

Introduction to Computer Vision: Visual perception and generation

Natalia Neverova

Course Schedule

< Today > +

Day Week Month

	SUN 07	MON 08	TUE 09	WED 10	THU 11	FRI 12	SAT 13
10:00 AM		10:00 am – 12:00 pm Week 2 - CV - Lecture 1	10:00 am – 12:00 pm Week 2 - CV - Lecture 2	10:00 am – 12:00 pm Week 2 - CV - Lecture 3	10:00 am – 12:00 pm Week 2 - CV - Lecture 4	10:00 am – 12:00 pm Week 2 - CV - Exam	
11:00 AM							
12:00 PM							
1:00 PM							
2:00 PM							
3:00 PM		2:30 pm – 4:30 pm Week 2 - CV - Lab 1	2:30 pm – 4:30 pm Week 2 - CV - Lab 2	2:30 pm – 4:30 pm Week 2 - CV - Lab 3	2:30 pm – 4:00 pm Week 2 - CV - Q&A		
4:00 PM							

Course Setup

- All course information is on CampusWire
- Labs: https://github.com/nrvr/ammi_cv_object_detection_labs
- Deadline for lab reports: **Sunday, May 14th** at midnight
- Project: Kaggle competition (4 weeks) on cassava disease recognition
- Grading: Labs + Exam + Project report (1 page)
- *Credits for materials of this course: Georgia Gkioxari, Ross Girshick, Justin Johnson, Aryan Jain*

Last week's recap

- Empirical risk minimization
- Stochastic gradient descent
- Multilayer neural networks & backprop
- Convolutional neural networks
- Important network building blocks and design motifs
- Image classification
- Vision transformers

This week

- Lectures 1-3: Perception (+ 3 labs)
 - Object detection – Tasks, evaluation & datasets
 - Object detection as an ML problem
 - Models and architectures for visual perception
- Lecture 4: Generative models for visual data

Introduction to Computer Vision: Object detection (part 1)

Natalia Neverova

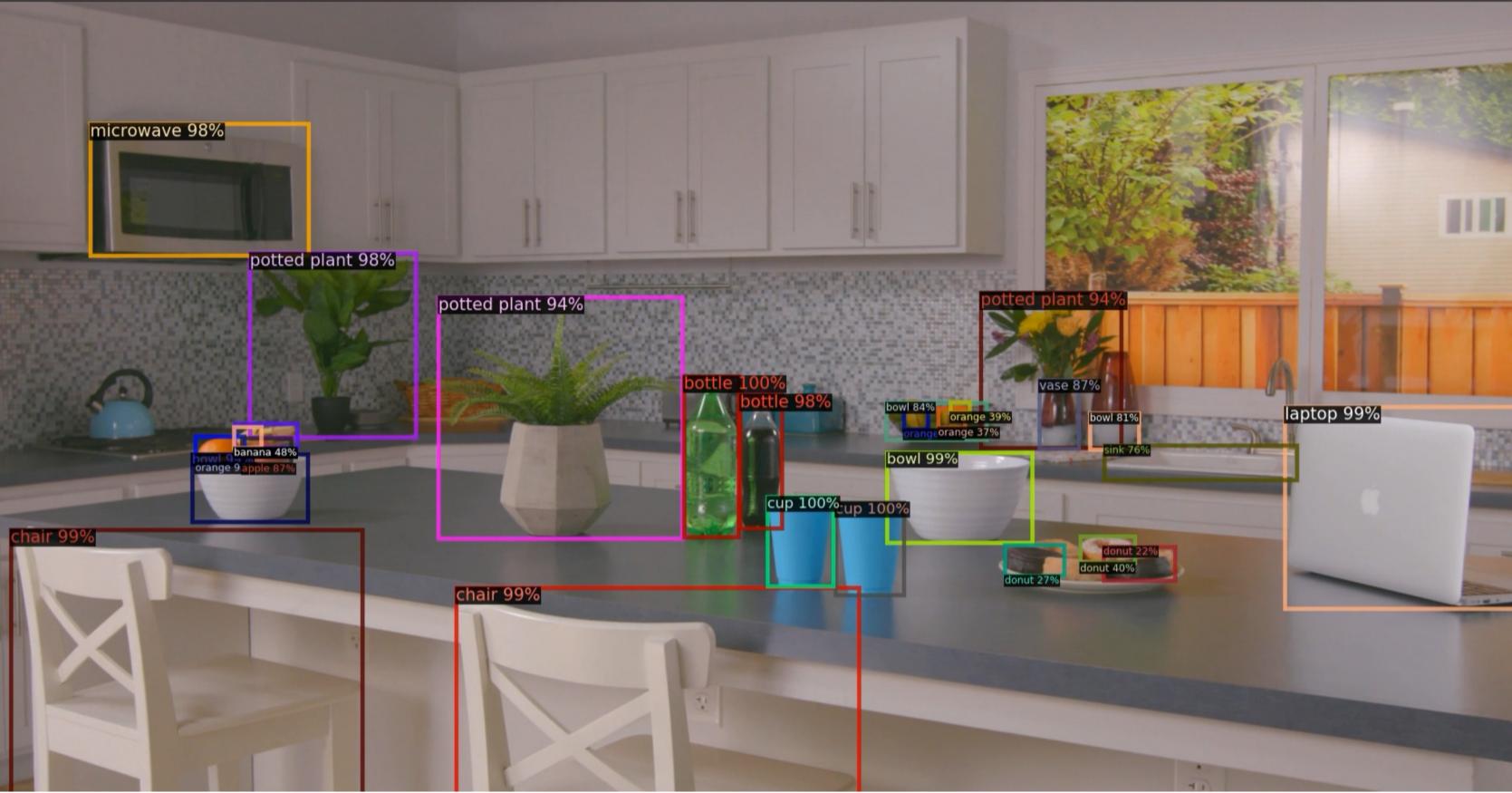
What we've seen so far: image classification



Label: kitchen

Visual perception – instance-level...

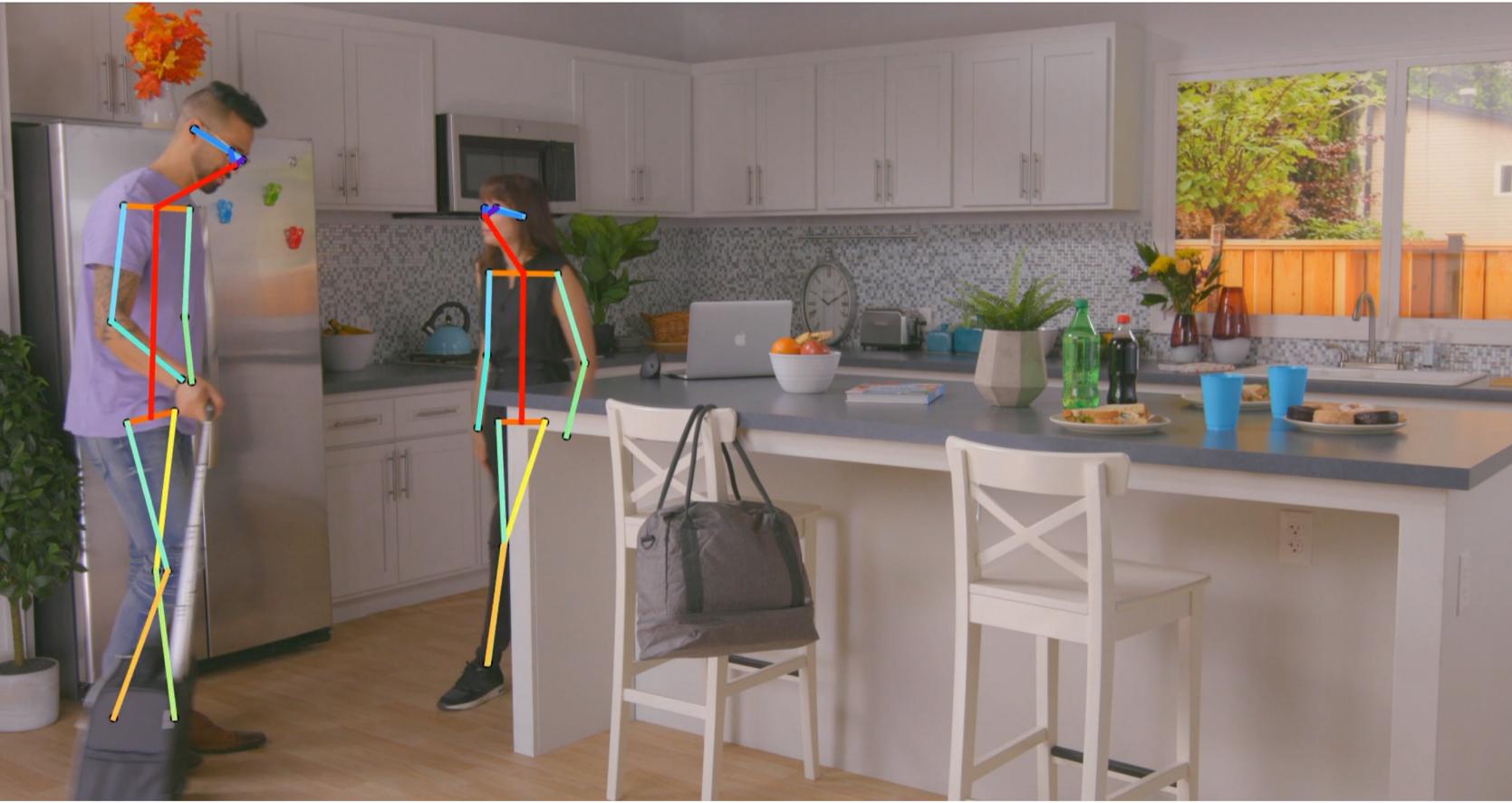
bounding box
detection



instance
segmentation



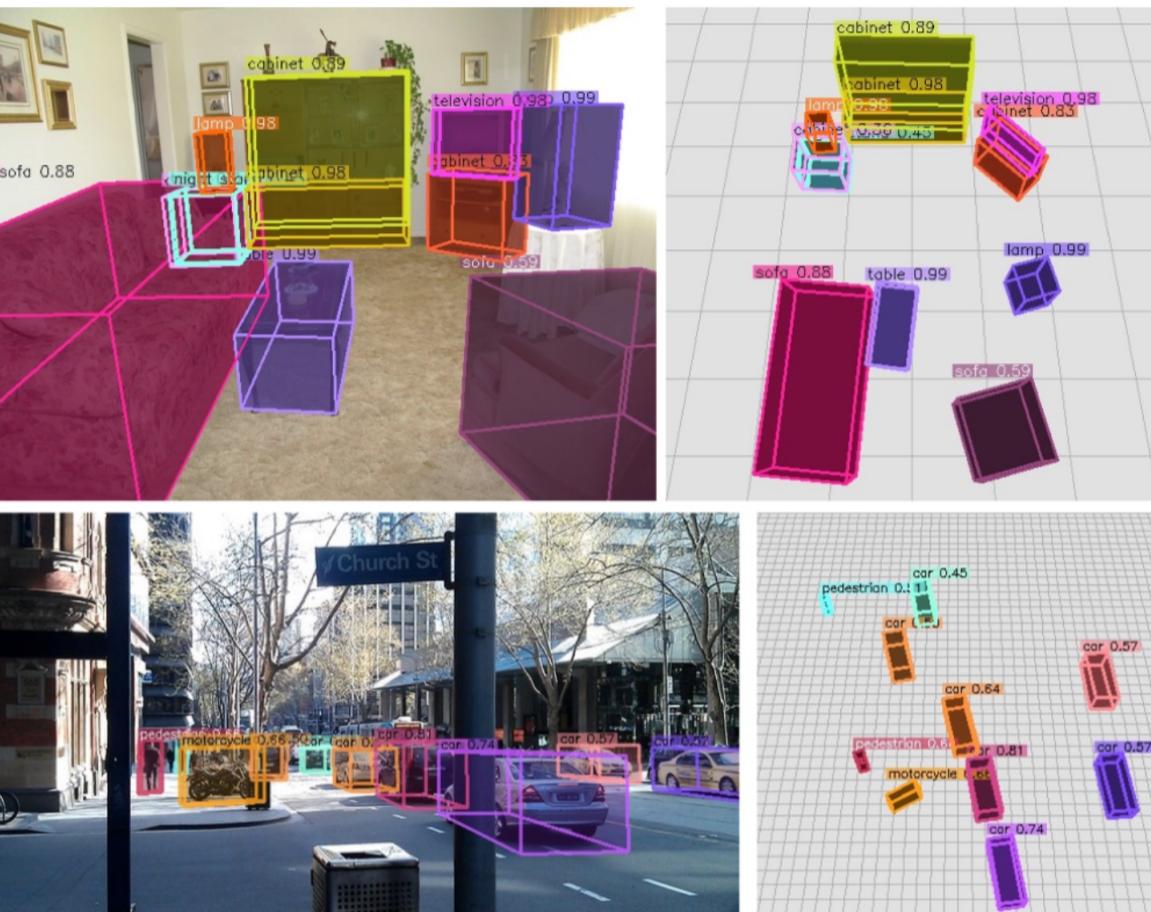
pose
estimation



dense
correspondences



Visual perception – instance-level...



3D detection



3D reconstruction

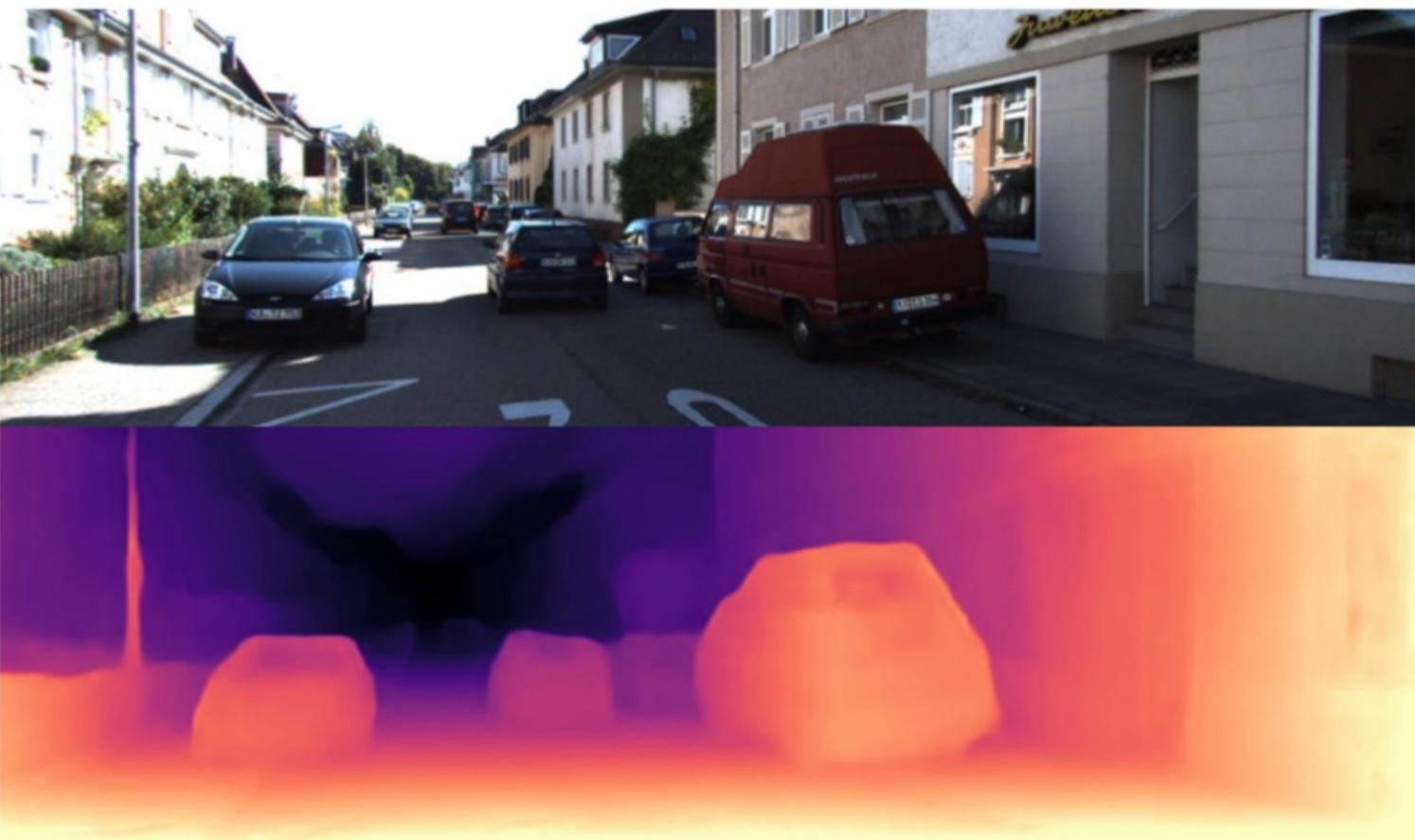


3D pose and
shape
estimation

Visual perception – or dense



semantic
segmentation

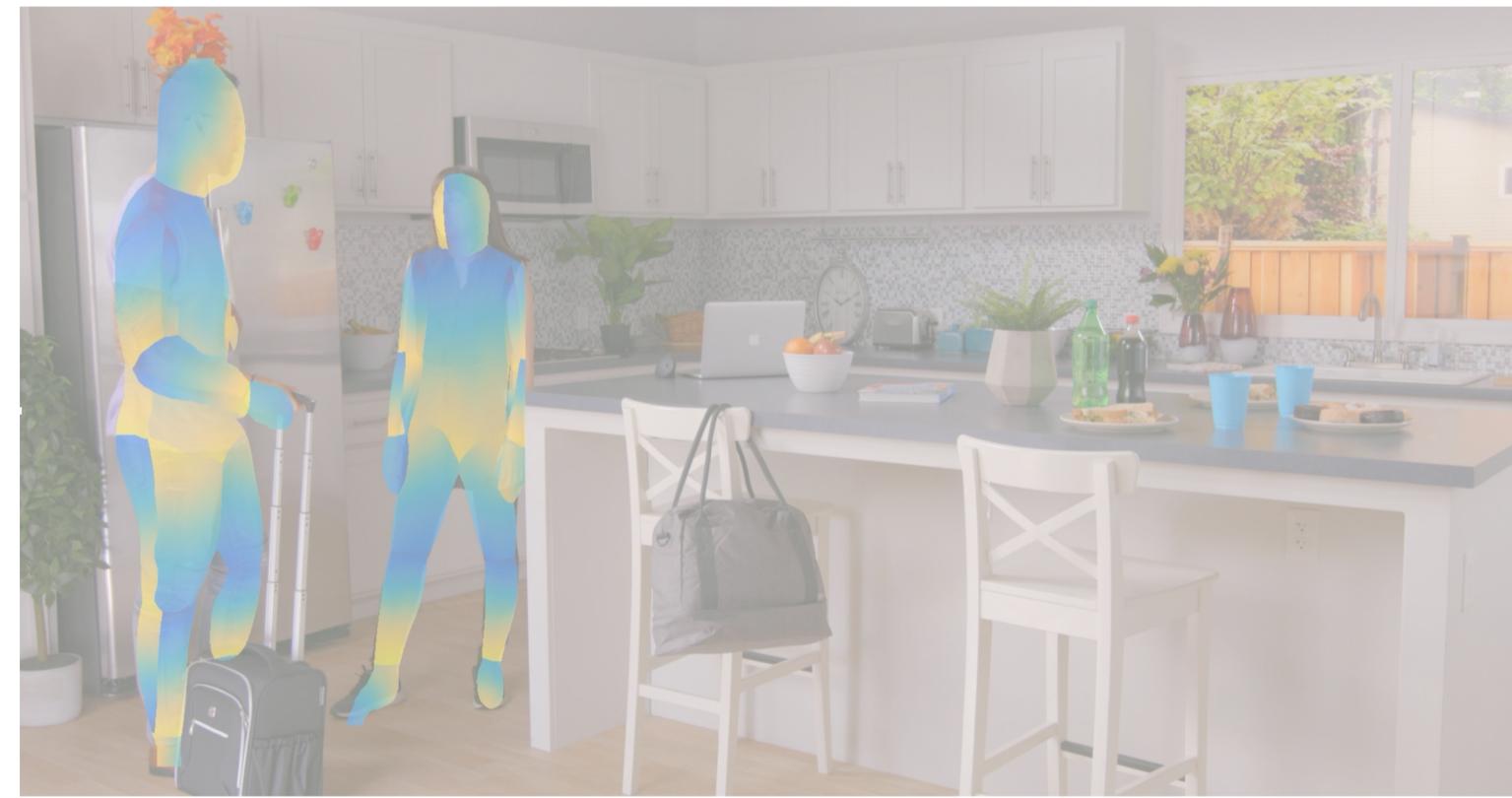
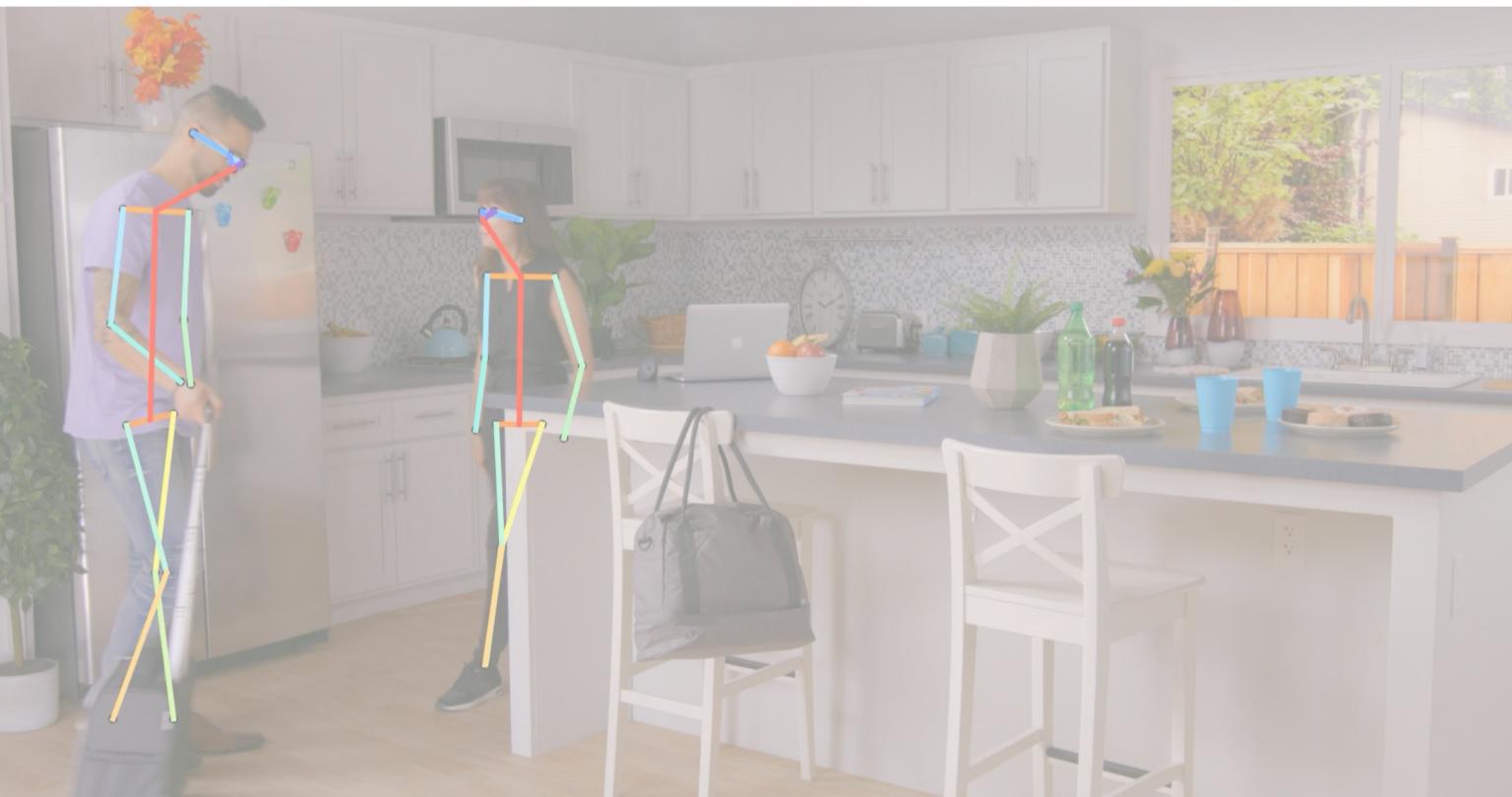


depth
estimation



optical flow
estimation

Today's focus



Today's focus

- What is the object detection task?
- How is this task formalized?
- How is this task evaluated?
- Why is the task difficult?
- What are important datasets to know about?



What is the object detection task?

- **What** objects are in the image and **where** are they?



What is the object detection task?

- **What** objects are in the image and **where** are they?

Answer: What and Where?

