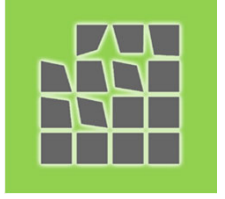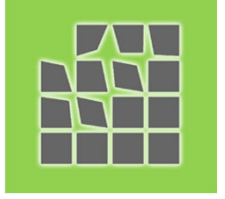# Image Segmentation
# in the Era of Deep Learning

Institute of Information Theory
and Automation of the AS CR

# Outline
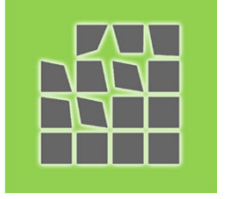
- Segmentation tasks

- Supervised training – datasets, loss

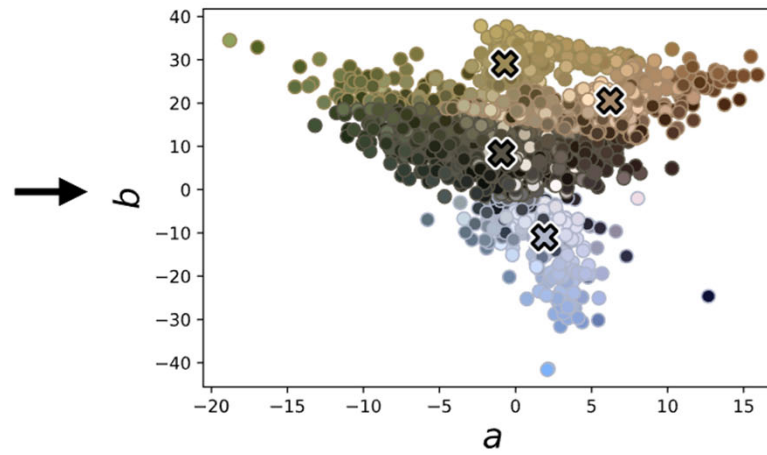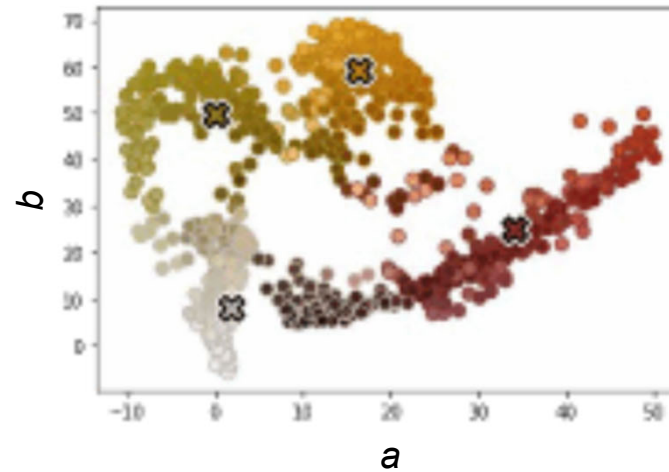- Architectures
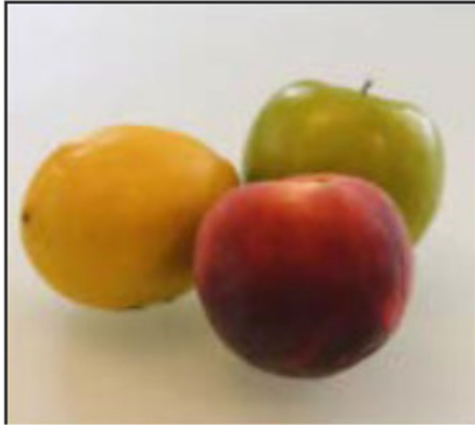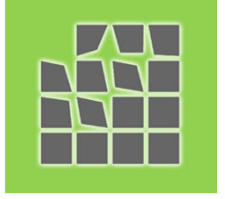
- Evaluations metrics

# Segmentation Tasks

- Image classification

- Object detection

- Segmentation

– Semantic - SS

– Instance - IS
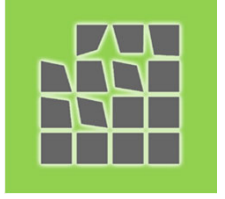
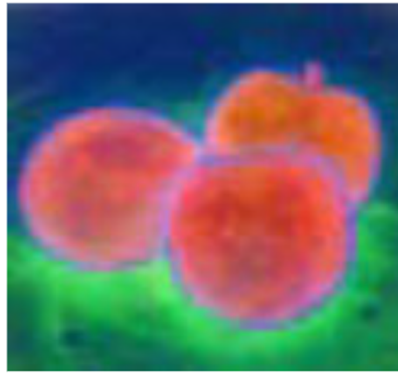– Panoptic (objects, stuff) - PS

# What is segmentation?



(a)

(b)

(c)

