# CS 190I
# Deep Learning
# Object Detection

Lei Li (leili@cs)

UCSB

Acknowledgement: Slides borrowed from Bhiksha Raj's 11485 and Mu Li & Alex Smola's 157 courses on Deep Learning, with modification

# Survey Result

## Lecture Topic:

| | | |
|---|---|---|
| A | VAE | 17% |
| B | GAN | 35% |
| C | Diffusion Generative Models | 37% |
| D | Deep Learning for Recommender Systems | 30% |
| E | Graph Neural Network | 42% |
| F | ChatGPT | 65% |

## Industrial Lecture:

| | | |
|---|---|---|
| A | Computer Vision applications in industry | 38% |
| B | NLP applications in industry | 71% |
| C | Recommender systems in industry | 30% |

# NLP Seminar

## Why Task Centricity Matters When You Build for Industry Applications

**Sameena Shah, Ph.D.**

Managing Director, J.P. Morgan Artificial Intelligence Research

Friday, February 24th, 2023

12:00 pm - 1:00 pm

**HH 1010**

**Host:** William Wang

**Abstract:**

The use of AI in finance is gaining traction as organizations realize the advantages of using algorithms to streamline and improve the accuracy of financial tasks. The data, models, algorithms, and solutions are all evolving over time but the task often stays for a longer time. We need to develop a task-centric view of AI. Step through use cases that examine how a task-centric view of AI can be used to minimize financial risk, maximize financial returns, optimize venture capital funding by connecting entrepreneurs to the right investors.

**Bio:**

Sameena Shah is a Managing Director, Artificial Intelligence Research in Digital & Platform Services, where she and the team work across the firm to create Artificial Intelligence technologies for business transformation and growth. She is a highly accomplished leader with over 20 years of educational and industry experience in AI, engineering, data. Her leadership has resulted in award-winning AI technologies that have transformed products and businesses. Read more about Sameena and her accomplishments in the flyer attached below.

# Recap

- Gradient descent can be sped up by incremental updates
  - Convergence is guaranteed under most conditions
    - Learning rate must shrink with time for convergence
  - Stochastic gradient descent: update after each observation. Can be much faster than batch learning
  - Mini-batch updates: update after batches. Can be more efficient than SGD

- Convergence can be improved using smoothed updates
  - AdaGrad, RMSprop, Adam and more advanced techniques

# AdaGrad

- AdaGrad (Duchi, Hazan, and Singer 2010) very popular adaptive method.

$$G_{t+1} = G_t + \nabla \ell(x_t)^2$$

$$x_{t+1} = x_t - \eta \frac{1}{\sqrt{G_{t+1} + \epsilon}} \nabla \ell(x_t)$$

element-wise

- Benefits:
  - AdaGrad does not require tuning learning rate $\eta$
  - Actual learning rate will decrease
  - Can drastically improve over SGD

# Image classification

This lecture

Object Detection
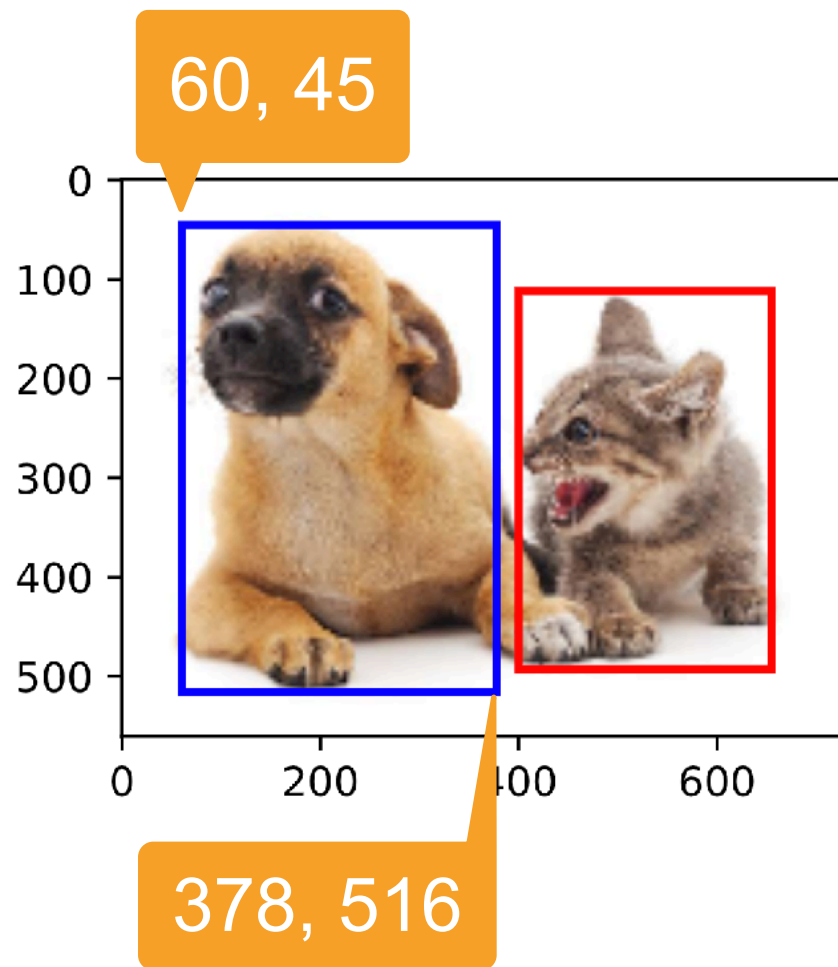
Dog



Critical in applications: autonomous driving

# Locating the Object: Bounding Box

- A bounding box can be defined by 4 numbers,
  - (top-left x, top-left y, bottom-right x, bottom-right y)
  - (top-left x, top-right y, width, height)

60, 45

378, 516

# Object Detection Dataset



- Each row present an object
  - Image_filename, object_category, bounding box
- PASCAL VOC
  - 11.53k images, 27.45k objects, 20 classes
- COCO ([cocodataset.org](cocodataset.org))
  - 80 object classes
  - 330K images
  - 1.5M objects

# Object Detection Dataset

- Open Image (v6):
  - 9M images,
  - 1.9M images with 16M bounding boxes, 600 classes
  - Includes 3.3M visual relations (of 1466 types)
  - [https://storage.googleapis.com/openimages/web/factsfigures.html](https://storage.googleapis.com/openimages/web/factsfigures.html)
- BDD100k
  - 100k videos in driving scenario
  - [https://github.com/bdd100k/bdd100k](https://github.com/bdd100k/bdd100k)



BDD100K: A Diverse Driving Dataset for Heterogeneous Multitask Learning

# Anchor Boxes

- A detection algorithm often

  – Proposes multiple regions, called anchor boxes

  – Predict if an anchor box contains an object

  – If yes, predict the offset from the anchor box to the ground truth bounding box

# IoU - Intersection over Union

- IoU measures the similarity between two boxes
  - 0 means no-overlapping
  - 1 means identical

- It's an especial case of Jacquard index
  - Given sets $A$ and $B$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$IoU = \frac{\text{Intersection}}{\text{Union}}$$

# Assign Labels to Anchor Boxes

- Each anchor box is a training example

- Label each anchor box with
  - Background
  - Associate with a bounding box

- We may generate a large amount of anchor boxes
  - Leads to a large portion of negative examples

# Assign Labels to Anchor Boxes

# Output with non-maximum suppression (NMS)

- Each anchor box generates one bounding box prediction

- Select the one with the highest score (not background)

- Remove all other predictions with IoU > $\theta$ compared to the selected one

- Repeat until all are selected or removed

# Region-based CNNs

# R-CNN



1. Input image   2. Extract region proposals (~2k)   3. Compute CNN features   4. Classify regions

- Select anchor boxes with a heuristic algorithm
- Use a pre-trained networks to extract features for each anchor box
  - Adding classifier layer
  - and regression layer to predict bounding boxes

# Region of Interest (RoI) Pooling

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

2 x 2 RoI Pooling

| 5 | 6 |
|---|---|
| 9 | 10 |

- Given an anchor box, uniformly cuts it into $n \times m$ blocks, output the maximal value in each block

- Returns $nm$ values for each anchor box

- A special case of maxpooling

# Fast RCNN



- A CNN to extract features

- Siding windows on the feature maps

- RoI pooling returns fixed length feature for each anchor box

# Faster R-CNN



- Use a region proposal network to replace select search for high quality anchor boxes

# Faster R-CNN



- Use a region proposal network to replace select search for high quality anchor boxes

