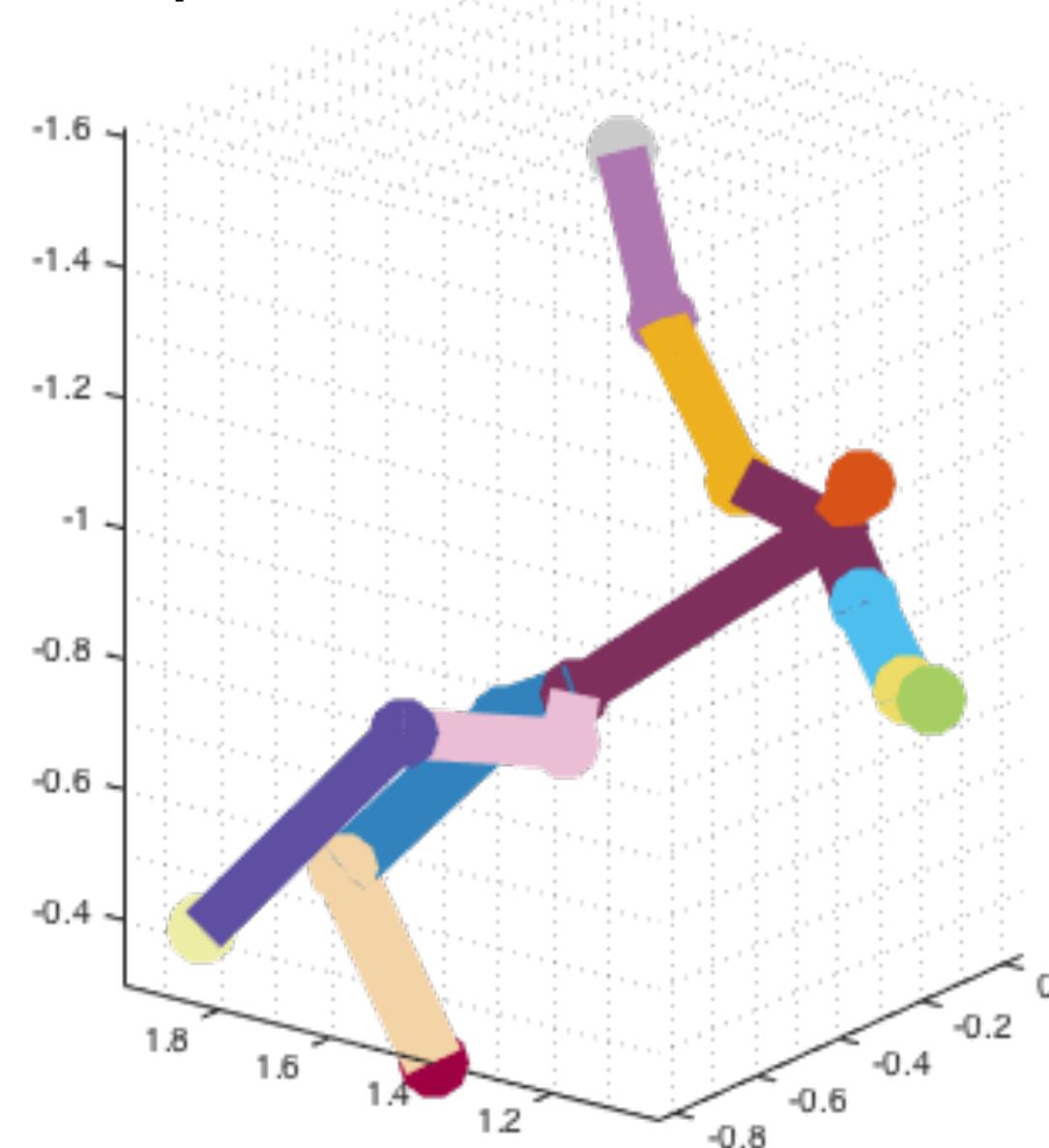
Human body analysis