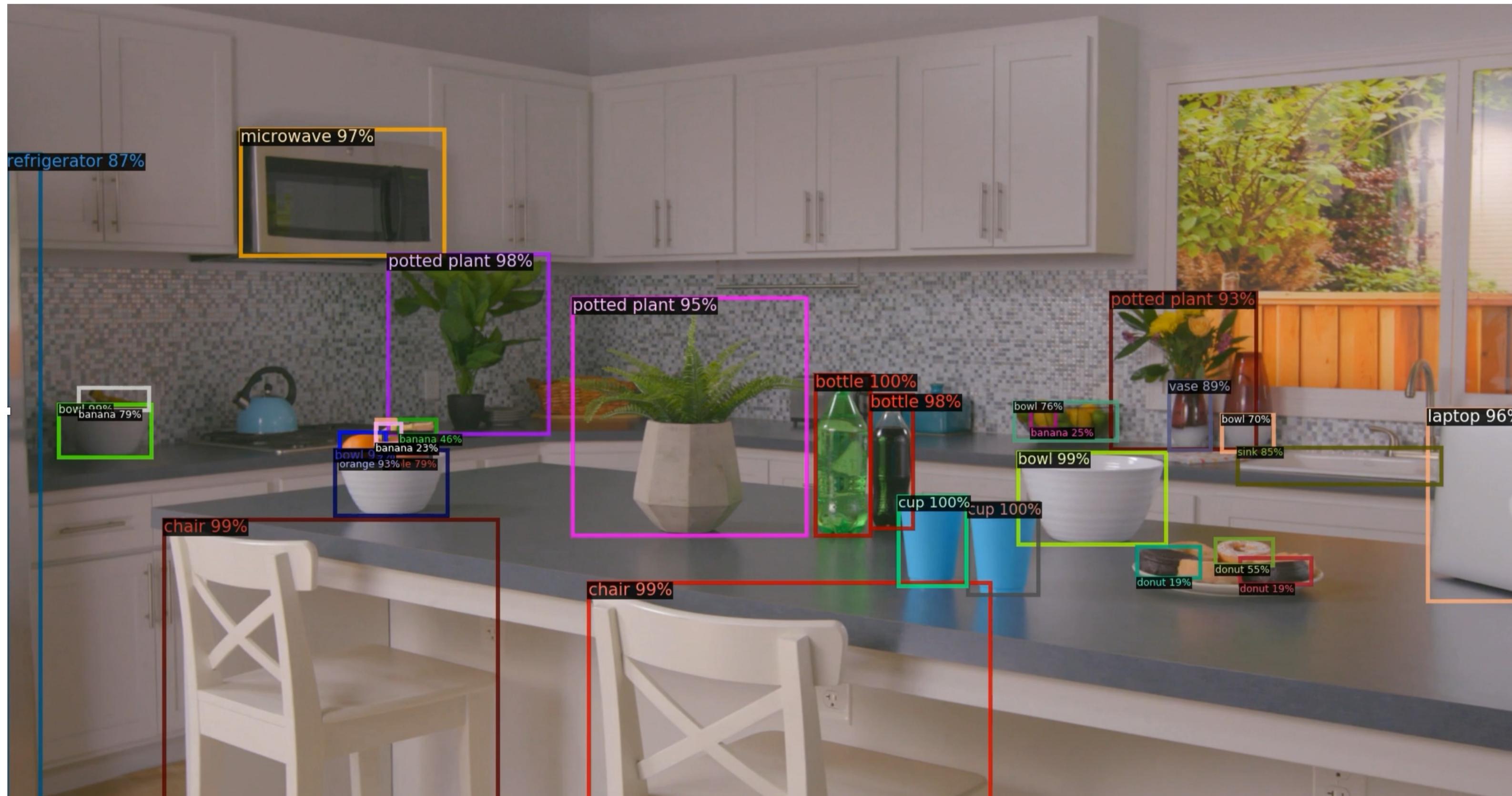


Task formulation (testing) - input



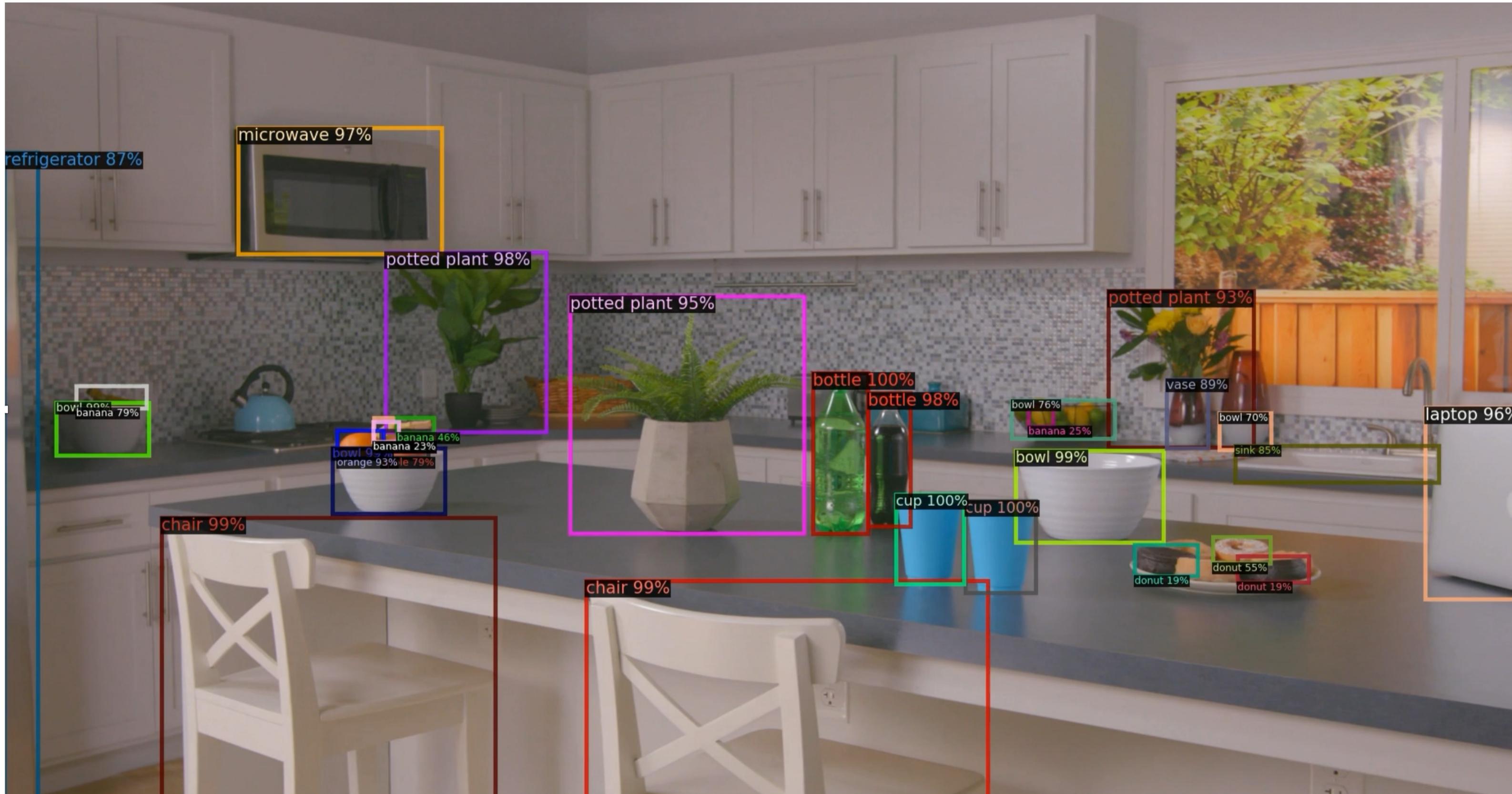
Input: an image (not seen during training)

Task formulation (testing) - output



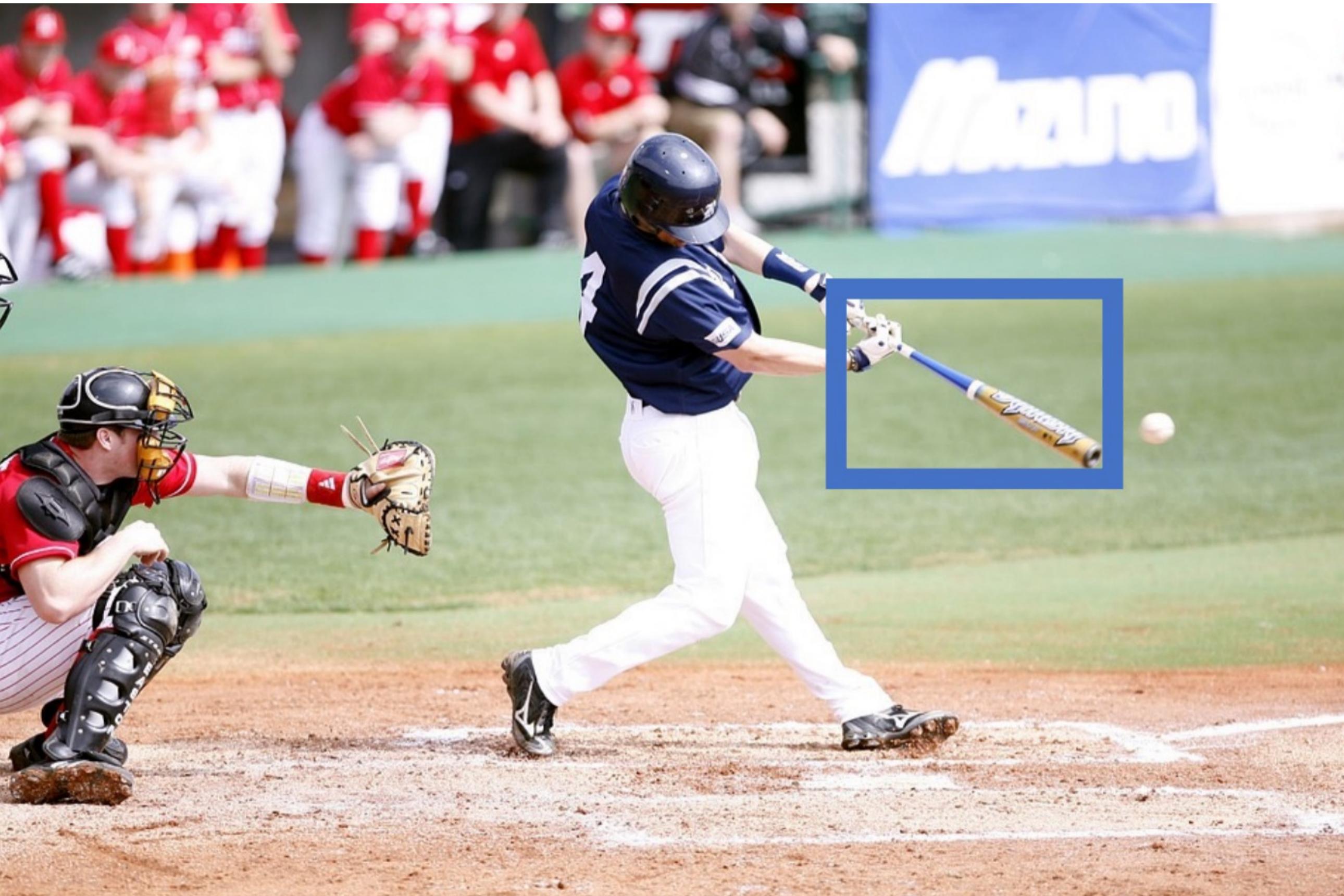
Output: boxes, categories and confidence scores

Where means bounding boxes



Output: boxes, categories and confidence scores

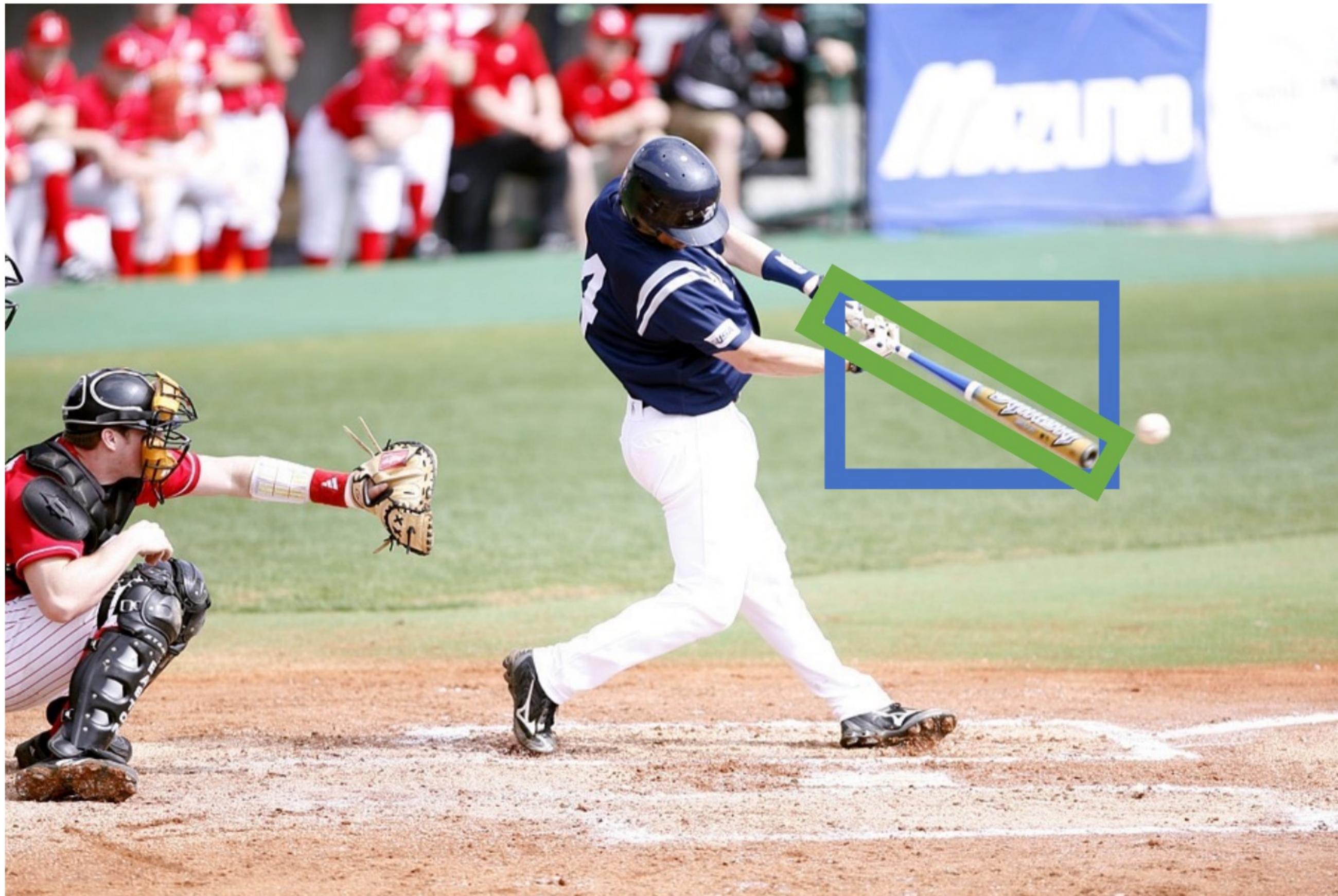
Where means bounding boxes



Bounding boxes are typically axis-aligned



Where means bounding boxes



Oriented boxes are much less common

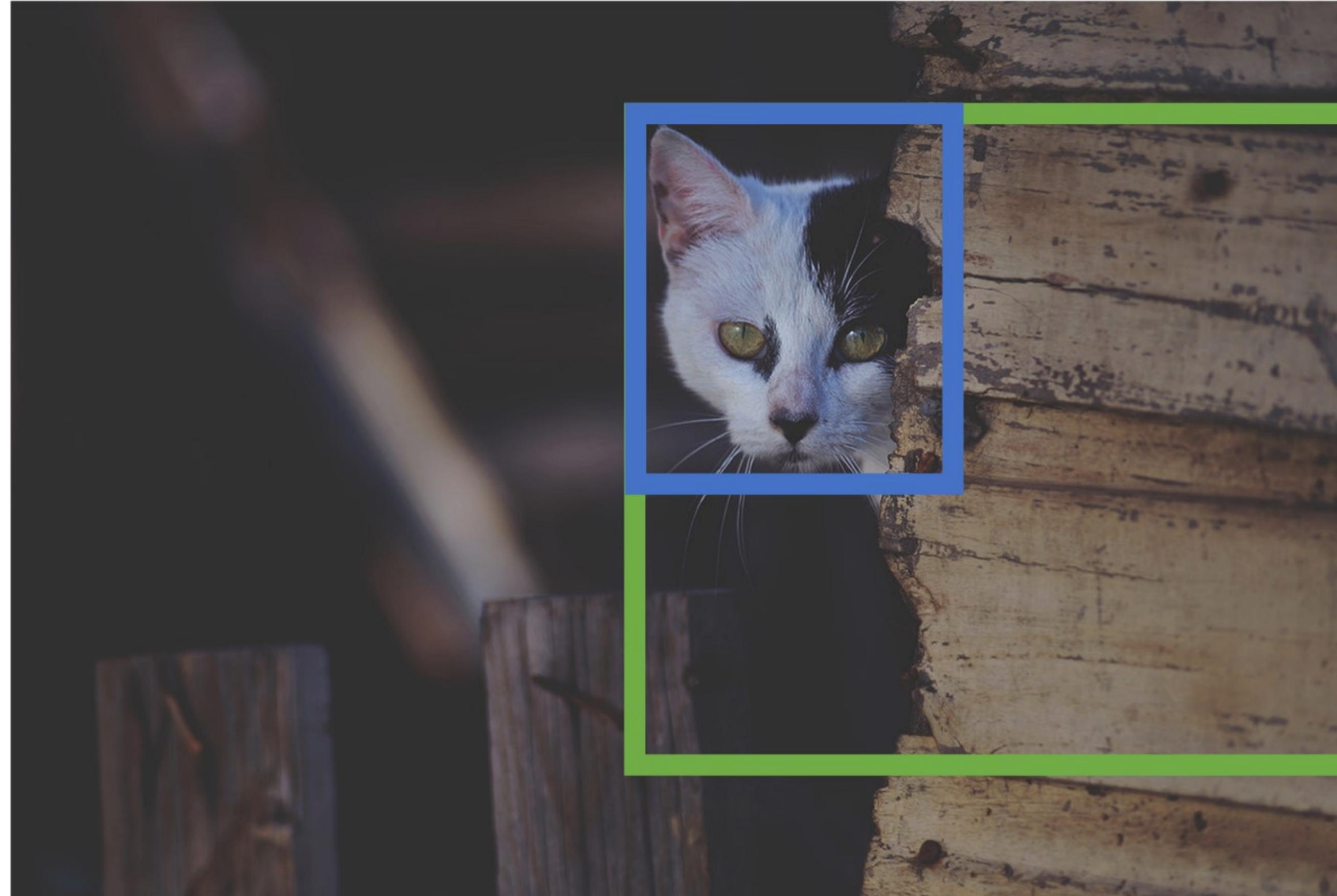


Where means bounding boxes



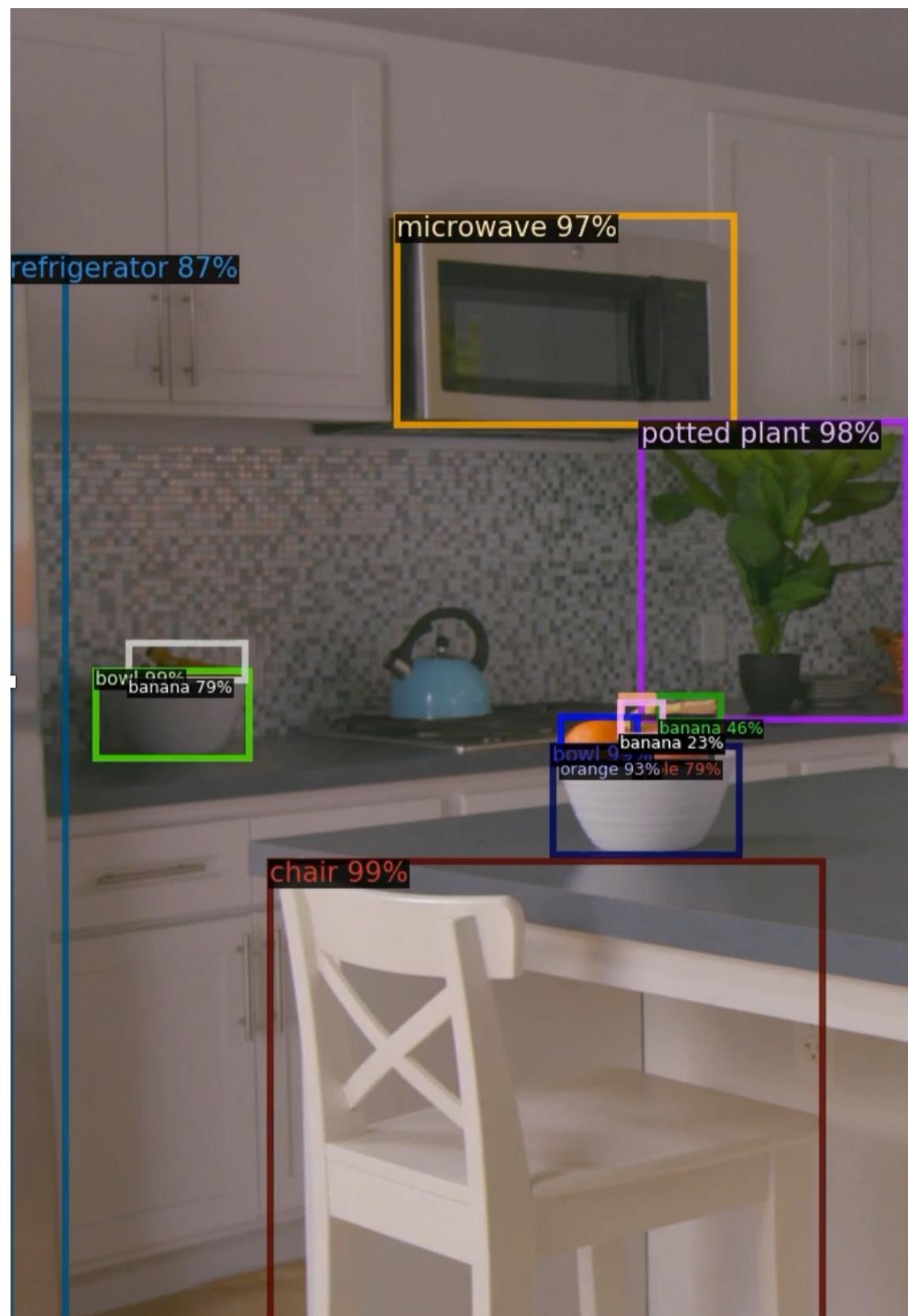
Bounding boxes (usually) cover only the visible portion of the object

Where means bounding boxes



Amodal detection: box covers the entire extent of the object, even occluded parts
(not considered in this course)

Task formulation (testing) - output



Input: an RGB image

Output:

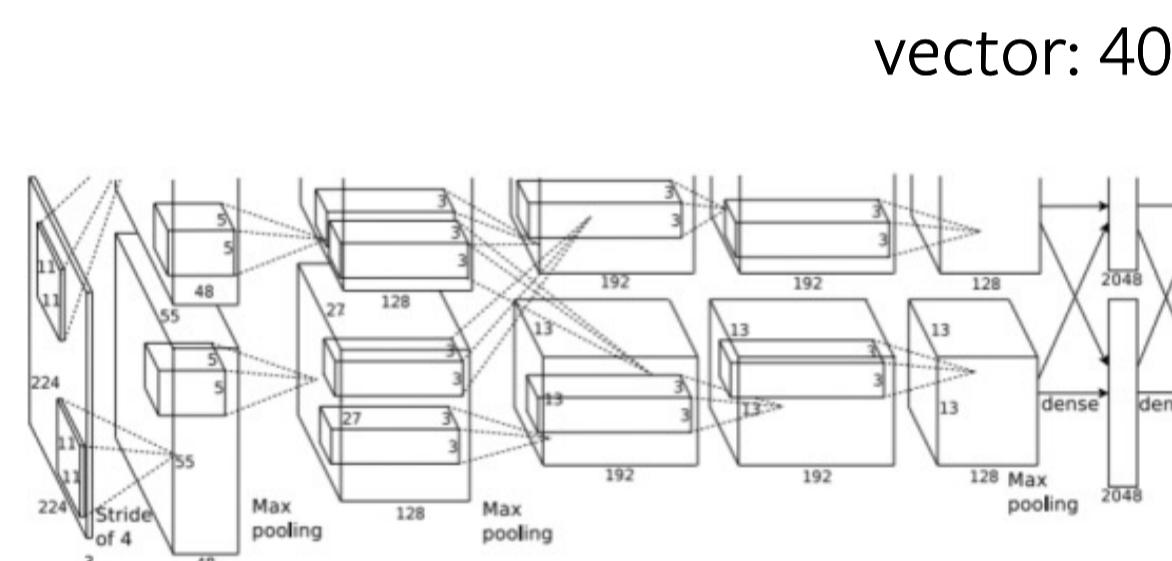
- **set of boxes** – each test image may have a different number of boxes;
- **category (class) per box** – comes from a pre-defined set of object categories (given by the dataset);
- **confidence score per box** – a number where higher indicates that the model “thinks” that the output is more likely correct

Object detection: challenges



- **Multiple outputs:** need to output variable numbers of objects per image;
- **Multiple types of output:** need to predict “what” (category label) as well as “where” (bounding box);
- **Large images:** classification works at 224x224; need higher resolution for detection, often 800x600

So far: image classification



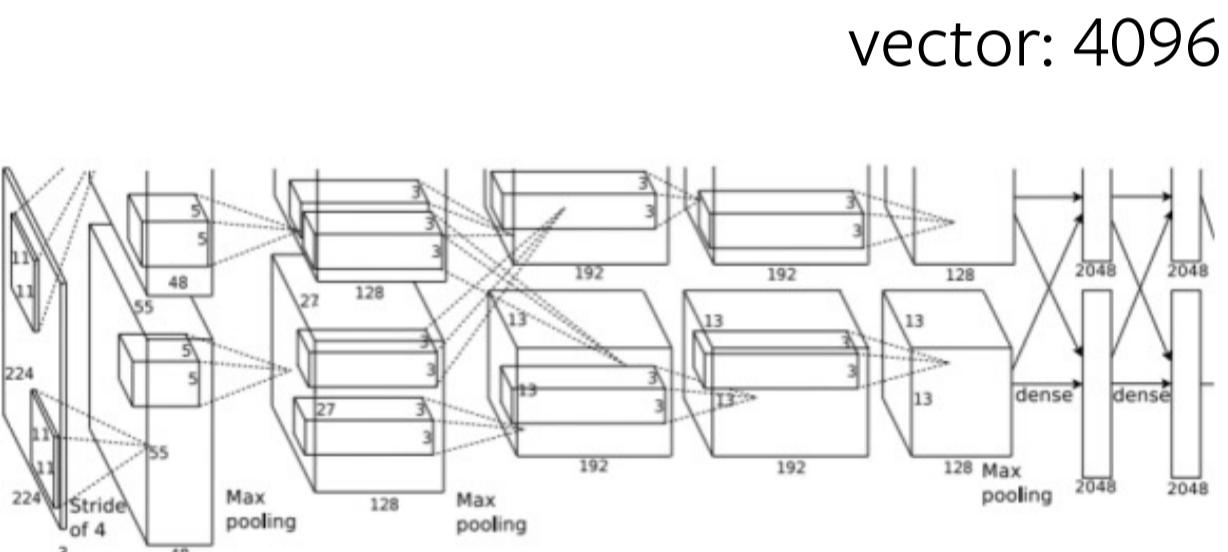
vector: 4096

fully connected:
4096 to 1000

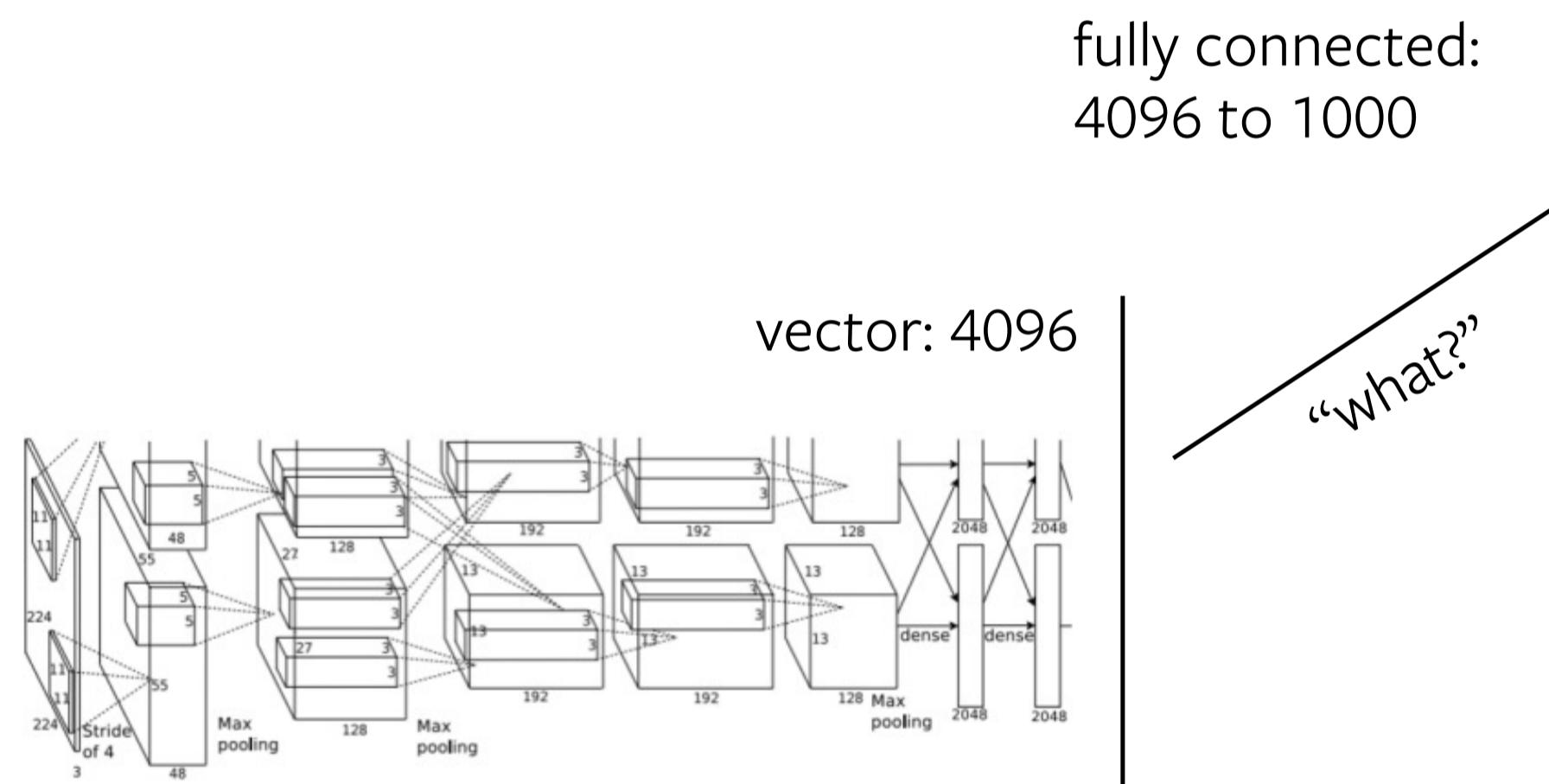
class scores:
cat: 0.9
dog: 0.05 →
car: 0.01
...

correct label:
cat
↓
softmax loss

Detecting a single object



Detecting a single object



fully connected:
4096 to 1000

vector: 4096

“what?”

class scores:
cat: 0.9
dog: 0.05 →
car: 0.01
...