**Faster RCNN**

https://gluon-cv.mxnet.io/model_zoo/detection.html

# Single Shot Multibox Detection (SSD)

# Generate Anchor Boxes

- For each pixel, generate multiple anchor boxes centered at this pixel

- Given $n$ sizes $s_1, \ldots, s_n$ and $m$ ratios $r_1, \ldots, r_m$, generate $n+m-1$ anchor boxes

$$(s_1, r_1), (s_2, r_1), \ldots, (s_n, r_1), (s_1, r_2), \ldots, (s_1, r_m)$$

# SSD Model

- A base network to extract feature, followed by conv-blocks to halve width and height

- Generate anchor boxes at each sale
  - Bottom for small objects and top for large objects

- Predict class and bounding box for each anchor box

# Quiz

- https://edstem.org/us/courses/31035/lessons/56184/slides/318357

# You Only Look Once (YOLO)

Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi

CVPR 2016

# YOLO

- Anchor boxes are highly overlapped in SSD
- YOLO cuts the input image uniformly into $S$ x $S$ anchor boxes
- Each anchor box predicts $B$ bounding boxes



S × S grid on input

# Each cell predicts boxes and confidences: P(Object)



Redmon. et al. 2016

# Each cell predicts boxes and confidences: P(Object)



normalize
(x,y,w,h)
P

# Each cell predicts boxes and confidences: P(Object)



(x,y,w,h)
P

Redmon. et al.
2016

# Each cell predicts boxes and confidences: P(Object)



Redmon. et al.
2016

# Each cell predicts boxes and confidences: P(Object)



(x,y,w,h)
P

(x,y,w,h)
P

Redmon. et al.
2016

34

# Each cell predicts boxes and confidences: P(Object)



Redmon. et al.
2016

# Each cell also predicts a class probability.



Redmon. et al. 2016

# Each cell also predicts a class probability.



Bicycle

Car
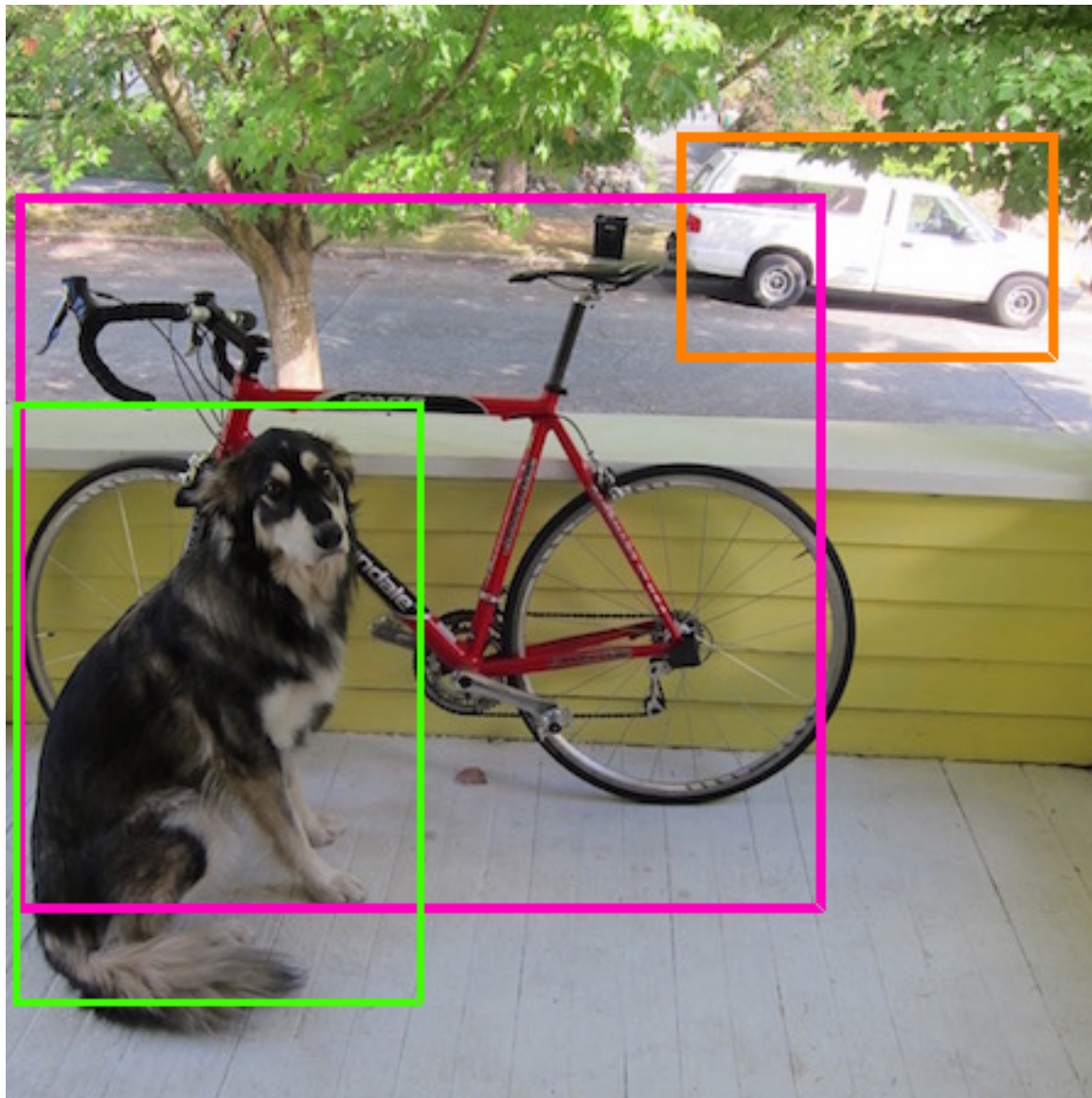
Dog

Dining Table

# Conditioned on object: P(Car I Object)

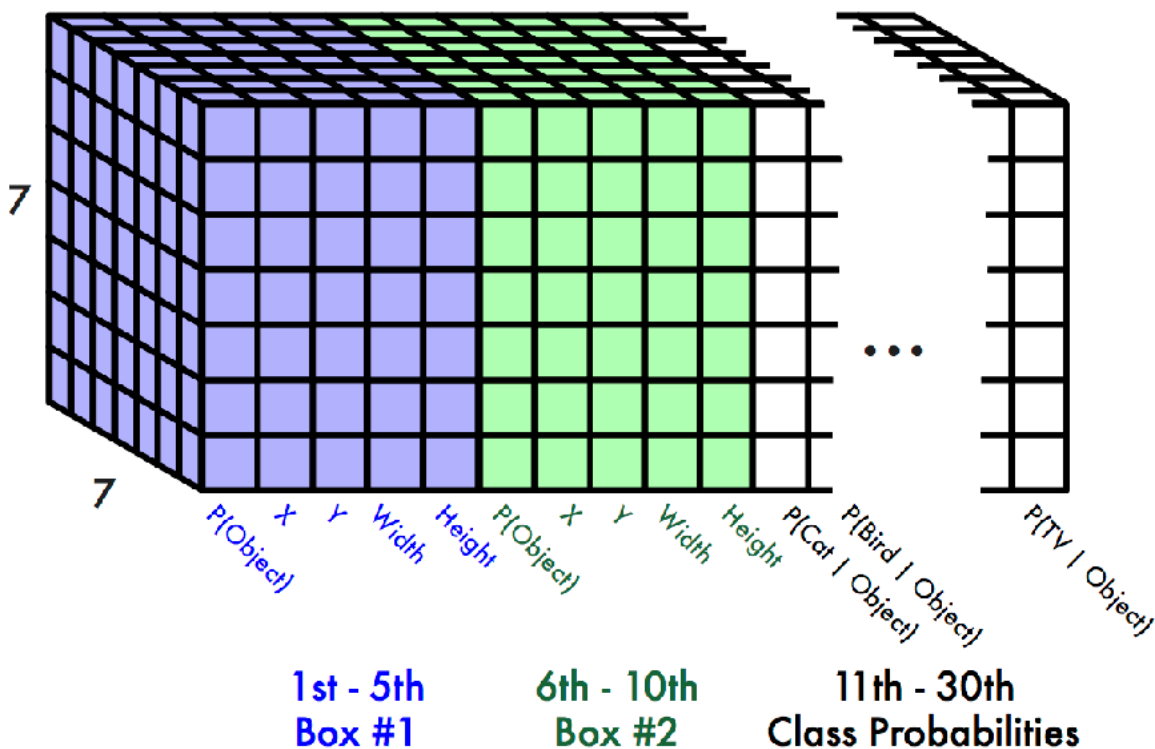# Then we combine the box and class predictions.