# Clustering Problem



k=4

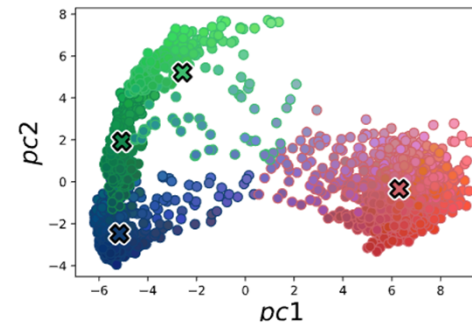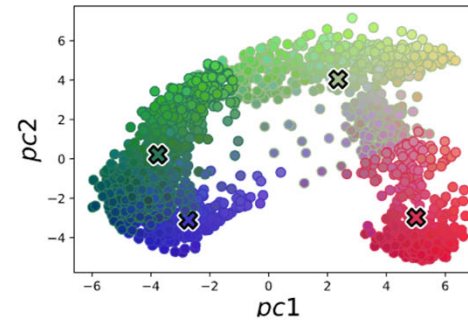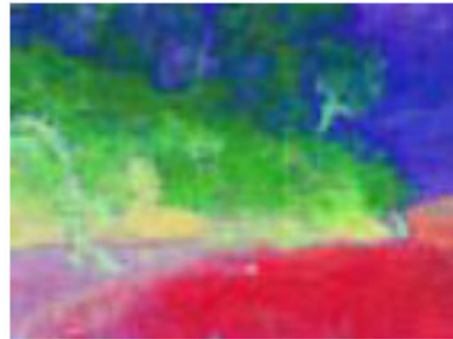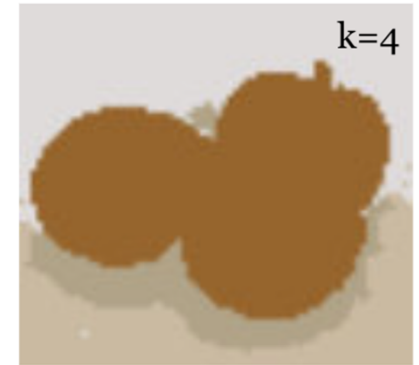k=4

# Clustering Problem



Input

Feature map
(DINO)

Feature space
(2 PCA components)

Segmentation
(k-means clustering)

k=4

k=4

# Image Classification

- Image -> class label

# Object Detection

- Image -> { label, bounding_box (x,y,w,h) }

# Segmentation

- Semantic

label and mask



- Instance (countable objects = things)

label and mask per instance



- Panoptic (semantic + instance)
- things (person) : label and mask per instance
- stuff (water): label and mask

# Datasets for SS,IS,PS

• **Pascal Visual Object Classes** (Pascal VOC): many different object classes, bounding boxes and robust segmentation maps. (20 object categories, only for SS)

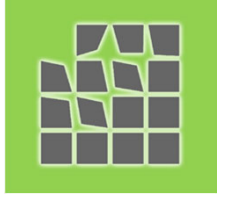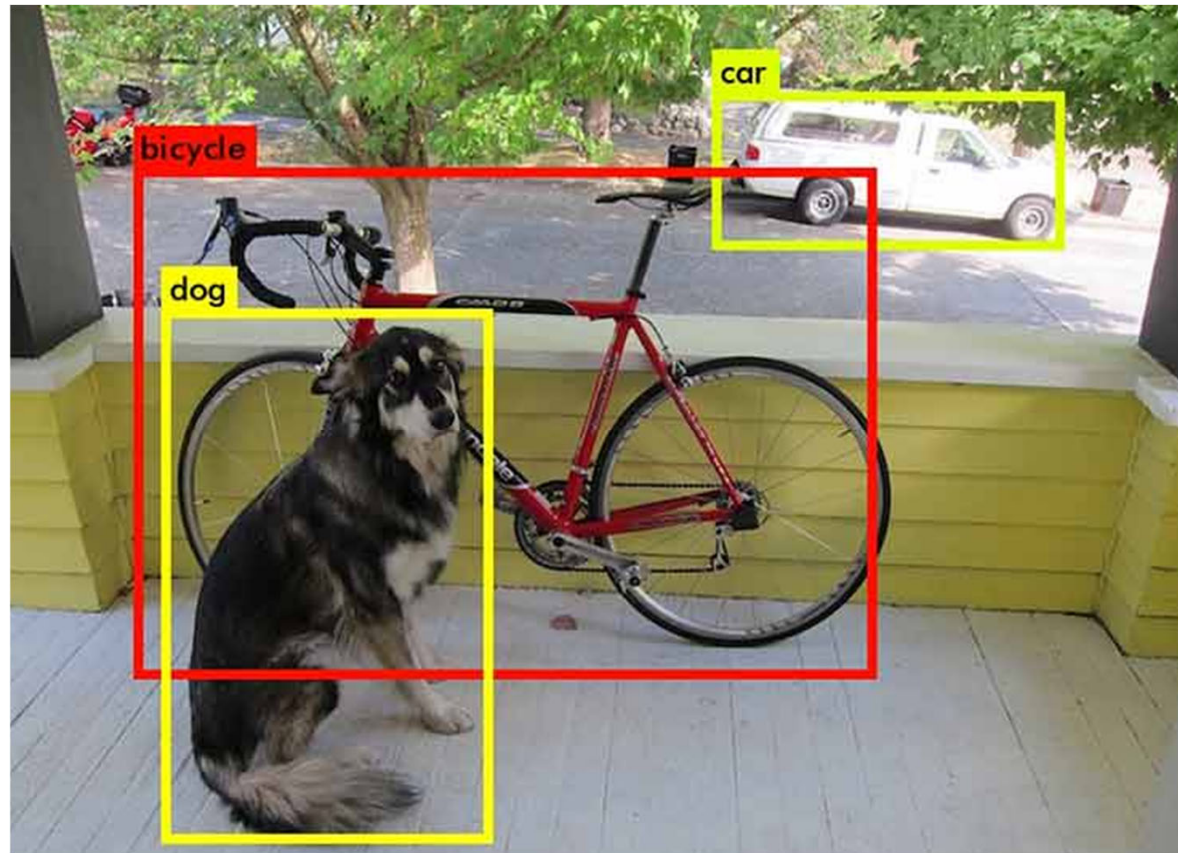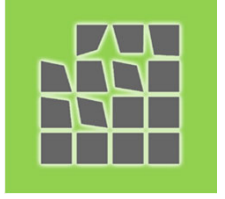• **MS COCO**: 330k images and annotations for many tasks including image captioning (80 thing and 91 stuff categories)

• **Cityscapes**: data from urban environments made up of 5k images with 20k annotations and 30 class labels.

