Topics:

- Advanced Architectures: Segmentation and Detection
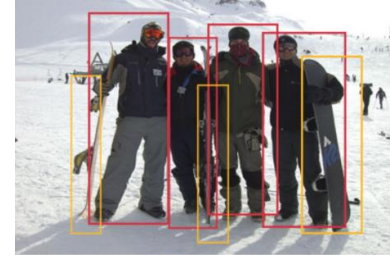
# CS 4644-DL / 7643-A
# ZSOLT KIRA

- **Assignment 3**
  - Due **March 8th 11:59pm EST**

- **Projects**
  - Project check-in due **March 14th**

- Meta office hours Friday 3pm ET on attention models

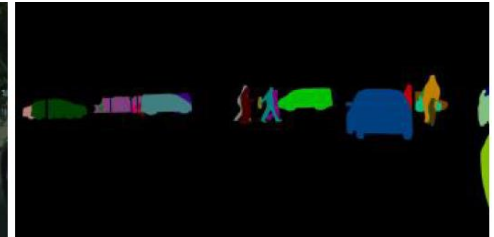**Classification**
(Class distribution per image)

**Object Detection**
(List of bounding boxes with class distribution per box)
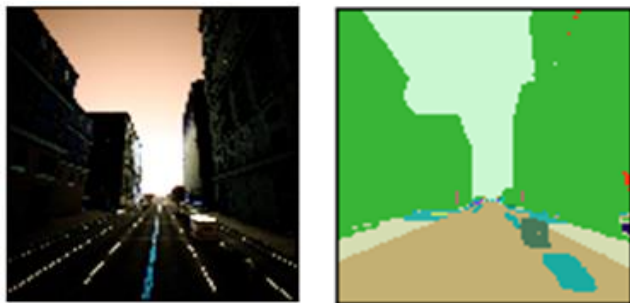
**Semantic Segmentation**
(Class distribution per pixel)

**Instance Segmentation**
(Class distribution per pixel with unique ID)

**Computer Vision Tasks**

Georgia Tech

# Given an image, output another image

- Each output contains class distribution per pixel
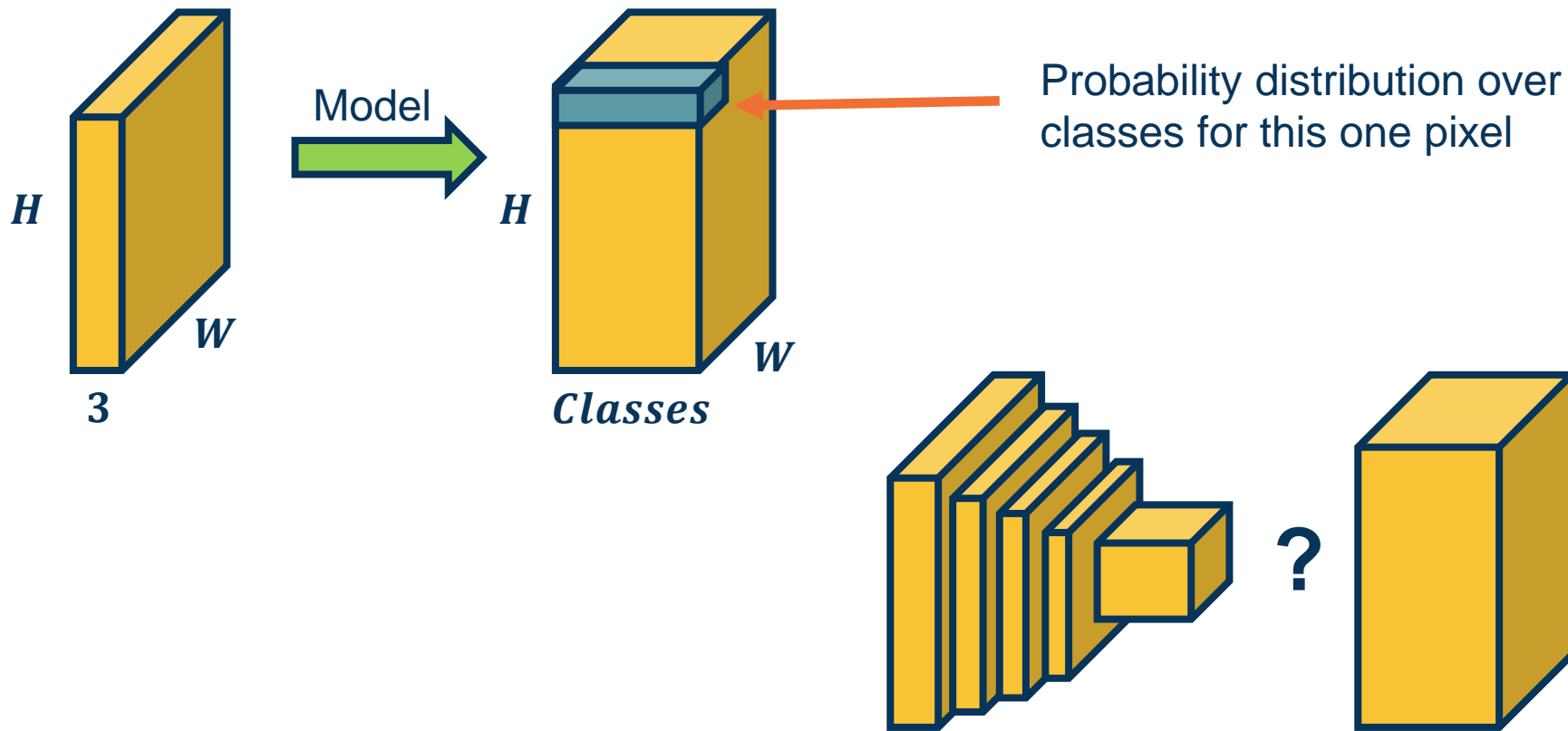
- More generally an image-to-image problem
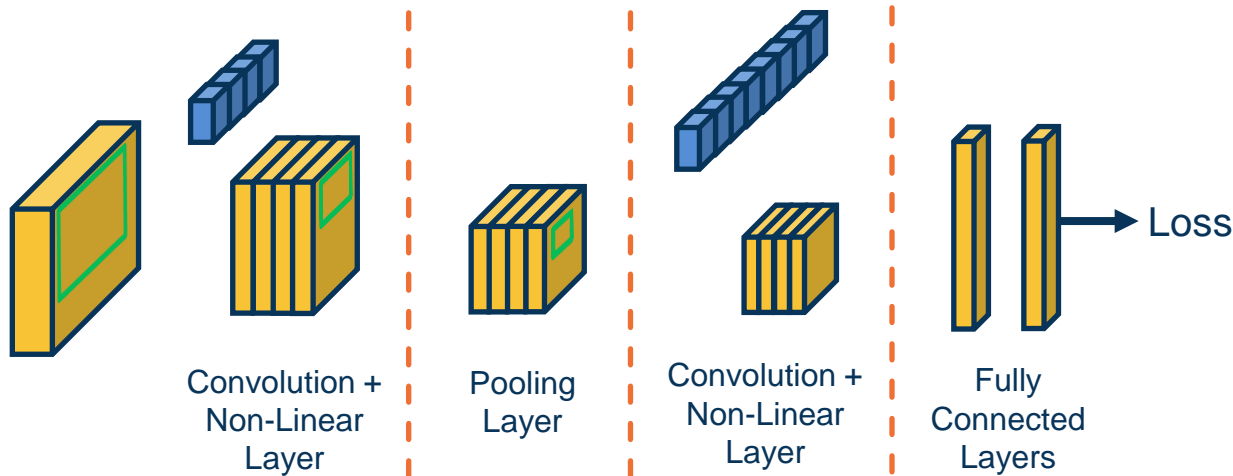