# Finally we do NMS and threshold detections

# The output

Each cell predicts:

- For each bounding box:
    - 4 coordinates (x, y, w, h)
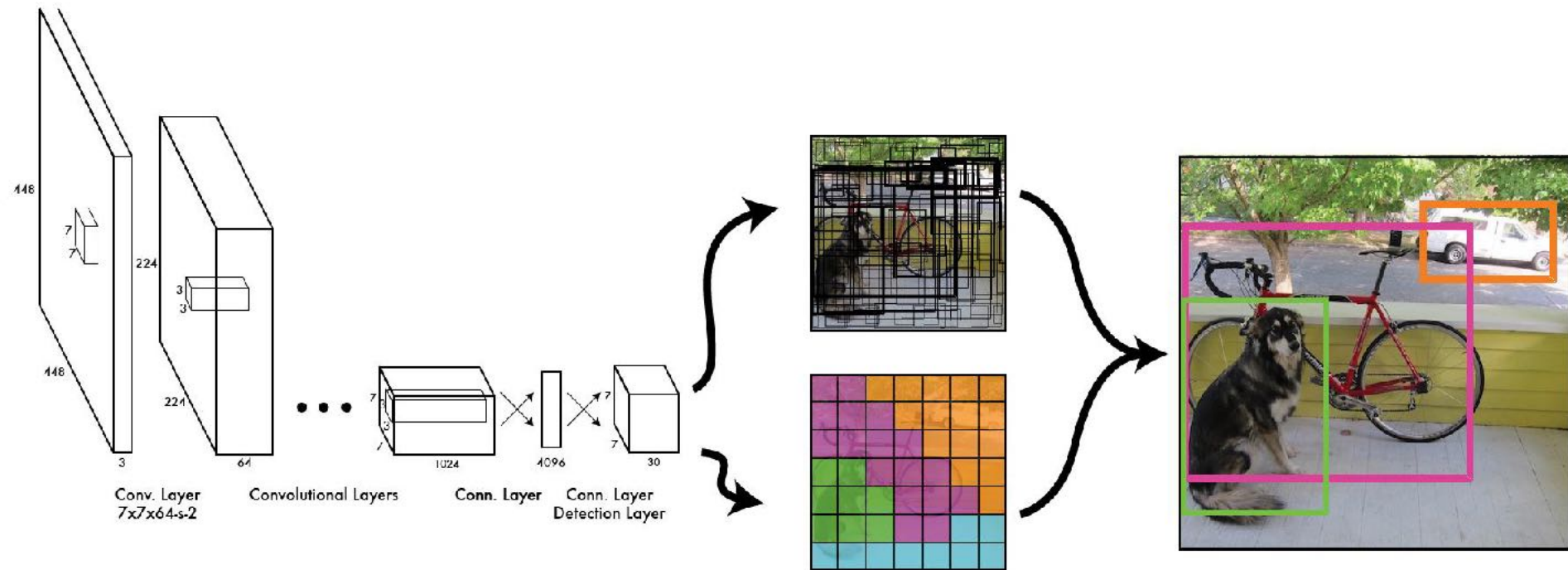    - 1 confidence value
- Some number of class probabilities

For Pascal VOC:
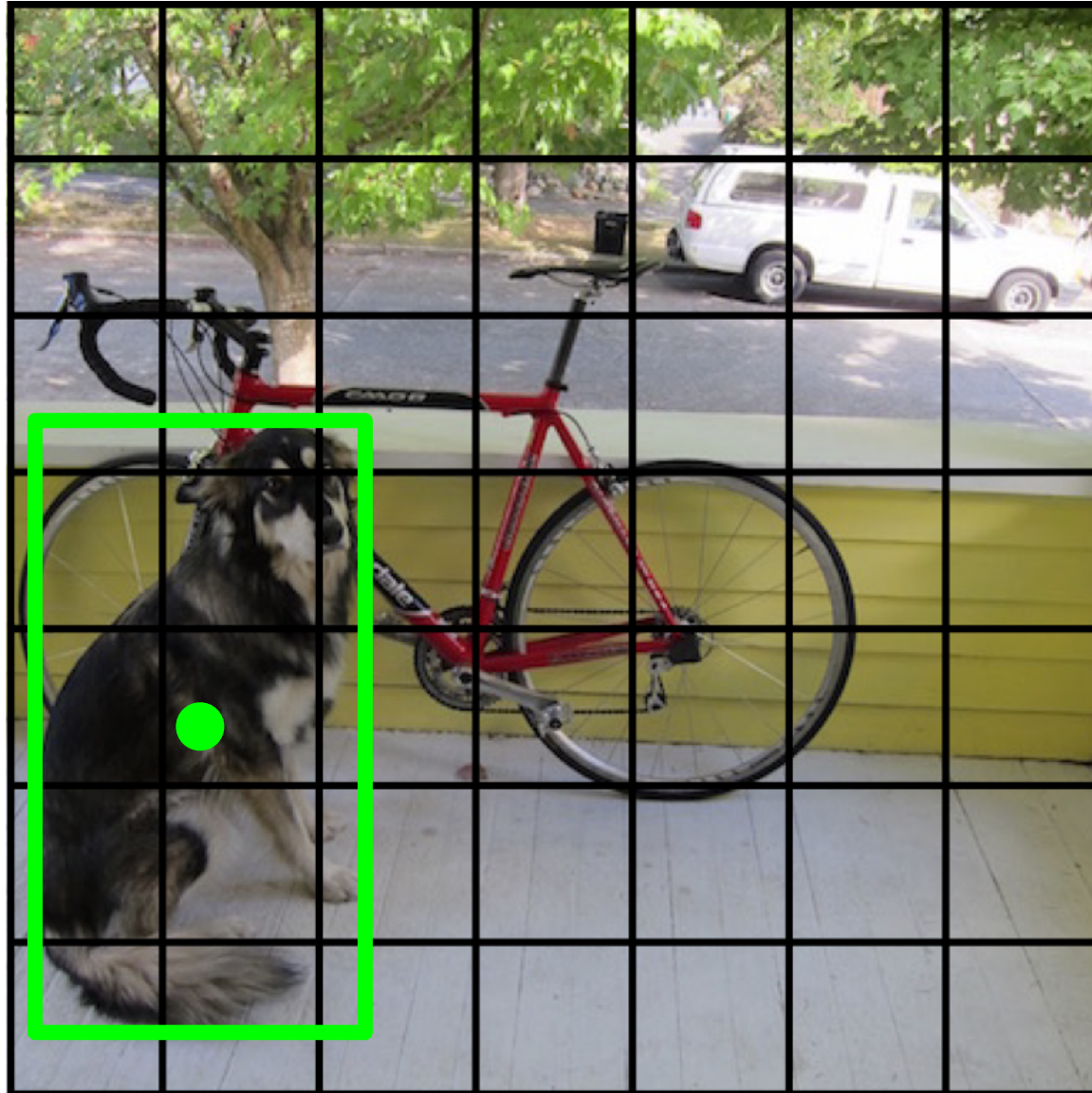
- 7x7 grid
- 2 bounding boxes / cell
- 20 classes