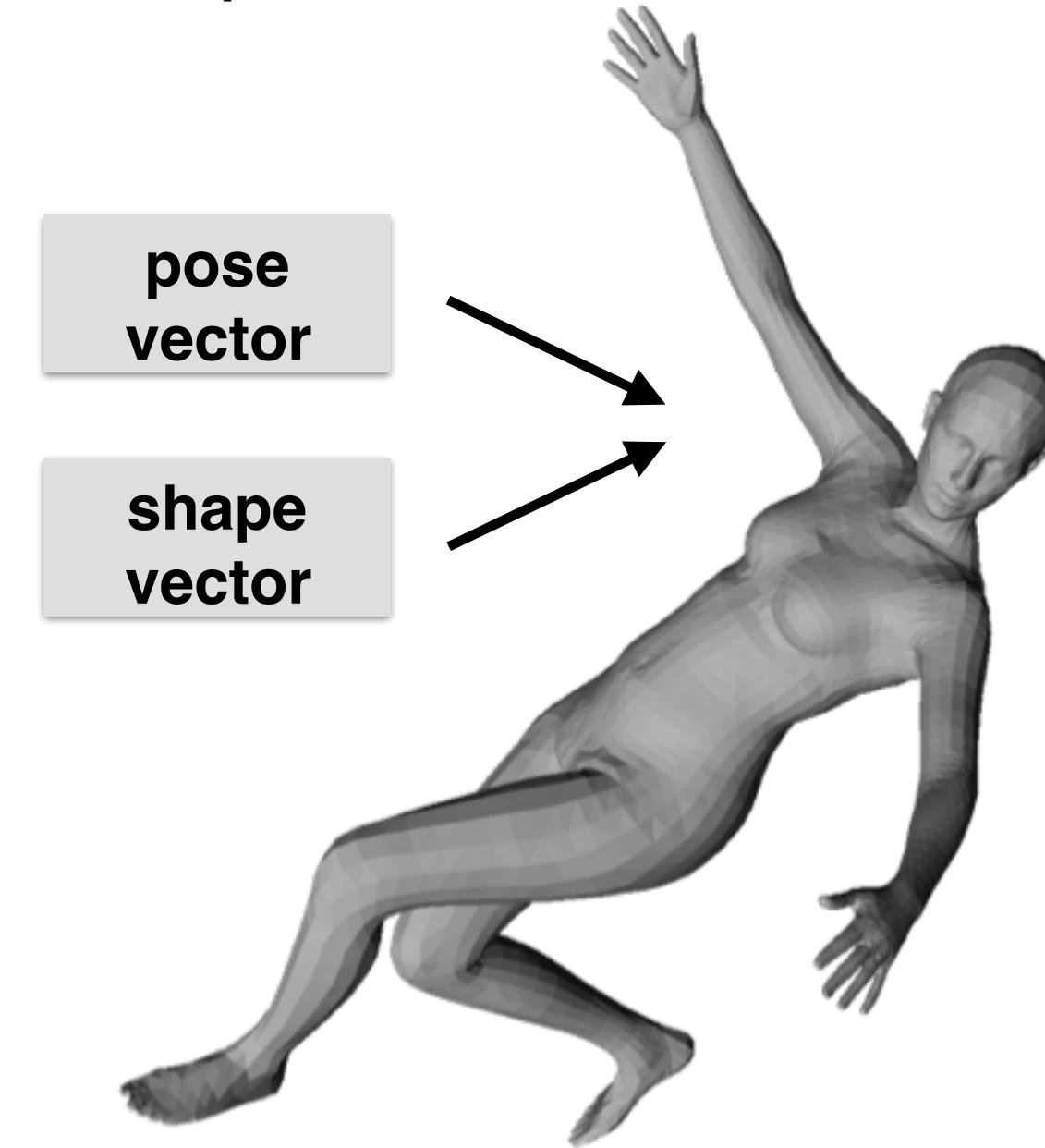
Challenges

How to model the body shape?