**Semantic Segmentation**
(Class distribution per pixel)



**Instance Segmentation**
(Class distribution per pixel with unique ID)

**Segmentation Tasks**

Georgia Tech

Model

H

W

3

H

W

*Classes*

Probability distribution over classes for this one pixel
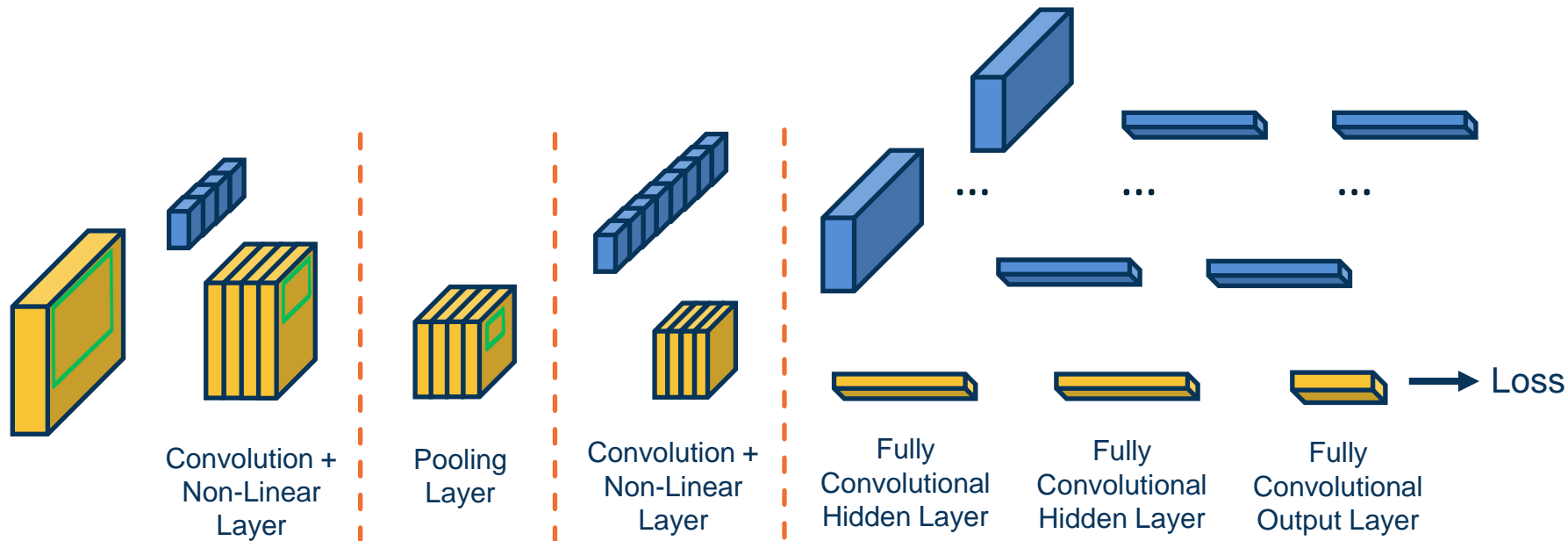
?

Georgia Tech

Fully connected layers no longer explicitly retain spatial information (though the network can still learn to do so)

**Idea: Convert fully connected layer to convolution!**

**Idea 1: Fully-Convolutional Network**

Georgia Tech

Convolution + Non-Linear Layer

Pooling Layer

Convolution + Non-Linear Layer

Fully Convolutional Hidden Layer

Fully Convolutional Hidden Layer
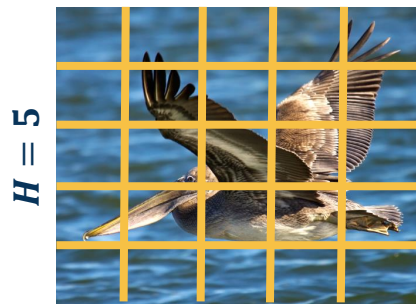
Fully Convolutional Output Layer

Loss

**Each kernel has the size of entire input! (output is 1 scalar)**

- This is equivalent to Wx+b!

- We have one kernel per output node

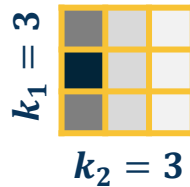**Converting FC Layers to Conv Layers**

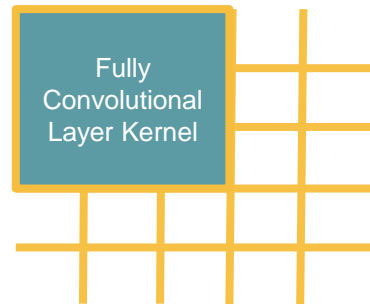Georgia Tech

**Original:**

$H = 5$

$W = 5$

Input

$k_1 = 3$

$k_2 = 3$

Conv Kernel

Fully Convolutional Layer Kernel

Output

**Larger:**

$H = 7$

$W = 7$

$k_1 = 3$

$k_2 = 3$

Fully Convolutional Layer Kernel

Fully Convolutional Layer Kernel

**Same Kernel, Larger Input**

Georgia Tech

# Why does this matter?

- We can stride the "fully connected" classifier across larger inputs!
- Convolutions work on arbitrary input sizes (because of striding)

**Original sized image**

**Larger Image**

**Larger Output Maps**

**Larger Output Size!**

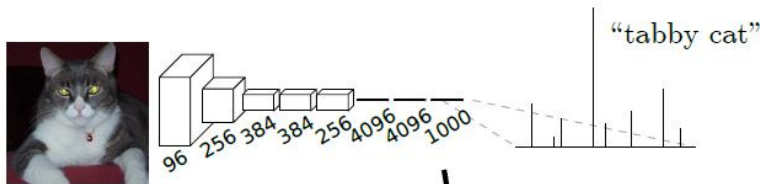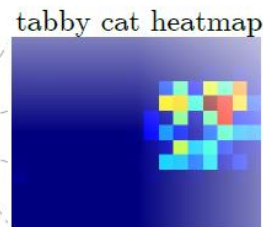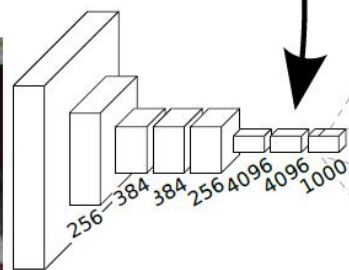*Long, et al., "Fully Convolutional Networks for Semantic Segmentation", 2015*

Convolutional Neural Network (CNN)

Image

Convolution + Non-Linear Layer

Pooling Layer

Convolution + Non-Linear Layer

Useful, lower-dimensional features

Encoder

We can develop learnable or non-learnable upsampling layers!

Decoder

(De)Convolution + Non-Linear Layer

(Un)Pooling Layer

(De)Convolution + Non-Linear Layer

"Image"

Useful, lower-dimensional features

**Idea 2: "De"Convolution and UnPooling**

Georgia Tech

**Example :** Max pooling

⬡ Stride window across image but perform per-patch **max operation**

$$X(0:1, 0:1) = \begin{bmatrix} 100 & 150 \\ 100 & 200 \end{bmatrix} \quad \Rightarrow \quad max(0:1,0:1) = 200$$

Copy value to position chosen as max in encoder, fill reset of this window with zeros



$H = 5$

$W = 5$

**Pooling**

**UnPooling**

$W = 5$

$H = 5$

**Idea:** Remember max elements in encoder! Copy value from equivalent position, rest are zeros

**Max Unpooling**

$$X = \begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \end{bmatrix}$$

**2x2 max pool**

$$Y = \begin{bmatrix} 150 & 150 \\ 100 & 110 \end{bmatrix}$$

**Encoder**

**Decoder**

**2x2 max unpool**

$$X = \begin{bmatrix} 300 & 450 \\ 100 & 250 \end{bmatrix}$$

$$Y = \begin{bmatrix} 0 & 300 & - \\ 0 & 0 & - \\ - & - & - \end{bmatrix}$$

**Max Unpooling Example (one window)**

Georgia Tech

$$X_{enc} = \begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \end{bmatrix}$$ →(2x2 max pool) $$Y_{enc} = \begin{bmatrix} 150 & 150 \\ 100 & 110 \end{bmatrix}$$

**Contributions from multiple windows are summed**

**Encoder**

**Decoder**

2x2 max unpool

$$X_{dec} = \begin{bmatrix} 300 & 450 \\ 100 & 250 \end{bmatrix}$$ → $$Y_{dec} = \begin{bmatrix} 0 & 300+450 & 0 \\ 100 & 0 & 250 \\ 0 & 0 & 0 \end{bmatrix}$$

**Max Unpooling Example**

Georgia Tech

Symmetry in Encoder/Decoder

# How can we *upsample* using convolutions and learnable kernel?

**Normal Convolution**



$H = 5$

$W = 5$

$k_1 = 3$

$k_2 = 3$

$H - k_1 + 1$

$W - k_2 + 1$

**Transposed Convolution (also known as "deconvolution", fractionally strided conv)**

Idea: Take each input pixel, multiply by learnable kernel, "stamp" it on output



$k_1 = 3$

$k_2 = 3$

$H = 5$

Georgia Tech

$$X = \begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \end{bmatrix} \qquad K = \begin{bmatrix} 1 & -1 \\ 2 & -2 \end{bmatrix}$$
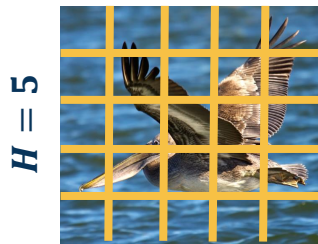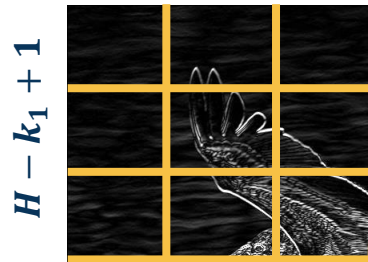
**Contributions from multiple windows are summed**

$$\begin{bmatrix} 120 & -120 & 0 & 0 \\ 240 & -240 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

**Incorporate**
$X(0,0)$

$$\begin{bmatrix} 120 & -120+150 & -150 & 0 \\ 240 & -240+300 & -300 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

**Incorporate**
$X(1,0)$

**Transposed Convolution Example**

Georgia Tech

**Convolutional Neural Network (CNN)**

We can either learn the kernels, or take corresponding encoder kernel and rotate 180 degrees (no decoder learning)

Useful, lower-dimensional features

Image

Convolution + Non-Linear Layer

Pooling Layer

Convolution + Non-Linear Layer

**Decoder**

**Encoder**

(De)Convolution + Non-Linear Layer

(Un)Pooling Layer

(De)Convolution + Non-Linear Layer

"Image"

Useful, lower-dimensional features

**Symmetry in Encoder/Decoder**

Georgia Tech

Input Image → CNN → Predictions

We can start with a pre-trained trunk/backbone (e.g. network pretrained on ImageNet)!