7 x 7 x (2 x 5 + 20) = 7 x 7 x 30 tensor = **1470 outputs**

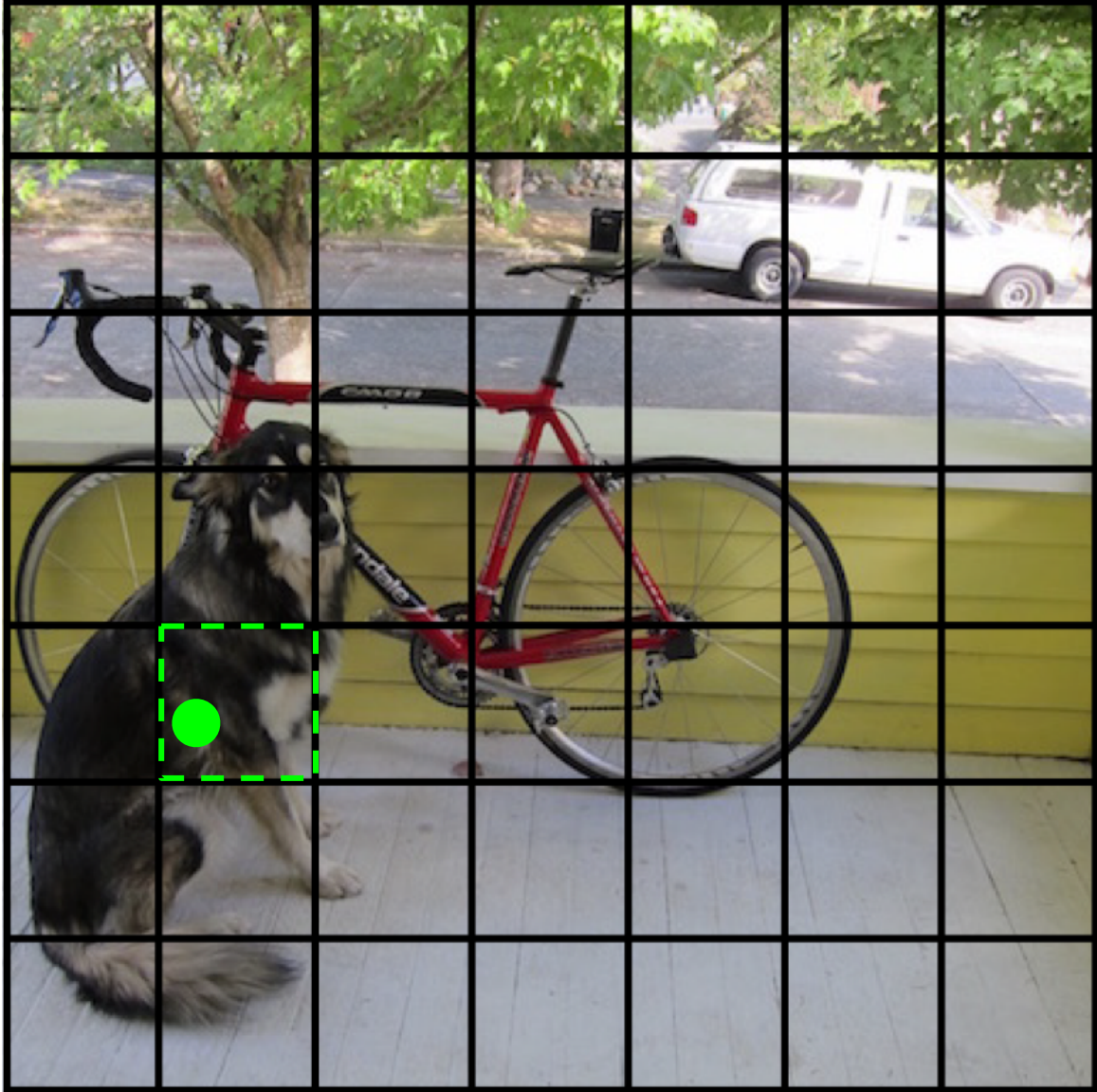# A single pipeline for detection



backbone network: VGG16, ResNet101, …

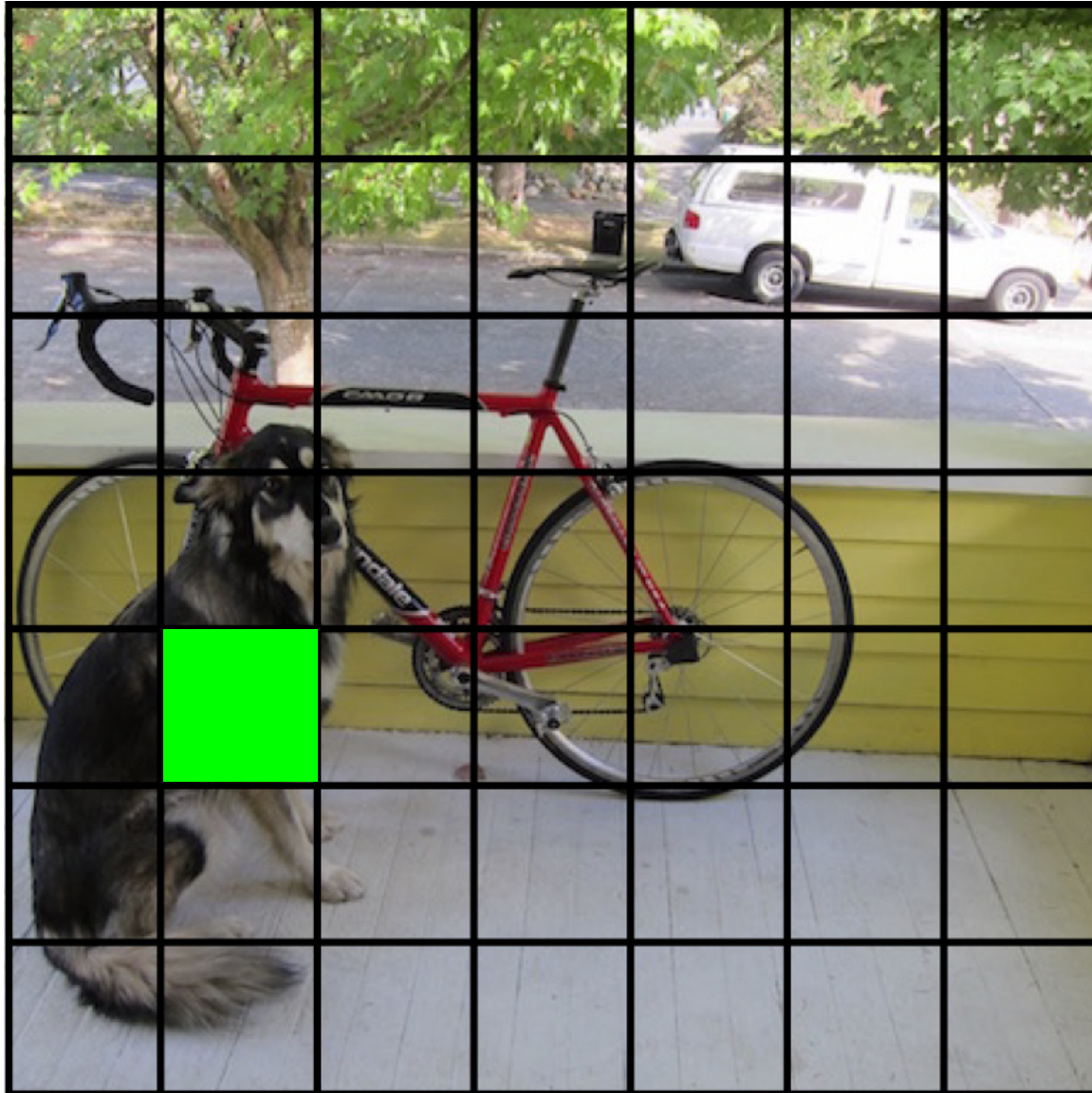# During training, match example to the right cell