(a) Skeleton representation

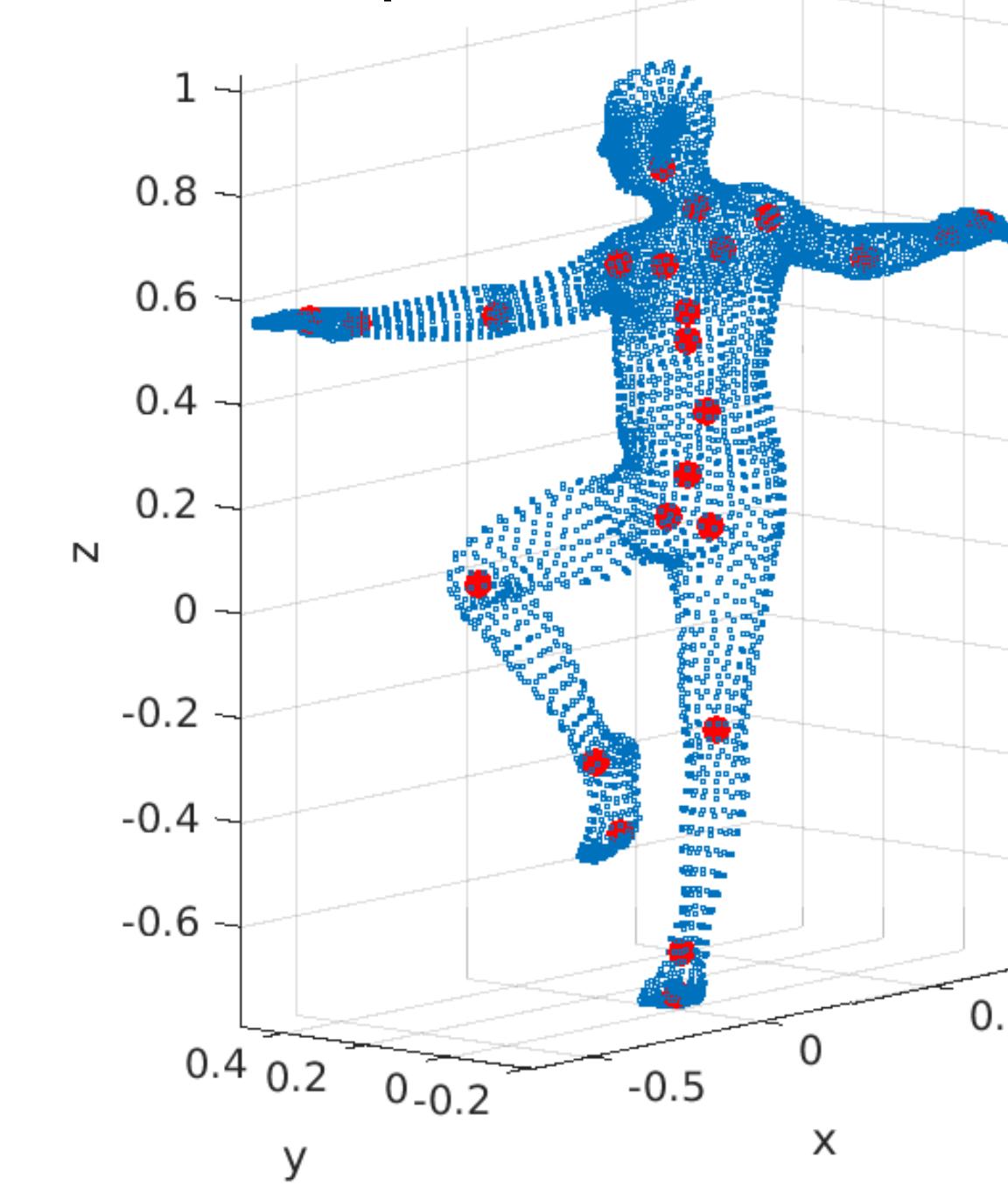