Transfer Learning

Georgia Tech

# U-Net

**You can have skip connections to bypass bottleneck!**



*Ronneberger, et al., "U-Net: Convolutional Networks for Biomedical Image Segmentation", 2015*

Georgia Tech

# Summary

- Various ways to get **image-like outputs,** for example to predict segmentations of input images

- Fully convolutional layers essentially apply the striding idea to the output classifiers, supporting arbitrary input sizes
  - (without output size depending on what the input size is)

- We can have various upsampling layers that actually increase the size

- Encoder/decoder architectures are popular ways to leverage these to perform general image-to-image tasks
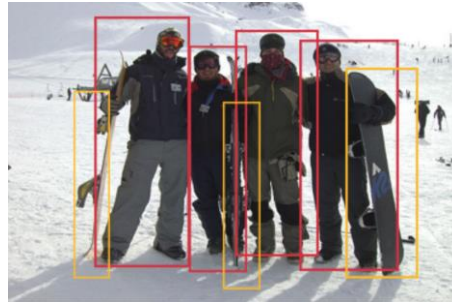
Georgia Tech

Single-Stage Object Detection

# Given an image, output a list of bounding boxes with probability distribution over classes per box
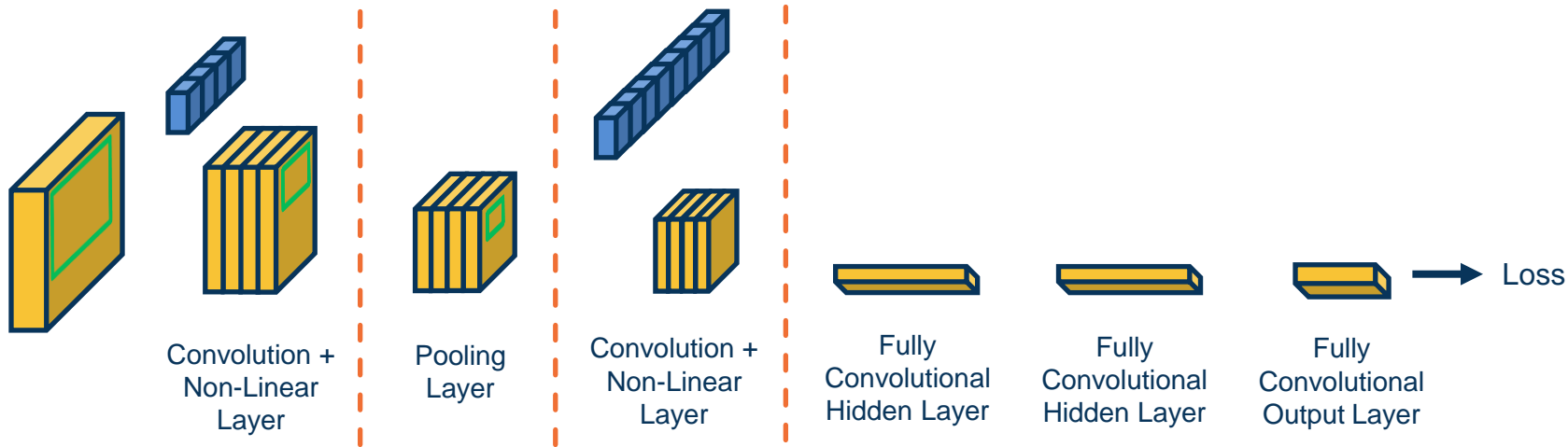
Problems:

- Variable number of boxes!

- Need to determine candidate regions (position and scale) first



**Object Detection**
(List of bounding boxes with class distribution per box)

Georgia Tech

Convolution +
Non-Linear
Layer

Pooling
Layer

Convolution +
Non-Linear
Layer

Fully
Convolutional
Hidden Layer

Fully
Convolutional
Hidden Layer

Fully
Convolutional
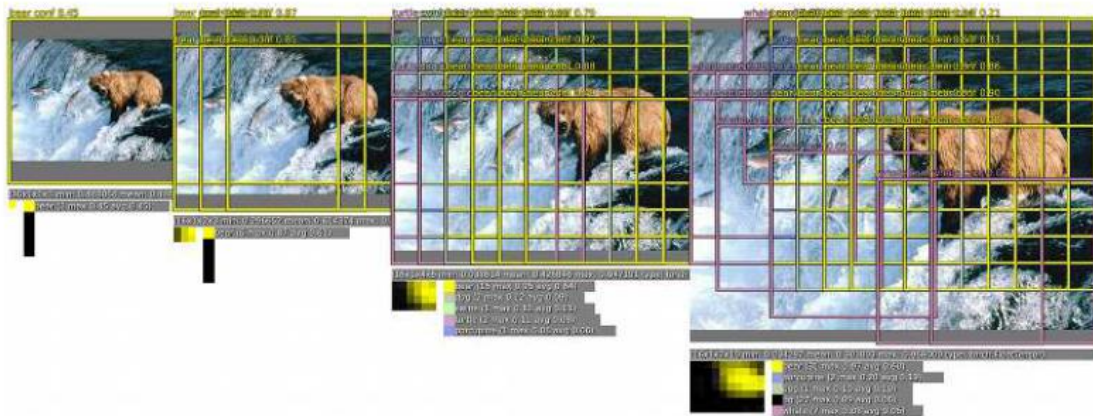Output Layer

Loss

**We can use the same idea of fully-convolutional networks**

- Use ImageNet pre-trained model as backbone (e.g. taking in 224x224 image)

- Feed in larger image and get classifications for different windows in image

**Object Detection Tasks**

Georgia Tech

Convolution + Non-Linear Layer — Pooling Layer — Convolution + Non-Linear Layer — Fully Convolutional Hidden Layer — Fully Convolutional Hidden Layer — Fully Convolutional Output Layer — Cross-Entropy Loss — Mean Squared Error (MSE)

**We can have a *multi-headed architecture***

- One part predicting distribution over class labels (classification)
- One part predicting a bounding box for each image region (regression)
  - Refinement to fit the object better (outputs 4 numbers)
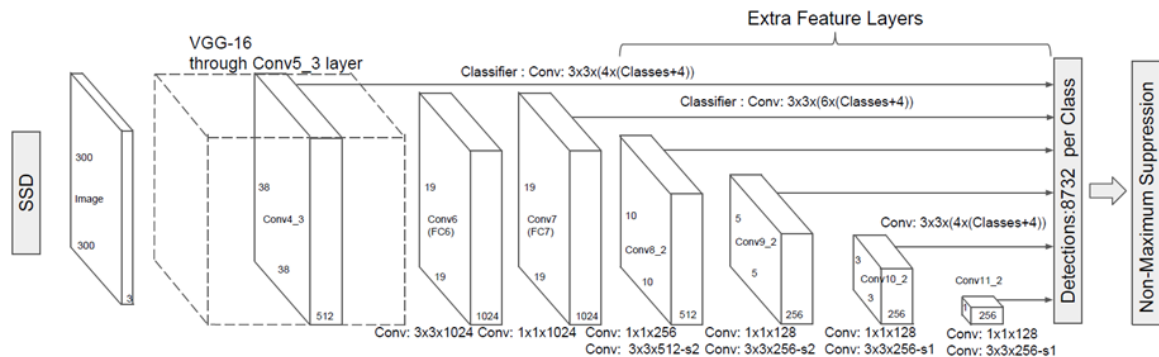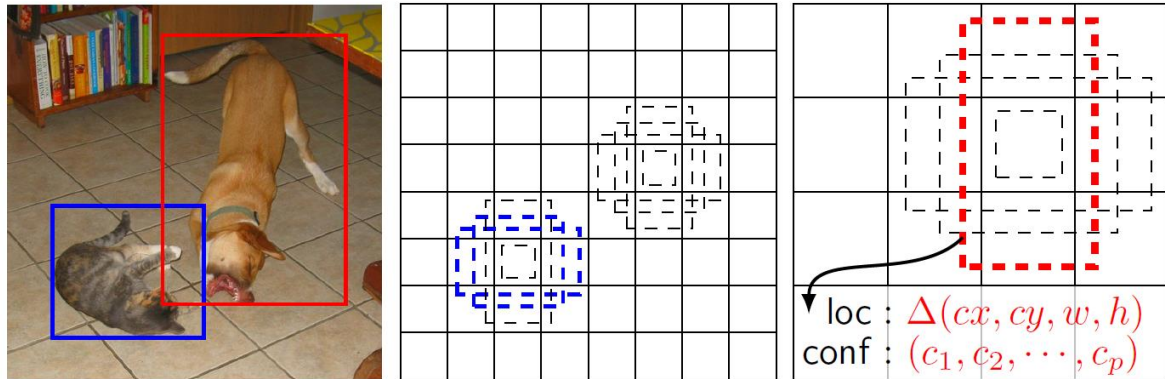- Both heads *share features*! Jointly optimized (summing gradients)