• **ADE20k**: 20k annotated images of scene categories from the SUN and Places database.

• **Cell Tracking Challenge**: 2D and 3D time-lapse images (for SS and IS)

# Supervised Training

Annotation: {x,y} = {"input image", "output segmentation"}



$$x \qquad f_\theta(x) \qquad\qquad y$$

Training: $\min_\theta \sum_m \mathcal{L}(f_\theta(x^{(m)}), y^{(m)})$  **Must be differentiable!**

# Cross-Entropy Loss

- Training set: $(x^{(i)}, y^{(i)})$  i … image (pixel) index

$$y^{(i)} \in \mathbb{R}^C$$
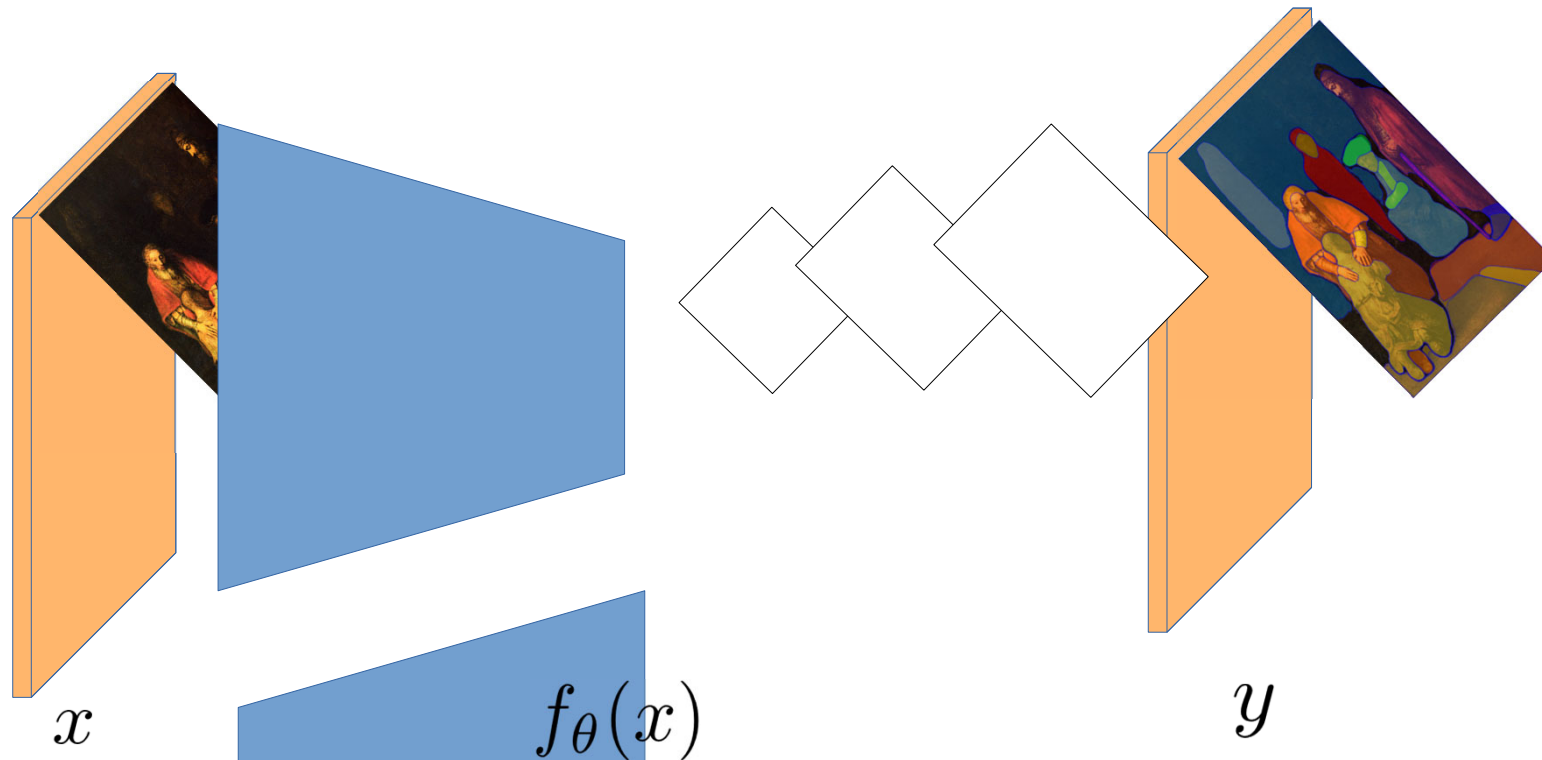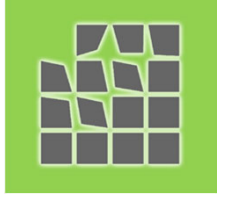
- one-hot vector

$$y_c^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is from the c-th class} \\ 0 & \text{elsewhere} \end{cases}$$

$$y^{(i)} = [0, \quad 1, \quad 0, \quad \ldots, \quad 0]$$

1. class (car)    2. class (bird)      C-th class
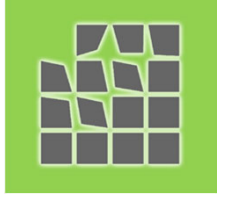
$$\hat{y}^{(i)}$$