(b) Parametric representation

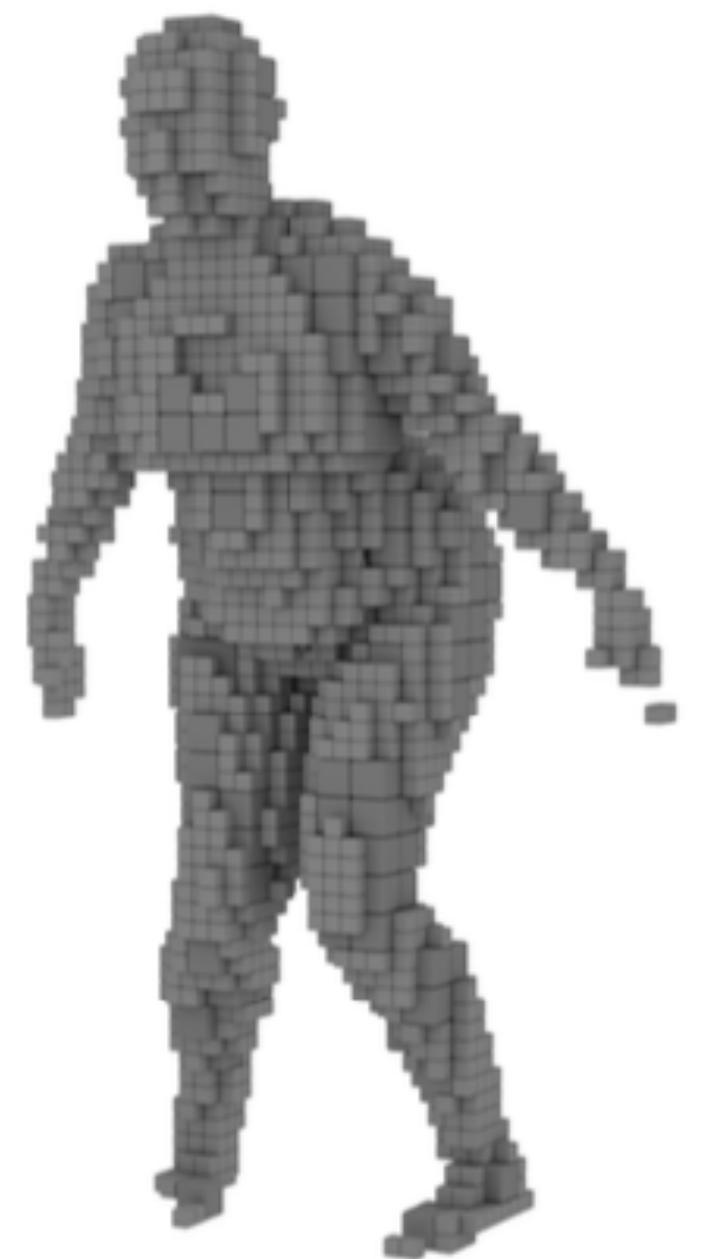


Loper et al. [1]