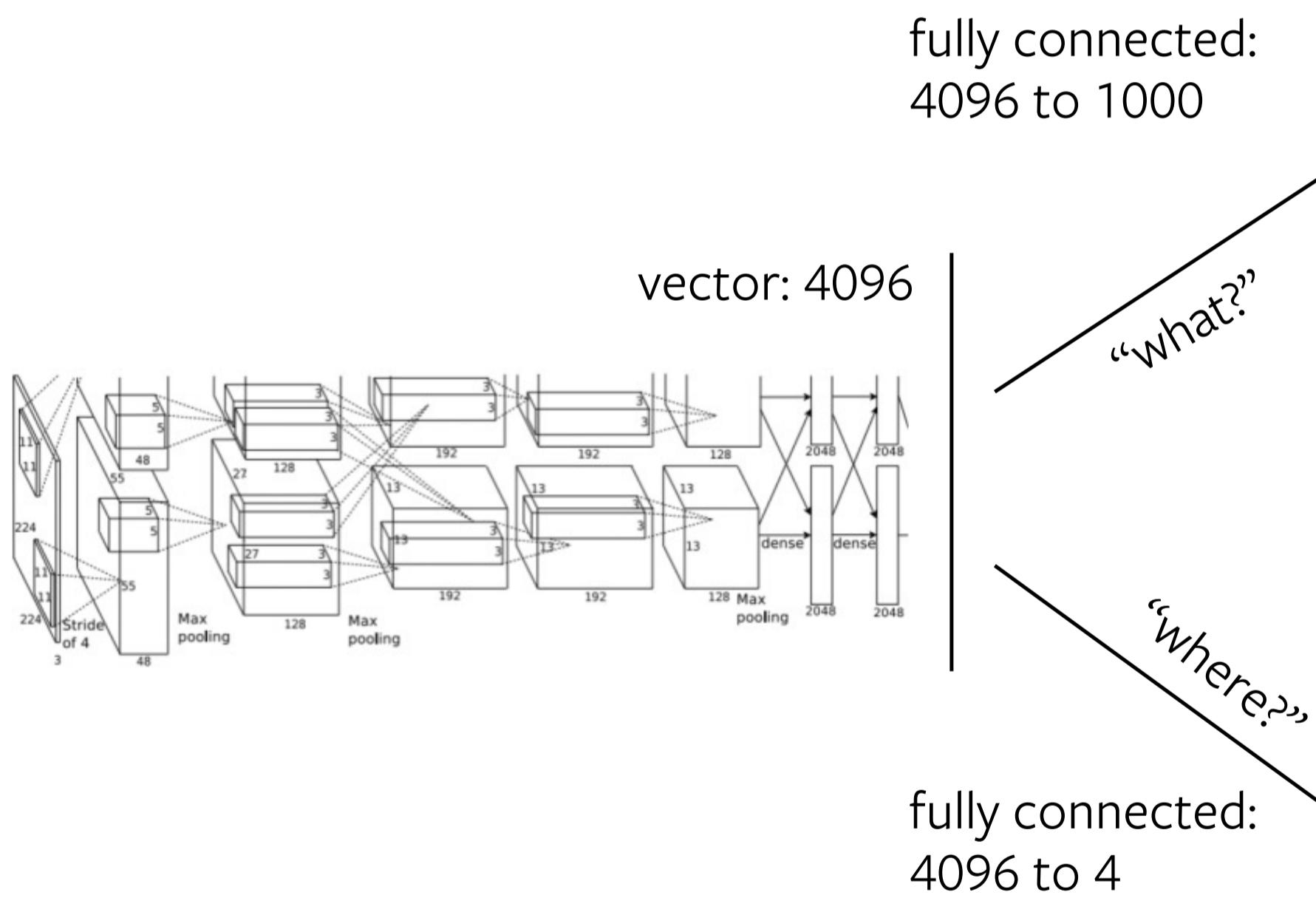
correct label:
cat
↓
softmax loss

treat localization as a regression problem!

Detecting a single object



treat localization as a regression problem!



class scores:
cat: 0.9
dog: 0.05 →
car: 0.01
...

box coordinates:
(x,y,w,h) →

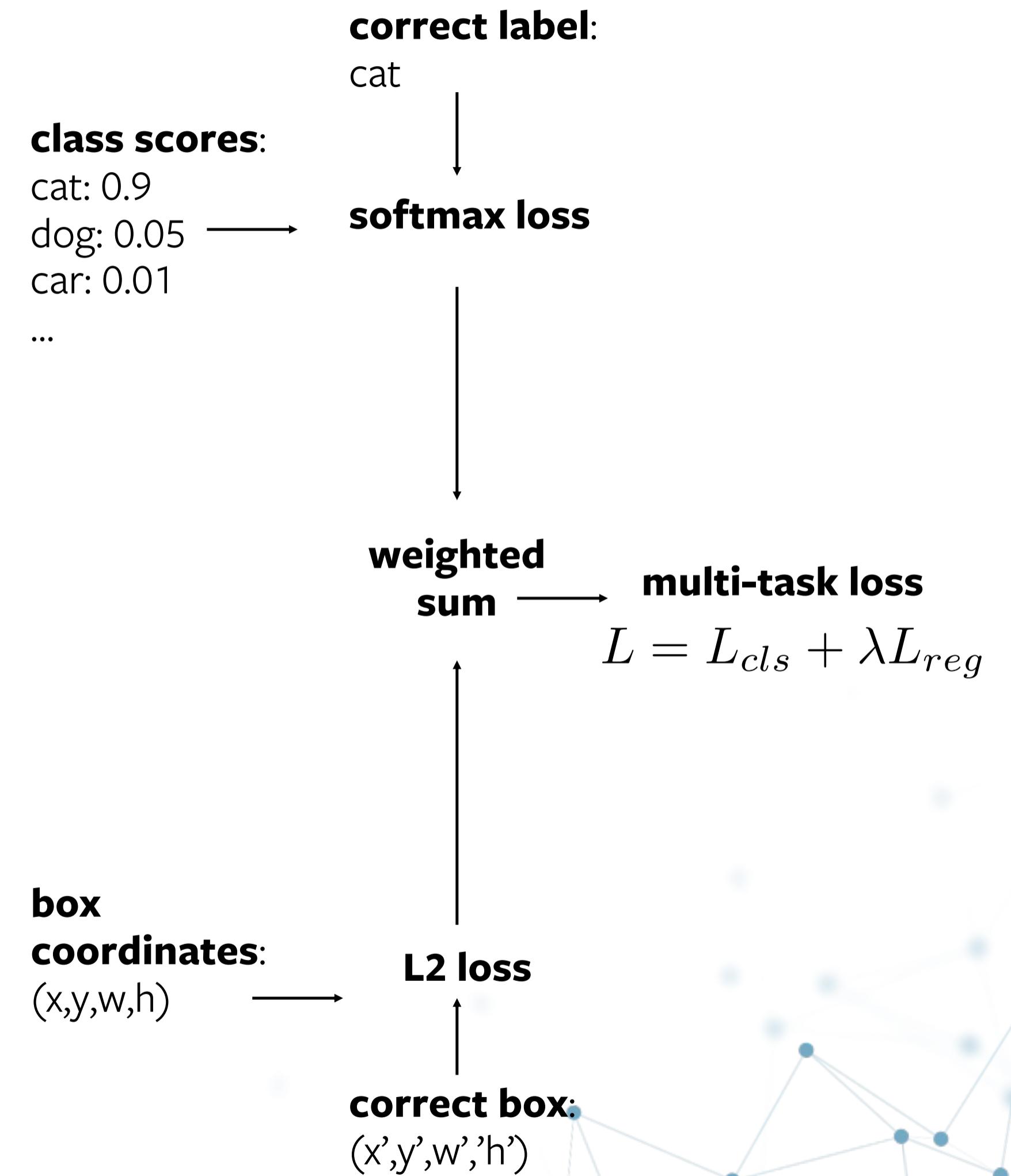
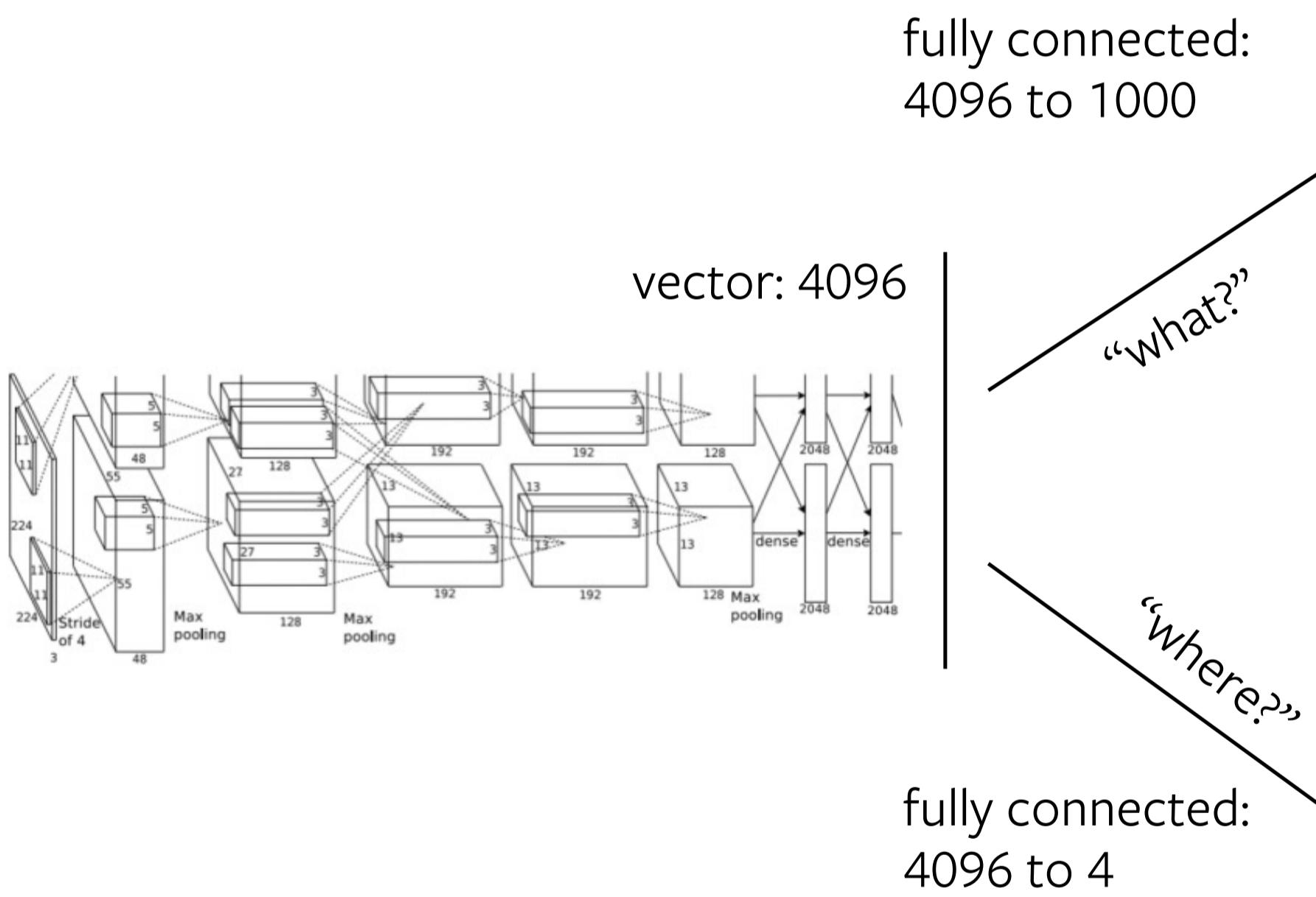
L2 loss
correct box:
(x',y',w',h')

correct label:
cat
↓
softmax loss

Detecting a single object



treat localization as a regression problem!

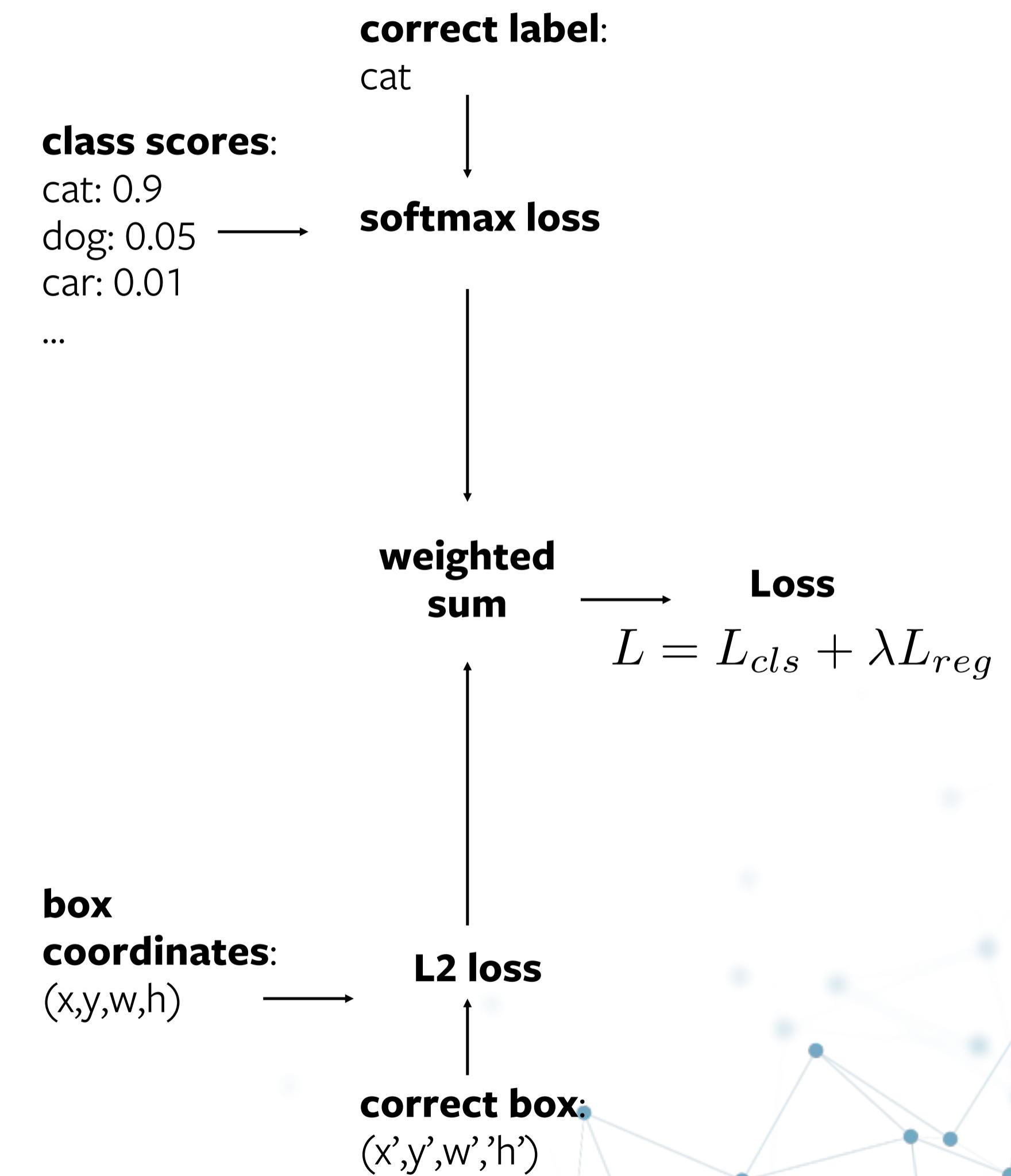
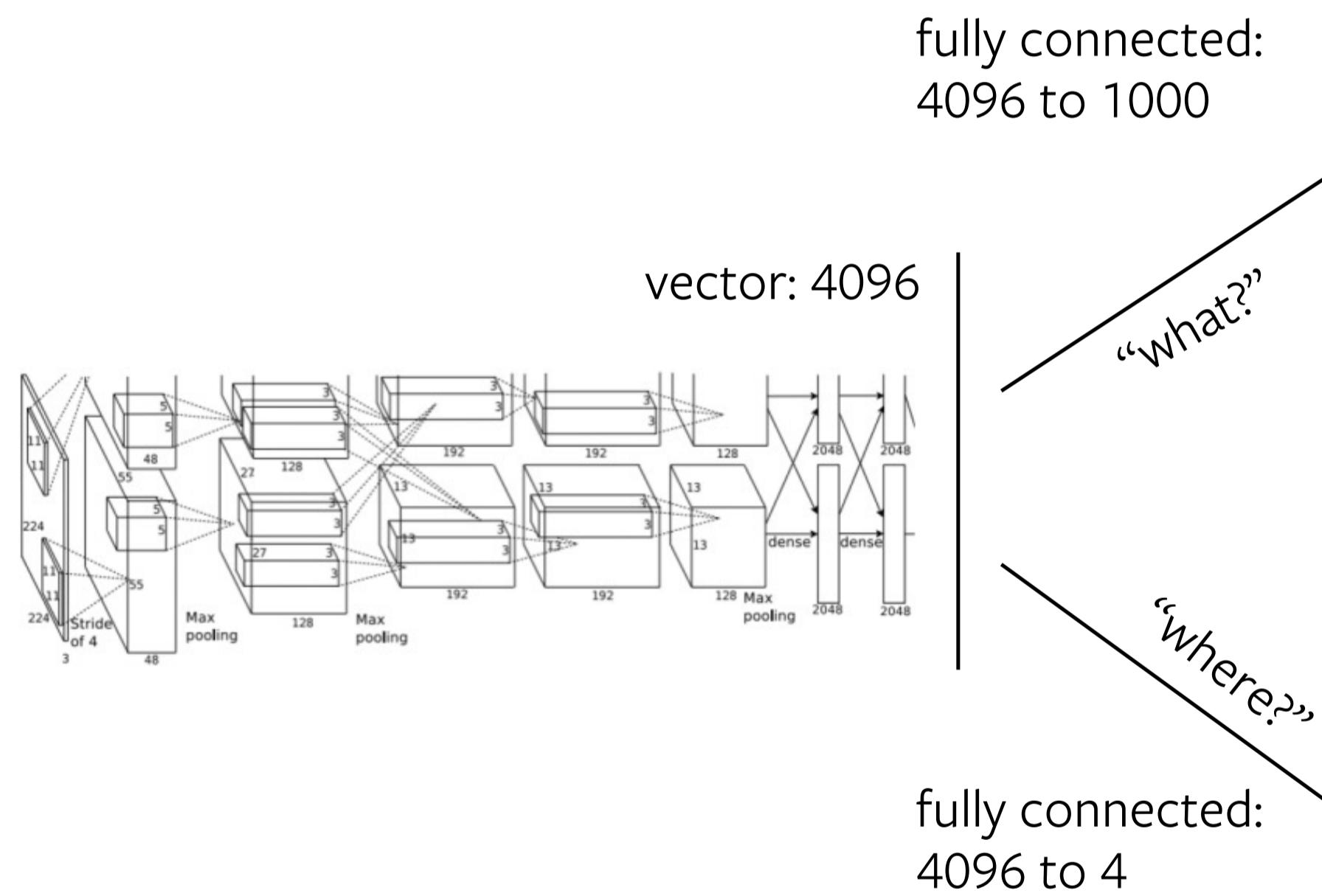


Detecting a single object



treat localization as a regression problem!

**things get much more complicated with multiple objects!
(next lecture)**



Recap: Task formulation (testing)

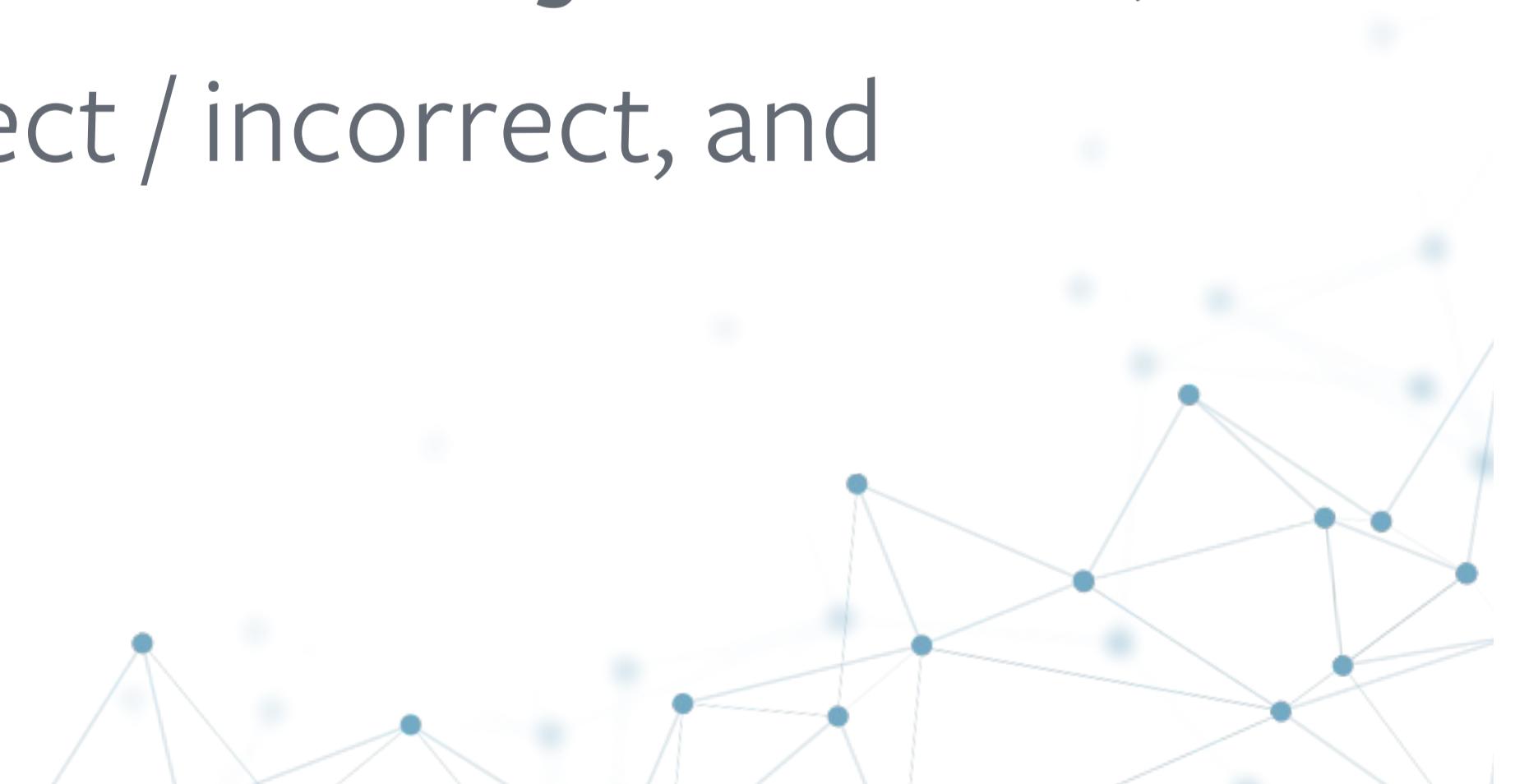
The test-time task formulation is deceptively simple:

- image in – boxes, categories and scores out. That's it!
- that's actually a hard problem.