- Prediction:

- CE Loss:

$$L = -\sum_{i \in \mathbb{B}} \sum_{c=1}^{C} y_c^{(i)} \log(\hat{y}_c^{(i)})$$

- Binary CE: $C=2$

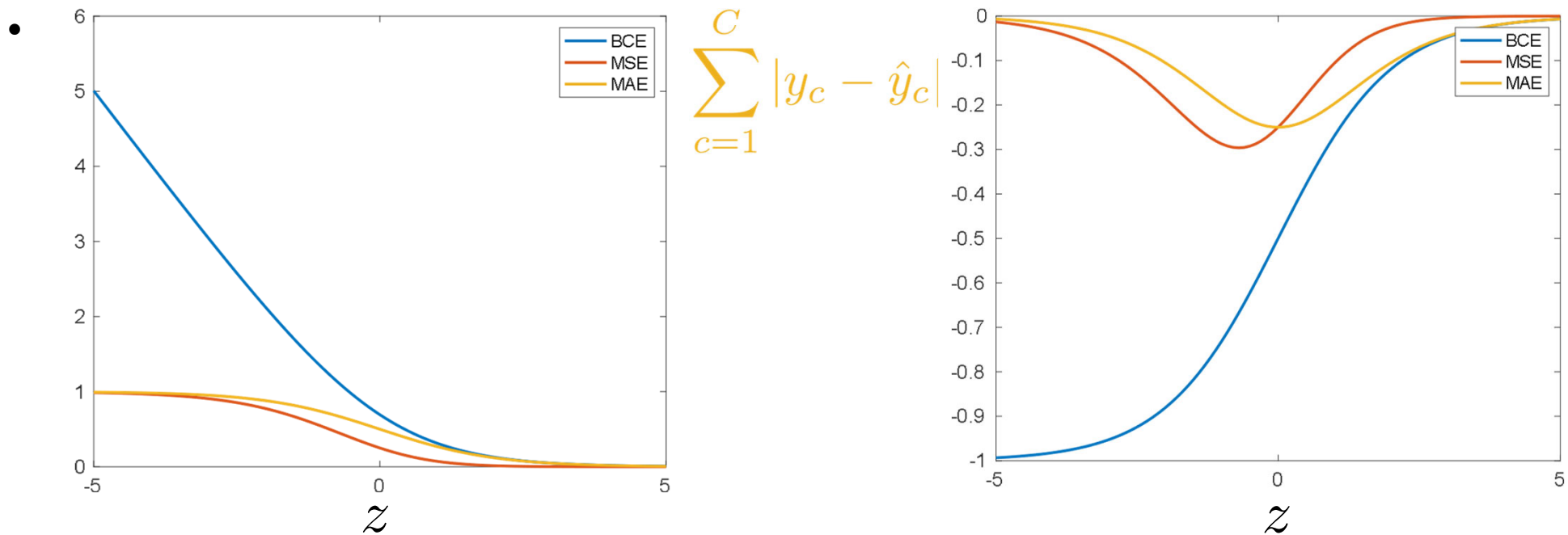$$L = -\sum_{i} y_1^{(i)} \log(\hat{y}_1^{(i)}) + (1 - y_1^{(i)}) \log(1 - \hat{y}_1^{(i)})$$

# $L_p$ norm

- GT: $\quad y \in \mathbb{R}^C$
- Prediction: $\quad \hat{y} = \sigma(z)$

- $L_2$: $\quad \displaystyle\sum_{c=1}^{C} (y_c - \hat{y}_c)^2$ $\qquad\qquad\qquad\qquad -\displaystyle\sum_{c=1}^{C} y_c \log(\hat{y}_c)$ CE:

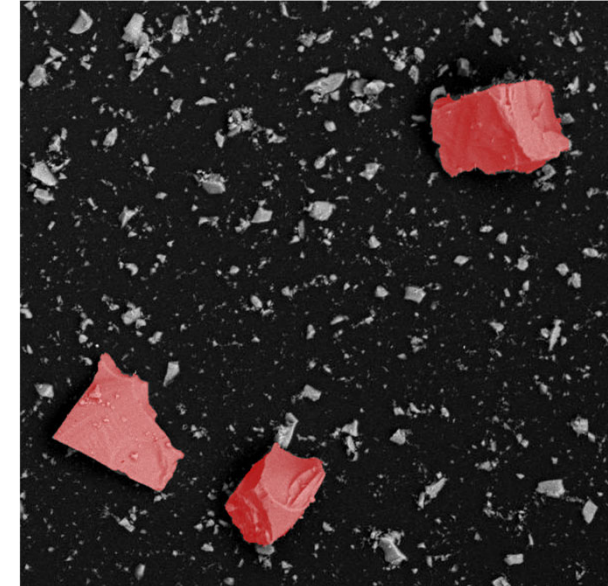-  $\displaystyle\sum_{c=1}^{C} |y_c - \hat{y}_c|$
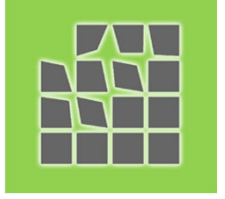
# Unbalanced classes