center of object

# During training, match example to the right cell
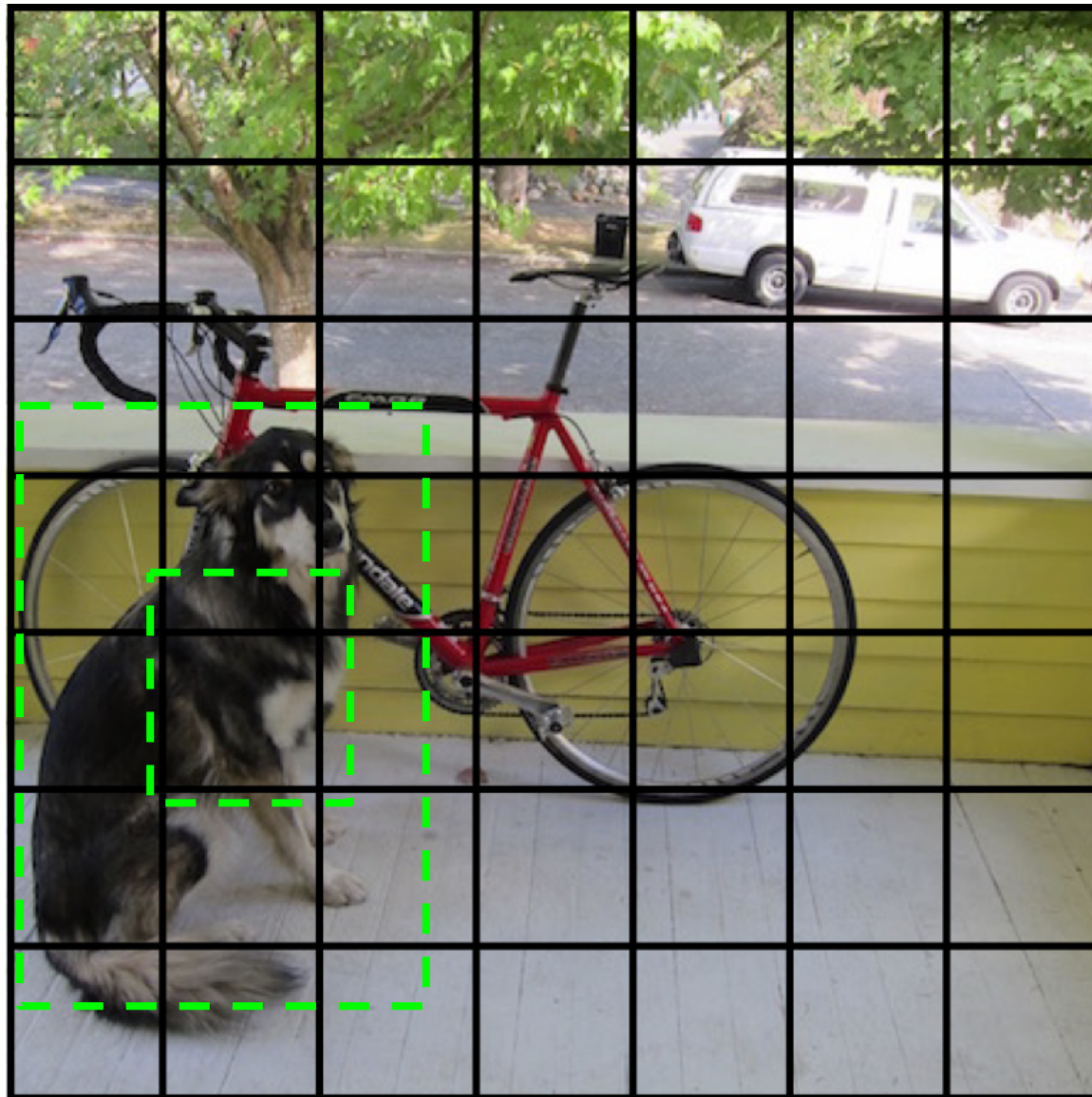
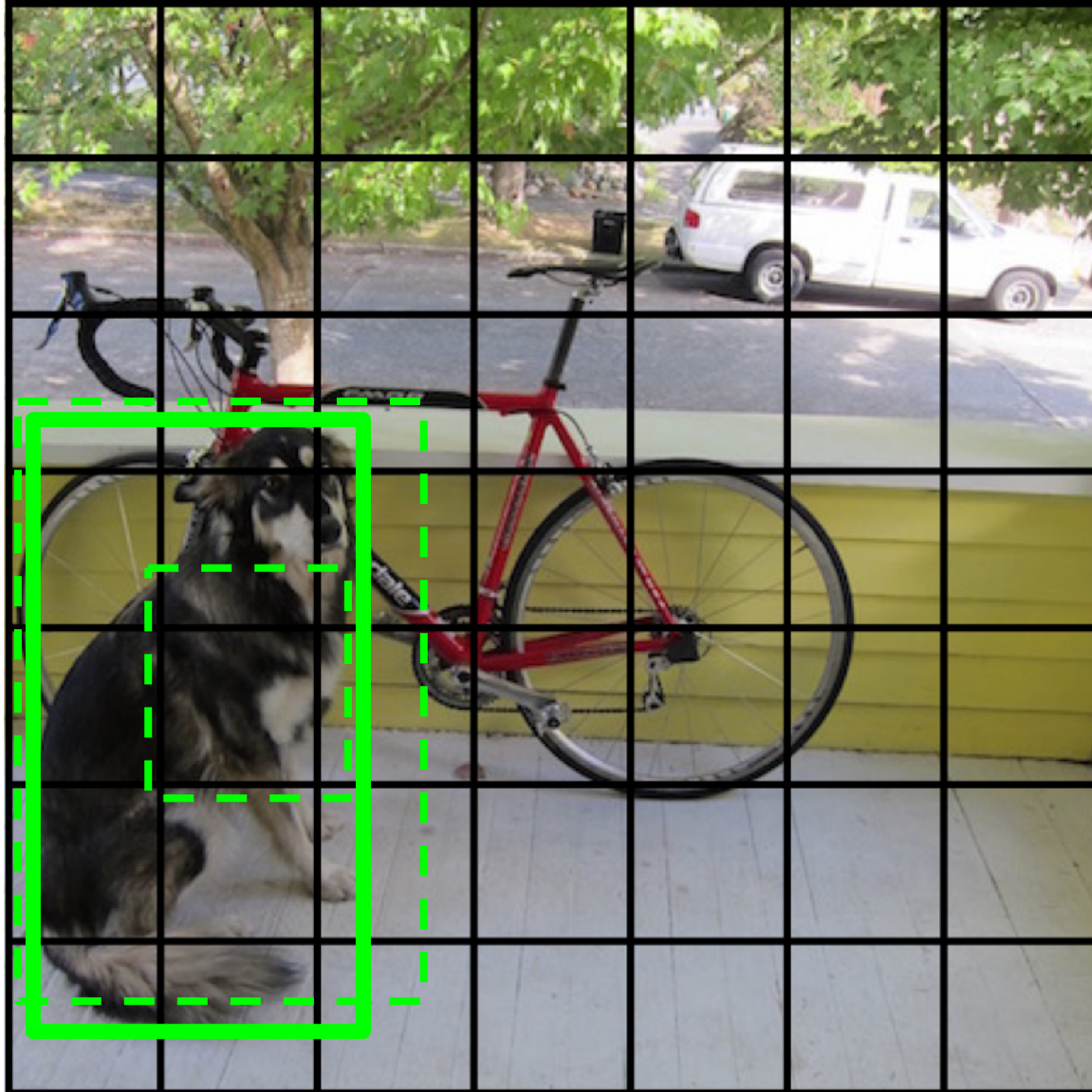# Adjust that cell's class prediction
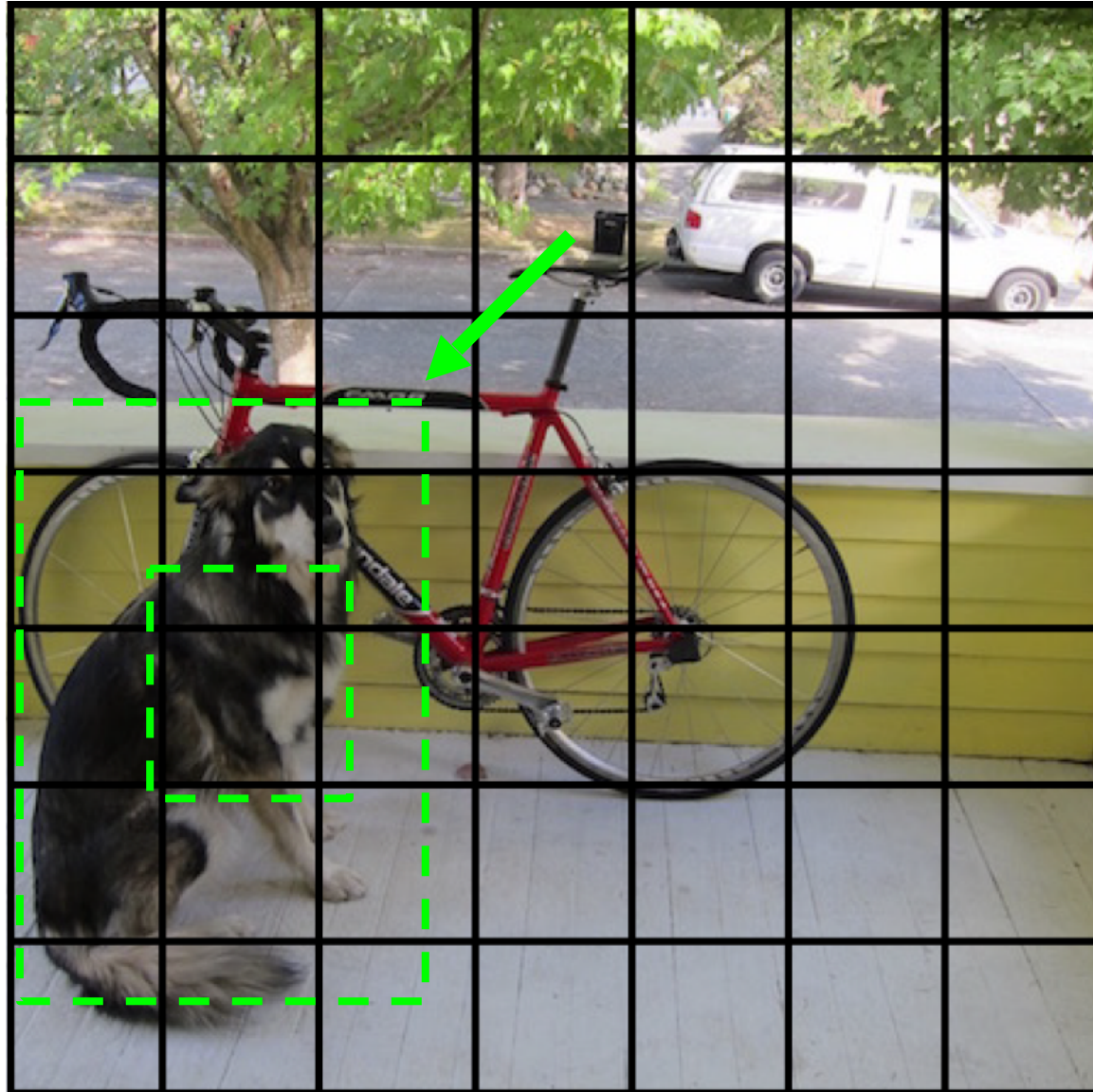