Georgia Tech

Can also do this at multiple scales to result in a large number of detections

- Various tricks used to increase the resolution (decrease subsampling ratio)

- Redundant boxes are combined through **Non-Maximal Suppression (NMS)**

*Sermanet, et al., "OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks", 2013*

**Object Detection Tasks**

Georgia Tech

Single-shot detectors use an idea of **grids** as anchors, with different scales and aspect ratios around them

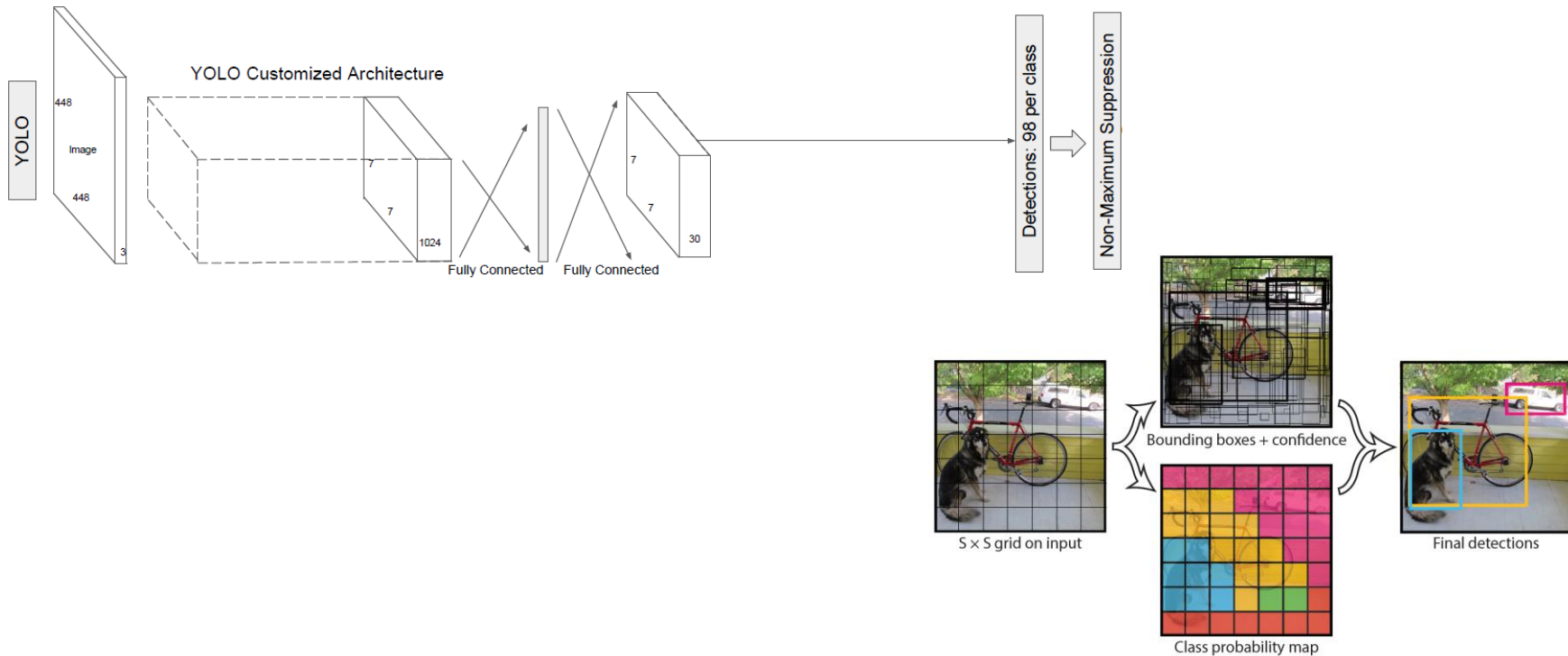◆ Various tricks used to increase the resolution (decrease subsampling ratio)



$$loc : \Delta(cx, cy, w, h)$$
$$conf : (c_1, c_2, \cdots, c_p)$$

*Liu, et al., "SSD: Single Shot MultiBox Detector", 2015*

Georgia Tech

# Similar network architecture but single-scale (and hence faster for same size)



*Redmon, et al., "You Only Look Once:Unified, Real-Time Object Detection", 2016*
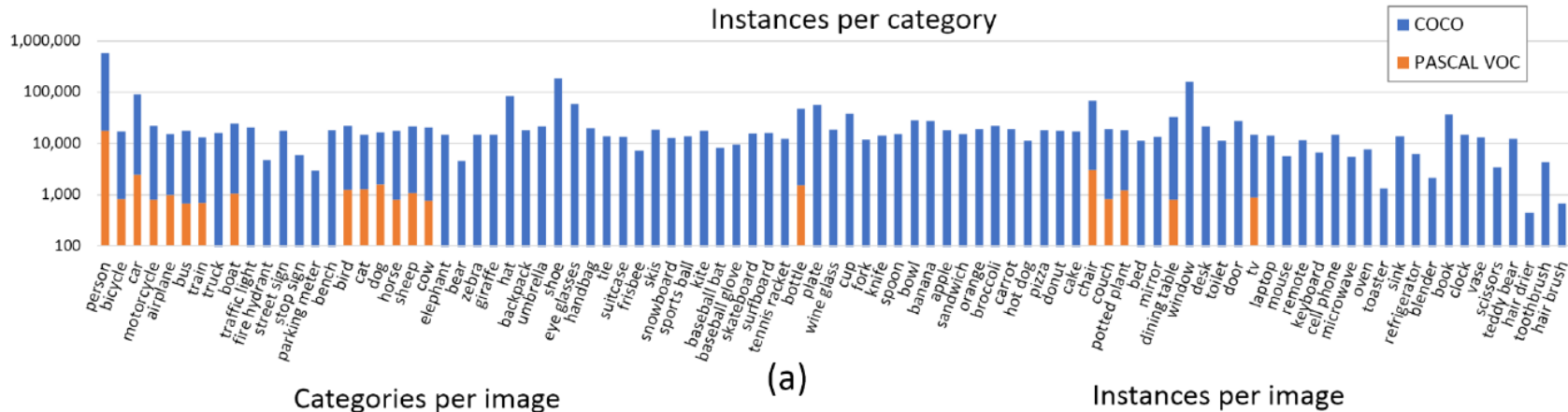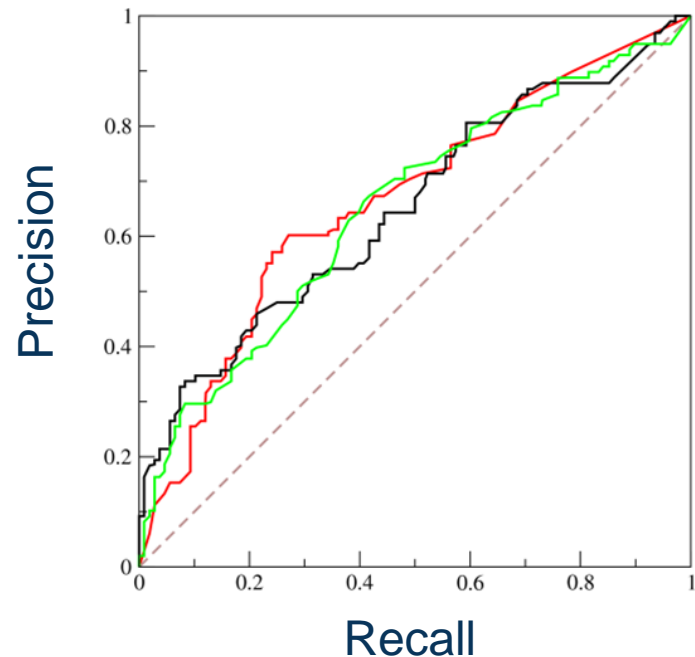
## You Only Look Once (YOLO)

What is COCO?