- CE: larger areas contribute to the error more

- weighted CE:

$$L = -\sum_{i \in \mathbb{B}} \sum_{c=1}^{C} w_c y_c^{(i)} \log(\hat{y}_c^{(i)})$$
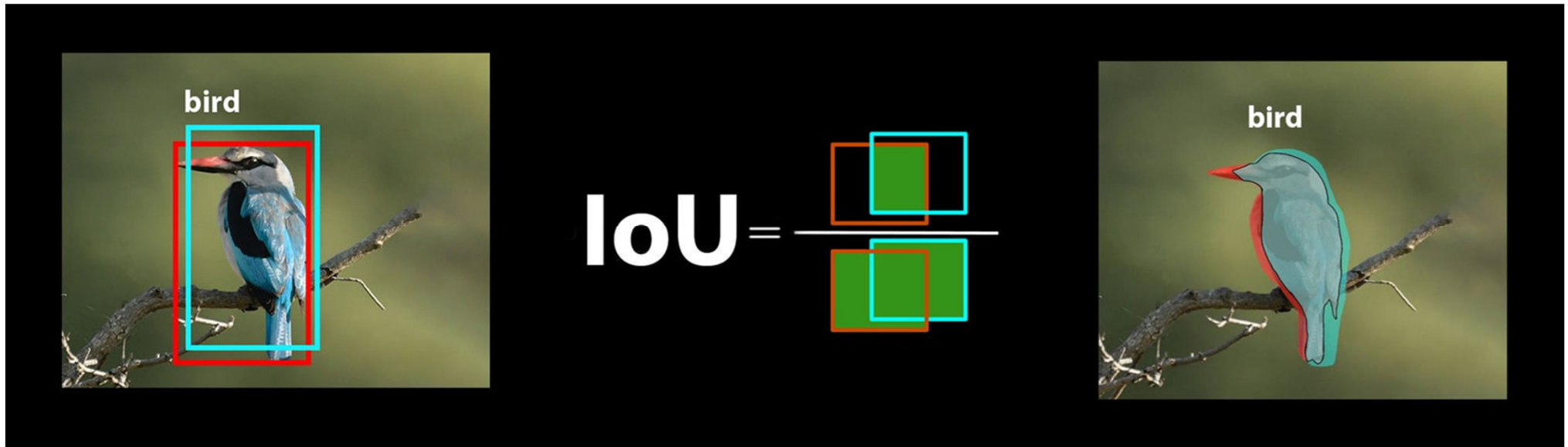


- IoU

# IoU - Intersection over Union
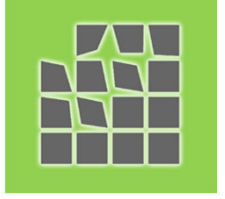
- Jaccard index



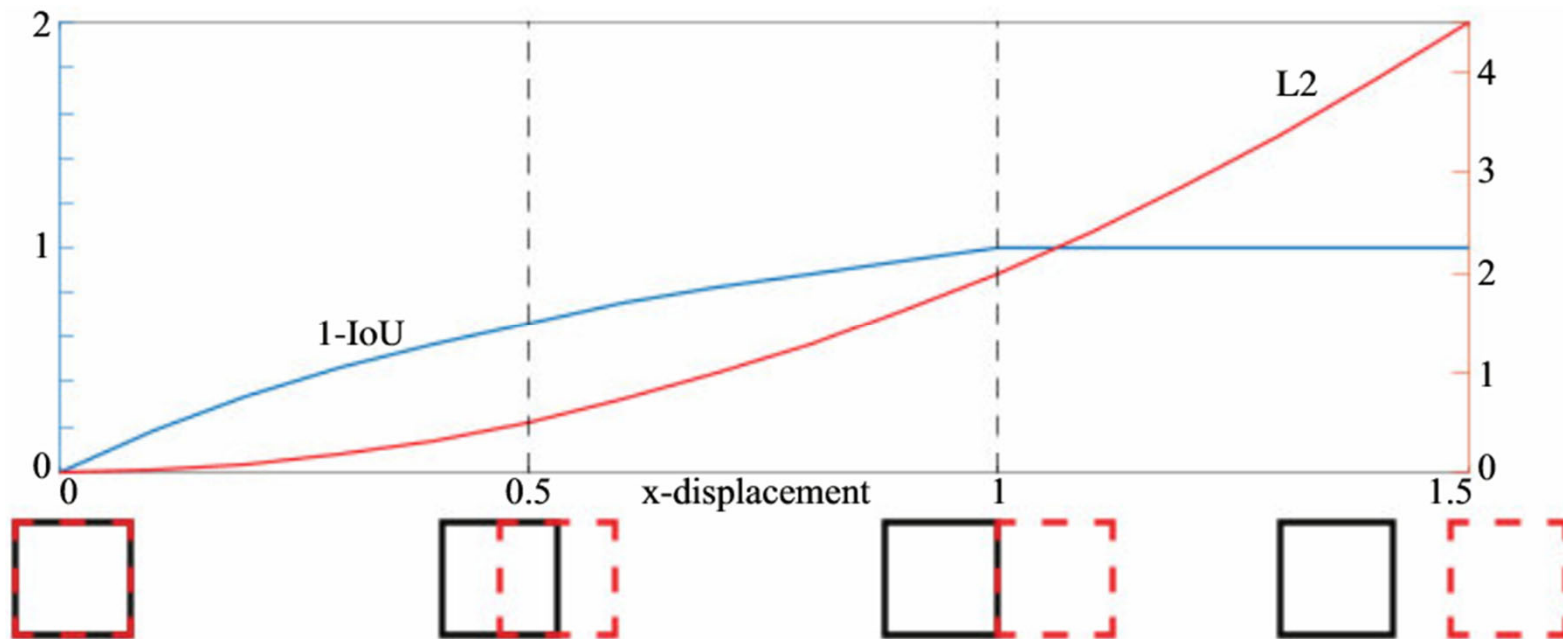$$\text{IoU} \in \langle 0, 1 \rangle \qquad L = 1 - \text{IoU}$$