(c) Point cloud representation



(d) Voxel representation

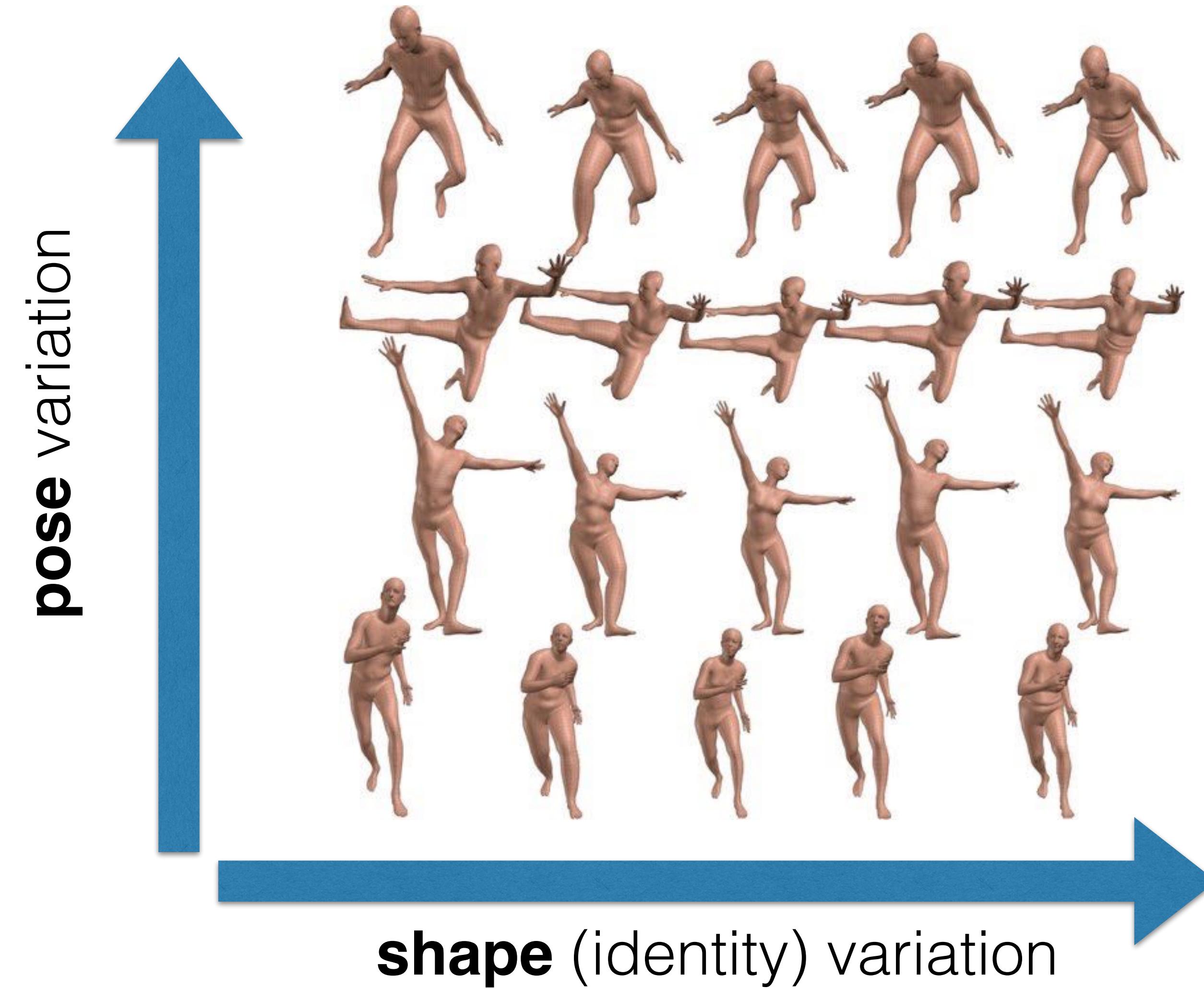


Tatarchenko et al. [2]

[1] Loper et al. SMPL: A Skinned Multi-Person Linear Model, SIGGRAPH Asia 2015

[2] Tatarchenko et al. Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs, ICCV 2017

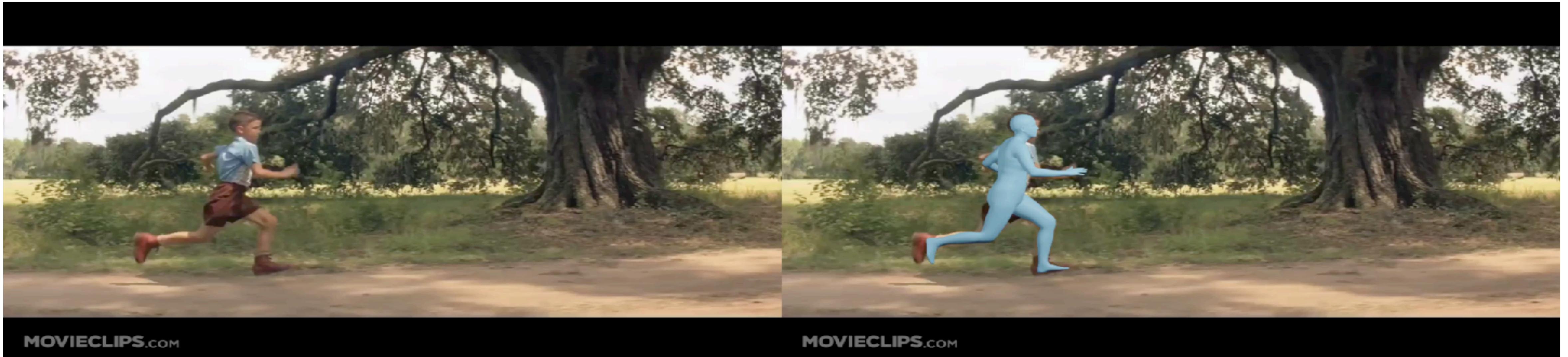
SMPL parametric body model: surface & joints