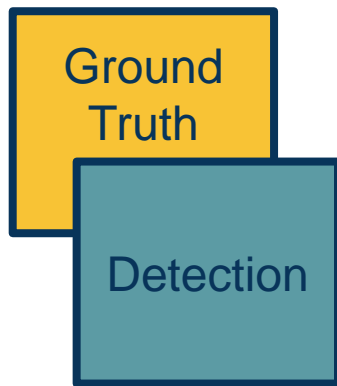COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✔ Object segmentation
- ✔ Recognition in context
- ✔ Superpixel stuff segmentation
- ✔ 330K images (>200K labeled)
- ✔ 1.5 million object instances
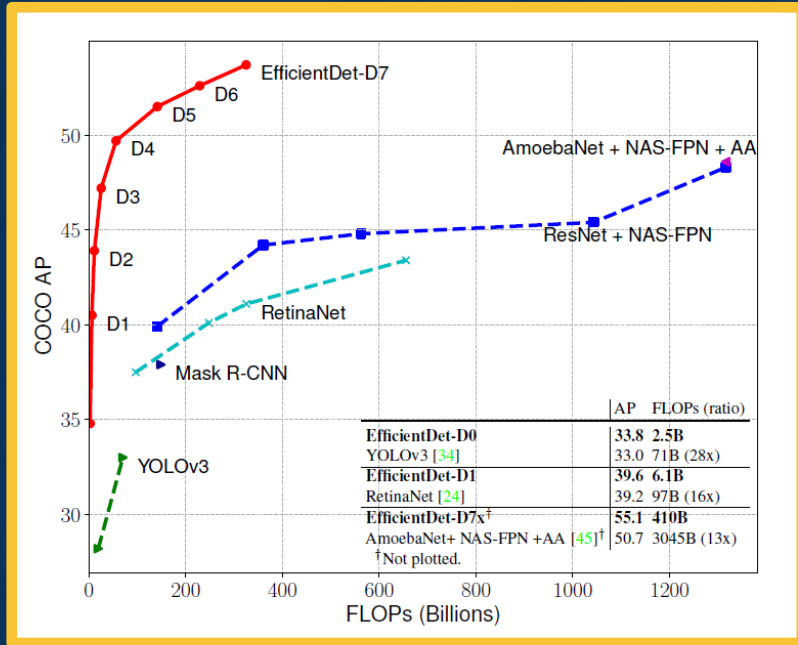- ✔ 80 object categories

Instances per category

*Lin, et al., "Microsoft COCO: Common Objects in Context", 2015. https://cocodataset.org/#explore*

**Datasets**

Georgia Tech

1. For each bounding box, calculate intersection over union (IoU)

2. Keep only those with IoU > threshold (e.g. 0.5)

3. Calculate precision/recall curve across classification probability threshold

4. Calculate **average precision (AP)** over recall of [0, 0.1, 0.2, …, 1.0]

5. Average over all categories to get mean Average Precision (mAP)



Ground Truth

Detection

Precision

Recall

$$AP = \frac{1}{11} \sum_{i \in [0, 0.1, \ldots 1.0]} AP_i$$

**Evaluation – Mean Average Precision (mAP)**

Georgia Tech
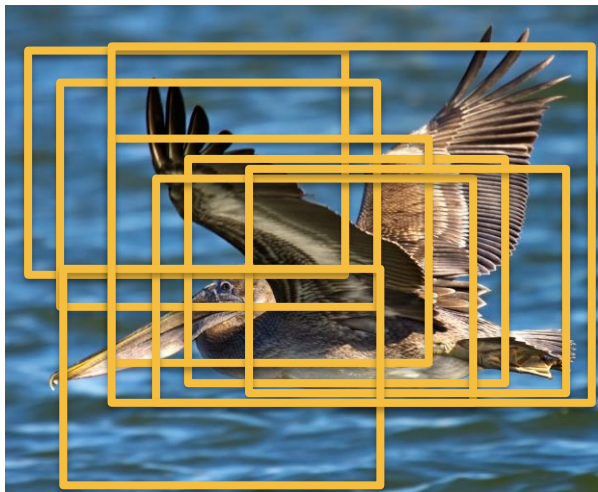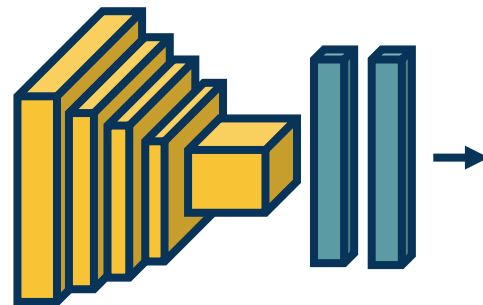
# Results



**EfficientDet**



**PP-YOLO**

*Tan, et al., "EfficientDet: Scalable and Efficient Object Detection", 2020*
*Long et al., "PP-YOLO: An Effective and Efficient Implementation of Object Detector", 2020*

For each crop, Resize

Instead of making dense predictions across an image, we can decompose the problem:

- ⬡ Find regions of interest (ROIs) with object-like things
- ⬡ Classifier those regions (and refine their bounding boxes)

*Girshick, et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", 2014*

**R-CNN**

Georgia Tech

We can use **unsupervised (non-learned!) algorithms** for finding candidates

**Downsides:**

⬡ Takes 1+ second per image

⬡ Return thousands of (mostly background) boxes

**Resize each candidate** to full input size and classify



*Uijlings, et al., "Selective Search for Object Recognition", 2012*

Georgia Tech

# What is the problem with this?
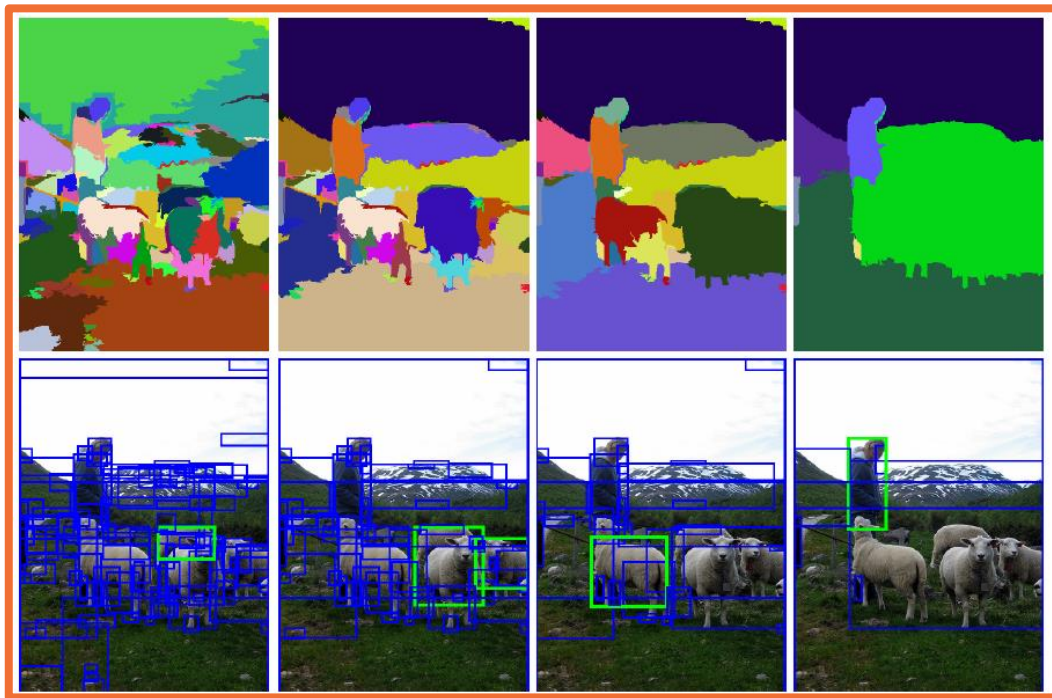


**Computation for convolutions re-done for each image patch, even if overlapping!**

*Girshick, et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", 2014*

**Inefficiency of R-CNN**

Georgia Tech

Map each ROI in image to corresponding region in feature maps



Extract Feature Map Region

?

**Idea:** *Reuse* computation by finding regions in *feature maps*

⬡ Feature extraction only done once per image now!

⬡ Problem: Variable input size to FC layers (different feature map sizes)

*Girshick, "Fast R-CNN", 2015*

**Fast R-CNN**

Georgia Tech

For each grid element, max pool however many values there are to one scalar

$$\begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \\ 65 & 75 & 10 \end{bmatrix}$$

Given an arbitrarily-sized feature map, we can use **pooling** across a grid (ROI Pooling Layer) to convert to fixed-sized representation

**ROI Pooling**

Georgia Tech

Map each ROI in image to corresponding are in feature maps

Extract Feature Map Region

ROI Pooling

We can now train this model **end-to-end** (i.e. backpropagate through entire model including ROI Pooling)!

Georgia Tech

- **Idea:** Why not have the neural network *also* generate the proposals?