**Dog = 1**
Cat = 0
Bike = 0
...

# Look at that cell's predicted boxes

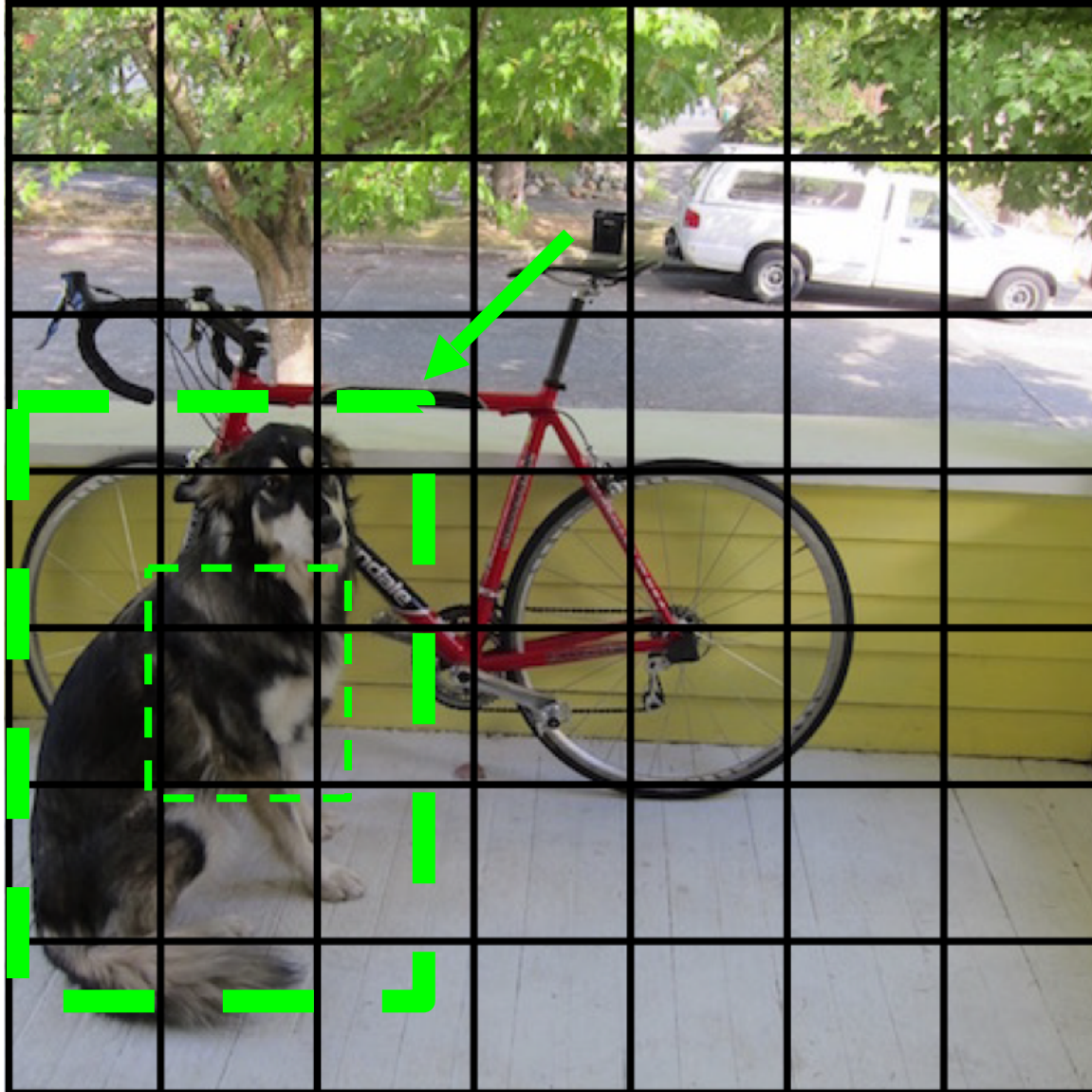# Find the best one, adjust it, increase the confidence