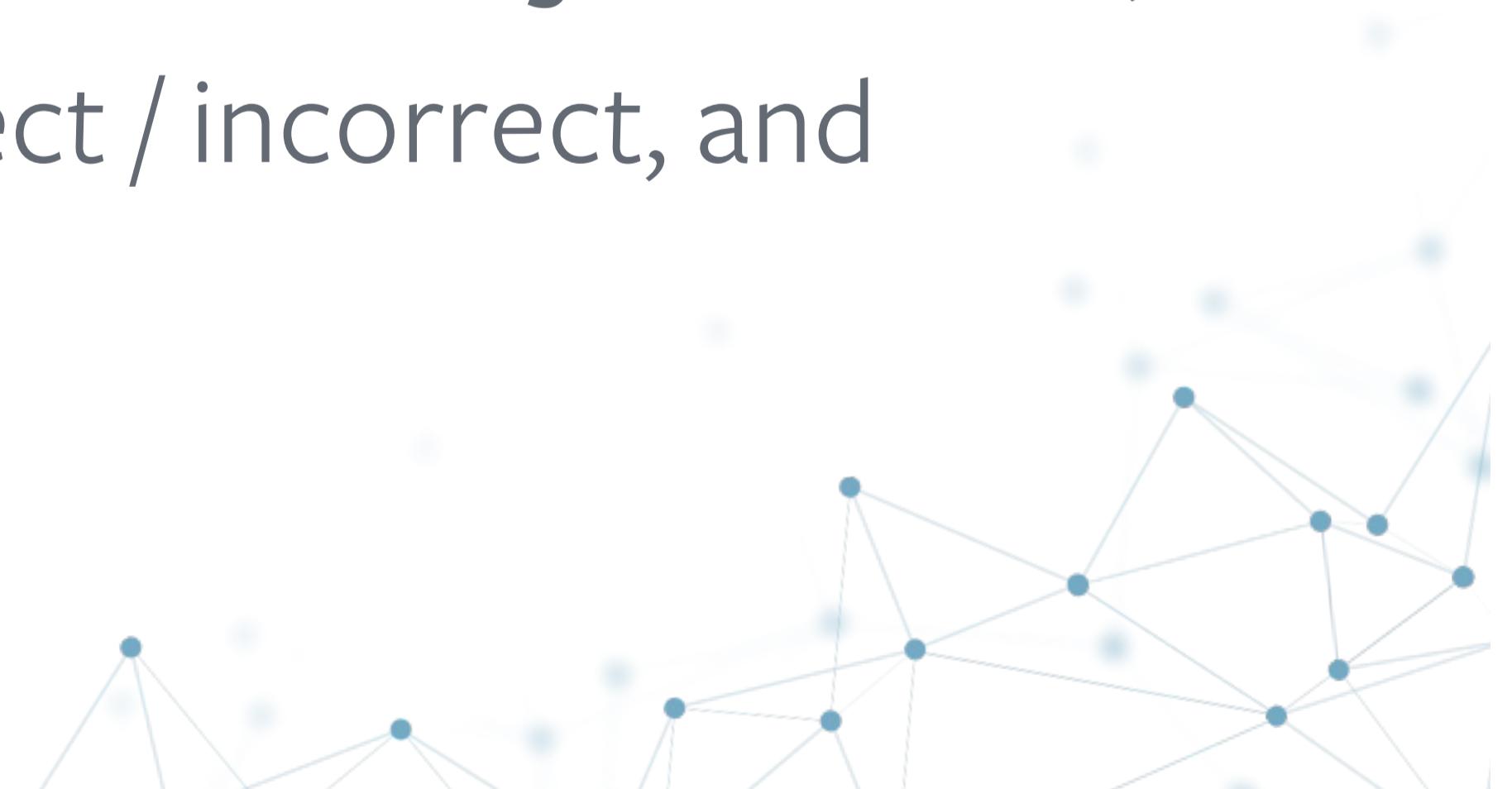
Evaluating an object detector

- Given the inputs, **how good are the outputs?**
- This problem is more nuanced than image classification:
need to evaluate both **what and where**;
ground truth and prediction are almost **never exactly the same**;
key questions: when is a prediction correct / incorrect, and
what are we missing?



Evaluating an object detector

- Given the inputs, **how good are the outputs?**
- This problem is more nuanced than image classification:
need to evaluate both **what and where**;
ground truth and prediction are almost **never exactly the same**;
key questions: when is a prediction correct / incorrect, and
what are we missing?



Motivational example



Input (1 out of N, typically O(1000), in the test set)

Motivational example



Kite, conf score = 0.97 (prediction)

Output of the model

Motivational example

Kite, conf score = 0.97 (prediction)



Output and ground truth

Motivational example

Kite, conf score = 0.97 (prediction)



Is this output good? Is it correct? Is it incorrect?

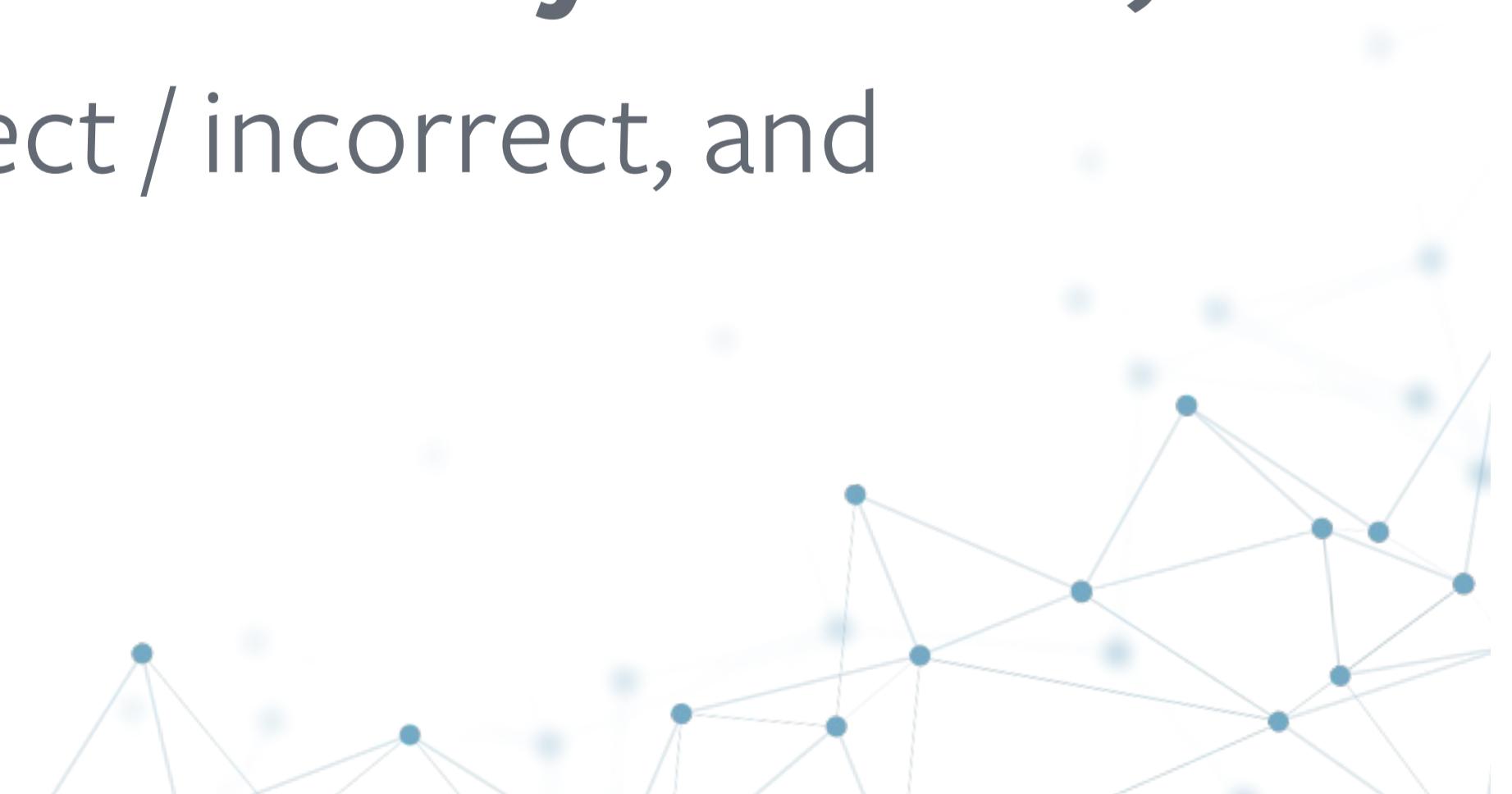
How should a good detector behave?

- Perfect localization is desired by not necessary for most applications
close is good enough;
- Names the object correctly:
don't call a cat, dog;
- Does not find an object more than once:
don't duplicate detections;
- Finds all instances of each object category:
don't miss anything.



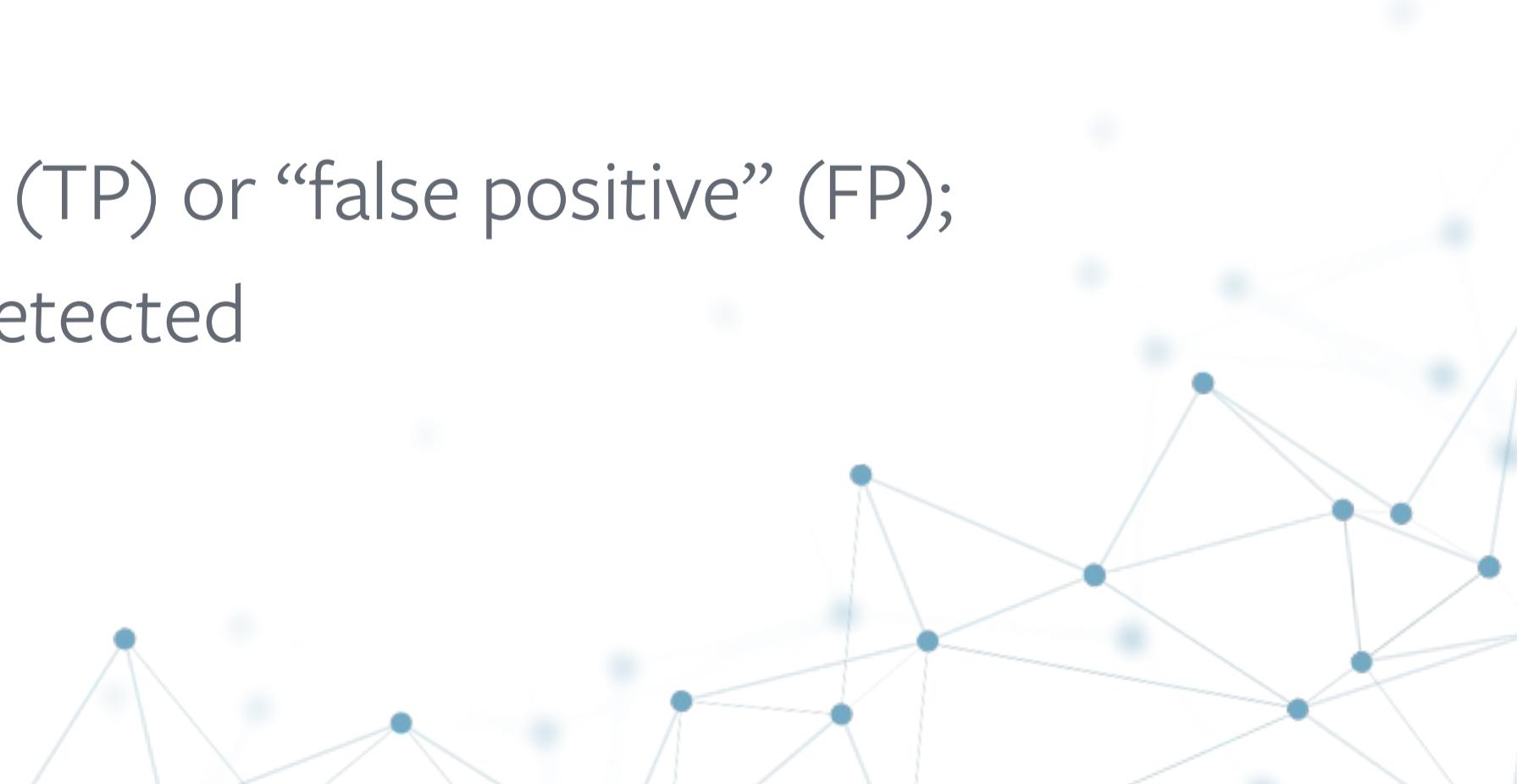
Evaluating an object detector

- Given the inputs, **how good are the outputs?**
- This problem is more nuanced than image classification:
need to evaluate both **what and where**;
ground truth and prediction are almost **never exactly the same**;
key questions: when is a prediction correct / incorrect, and
what are we missing?



Evaluation code

- The evaluation code has access to:
 - ground-truth annotations (boxes and classes) – privileged information;
 - predictions from the model (boxes, classes, scores)
- The code evaluates each category independently
- In particular, the evaluation code:
 - sorts the predictions from highest to lowest confidence;
 - processes each prediction one at a time in this order;
 - assigns a binary classification to each prediction: either “true positive” (TP) or “false positive” (FP);
 - counts the number of false negatives (FN), i.e. objects that were not detected



Terminology

- True positive (TP)

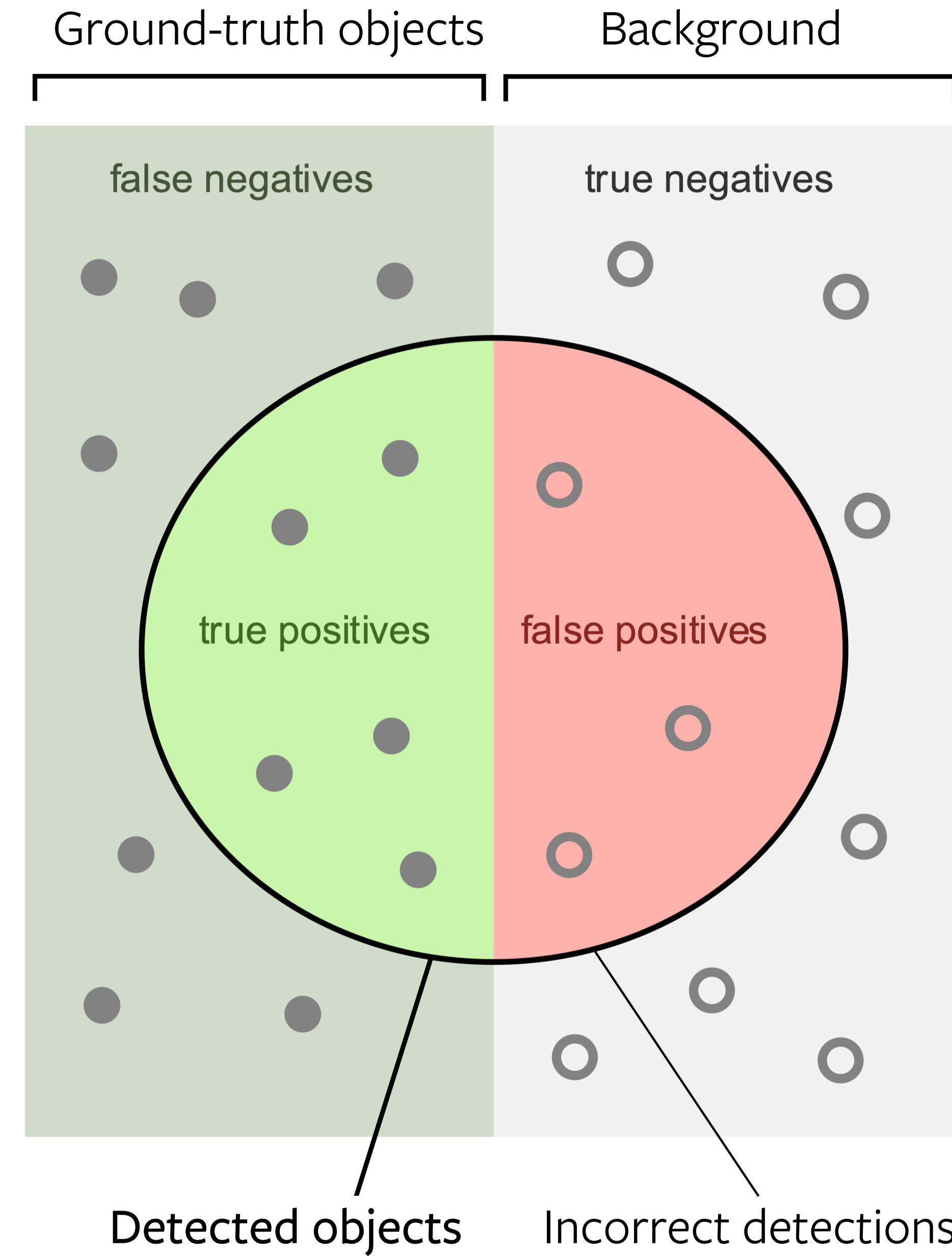
a prediction that is correct

- False positive (FP)

a prediction that is incorrect

- False negative (FN)

a ground truth object that the detector missed,
i.e. it is not in the predictions



Criteria that define a true positive (TP)

Prediction P is a true positive if there exists a ground truth G such that:

- (1) P 's box has **sufficient similarity** to G 's box
- (2) P and G have the **same category label**
- (3) G **was not matched** to a higher scoring prediction P'

If (1-3) are satisfied, then we say P and G are **matched**

(Subtle note: Ranking by score is used to define matches; other strategies, such as optimal bipartite matching are possible, but rarely used in practice.)



Defining similarity between two boxes

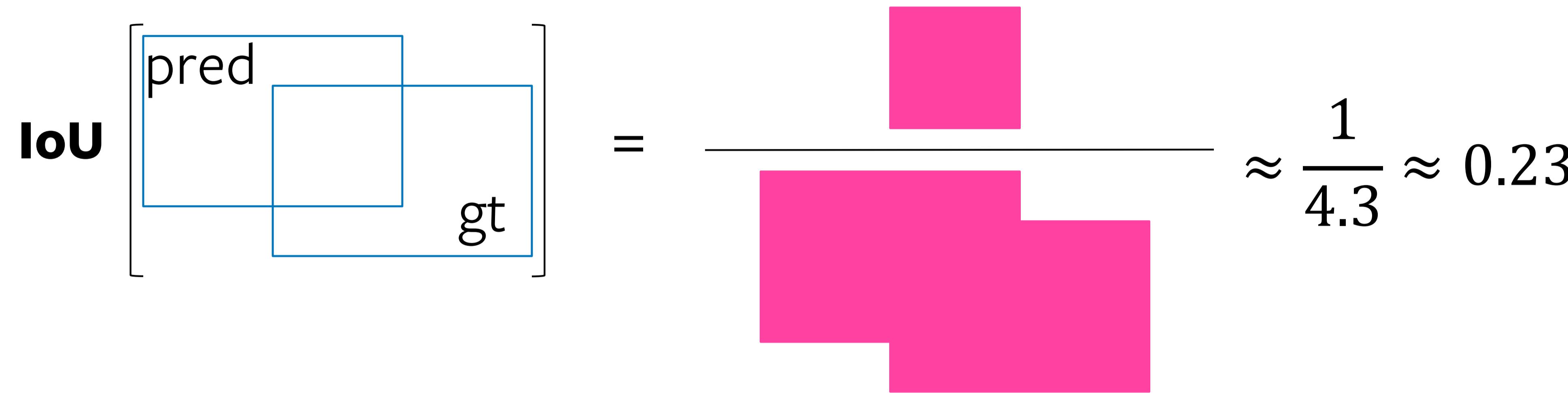
$$\text{IoU} = \frac{\text{Intersection}}{\text{Union}} \approx \frac{1}{4.3} \approx 0.23$$

The diagram illustrates the calculation of Intersection over Union (IoU) for two boxes. On the left, a blue bracket labeled "pred" contains a white rectangle representing the predicted bounding box. To its right, another blue bracket labeled "gt" contains a white rectangle representing the ground truth bounding box. The intersection of these two boxes is highlighted in pink. This intersection is divided by the union of the two boxes, which is represented by a larger pink rectangle covering both the predicted and ground truth boxes. The resulting ratio is approximately 0.23.

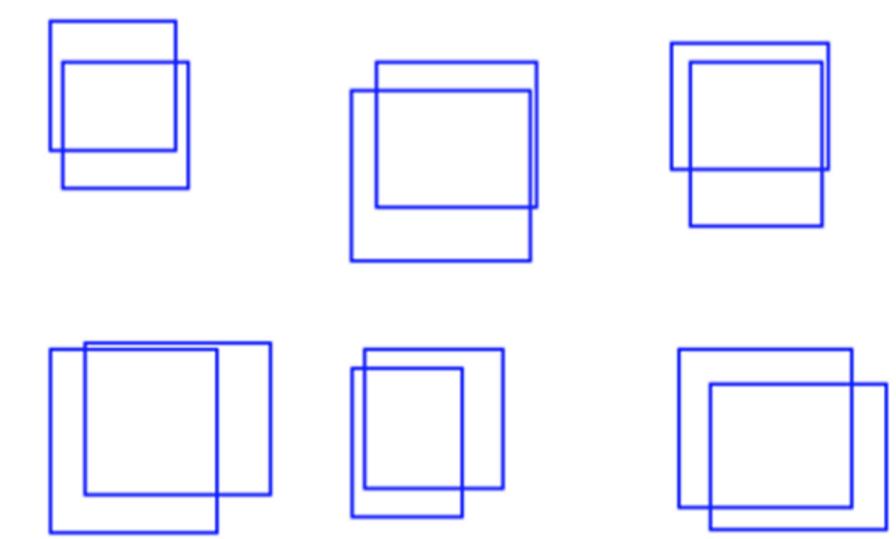
This definition of similarity is called **Intersection over Union (IoU)**

It's also called Jaccard Similarity

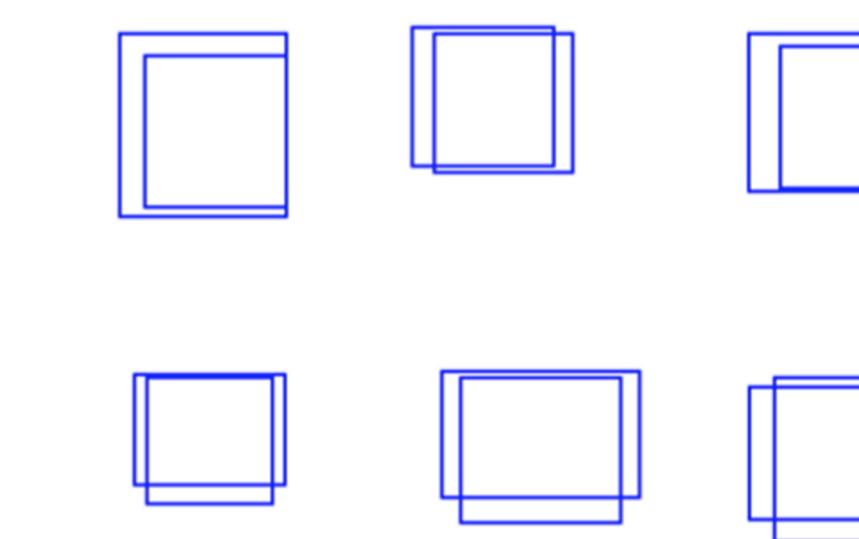
Defining similarity between two boxes



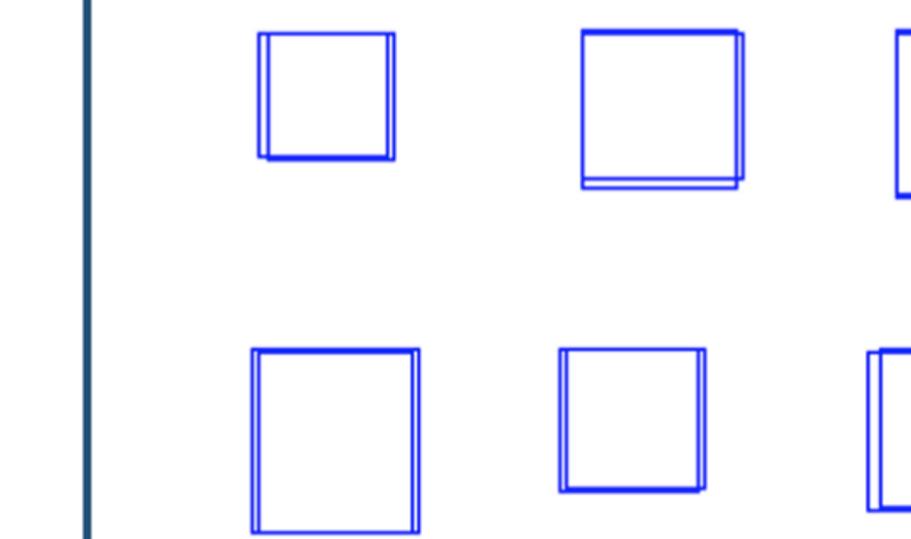
Random examples for your intuition:



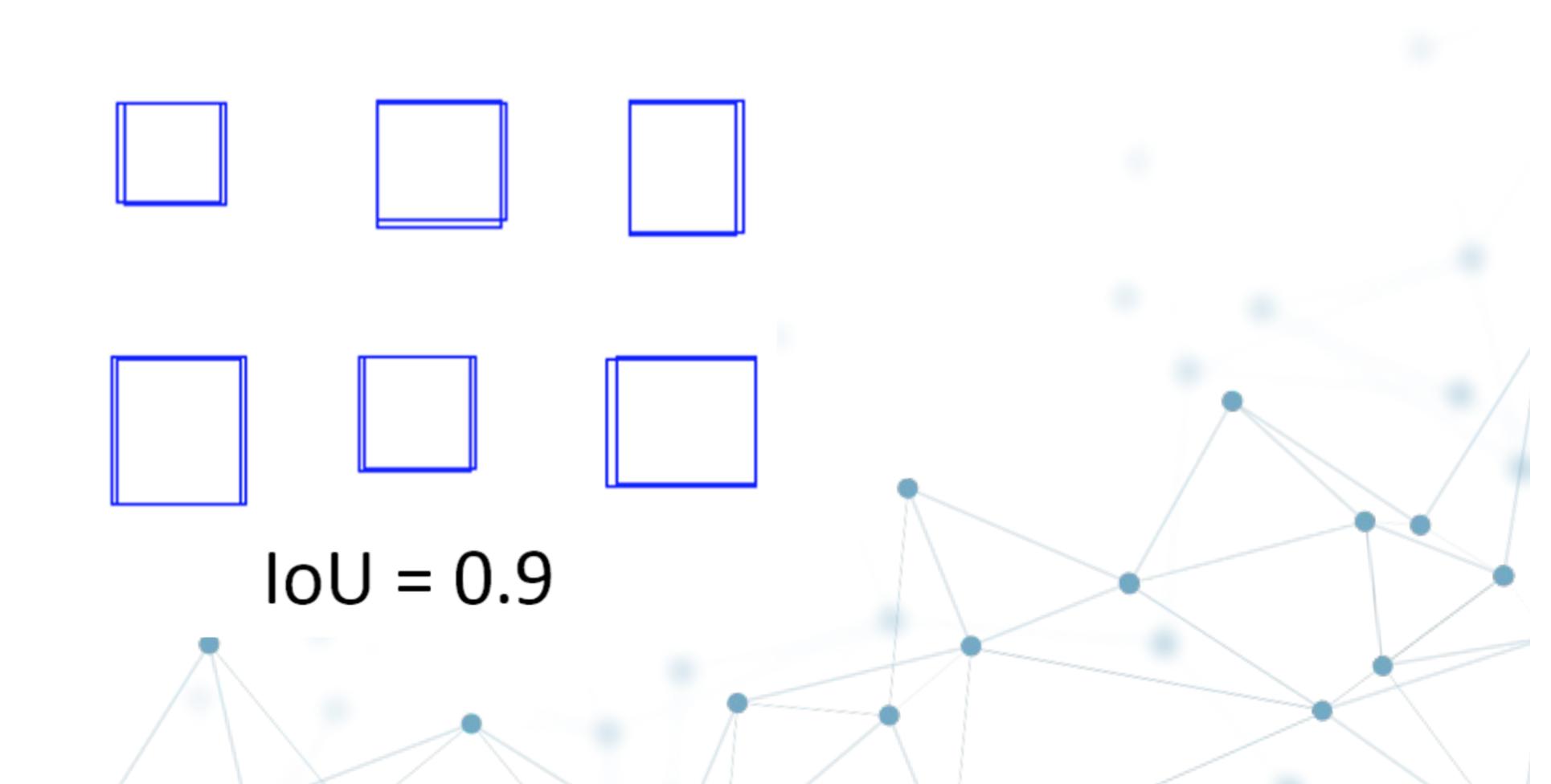
IoU = 0.5



IoU = 0.7

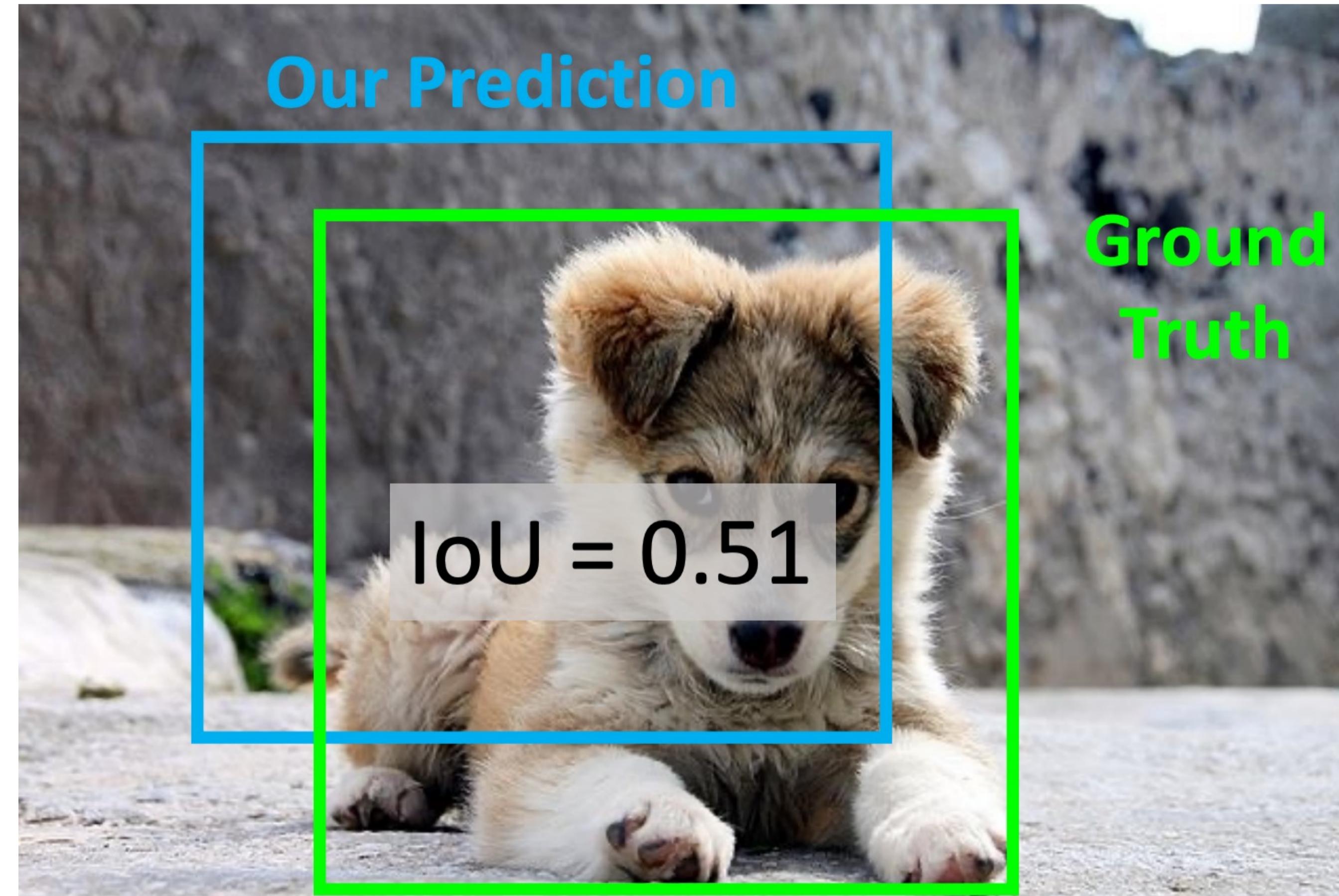


IoU = 0.9



Defining similarity between two boxes

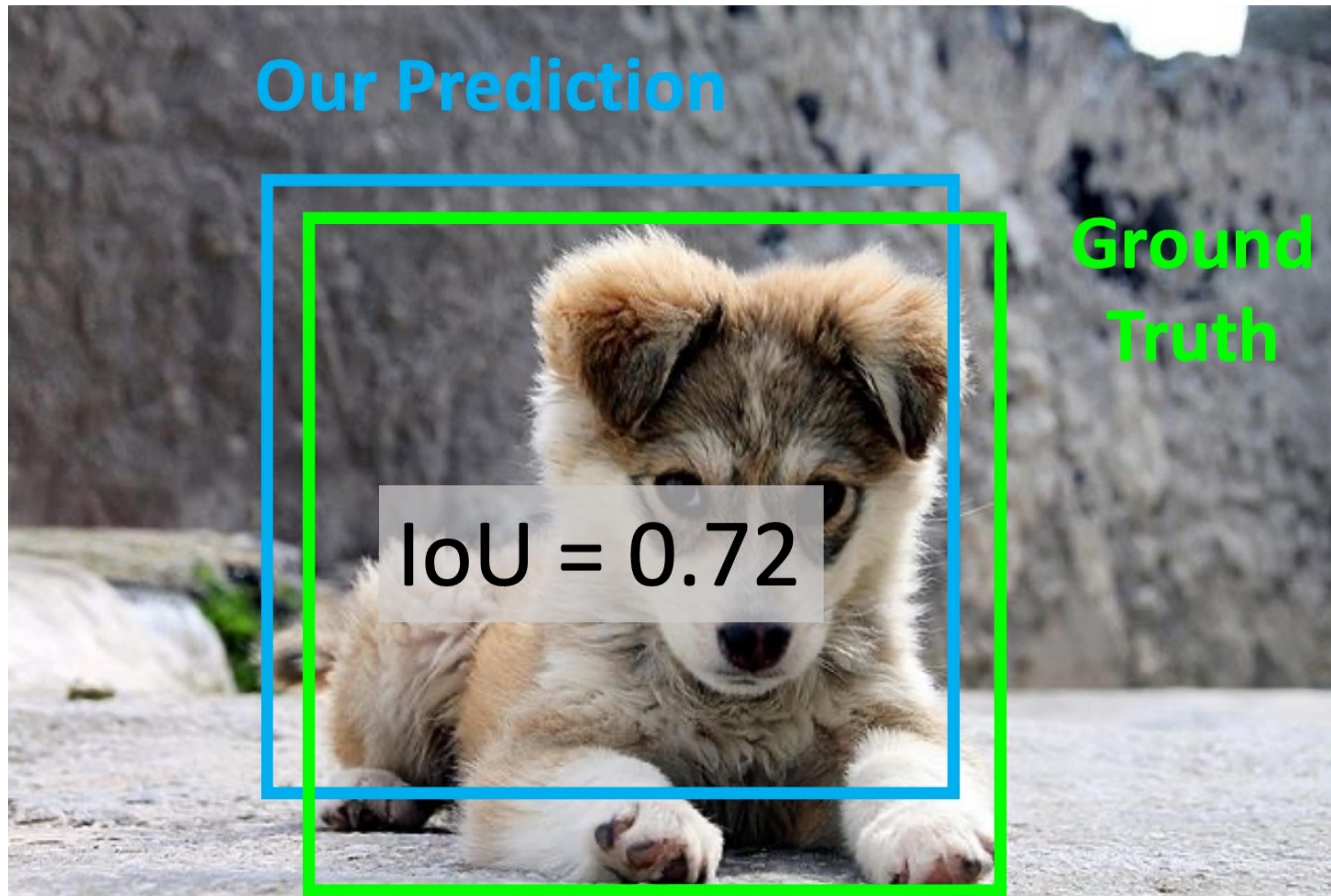
IoU > 0.5:
“decent”



Defining similarity between two boxes

$\text{IoU} > 0.7:$

“pretty good”



Defining similarity between two boxes

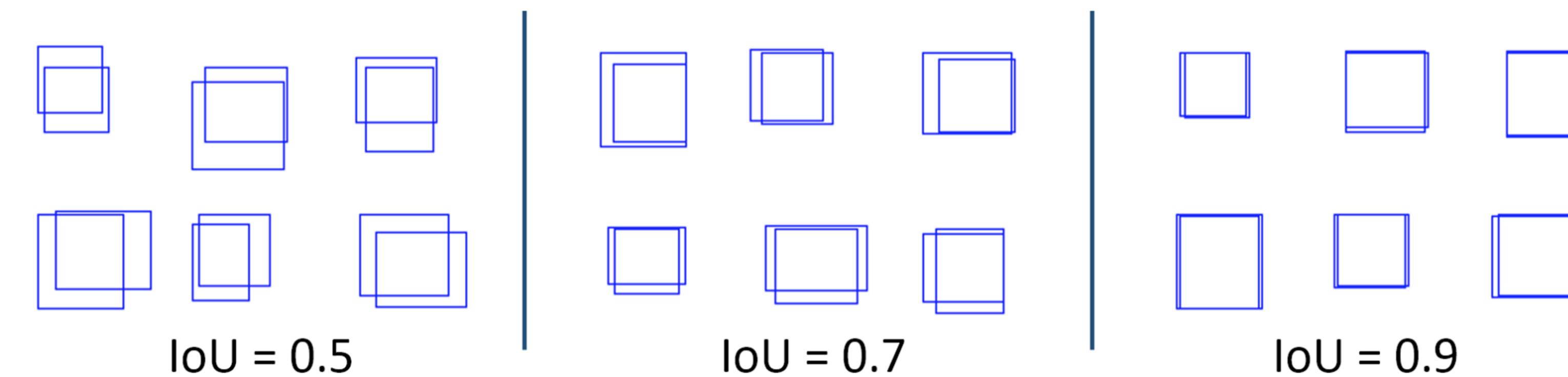
IoU > 0.9:

“almost perfect”



Defining sufficient similarity

How much IoU is good enough, i.e. **sufficient**?



Two schools of thought:

- (1) pick one threshold that seems reasonable (i.e., $\text{IoU}=0.5$)
- (2) consider multiple thresholds



Criterion that defines a false positive

Prediction P is a false positive if there is no matching ground truth



Criterion that defines a false negative

Ground-truth G is a false negative, if G has no matching prediction



Example evaluation for one image

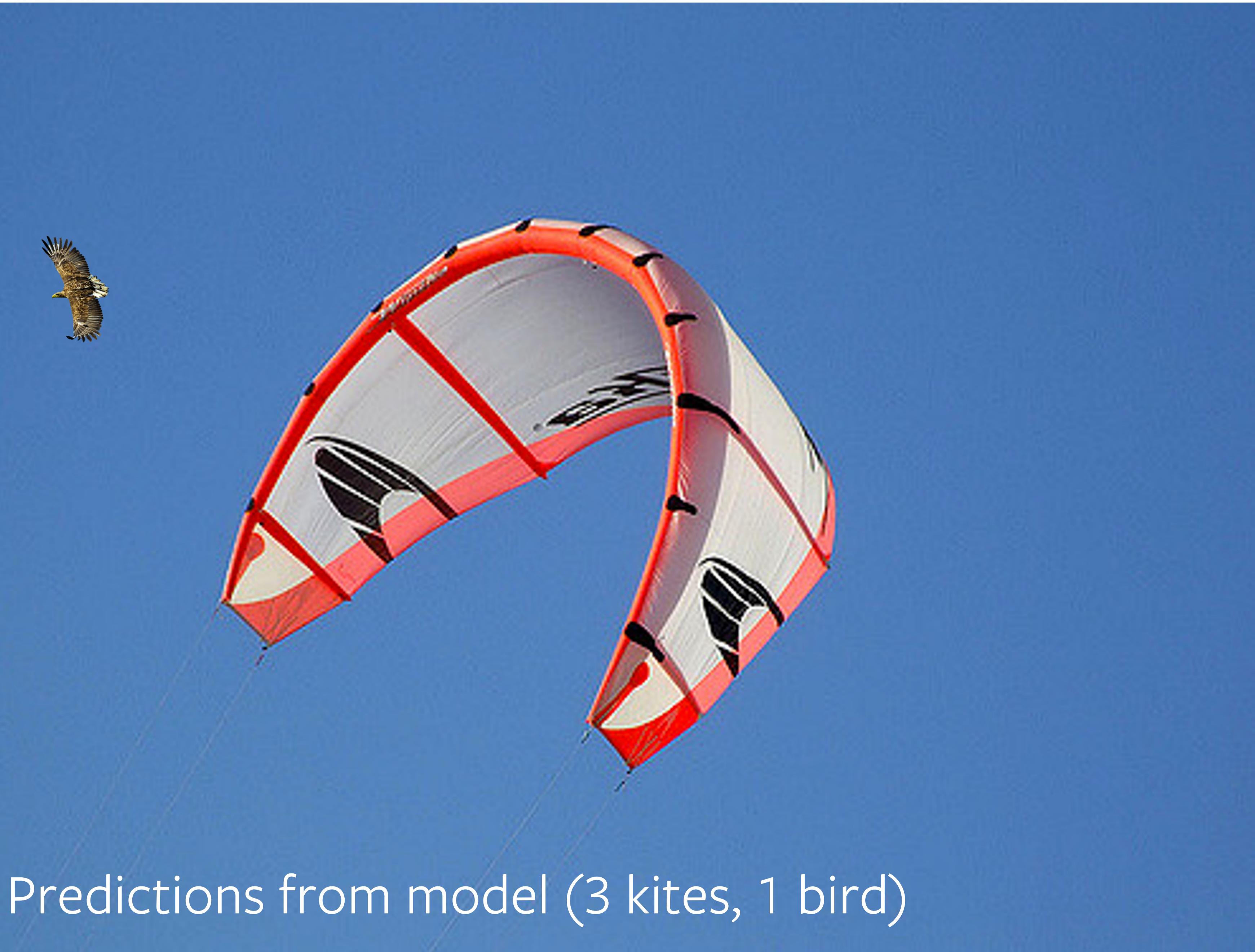
We will use IoU threshold 0.5

The dataset has two categories: {bird, kite}





Testing image



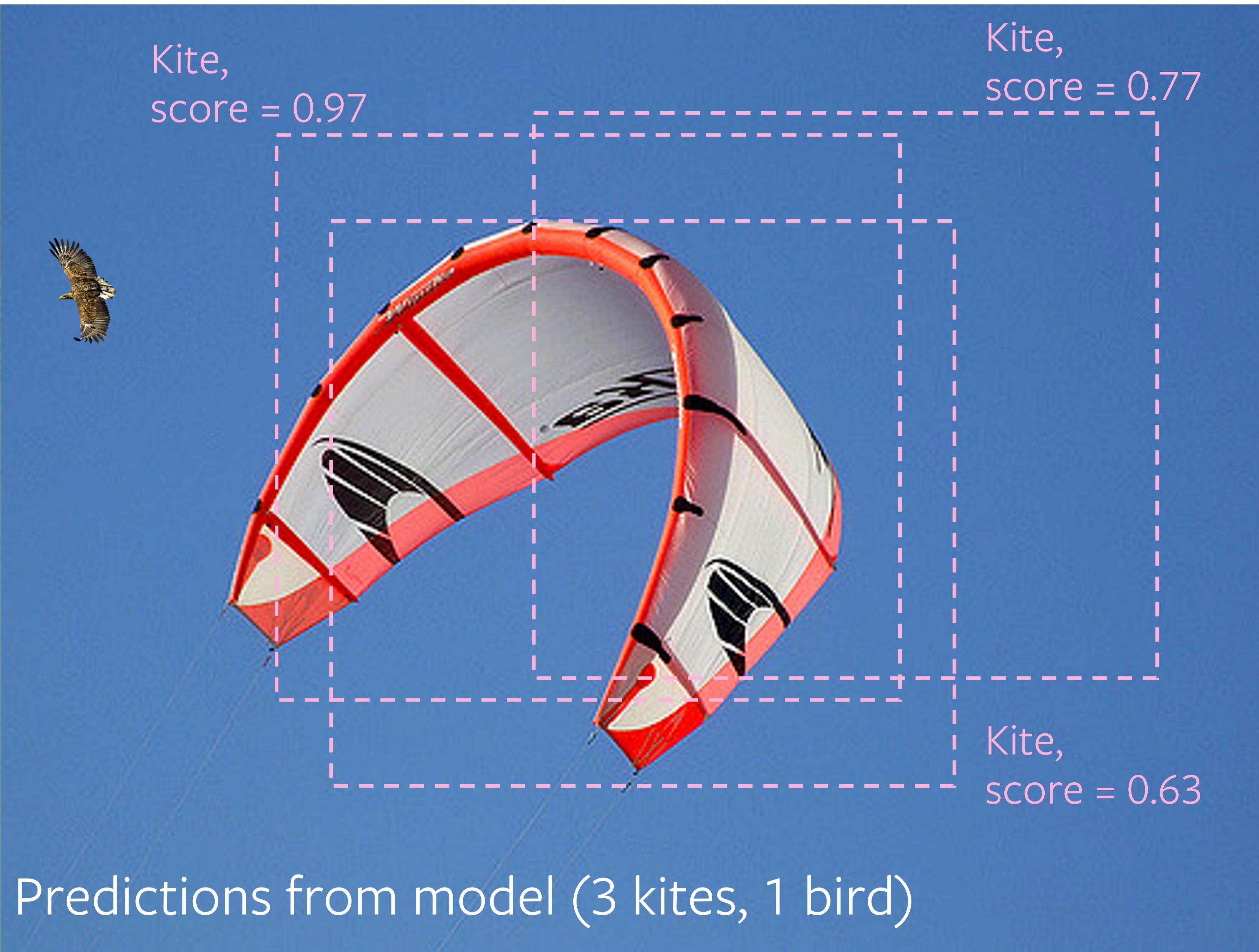
Predictions from model (3 kites, 1 bird)

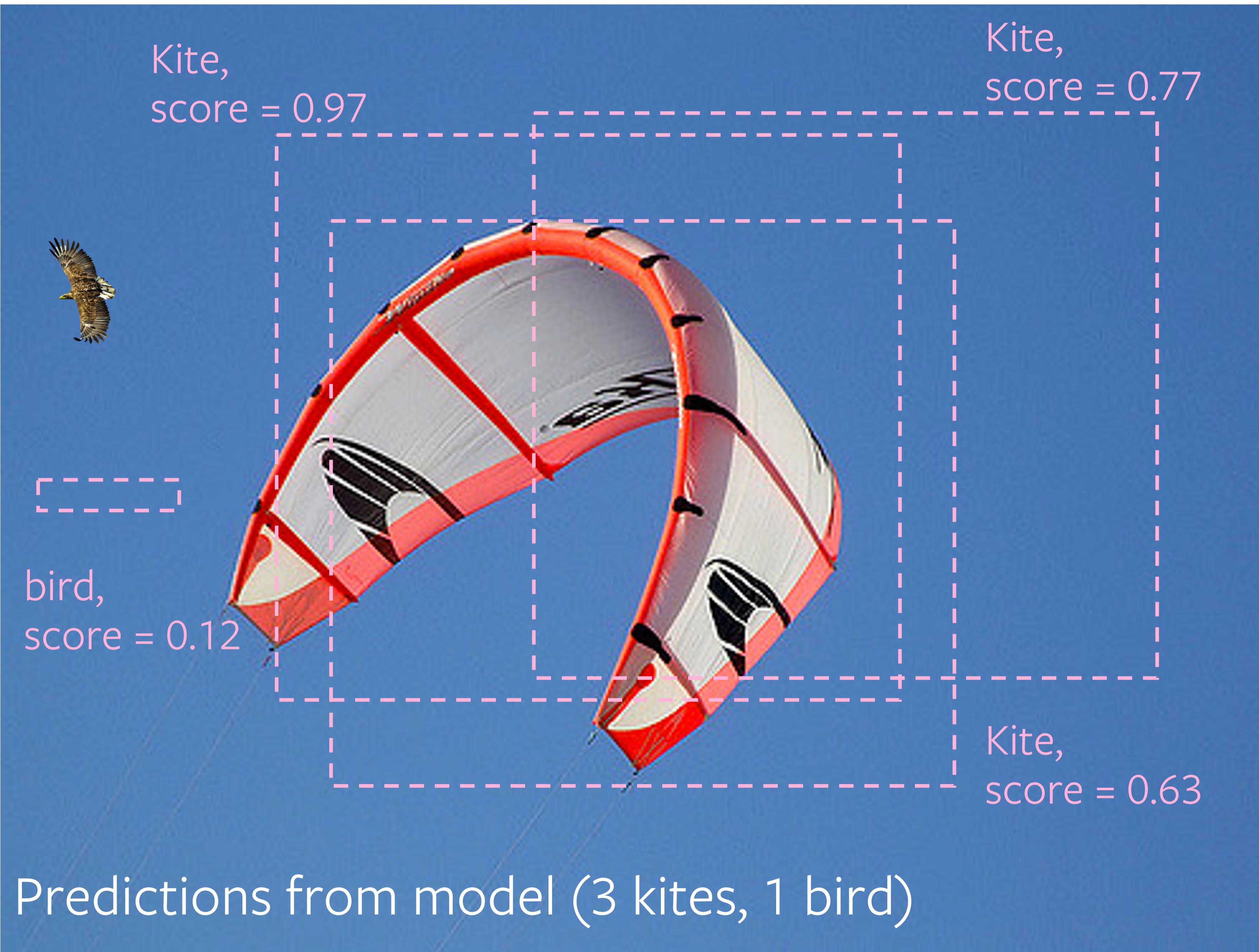
Kite,
score = 0.97



Predictions from model (3 kites, 1 bird)







bird
(ground truth)



Ground-truth available to evaluation code

Recall what the evaluation code does

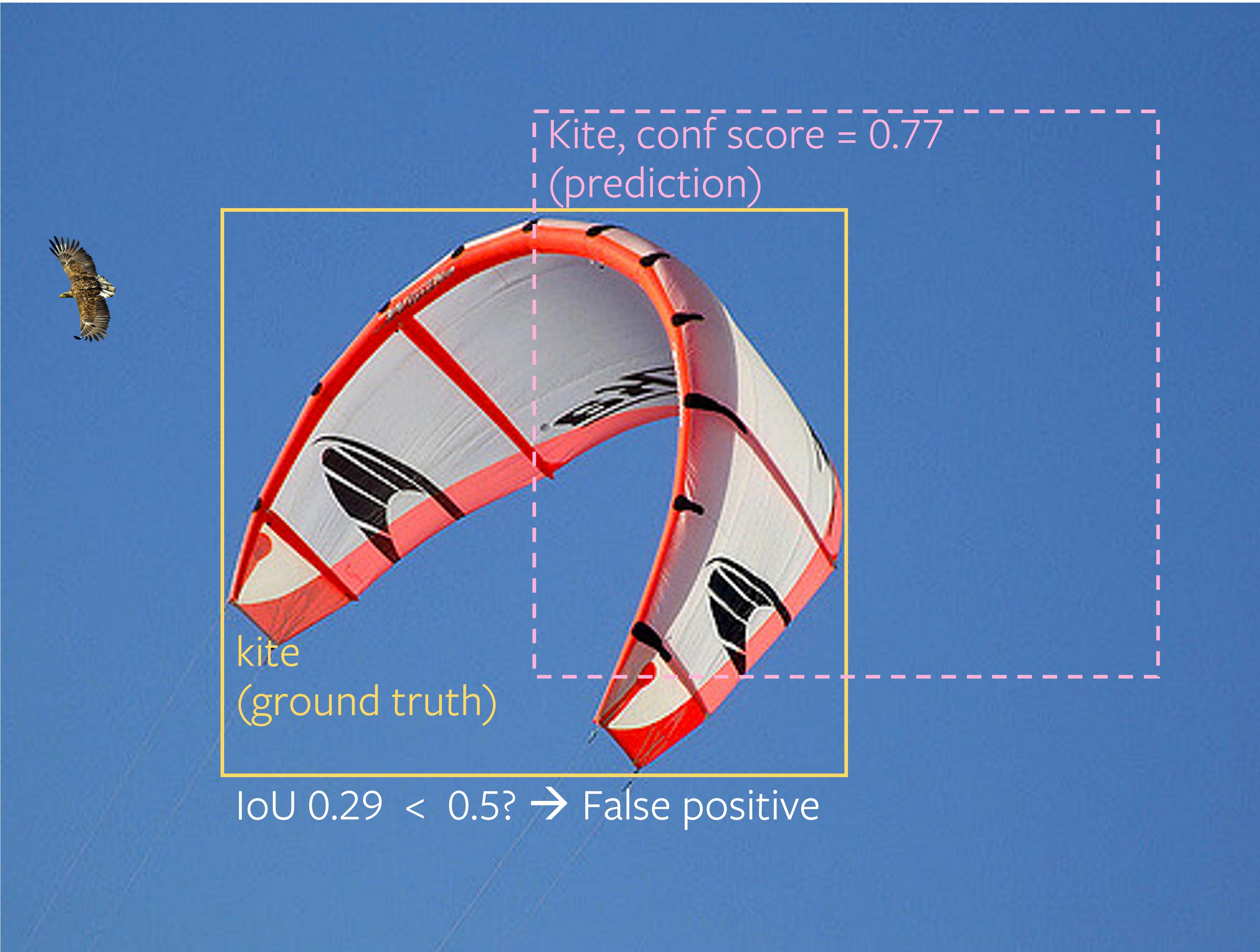
- The code evaluates each category independently
- In particular, the evaluation code:
 - sorts the predictions from highest to lowest confidence;
 - processes each prediction one at a time in this order;
 - assigns a binary classification to each prediction: either “true positive” (TP) or “false positive” (FP);
 - counts the number of false negatives (FN), i.e. objects that were not detected



Kite, conf score = 0.97 (prediction)



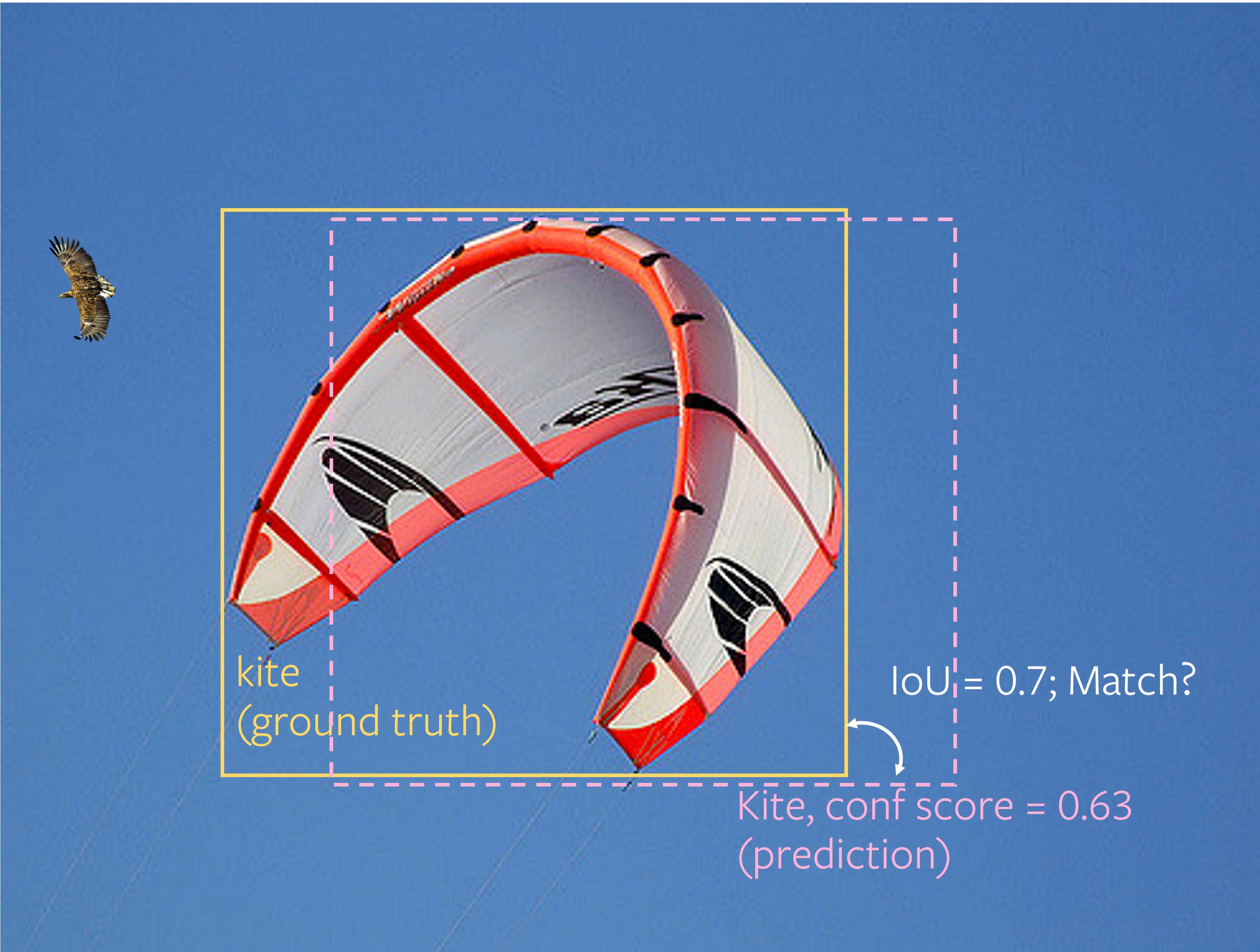
IoU $0.65 \geq 0.5$? (the IoU threshold) → True positive



Kite, conf score = 0.77
(prediction)

kite
(ground truth)

IoU 0.29 < 0.5? → False positive



Previously
matched

Kite, conf score = 0.97 (prediction)



IoU = 0.7; Match?