- For masks:

$$L = 1 - \sum_{c=1}^{C} \frac{\sum_i y_c^{(i)} \hat{y}_c^{(i)}}{\sum_i y_c^{(i)} + \hat{y}_c^{(i)} - y_c^{(i)} \hat{y}_c^{(i)}}$$
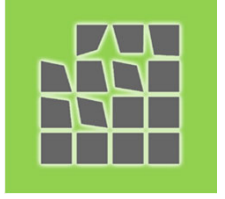
- Dice loss (3DV 2016)

# Bounding Box Loss

- Bounding box: [x,y,width,height]

- $L_p$ norm: same error irrespective of the area size

- IoU: no overlap has zero gradient



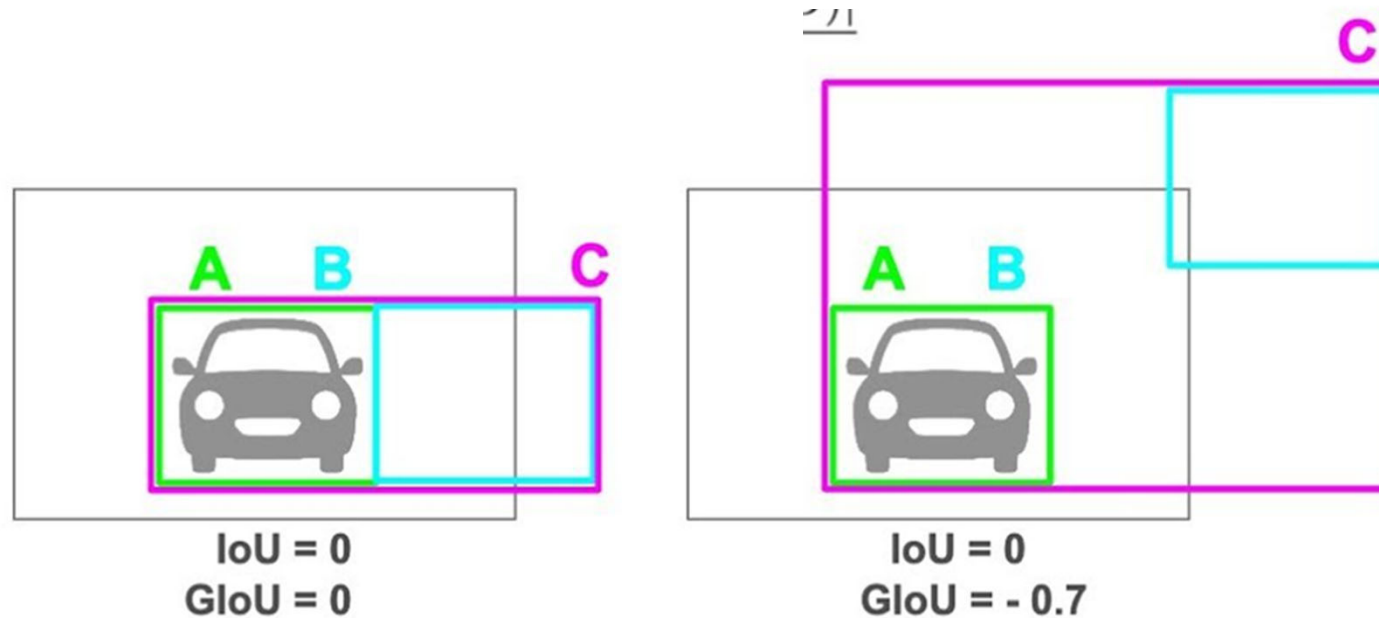- Loss = $L_p$ + ( 1-IoU )

# Generalized IoU
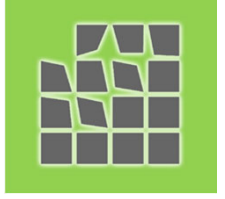
$$\mathrm{GIoU} = \mathrm{IoU} - \frac{|C \setminus (A \cup B)|}{|C|}$$



IoU = 0
GIoU = 0

IoU = 0
GIoU = - 0.7

$$\mathrm{GIoU} \in \langle -1, 1 \rangle$$

Defined only for bounding boxes.