  - Region Proposal Network (RPN) uses same features!
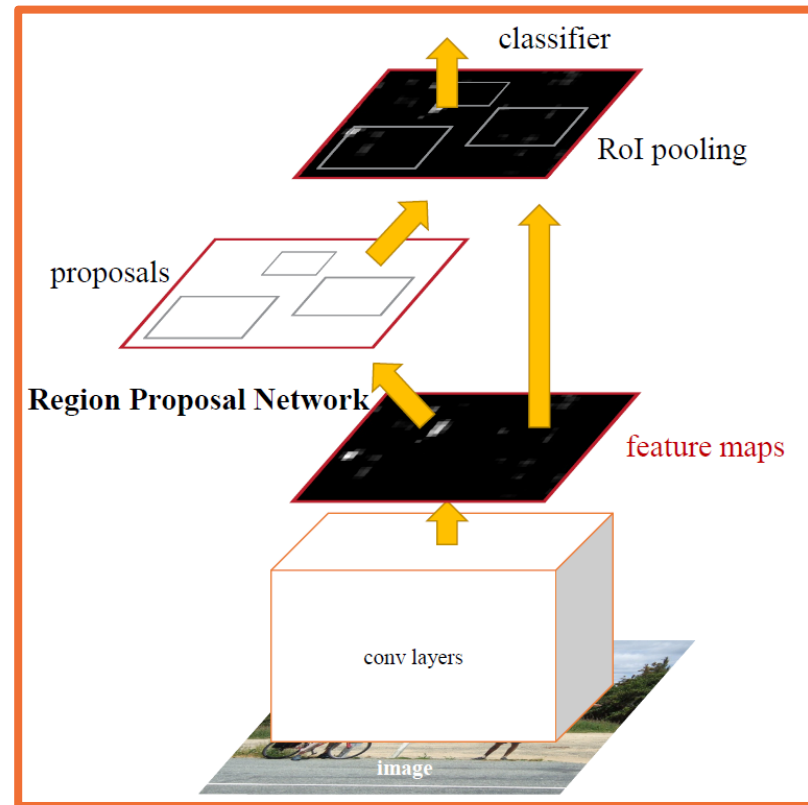
- Outputs ***objectness score*** and bounding box

- Top k selected for classification

- Note some parts (gradient w.r.t. bounding box coordinates) not differentiable so some complexity in implementation



*Ren, et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", 2016*

**Faster R-CNN**

Georgia Tech

**RPN** also uses notion of **anchors in a grid**

Boxes of various sizes and scales classified with objectness score and refined bounding boxes refined



*Ren, et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", 2016*

**Faster R-CNN**

Georgia Tech

Many new advancements have been made

For example, combining detection and segmentation

- Extract foreground (object) mask per bounding box



*He, et al., "Mask R-CNN", 2018*

https://paperswithcode.com/sota/object-detection-on-coco

**Mask R-CNN**

Georgia Tech

- A range of problems characterized by **density and type of output**

- **Semantic/instance segmentation:** Dense, spatial output

  - Leverage encoder/decoder architectures

- **Object detection:** Variable-length list of objects

  - Two-stage versus one-stage architectures

  - (Not covered): Anchor-based versus anchor-free methods

Georgia
Tech

# DETR
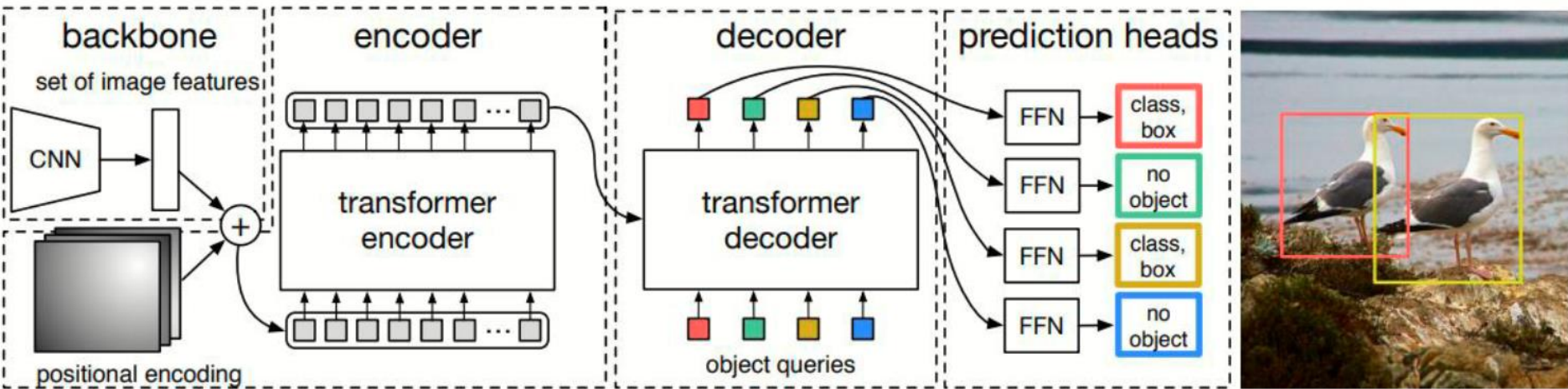## End-to-End Object Detection with Transformers

**End-to-End Object Detection with Transformers**

Nicolas Carion*, Francisco Massa*, Gabriel Synnaeve, Nicolas Usunier,
Alexander Kirillov, and Sergey Zagoruyko
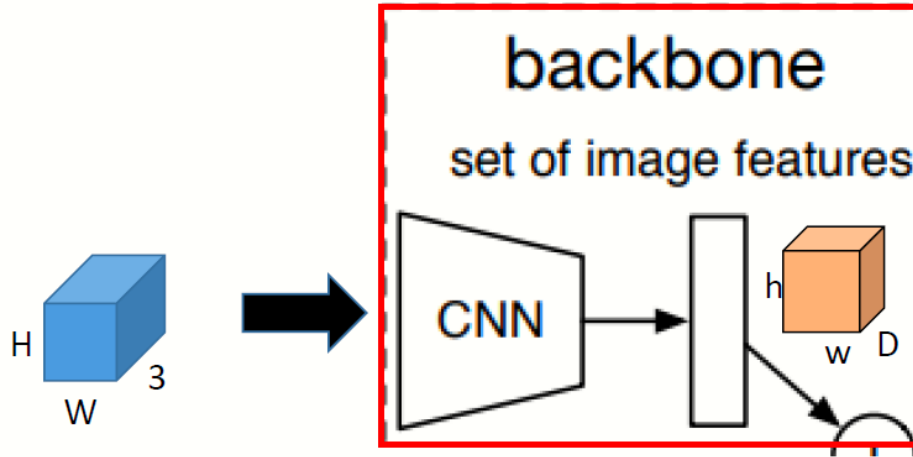
Facebook AI

**Abstract.** We present a new method that views object detection as a direct set prediction problem. Our approach streamlines the detection pipeline, effectively removing the need for many hand-designed components like a non-maximum suppression procedure or anchor generation that explicitly encode our prior knowledge about the task. The main ingredients of the new framework, called DEtection TRansformer or DETR, are a set-based global loss that forces unique predictions via bipartite matching, and a transformer encoder-decoder architecture. Given a fixed small set of learned object queries, DETR reasons about the relations of the objects and the global image context to directly output the final set of predictions in parallel. The new model is conceptually simple and does not require a specialized library, unlike many other modern detectors. DETR demonstrates accuracy and run-time performance on par with the well-established and highly-optimized Faster R-CNN baseline on the challenging COCO object detection dataset. Moreover, DETR can be easily generalized to produce panoptic segmentation in a unified manner. We show that it significantly outperforms competitive baselines. Training code and pretrained models are available at https://github.com/facebookresearch/detr.

arXiv:2005.12872v3 [cs.CV] 28 May 2020

**Slides by R. Q. FEITOSA**

Georgia Tech
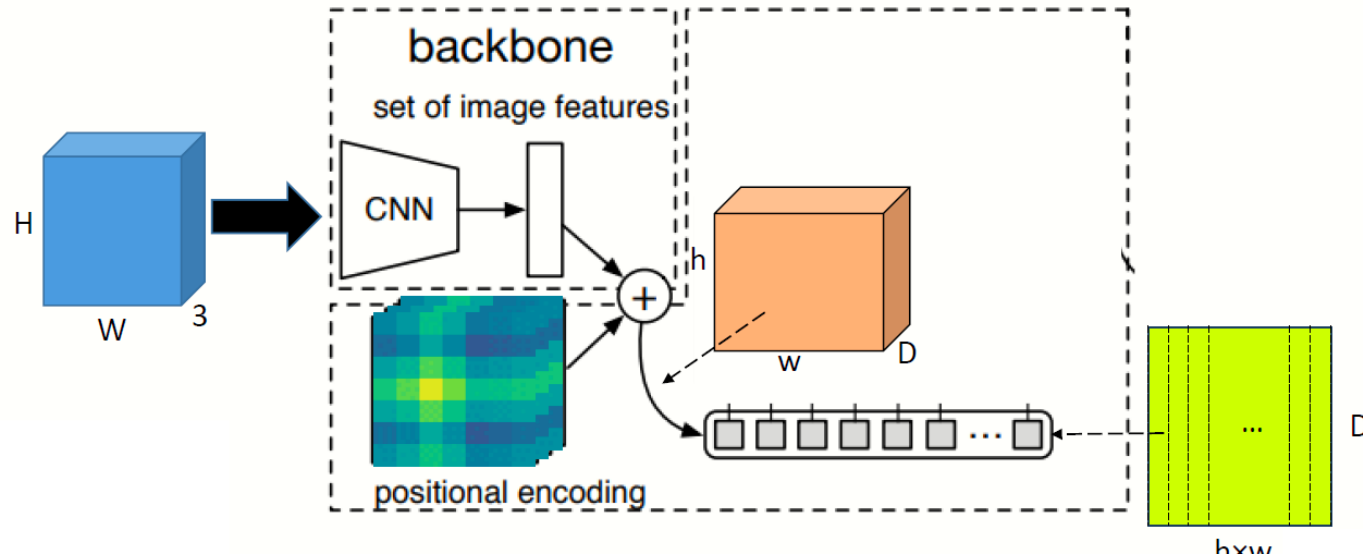
# DEtector TRansformer - DETR
## overview

Georgia Tech

# DEtector TRansformer - DETR
## backbone

A conventional CNN backbone to learn a 2D representation of an input image.