Kite, conf score = 0.63
(prediction)

Previously
matched

Kite, conf score = 0.97 (prediction)



Duplicate detection
→ False positive

IoU = 0.7; Match?
Nope!
Kite, conf score = 0.63
(prediction)

bird
(ground truth)



IoU 0.0 < 0.5?
→ False positive



bird,
score = 0.12



No matching prediction
bird → False negative
(ground truth)



Recap: evaluation code

- For each category C, a separate list of confidence ranked predictions
- For each prediction of C, a decision if it is a TP or a FP
- A set of FNs for C



How should a good detector behave?

- . ✓ Perfect localization is desired by not necessary for most applications
close is good enough -> **“close” is determined by IoU threshold** ;
- . ✓ Names the object correctly:
don’t call a cat, dog -> **exact category match**;
- . ✓ Does not find an object more than once:
don’t duplicate detections -> **confidence scores define ranking**;
- . ✓ Finds all instances of each object category:
don’t miss anything -> **penalizes missed objects**.

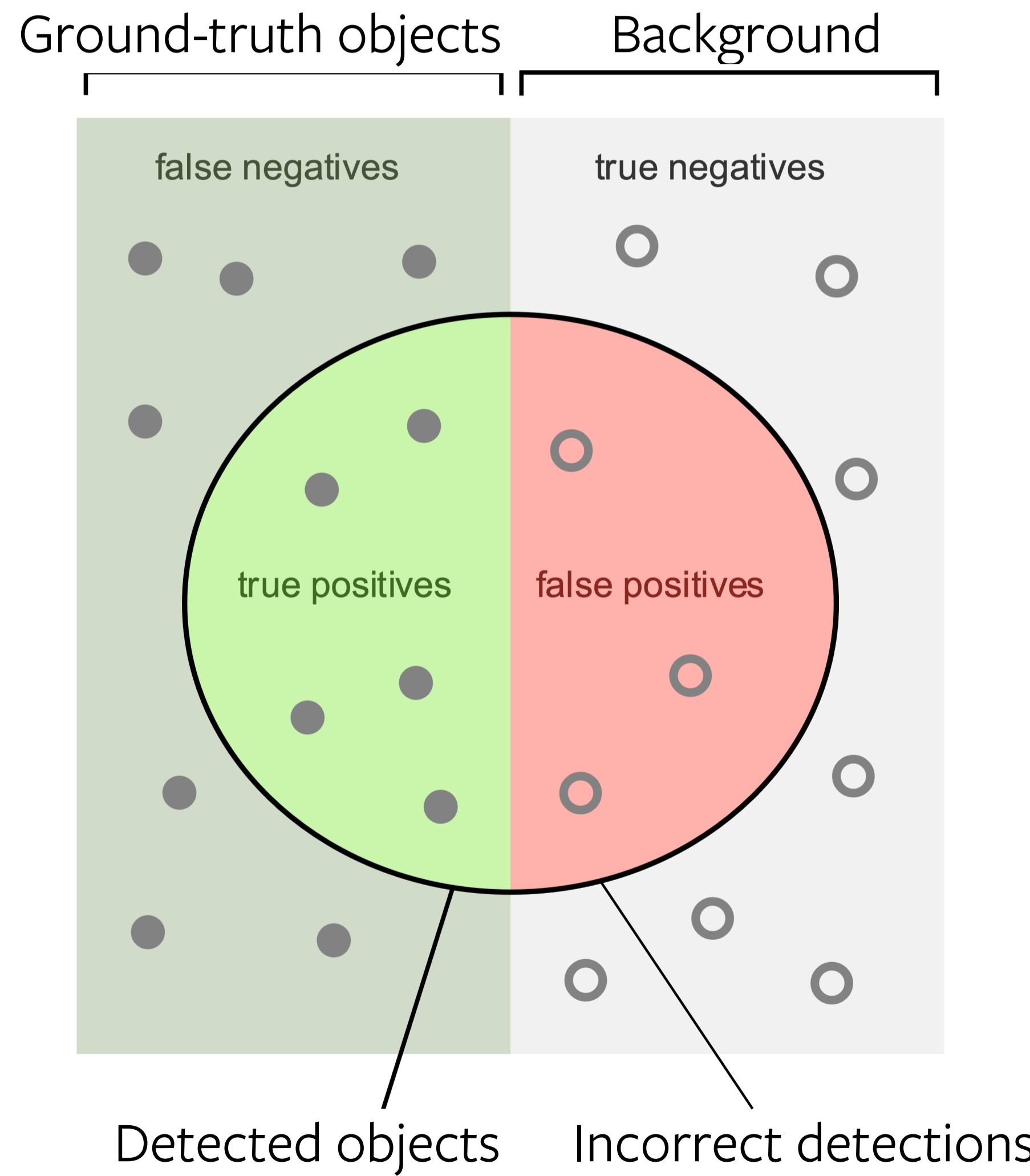


How do we quantify “goodness”?

- Thus far all we can do is categorize predictions (TP vs FP) and count false negatives (FNs);
- **What we want: an evaluation metric**
a single number that communicates how good the model is performing on the dataset (something like accuracy for image classification)
- Next, we'll build up the definition of Average Precision (AP)



Preliminaries: Precision and Recall



$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

“Of the detections, what fraction are correct?”

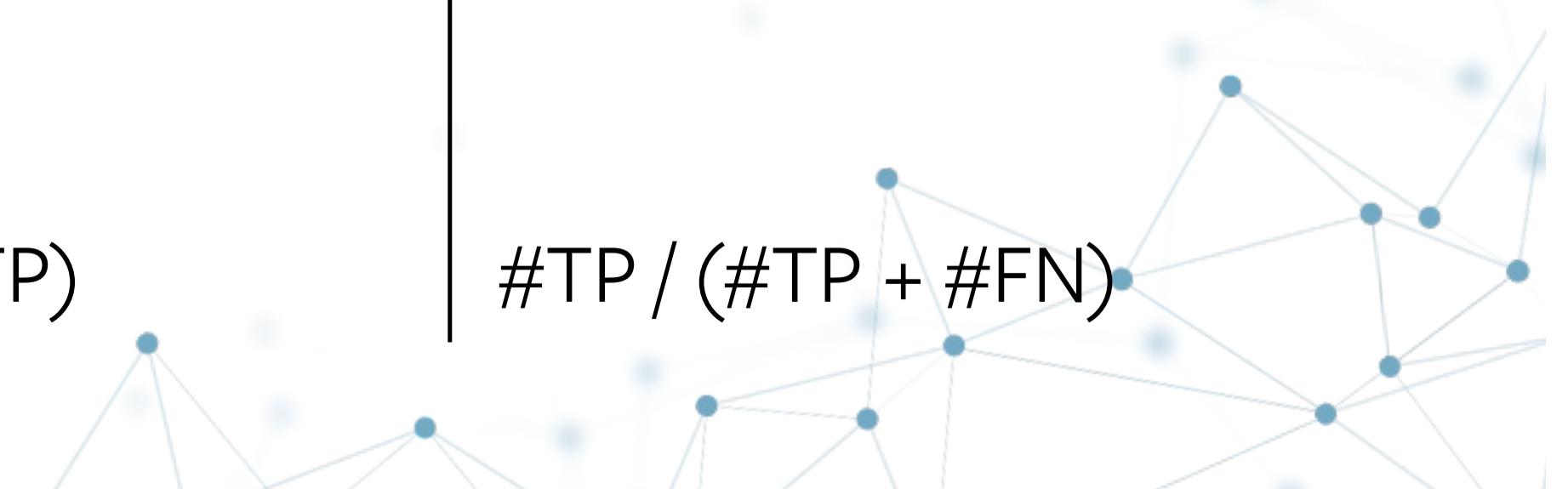
$$\#TP / (\#TP + \#FP)$$

73

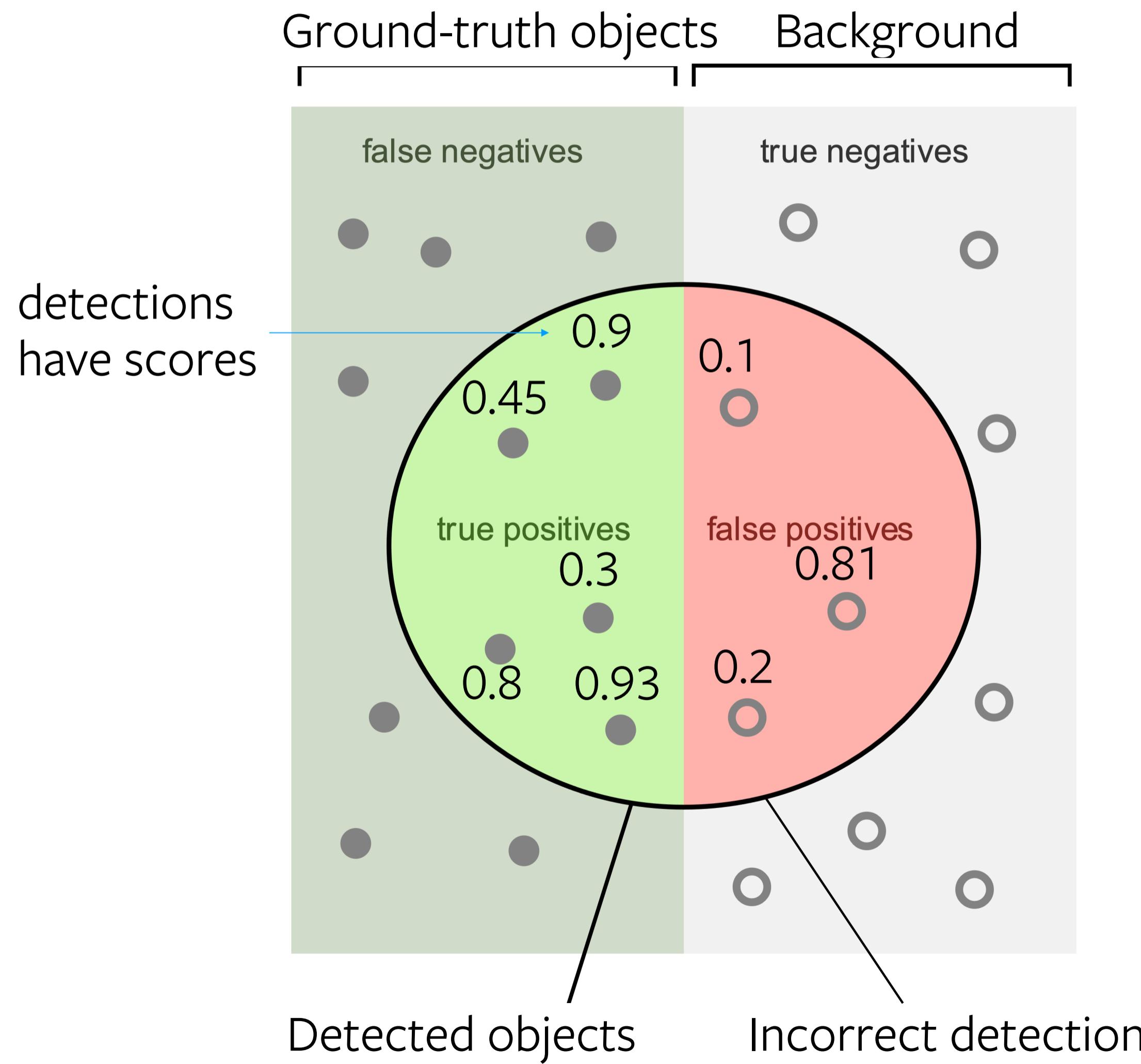
$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

“What fraction of the ground-truth was detected?”

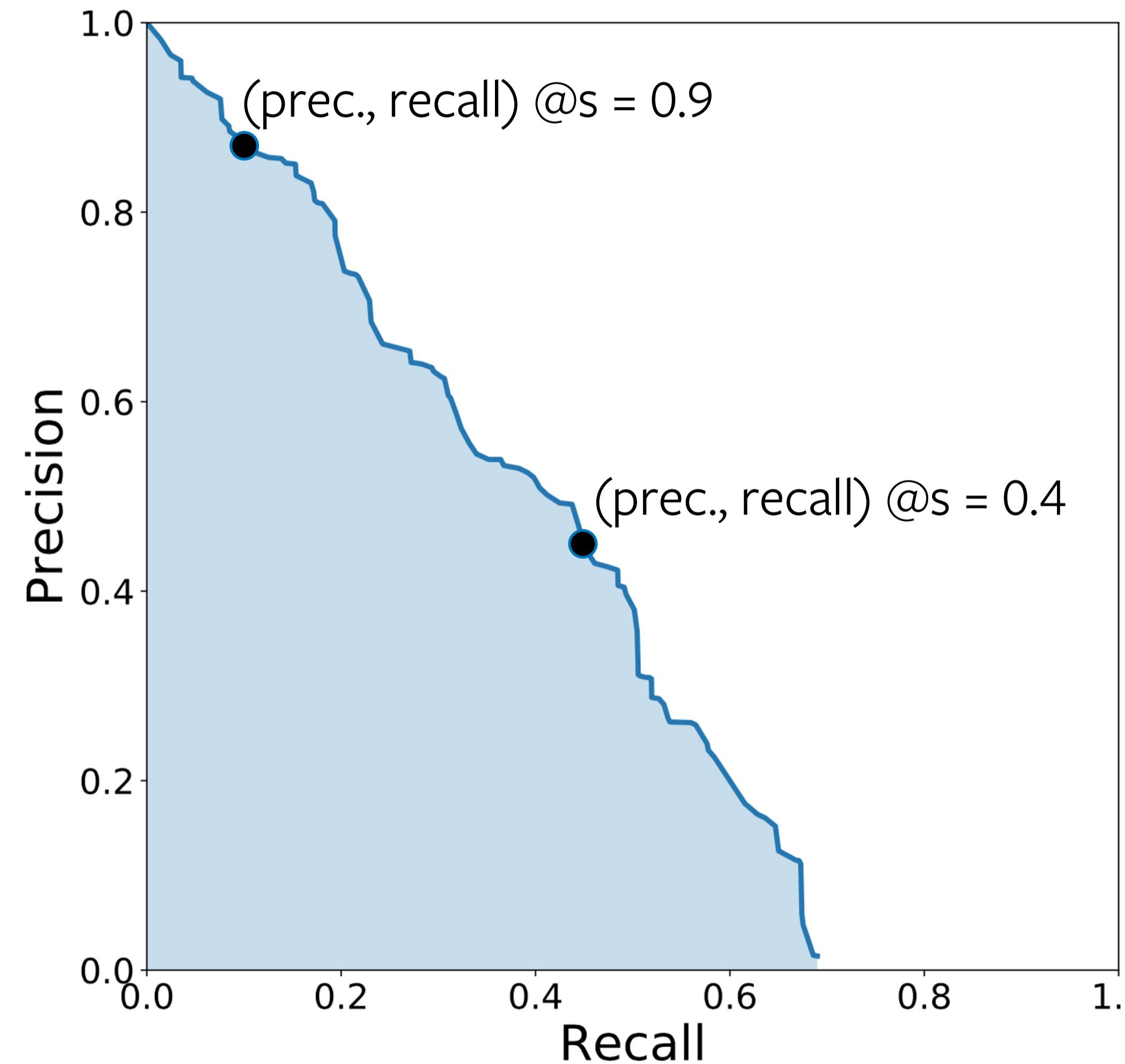
$$\#TP / (\#TP + \#FN)$$



Preliminaries: Precision and Recall “@s”



Precision-recall (PR) curve



A curve parameterized by score s

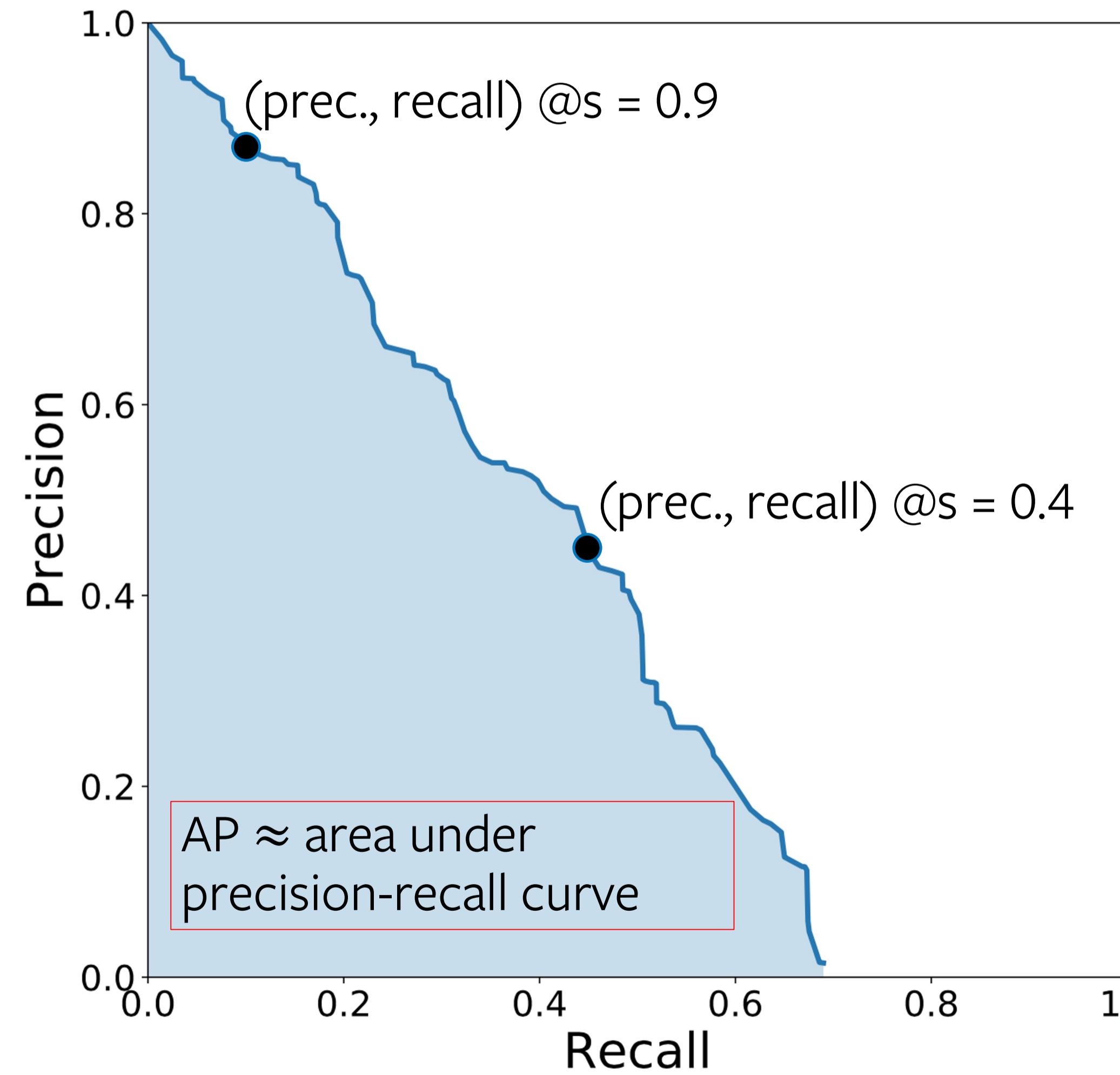
Sweeping $(\text{precision}, \text{recall})@s$ from high to low s traces out the PR curve.

Example points at $s = 0.9$ and $s = 0.4$.

(Note: this curve is for one category at a time; each category has a separate PR curve)



Average precision (AP) of the PR curve



Average precision (AP) := the area under the PR curve.

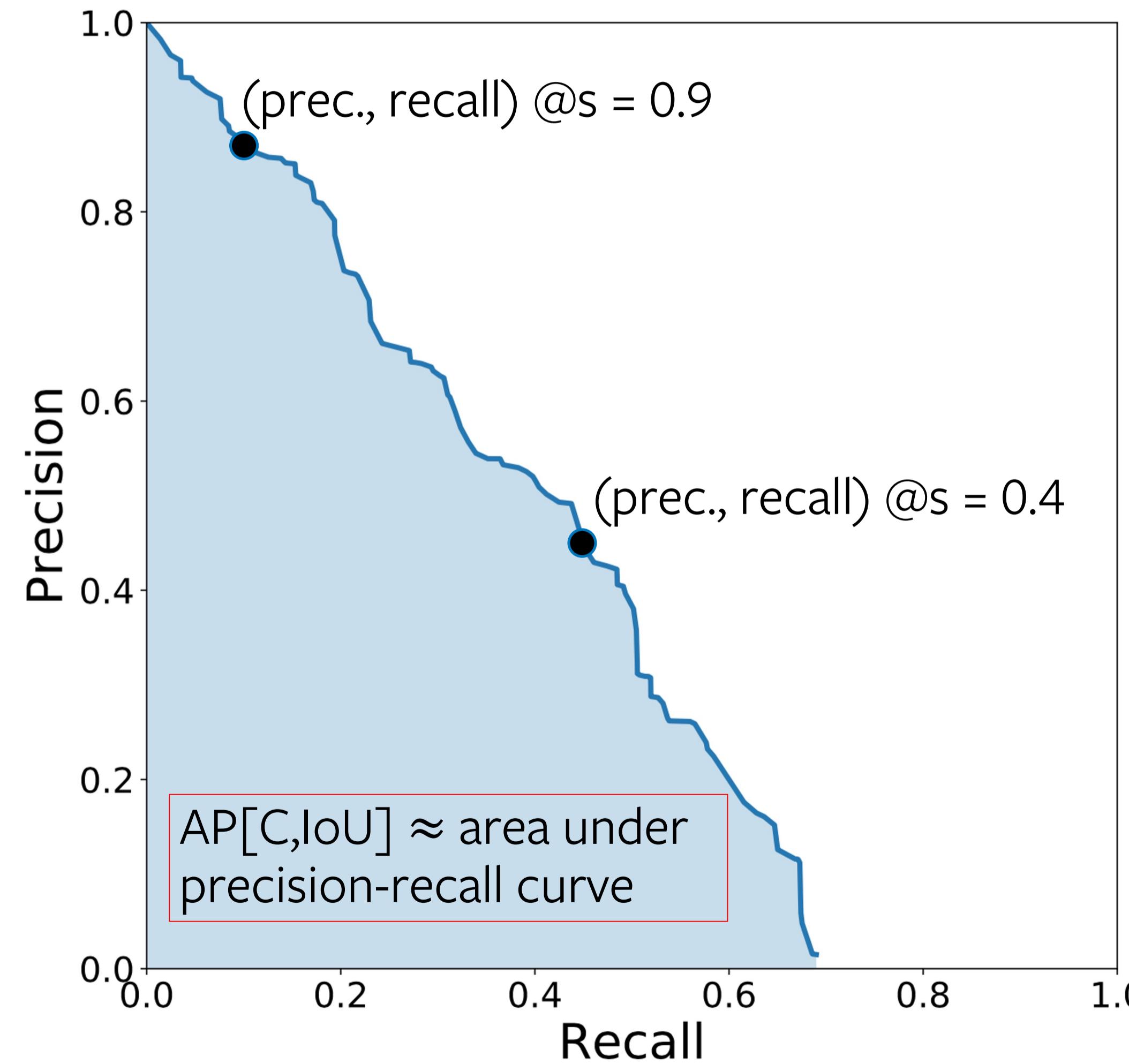
Average is over all levels of recall.

AP ranges from 0 to 1.

AP of 1 means perfect precision and perfect recall.