# Focal Loss

**CE:**

$$-\sum_{c=1}^{C} y_c \log(\hat{y}_c)$$

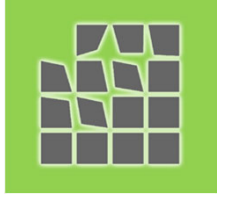**Focal:**

$$-\sum_{c=1}^{C} y_c (1-\hat{y}_c)^\gamma \log(\hat{y}_c)$$

# Other Semantic Loss Functions

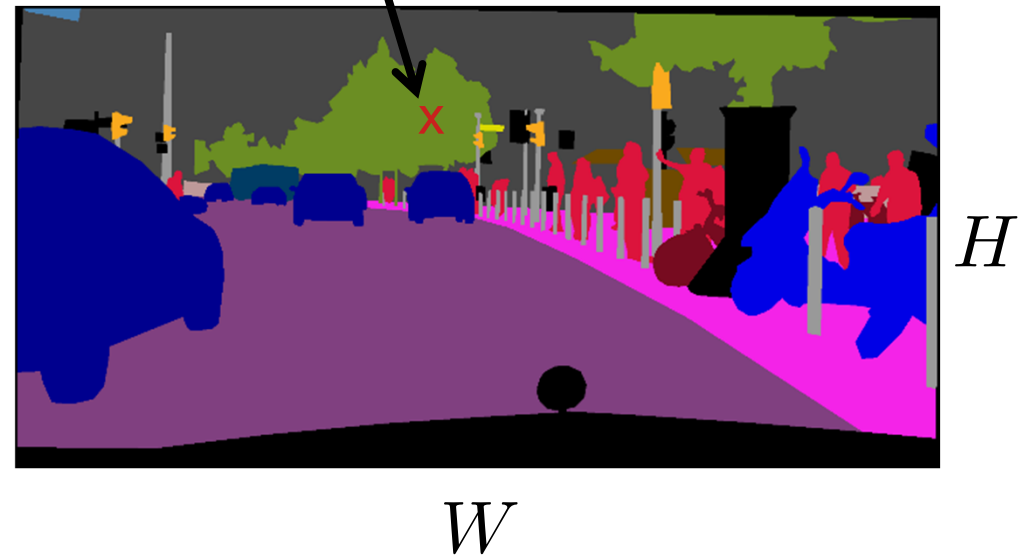- Dice loss (3DV 2016)

- Focal loss (ICCV 2017)

- …

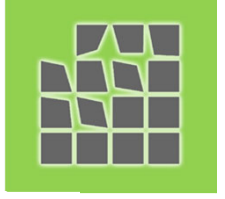- A survey of semantic losses (arXiV 2023)
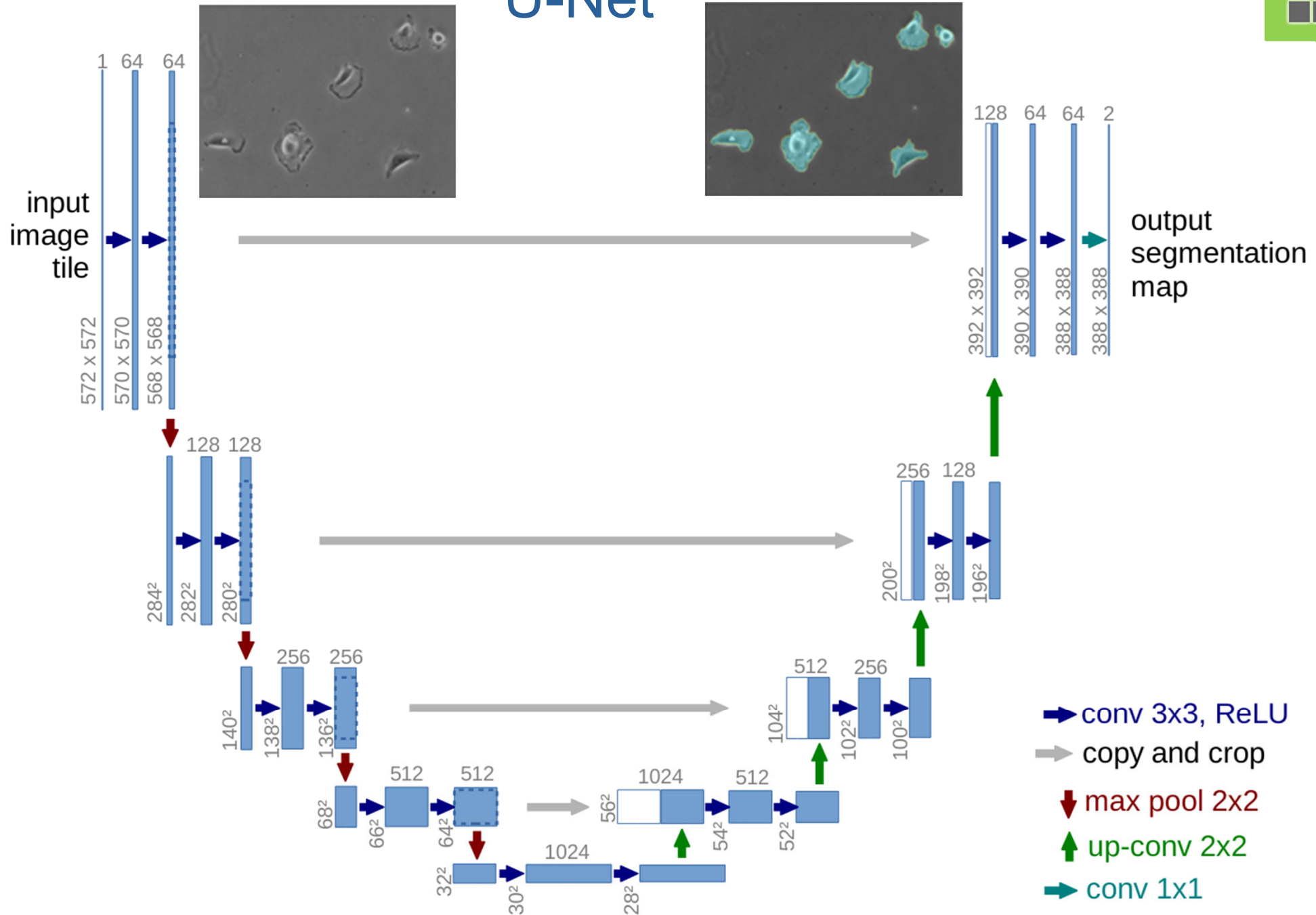
# Semantic Segmentation