**Find the best one, adjust it, increase the confidence**

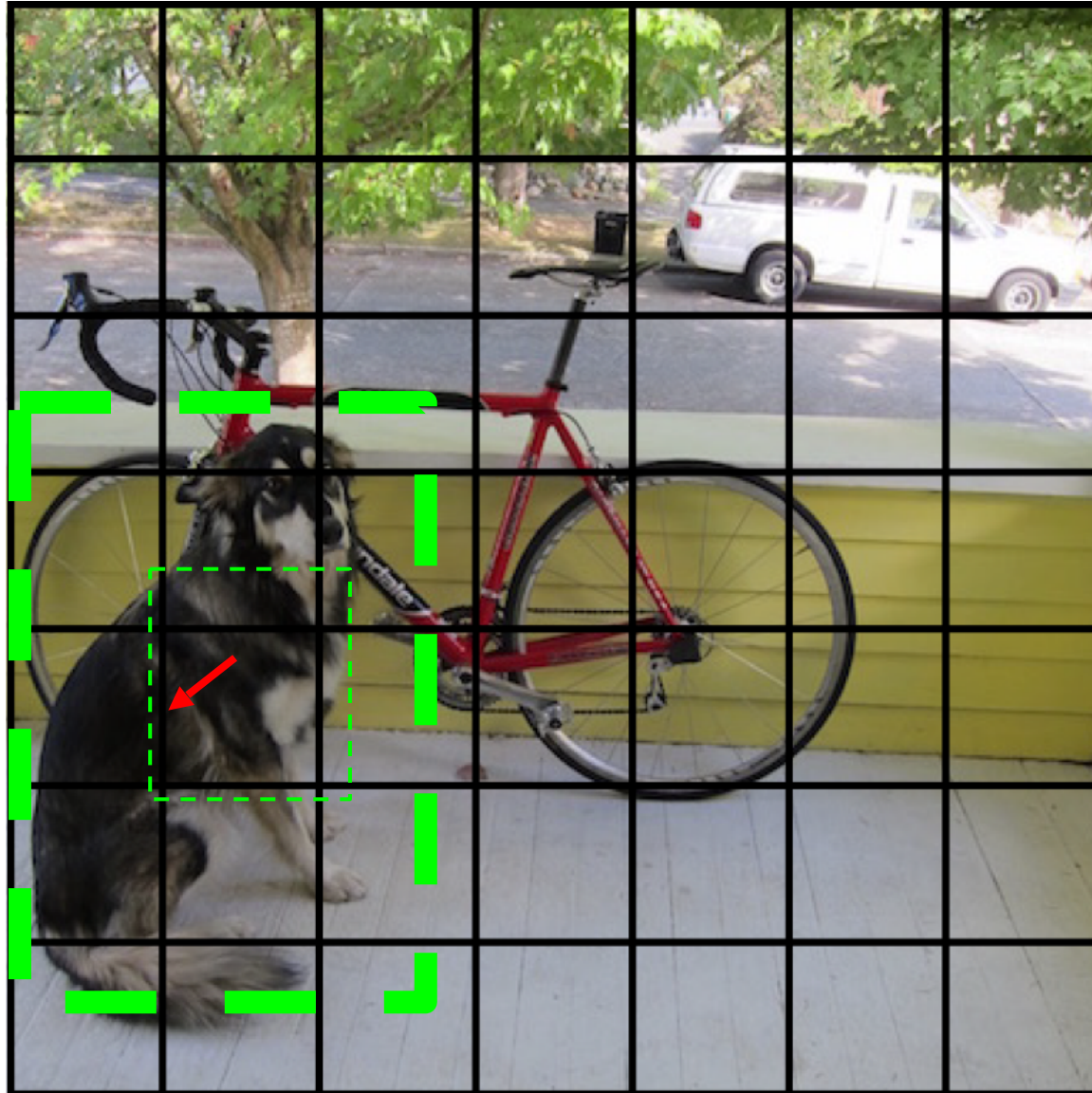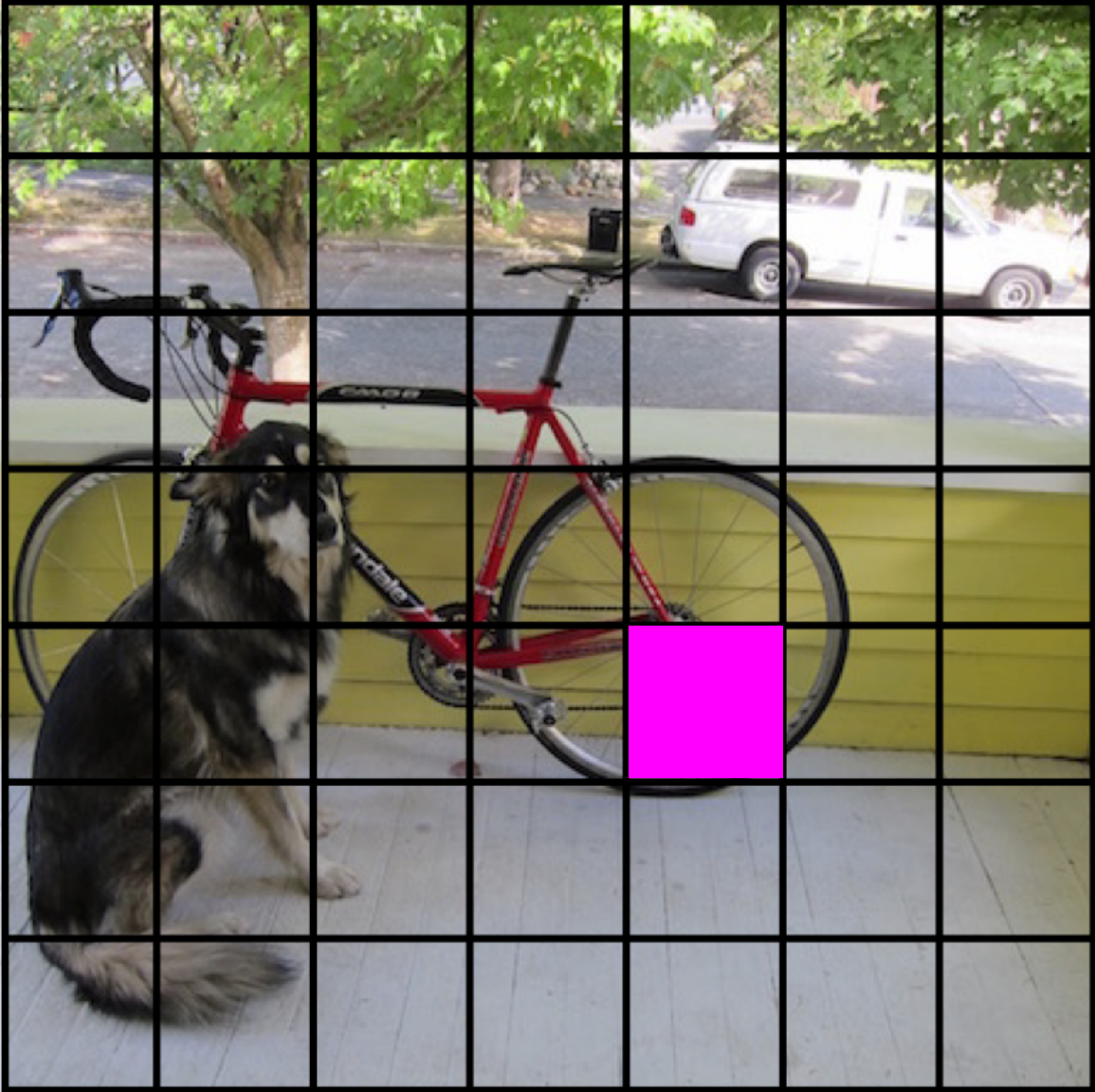# Find the best one, adjust it, increase the confidence
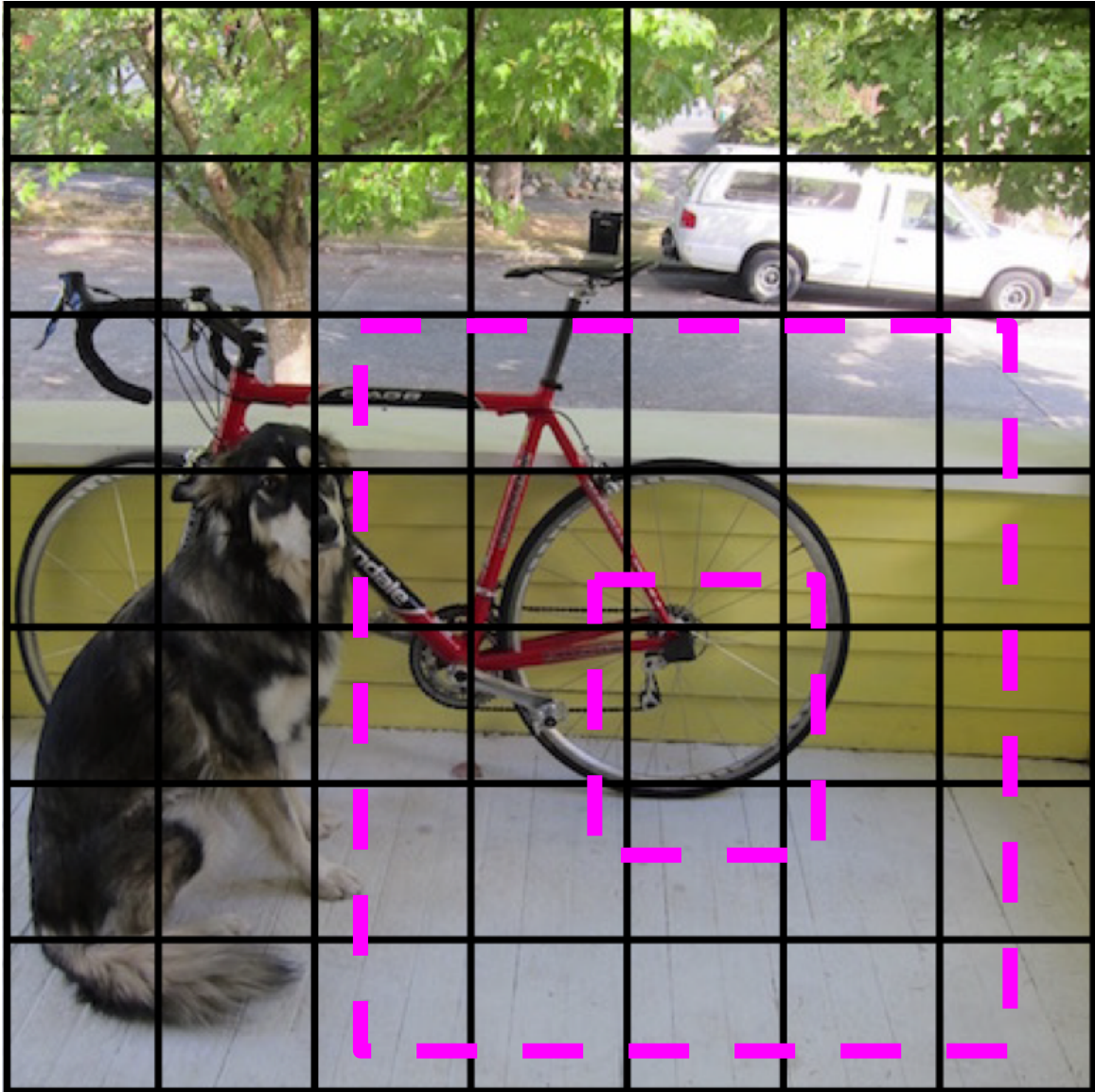
# Decrease the confidence of other boxes

# Decrease the confidence of other boxes

**Some cells don't have any ground truth detections!**

# Some cells don't have any ground truth detections!

# Decrease the confidence of these boxes

# Decrease the confidence of these boxes

# Don't adjust the class probabilities or coordinates

Conv. Layer
7x7x64-s-2
Maxpool Layer
2x2-s-2

Conv. Layer
3x3x192
Maxpool Layer
2x2-s-2

Conv. Layers
1x1x128
3x3x256
1x1x256
3x3x512
Maxpool Layer
2x2-s-2

Conv. Layers
$\left.\begin{array}{l} 1\text{x}1\text{x}256 \\ 3\text{x}3\text{x}512 \end{array}\right\}\times 4$
1x1x512
3x3x1024
Maxpool Layer
2x2-s-2

Conv. Layers
$\left.\begin{array}{l} 1\text{x}1\text{x}512 \\ 3\text{x}3\text{x}1024 \end{array}\right\}\times 2$
3x3x1024
3x3x1024-s-2

Conv. Layers
3x3x1024
3x3x1024

Conn. Layer

Conn. Layer

# Training YOLO

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

if i-th cell contain object and j-th box has max IoU

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2$$

if i-th cell contain object and j-th box has max IoU

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

# Other tricks

- Pretraining on Imagenet
- SGD with decreasing learning rate
- Extensive data augmentation