AP is computed numerically given finite (prec., recall)@s samples

Average precision (AP) of the PR curve



Notation: AP[C, IoU]

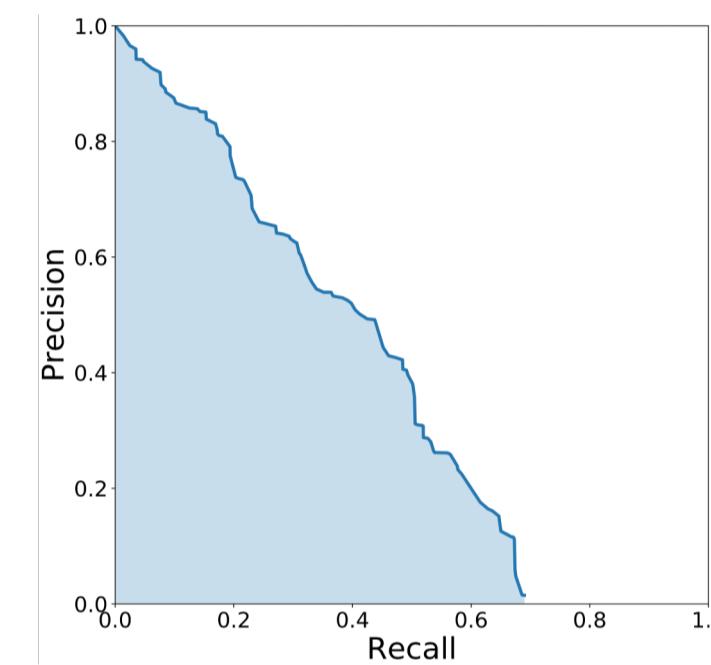
Category C

IoU threshold

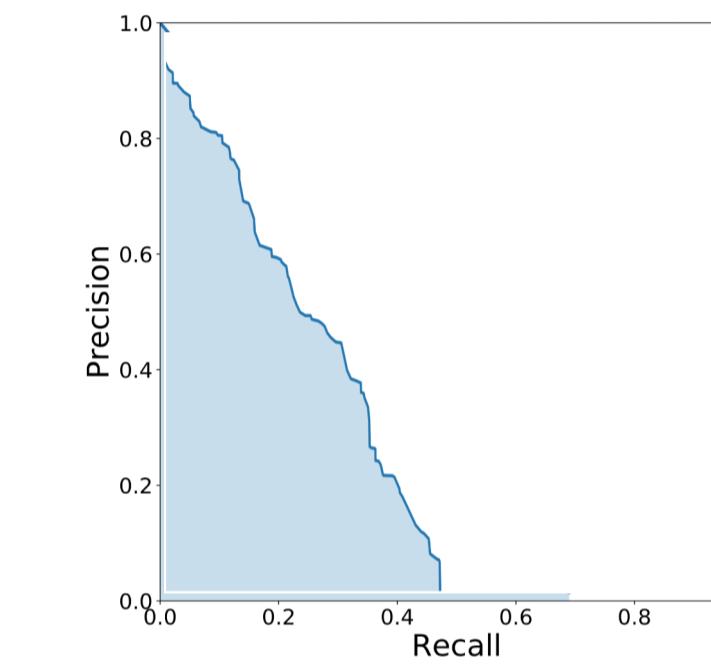
E.g., AP[“kite”, 0.5]



Extending AP to multiple categories



$\text{AP}[\text{"kite"}, 0.5]$



$\text{AP}[\text{"bird"}, 0.5]$

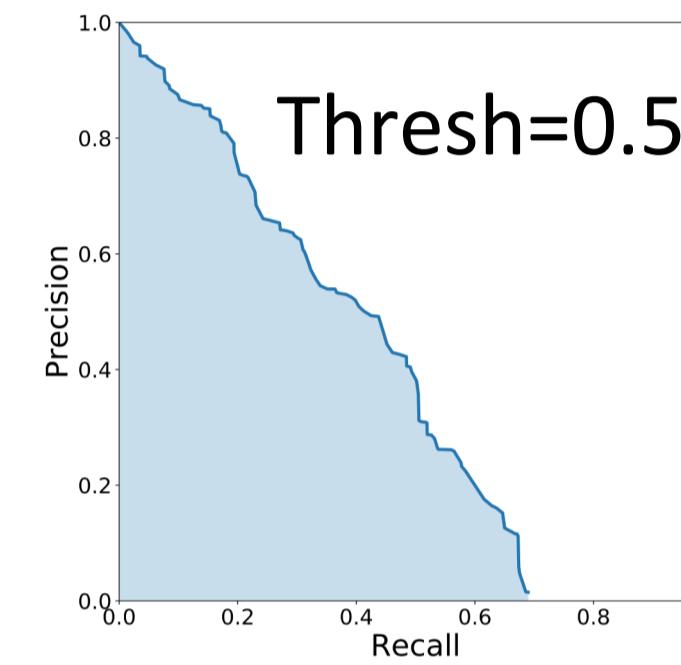
$C = \{\text{kite}, \text{cat}, \text{dog}, \dots, \text{bird}, \dots\}$

$$\text{AP}[0.5] = \frac{1}{|C|} \sum_{C \in C} \text{AP}[C, 0.5]$$

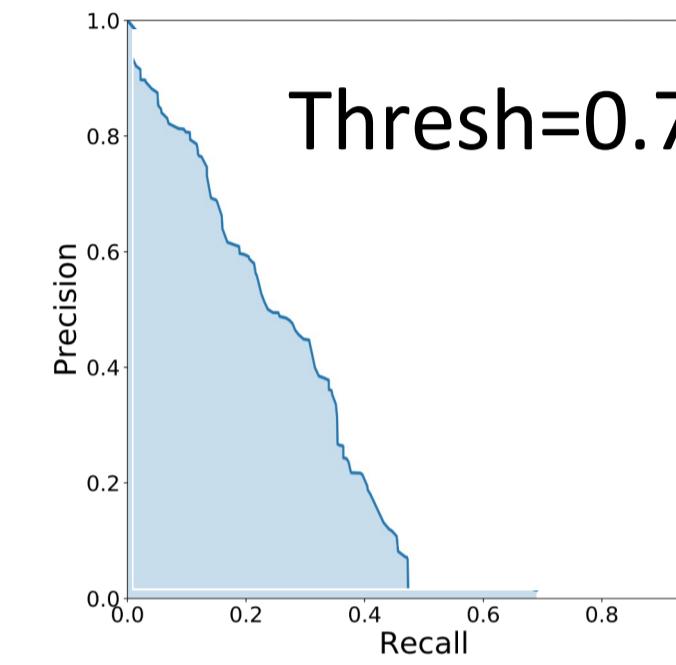
This is often called **mean average precision** or **mAP**



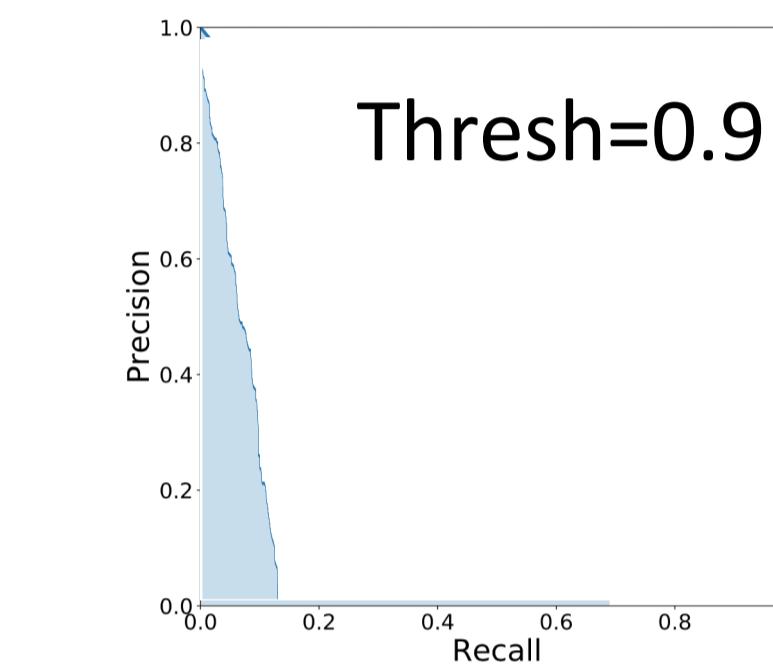
Extending AP to multiple IoU thresholds



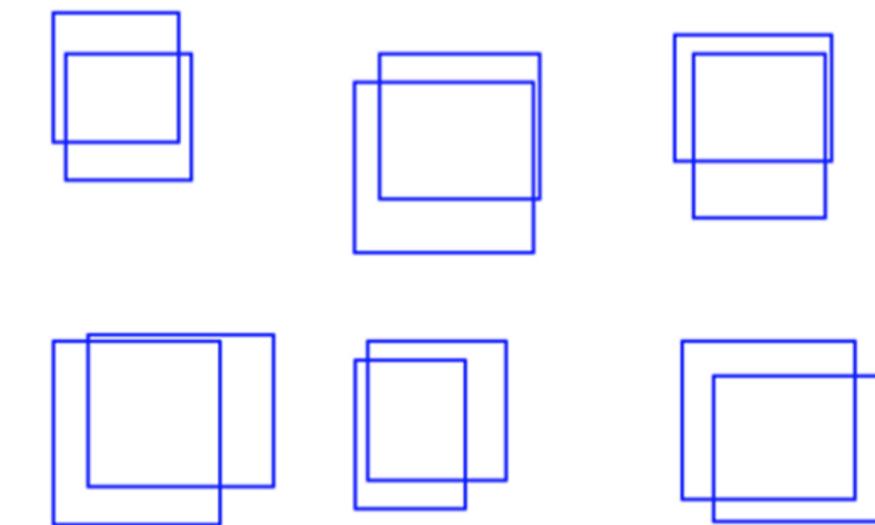
AP[“kite”, 0.5]



AP[“kite”, 0.7]

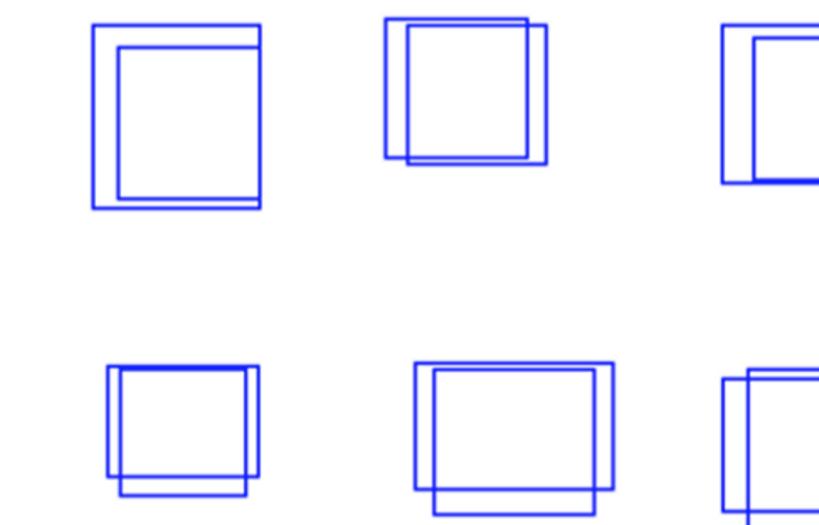


AP[“kite”, 0.9]

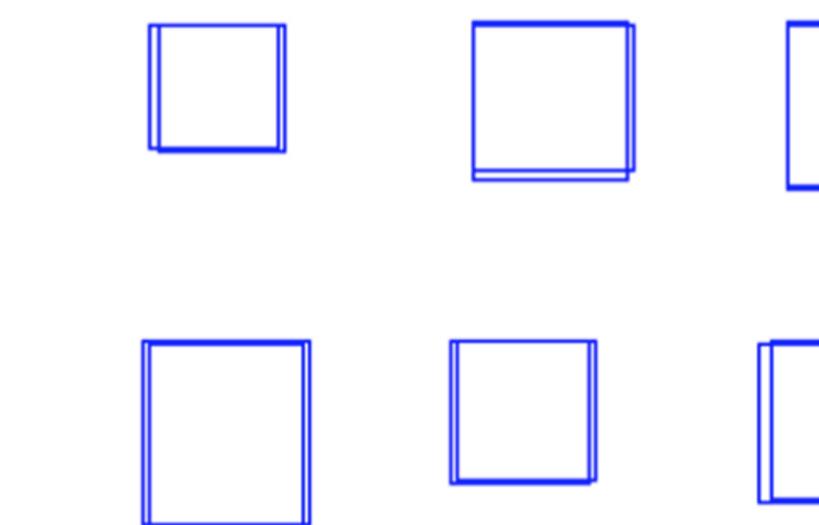


IoU = 0.5

Loose



IoU = 0.7

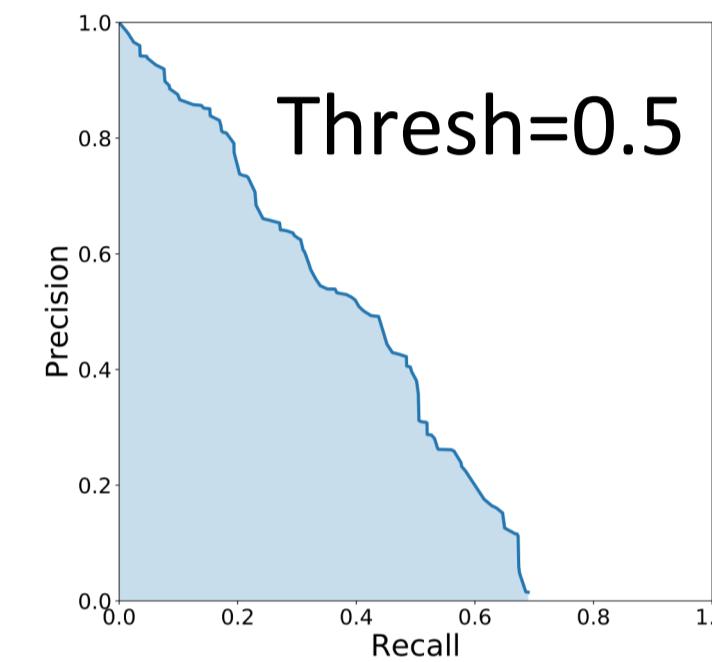


IoU = 0.9

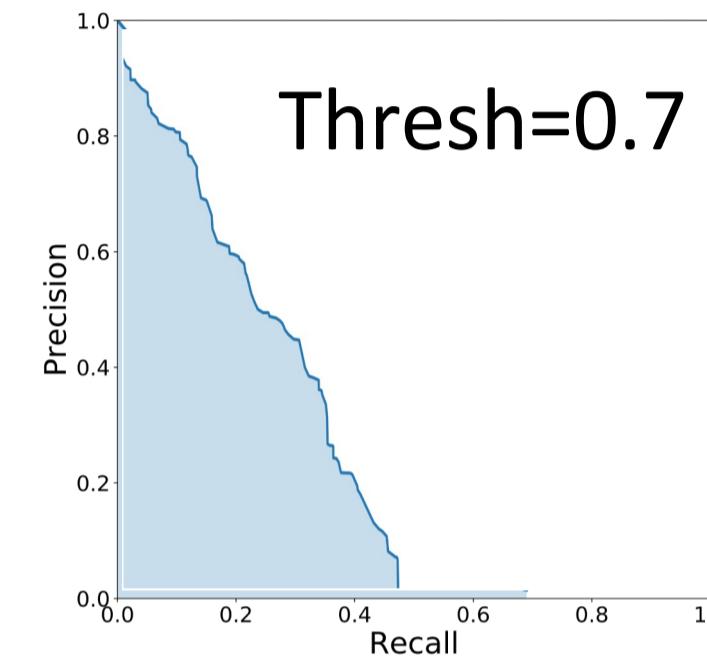
Tight



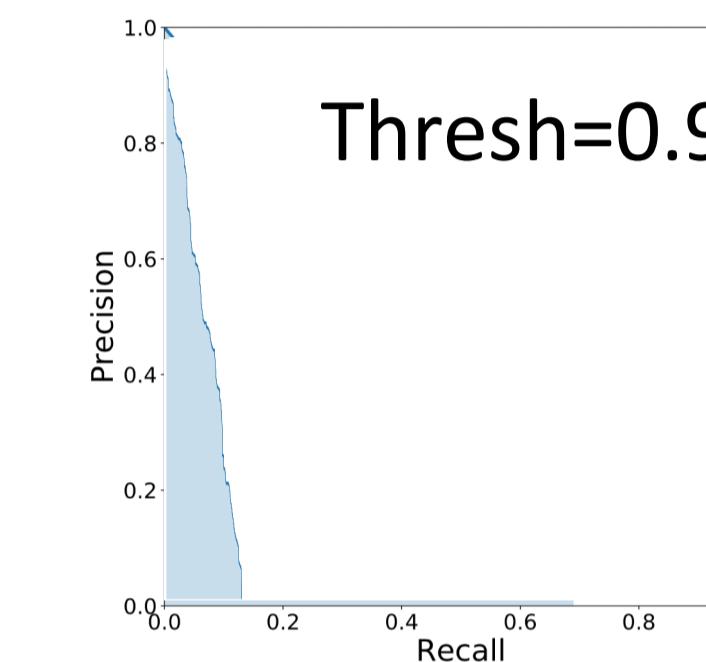
Extending AP to multiple IoU thresholds



AP["kite", 0.5]



AP["kite", 0.7]



AP["kite", 0.9]

$$T = \{0.50, 0.55, 0.60, 0.65, \dots, 0.95\}$$

$$AP["kite"] = \frac{1}{|T|} \sum_{t \in T} AP[kite, t]$$



Multiple categories and IoU thresholds

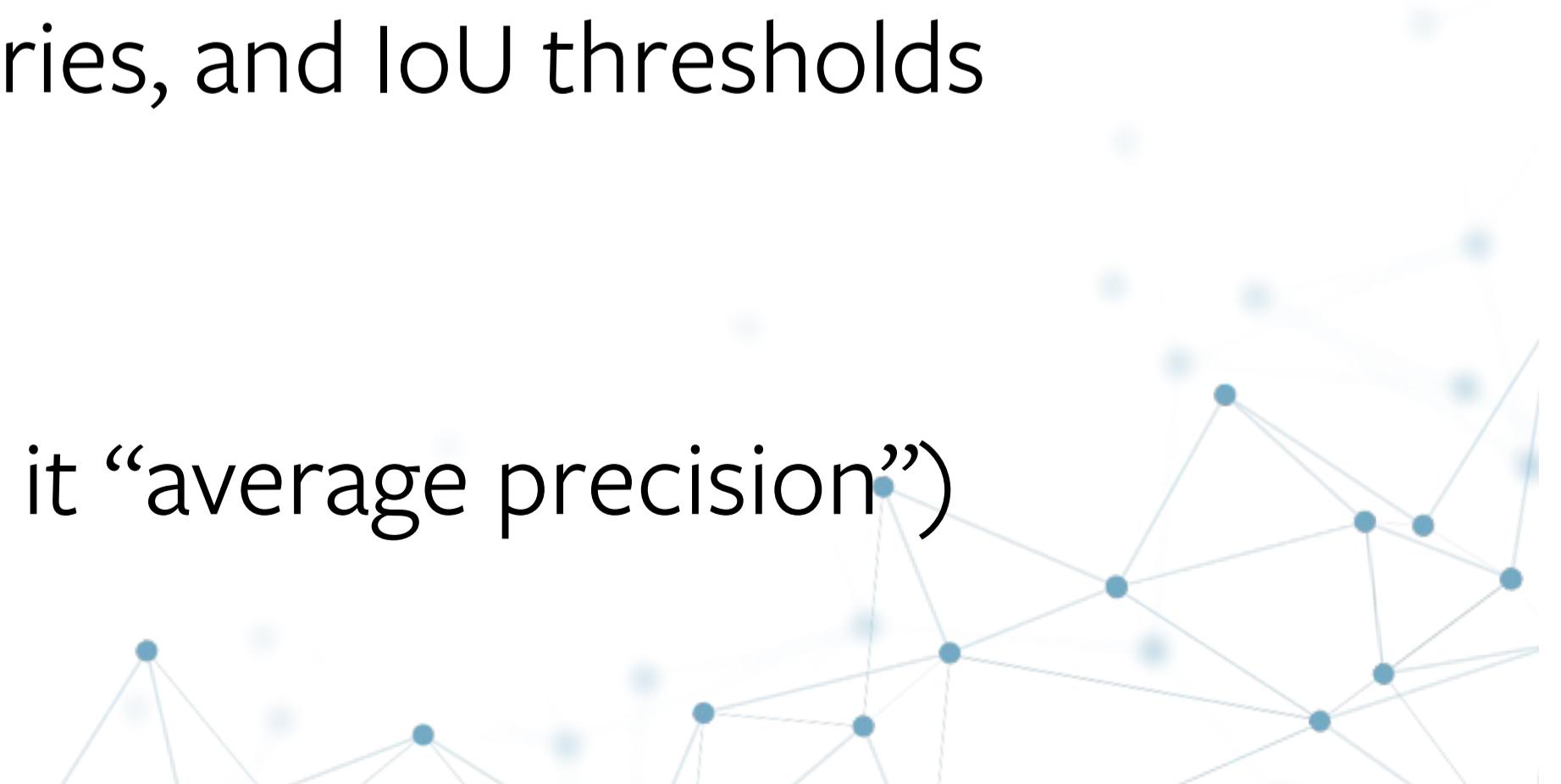
$$\mathbb{C} = \{\text{kite}, \dots, \text{bird}, \dots\}$$

$$\mathbb{T} = \{0.50, 0.55, 0.60, 0.65, \dots, 0.95\}$$

$$AP = \frac{1}{|\mathbb{T}| |\mathbb{C}|} \sum_{C \in \mathbb{C}, t \in \mathbb{T}} AP[C, t]$$

This is an average of precision over recall, categories, and IoU thresholds

(“mean mean average precision”? – Most just call it “average precision”)



Recap: Evaluation metric

- The most popular metric is “average precision” (AP)
- Precision values averaged over recall ranging from 0 to 1
- Caveat: there are multiple definitions of AP:
 - can be AP for a single category and single IoU threshold;
 - can be AP for a single category, averaged over multiple IoU thresholds;
 - low-level numerical computation details can differ;
 - can be difficult to disambiguate without context.

Rule of thumb: each dataset defines a metric (e.g., “COCO-style AP”, “PASCAL VOC-style 2007 AP”, “PASCAL VOC 2010-style AP”, etc.)



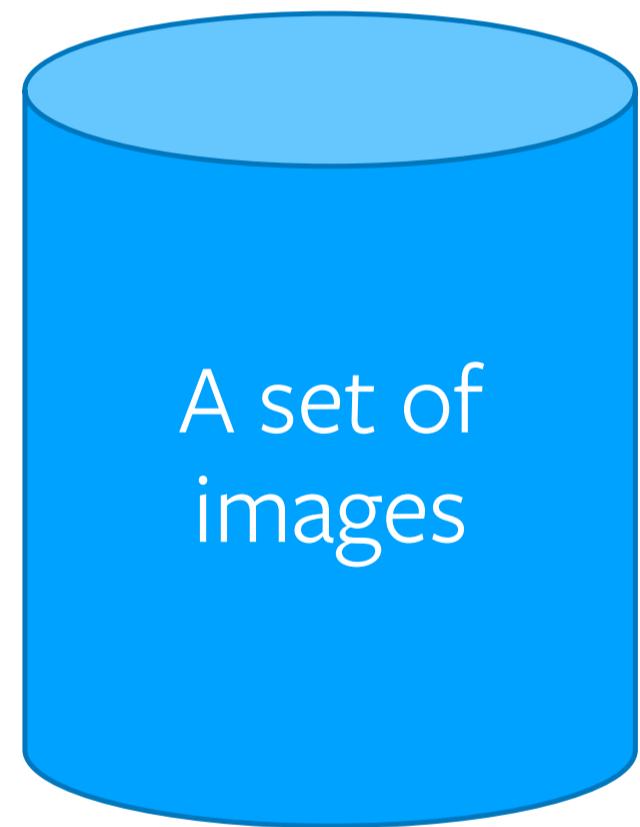
Which definition of AP to use?

- Largely a philosophical question
- Real-world applications are driven by real-world requirements
- AP is defined without any grounding in a particular application
 - the metric is generic -> doing well should be good for many applications
 - the goal of AP=1 is generically good (perfect recall, perfect precision)
- Early datasets (e.g., PASCAL VOC) used IoU threshold=0.5 – the detection problem was new and “too hard”. Larger datasets (e.g. COCO) extended to multiple IoU thresholds.
The detection problem is established, there’s good progress: aim higher!



Task formulation (training) - output

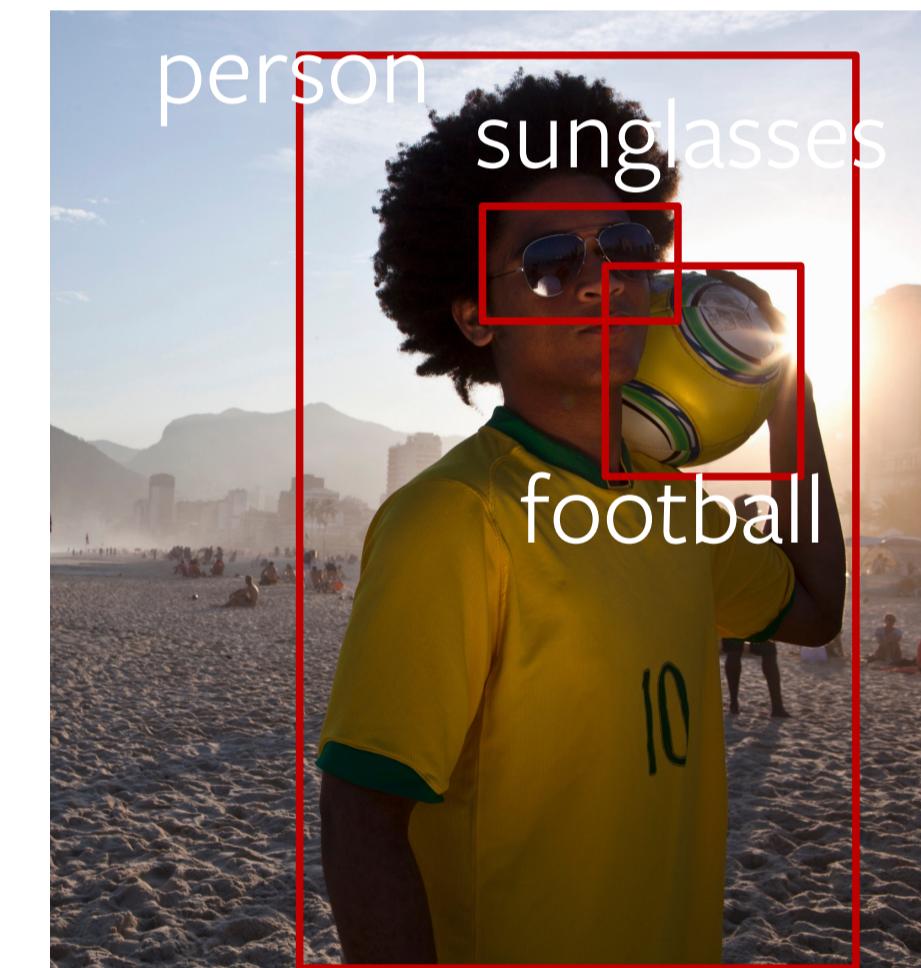
- A **dataset** composed of:



Typically, 10,000 – 100,000

Car,
Person,
Football,
Bottle,
Sunglasses,
Giraffe,
...

A pre-specified set of categories
Typically, 10 – 1,000



For each image:
boxes for each category

Recap: today's focus

- What is the object detection task?
- How is this task formalized?
- How is this task evaluated?
- Why is the task difficult?
- What are important datasets to know about?



Why is object detection difficult?

- Invariance
- Equivariance
- Imbalance



Invariance

A function f is invariant to a transformation g if

$$\forall x \ f(gx) = f(x)$$

$$f(R_\theta \begin{array}{c} \text{2} \\ \text{2} \end{array}) = f(\begin{array}{c} \text{2} \\ \text{2} \end{array}) = f(\begin{array}{c} \text{2} \\ \text{2} \end{array}) \\ = \text{"2"}$$

The “what” part of a detector (“ f ”) needs to be invariant to nuisance transformations – this is like image classification

Example: Photometric variation

Bias and gain $I_2 = \alpha I_1 + \beta$



negative bias



fractional gain

Photometric variations should not change the detector's output

Example: Intra-class variation

Appearance and structure



What makes a chair a chair?