# DEtector TRansformer - DETR
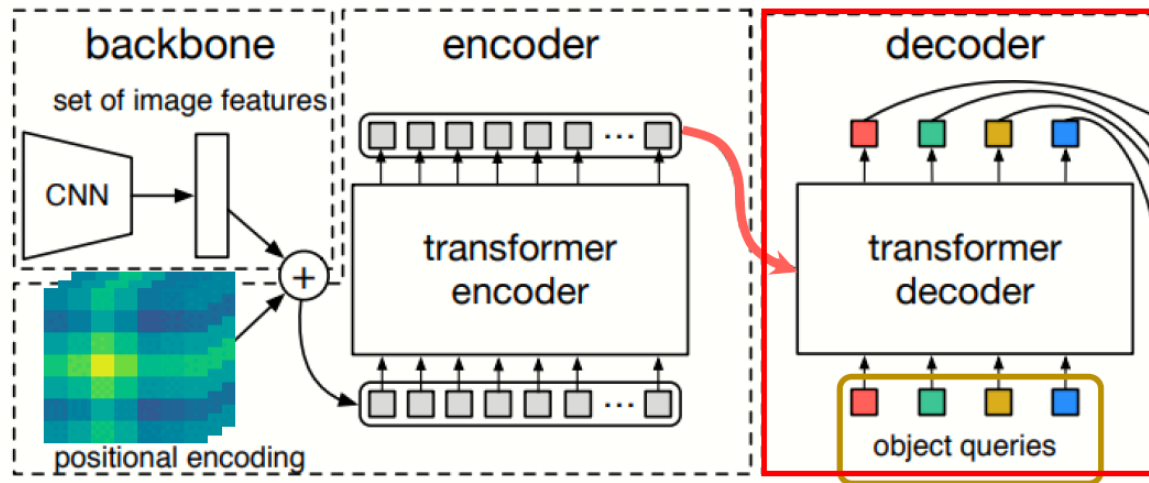## transformer encoder

DETR supplements the features with a positional encoding and flattens them before passing them into a transformer encoder.

# DEtector TRansformer - DETR
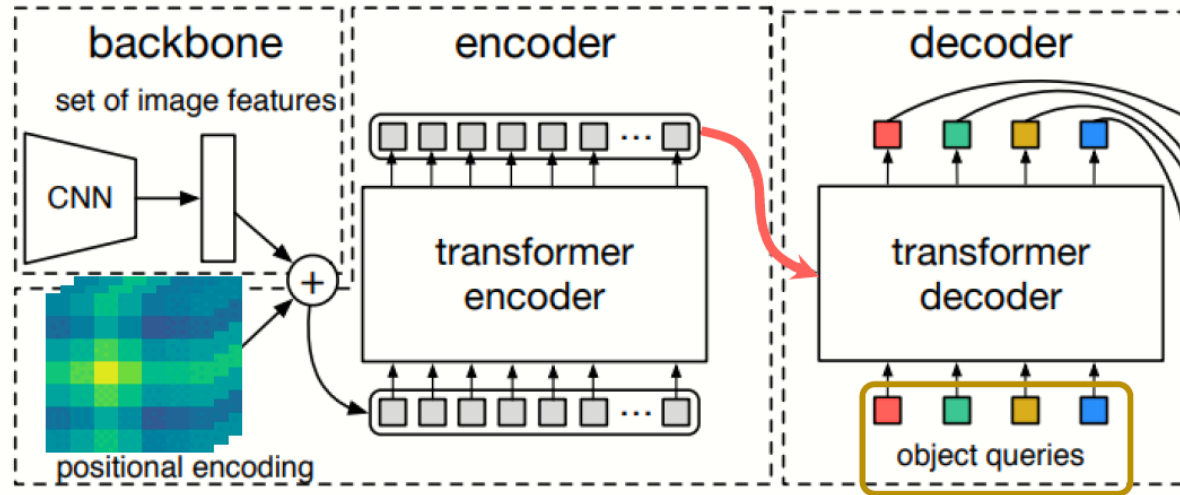## transformer decoder

The transformer decoder takes as input a small number of **learned positional** and content embeddings (object queries) and additionally attends to the encoder output.
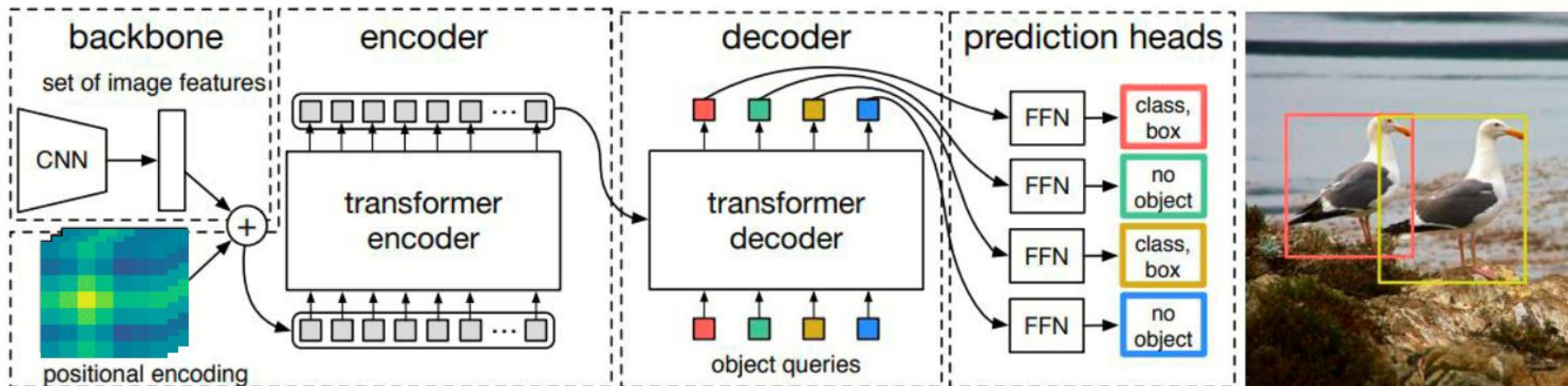
# DEtector TRansformer - DETR
## Object queries

- are randomly initialized embeddings,
- refined through the course of training, and
- then fixed for evaluation.
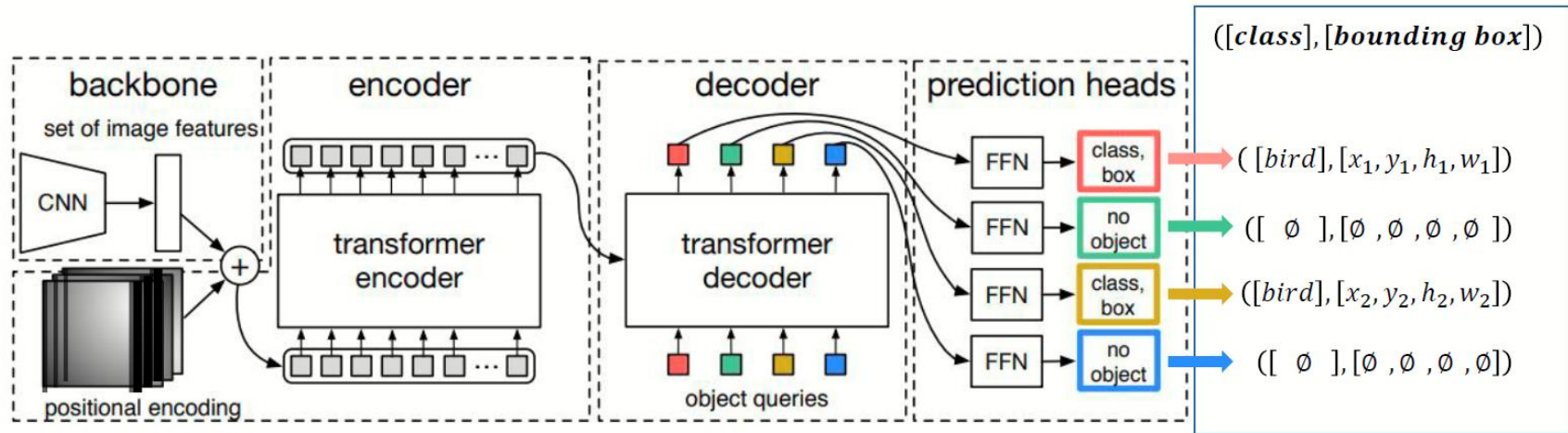
# DEtector TRansformer - DETR
## prediction heads

Each embedding at the decoder output feeds a shared feed-forward network (FFN) that predicts either a detection (class and bounding box) or a "no object" class.