[Loper et al. 2015]

Human pose estimation beyond 2D keypoints

- A rich literature also on 3D human pose & motion estimation



VIBE [Kocabas et al. CVPR 2020]

Human pose estimation beyond 2D keypoints

- A rich literature also on 3D human pose & motion estimation



4D Humans [Goel et al. ICCV 2023]

1. Beyond CNNs

- Attention & Transformer

- Vision Transformers

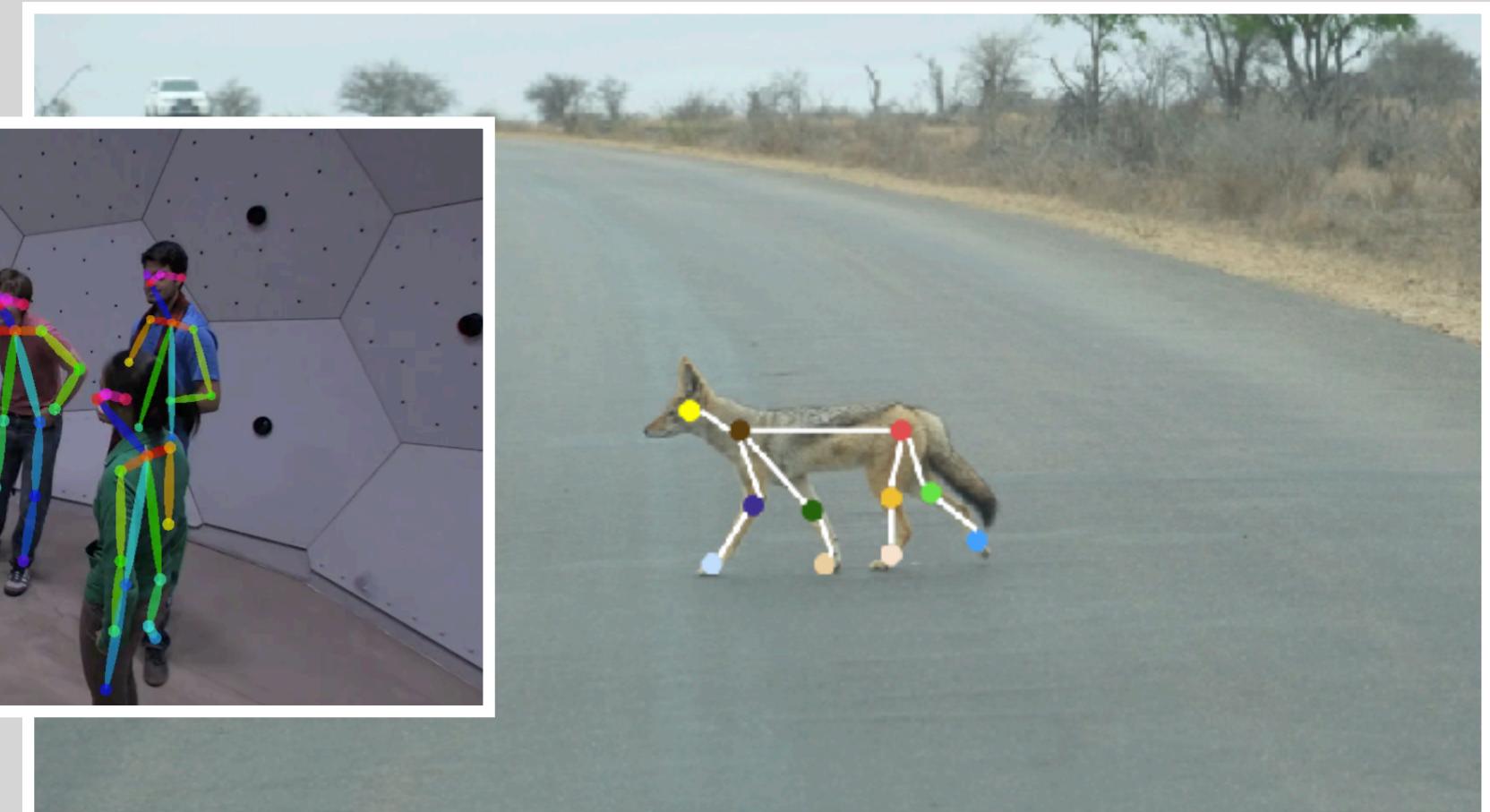
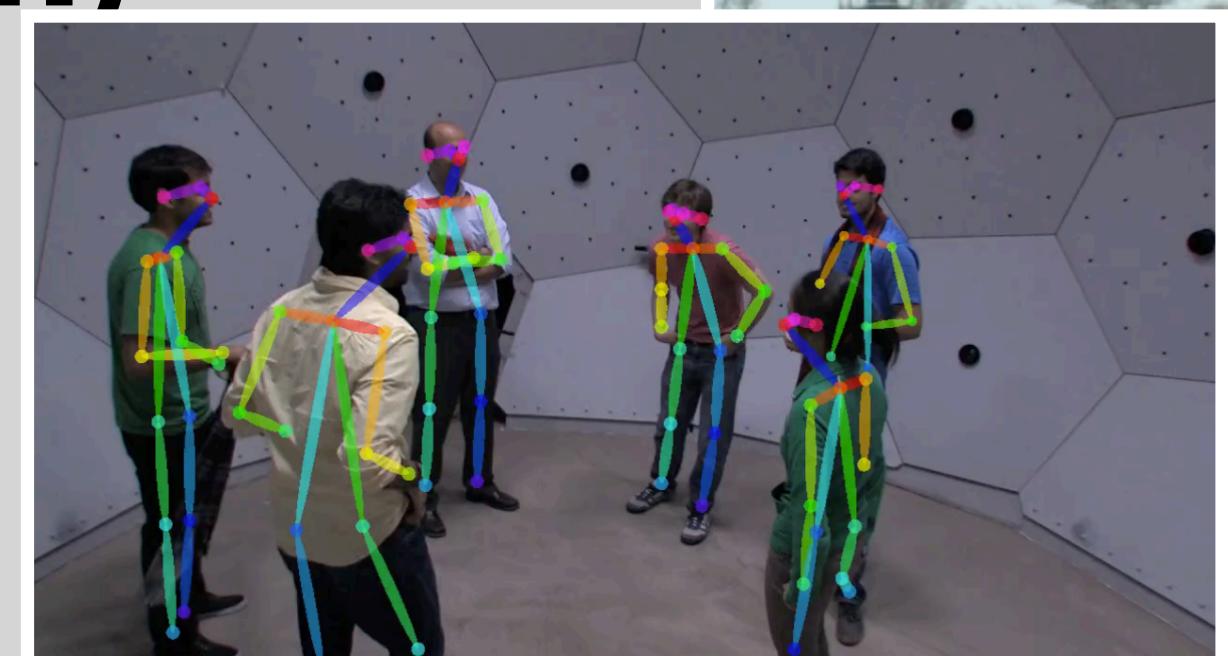
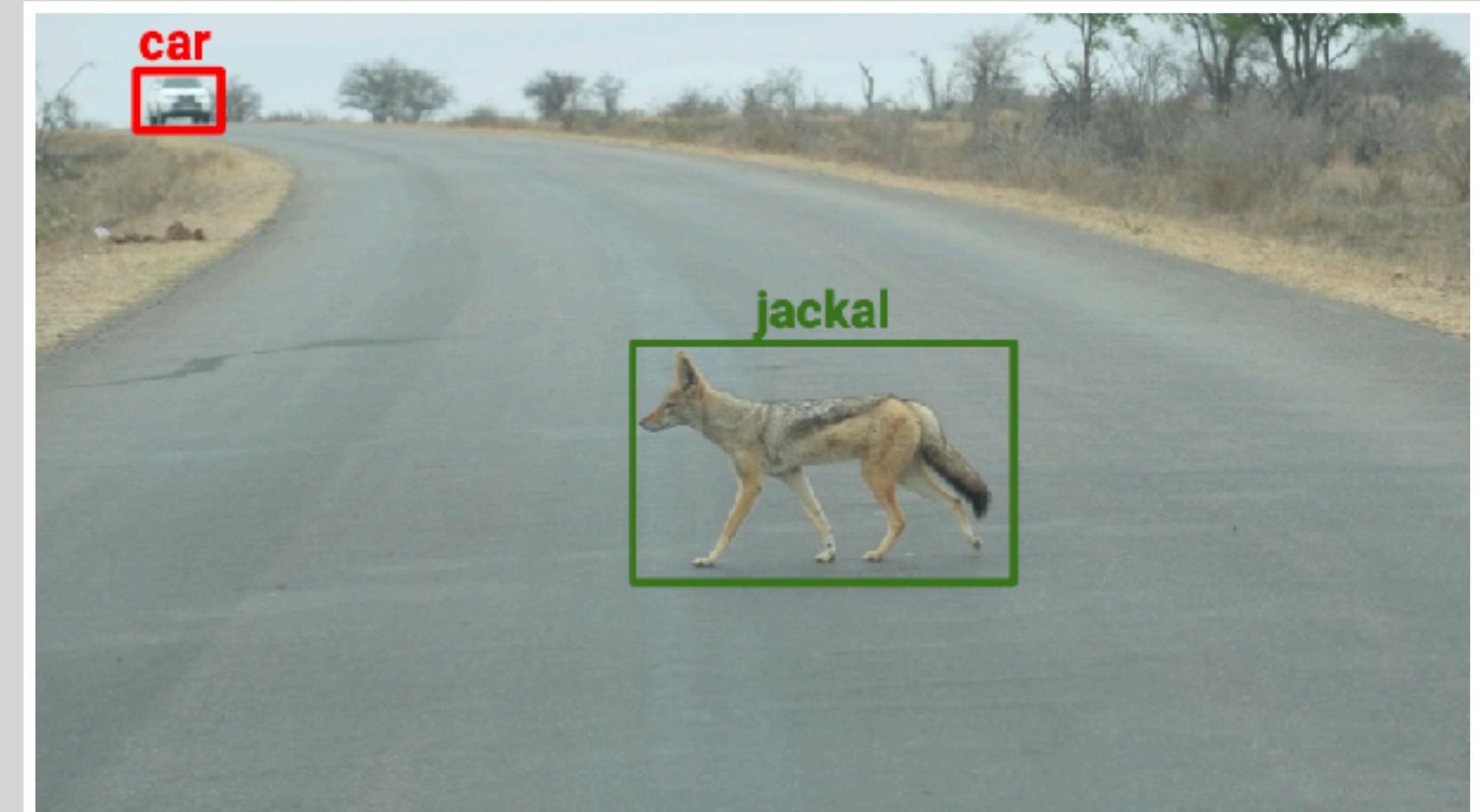
2. Beyond classification

- a. Intro to structured outputs

- b. Object detection (localization)

- c. Segmentation

- d. Human pose estimation



Key elements of DL for CV

Initially in CV

1 Model (i.e., architectural definition of connectivity and learnable parameters)

Next in CV

2 Loss

These days in CV !

3 Data

ML community

Optimization algorithm (i.e., variations of SGD)

Still many open questions

- 3D
- Videos
- Visual perception in robotics