(Unhelpful) answer: “A chair is a chair because it looks like a chair”

Variation within a category should not change the detector’s category prediction

Equivariance

A function f is equivariant with transformation g if

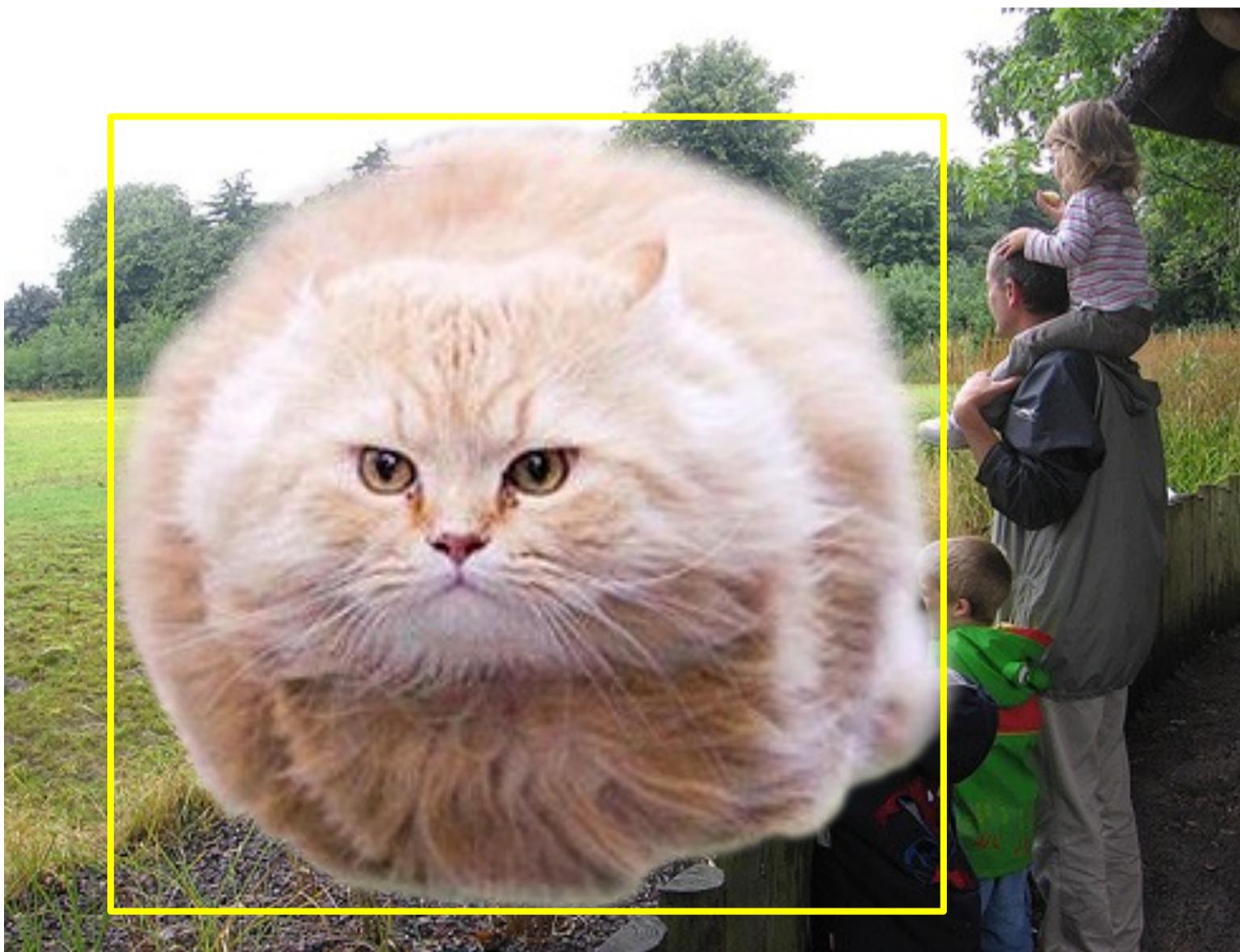
$$\forall x \exists M_g \quad f(gx) \approx M_g f(x)$$

$$f(S_\theta \text{ ↪}) = f(\text{ ↪}) = f(\text{ ↪}) + \theta$$

The “where” part of the detector needs to be equivariant with geometric transformations of objects – this is new compared to image classification

Example: Scale variation

Detection: “where” (i.e., the box) is equivariant with scale shift
(Contrast to classification: invariant to scale shift)



Cat

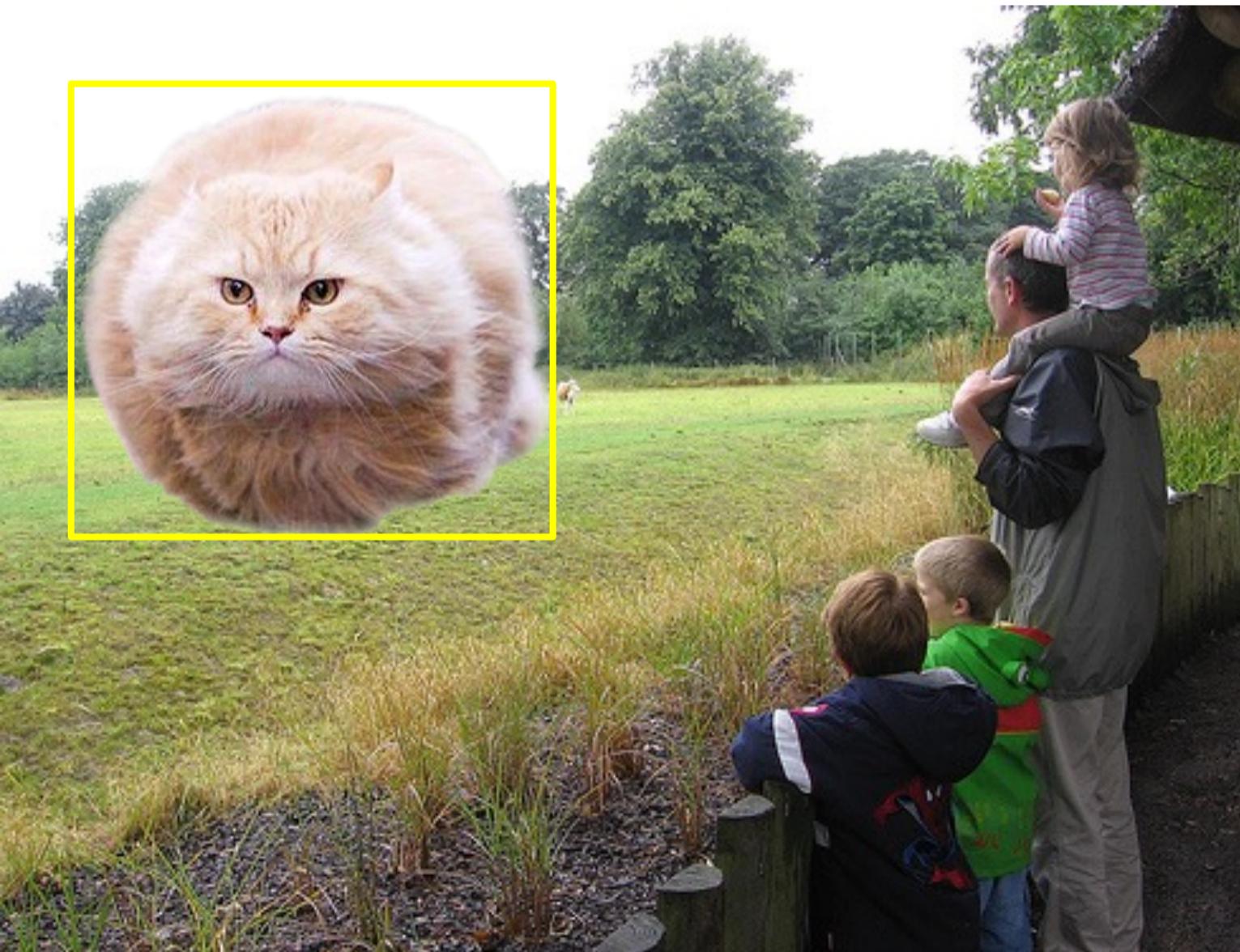


Cat



Example: Translation variation

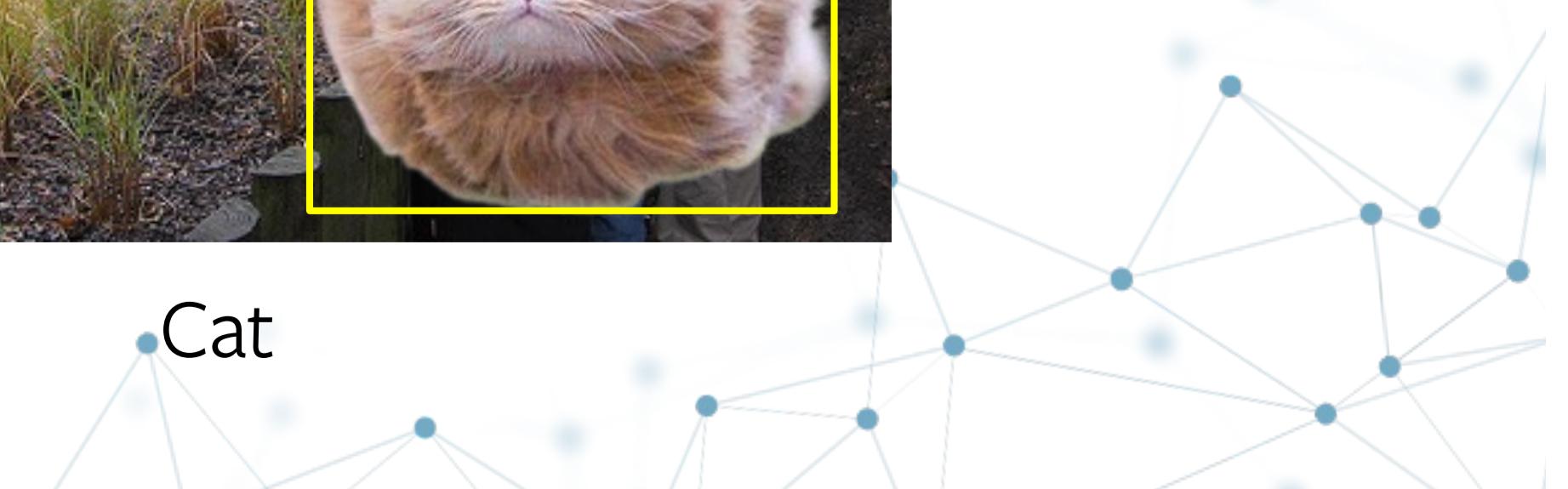
Detection: “where” (i.e., the box) is equivariant with translation
(Contrast to classification: invariant to translation)



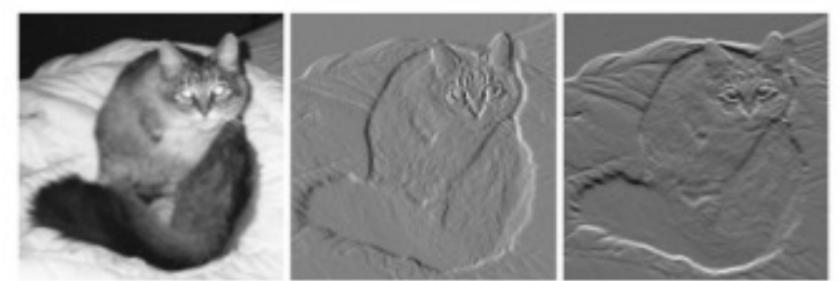
Cat



Cat



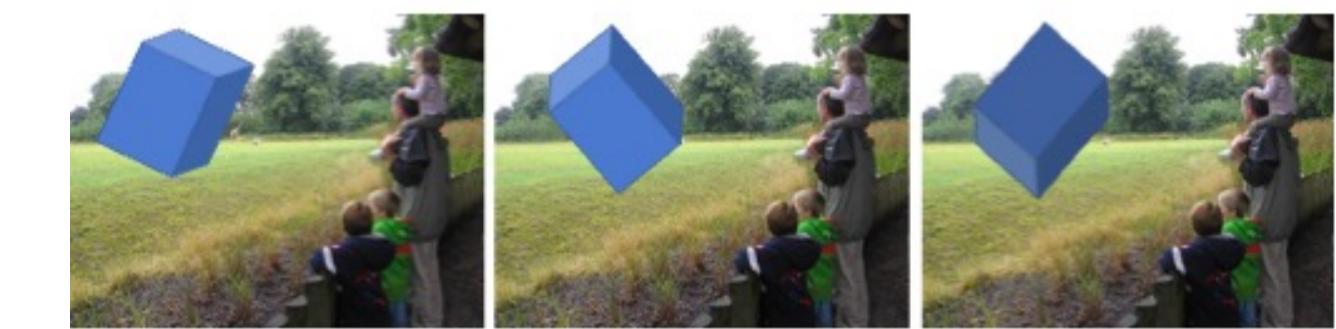
Sources of variation



Photometric



Translation



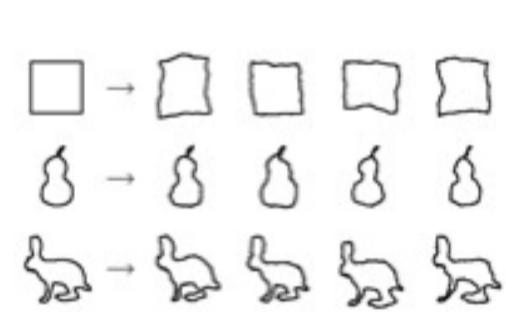
Rotation, rigid pose



Scale



Occlusion (self, other)



Deformation (local, articulated)



Intra-class
(appearance, structural)



Context (clutter,
spurious corr.)

And more...



Imbalance

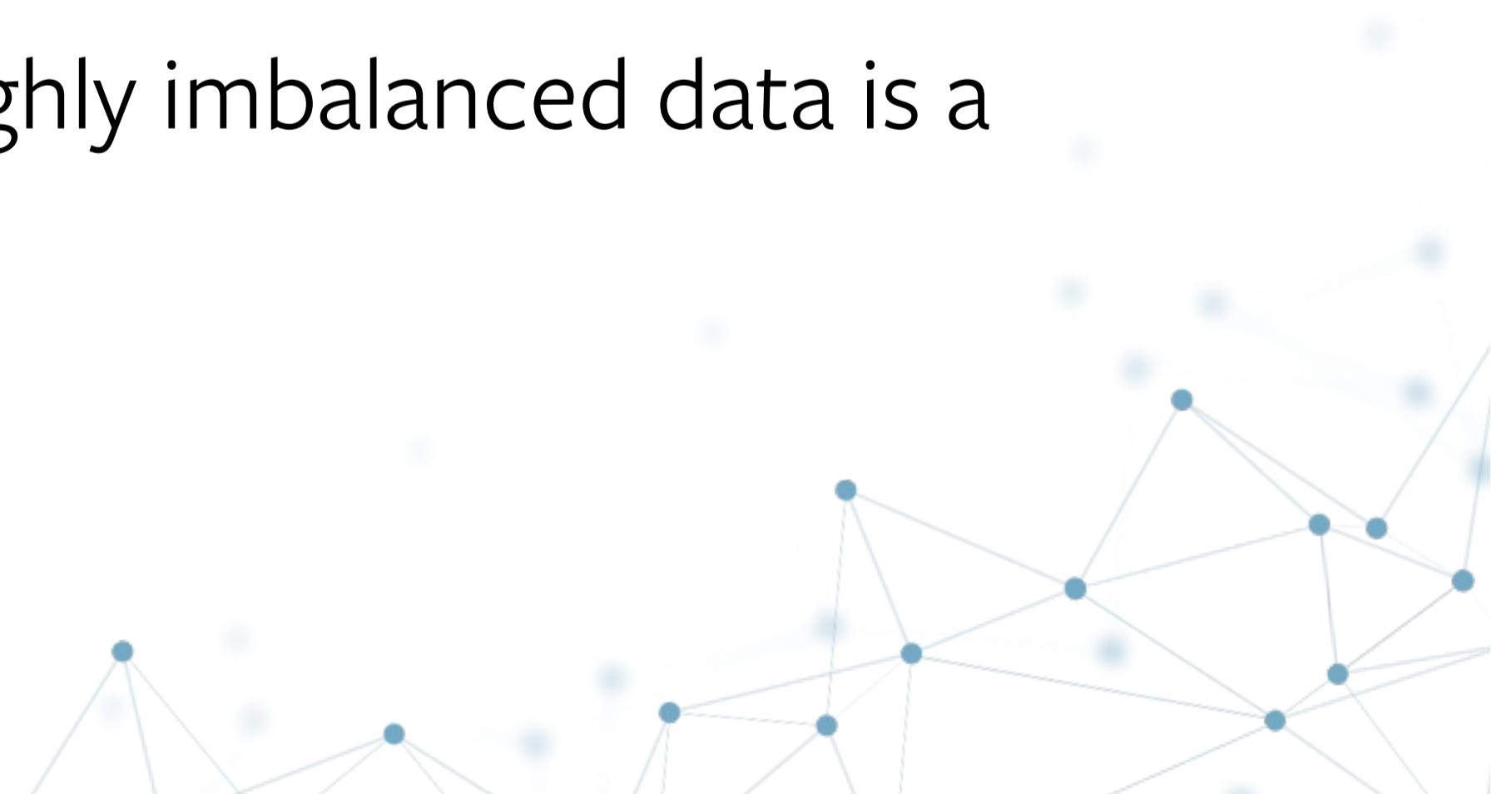
A test set will typically only contain $O(100) - O(1000)$ objects of each category C

There are vastly more possible boxes than this, $\sim O(1e15)$

These boxes belong to an implicit “background” category

Classification of category C is effectively C vs. “background” (and other categories)

From a machine learning perspective, learning with highly imbalanced data is a difficult problem



Recap: Why is object detection difficult?

- A good detector needs a combination of invariance and equivariance properties
 - “What” needs to be invariant to many challenging transformations
 - “Where” needs to be equivariant to many challenging transformations
 - Compared to image classification, which is only about invariance
- Object detection implicitly poses an imbalanced learning problem
 - $O(1e3)$ objects vs. $O(1e15)$ background boxes
 - Image classification is typically balanced

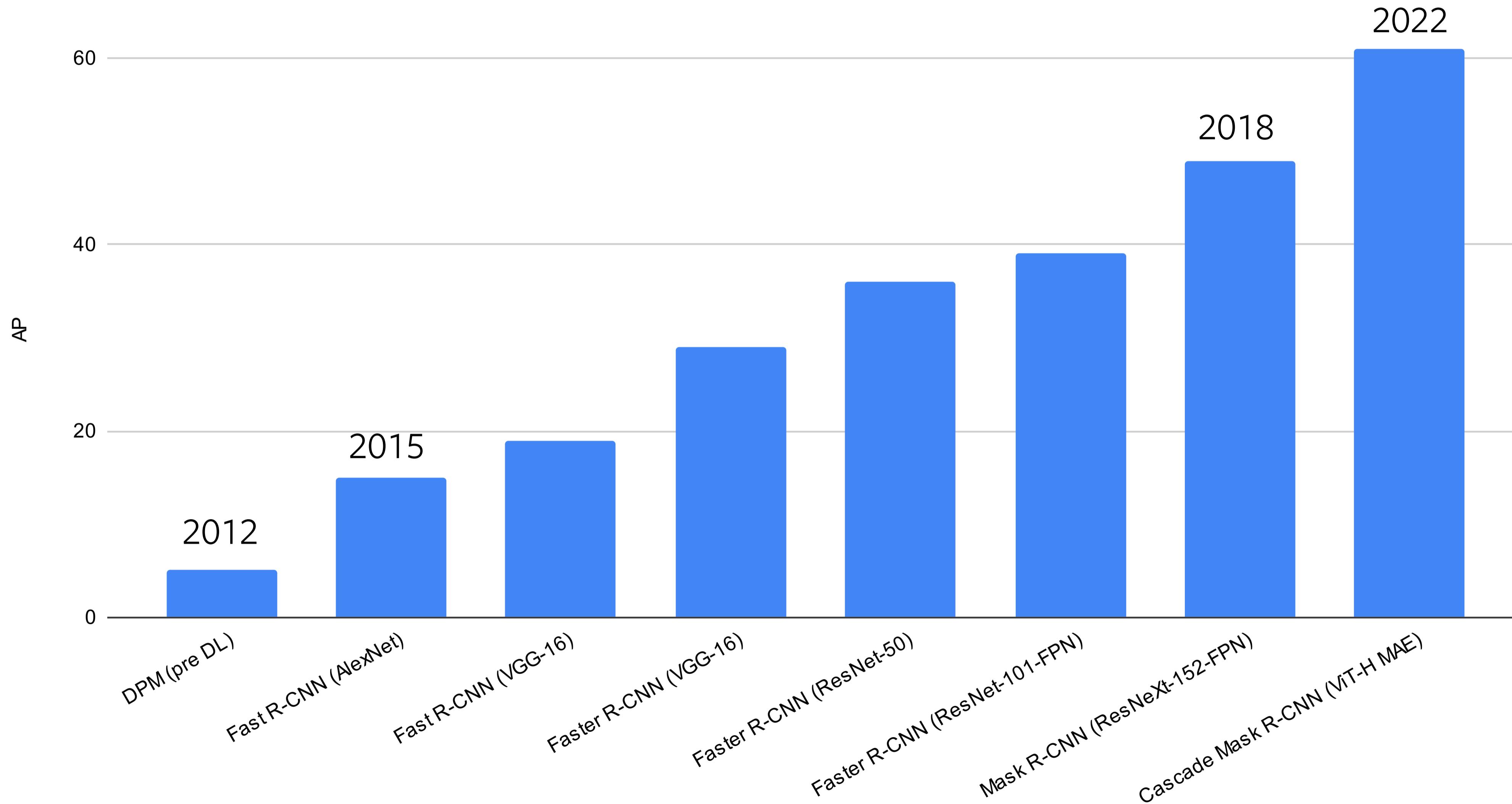


Popular object detection datasets

- PASCAL VOC (retired)
 - The original high-quality, large-scale (at the time) detection dataset
 - 21 object categories
 - ~10k – 20k images
- COCO
 - Leap forward in number of image; established the instance segmentation task
 - 80 object categories
 - ~160k images
- LVIS
 - Leap forward in number of categories
 - ~1200 object categories; natural long-tailed distribution
 - ~160k images (same images as COCO)
- Open Images
- Objects365



COCO average precision (%) over time



Recap: today's focus

- What is the object detection task?
- How is this task formalized?
- How is this task evaluated?
- Why is the task difficult?
- What are important datasets to know about?

Tomorrow: modern object detectors



This afternoon:

Detectron2 - generalized R-CNN Framework

A common framework for understanding

- R-CNN
- Fast R-CNN
- Faster R-CNN
- Feature Pyramid Networks (FPN) + Faster R-CNN
- Mask R-CNN
- ... and more (e.g., Cascade R-CNN, DensePose, Mesh R-CNN, ...)



Detectron2

PyTorch

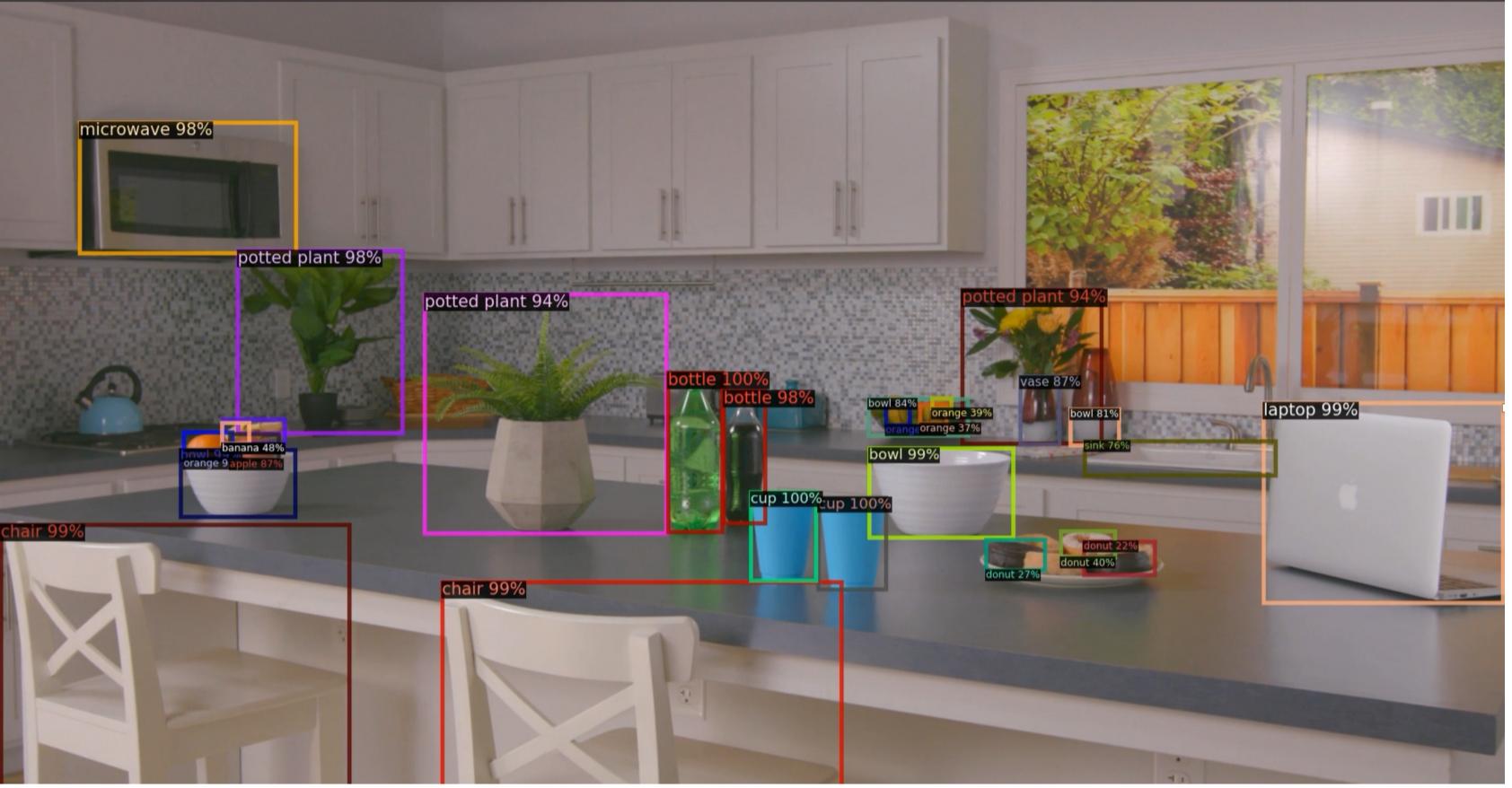
<https://github.com/facebookresearch/detectron2>



This afternoon:

Detectron2 - generalized R-CNN Framework

bounding box
detection



instance
segmentation,

semantic
segmentation



pose
estimation



dense
correspondences



Questions?