# DEtector TRansformer - DETR
## prediction heads

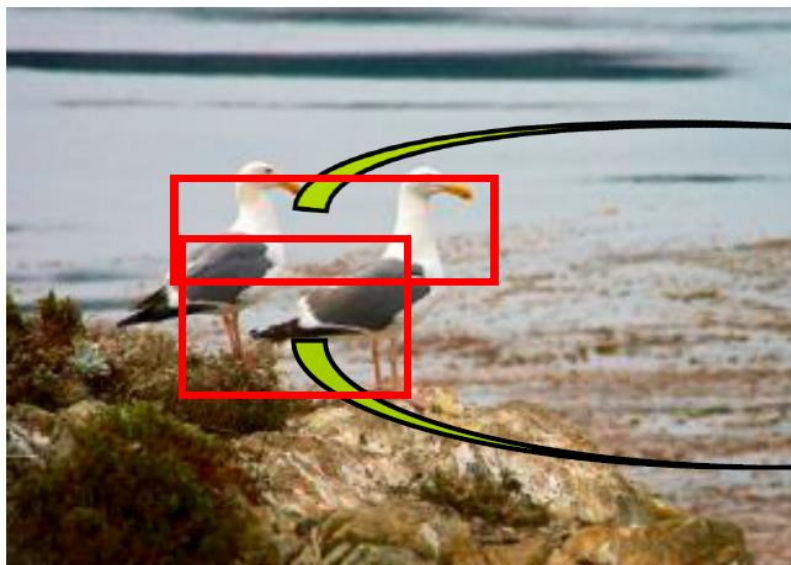Each output embedding of the decoder to a shared feed forward network (FFN) that predicts either a detection (class and bounding box) or a "no object" class.



$([class], [bounding\ box])$

$([bird], [x_1, y_1, h_1, w_1])$

$([\ \emptyset\ ], [\emptyset, \emptyset, \emptyset, \emptyset])$

$([bird], [x_2, y_2, h_2, w_2])$

$([\ \emptyset\ ], [\emptyset, \emptyset, \emptyset, \emptyset])$
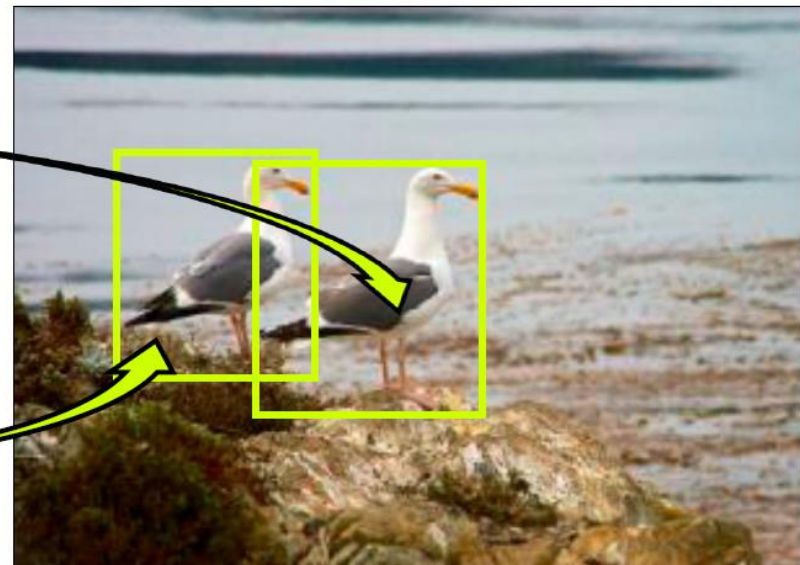
Georgia Tech

# Matching bounding boxes during training with the Reference

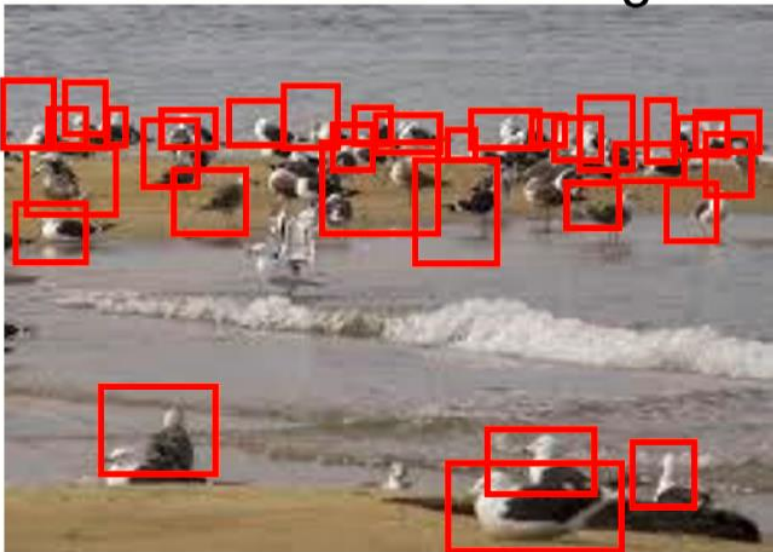**What is the target for each bounding box during training?**



**DETR while training**

**Reference – Ground Truth**

# The complexity of matching grows with *N*!



DETR while training

Reference – Ground Truth

*N!*

Georgia Tech
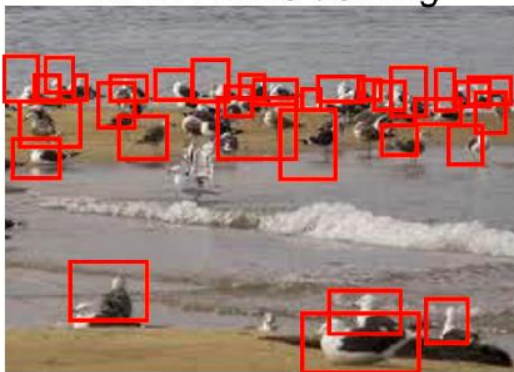
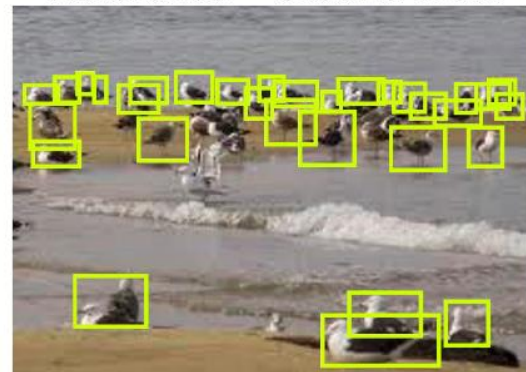# The Hungarian algorithm used for bipartite matching



DETR while training

Reference – Ground Truth

*N!*

The **Hungarian algorithm**[1] computes the optimal assignment efficiently. It considers both the class prediction and the similarity of predicted and ground truth boxes.

Georgia Tech

# DETR Demo

Leaderboard        Dataset

View [ box mAP ⌄ ] by [ Date ⌄ ] for [ All models ⌄ ]



FocalNet-H (DINO)

Swin-L (HTC++, multi-scale)

Cascade Mask R-CNN (Triple-ResNeXt152, multi-scale)

SNIPER (ResNet-101)

RetinaNet (ResNeXt-101-FPN)

Faster R-CNN (box refinement, context, multi-scale testing)

SSD512

Fast-RCNN

BOX MAP

80
60
40
20
0
−20

2016   2017   2018   2019   2020   2021   2022   2023   2024

● Other models    ━●━ Models with highest box mAP

Georgia Tech

View [mask AP ▾] by [Date ▾] for [All models ▾]



MASK AP

70

60    EVA    CBNetV2 (EVA02, single-scale)

Soft Teacher + Swin-L (HTC++, multi-scale)

Cascade Eff-B7 NAS-FPN (1280, self-training Copy Paste, single-scale)
50    Mask R-CNN (SpineNet-190, 1536x1536)
Cascade Mask R-CNN (ResNeXt152, CBNet)

PANet
40    Mask R-CNN (ResNeXt-101-FPN)

FCIS+++ +OHEM
30
MultiPath Network

20
        2017        2018        2019        2020        2021        2022        2023        2024        2025

● Other models    ●— Models with highest mask AP

Georgia Tech