# YOLO works across a variety of natural images

# It also generalizes well to new domains (like art)

# YOLO outperforms methods like DPM and R-CNN when generalizing to person detection in artwork



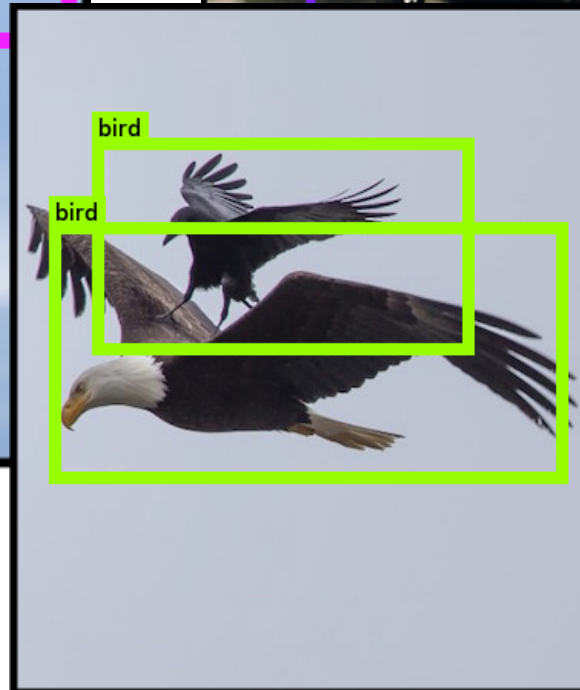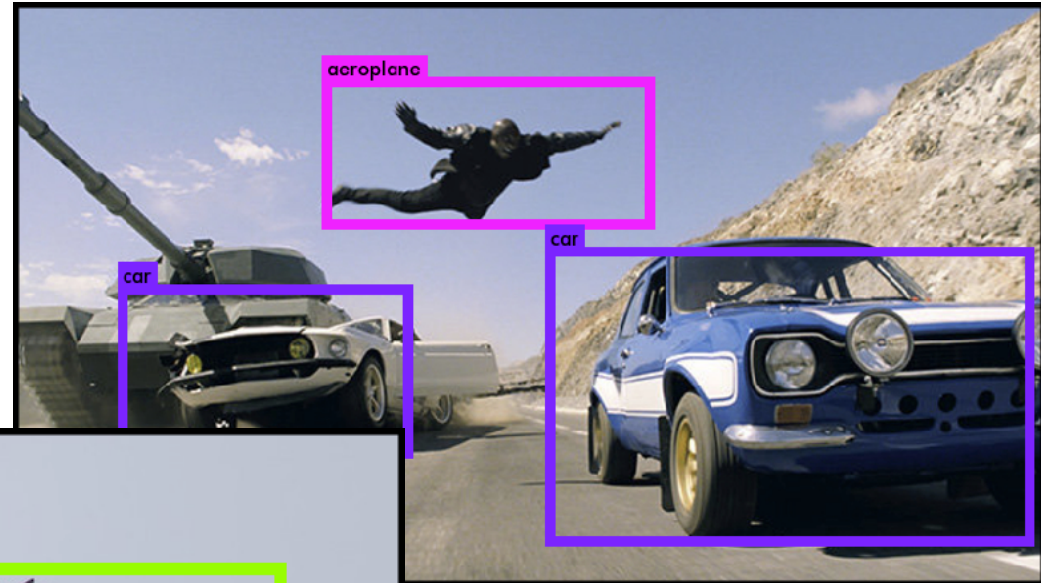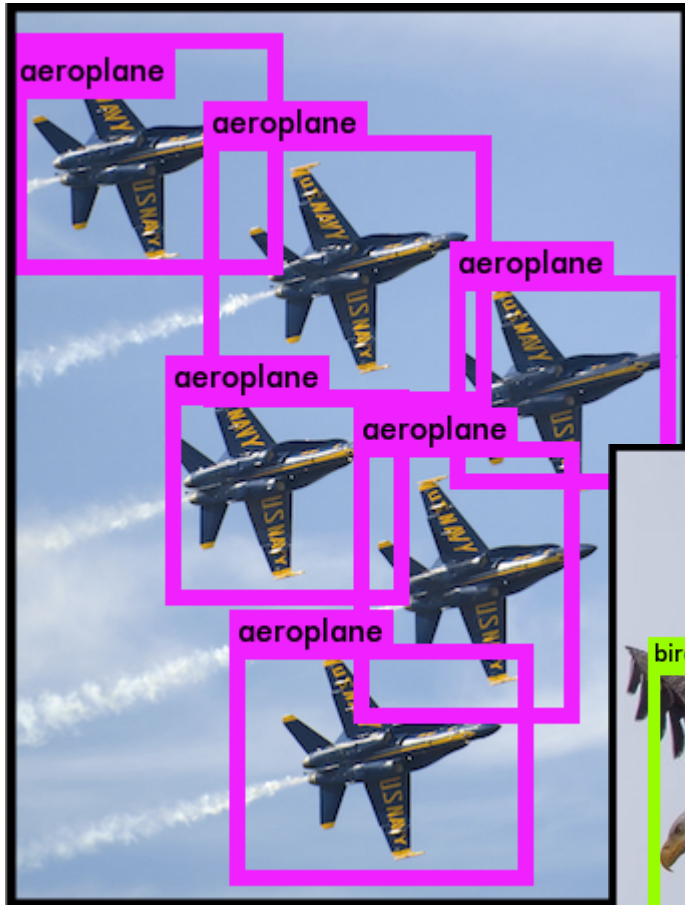| | VOC 2007 | Picasso | | People-Art |
|---|---|---|---|---|
| | AP | AP | Best $F_1$ | AP |
| **YOLO** | **59.2** | **53.3** | **0.590** | **45** |
| R-CNN | 54.2 | 10.4 | 0.226 | 26 |
| DPM | 43.2 | 37.8 | 0.458 | 32 |

*S. Ginosar, D. Haas, T. Brown, and J. Malik. Detecting people in cubist art. In Computer Vision-ECCV 2014 Workshops, pages 101–116. Springer, 2014.*

*H. Cai, Q. Wu, T. Corradi, and P. Hall. The cross-depiction problem: Computer vision algorithms for recognising objects in artwork and in photographs.*

# Results: Performance vs Speed

| | Pascal 2007 mAP | Speed | |
| --- | --- | --- | --- |
| DPM v5 | 33.7 | .07 FPS | 14 s/img |
| R-CNN | 66.0 | .05 FPS | 20 s/img |
| Fast R-CNN | 70.0 | .5 FPS | 2 s/img |
| Faster R-CNN | 73.2 | 7 FPS | 140 ms/img |
| YOLO | 69.0 | 45 FPS | 22 ms/img |

# **YOLO Series**

- YOLO
- YOLOv2 improves the detection of small objects in groups and the localization accuracy.
  - and adding batch norm
- YOLOv3,
  - 106 layer resnet
  - multi-scale detection (three scales)
- YOLOv4, …

# **Additional Tricks: Mixup**

- Apply to object detection as well

# Results for YOLOv3

| Incremental Tricks | mAP | Δ | Cumu Δ |
|---|---|---|---|
| - data augmentation | 64.26 | -15.99 | -15.99 |
| baseline | 80.25 | 0 | 0 |
| + synchronize BN | 80.81 | +0.56 | +0.56 |
| + random training shapes | 81.23 | +0.42 | +0.98 |
| + cosine lr schedule | 81.69 | +0.46 | +1.44 |
| + class label smoothing | 82.14 | +0.45 | +1.89 |
| + mixup | **83.68** | **+1.54** | **+3.43** |

Zhi et al, *Bag of Freebies for Training Object Detection Neural Networks*

# Summary

- Object Detection
  - RCNN
  - YOLO: single pipeline model (e2e) for object detection

# Next Up

- Recurrent neural networks

- Friday talk on NLP in industry

**NLP Seminar**

**Why Task Centricity Matters When You Build for Industry Applications**

**Sameena Shah, Ph.D.**

Managing Director, J.P. Morgan Artificial Intelligence Research

Friday, February 24th, 2023

12:00 pm - 1:00 pm

**HH 1010**

**Host:** William Wang

**Abstract:**

The use of AI in finance is gaining traction as organizations realize the advantages of using algorithms to streamline and improve the accuracy of financial tasks. The data, models, algorithms, and solutions are all evolving over time but the task often stays for a longer time. We need to develop a task-centric view of AI. Step through use cases that examine how a task-centric view of AI can be used to minimize financial risk, maximize financial returns, optimize venture capital funding by connecting entrepreneurs to the right investors.

**Bio:**

Sameena Shah is a Managing Director, Artificial Intelligence Research in Digital & Platform Services, where she and the team work across the firm to create Artificial Intelligence technologies for business transformation and growth. She is a highly accomplished leader with over 20 years of educational and industry experience in AI, engineering, data. Her leadership has resulted in award-winning AI technologies that have transformed products and businesses. Read more about Sameena and her accomplishments in the flyer attached below.