- Classify every pixel

- Output: $y \in \mathbb{R}^{W \times H \times C}$
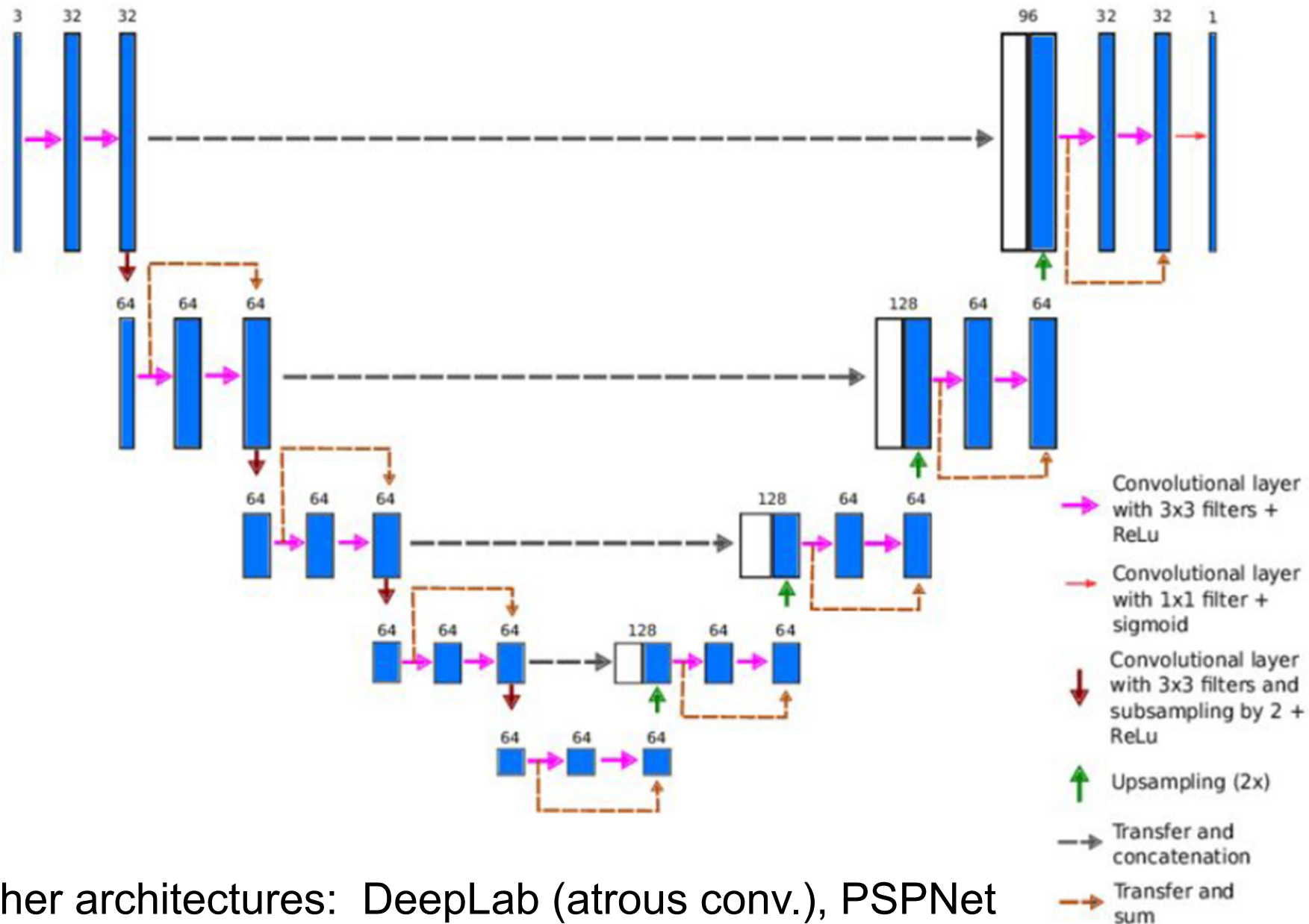
$$[p_1, p_2, \ldots, p_i, \ldots, p_{C-1}, p_C]$$
$$0 \quad 0 \quad 1 \quad 0 \quad 0$$



$H$

$W$

# U-Net



input image tile

output segmentation map

→ conv 3x3, ReLU
⇒ copy and crop
↓ max pool 2x2
↑ up-conv 2x2
→ conv 1x1

Ronneberger et al., U-Net, 2015

21

# Res U-Net



Convolutional layer with 3x3 filters + ReLu

Convolutional layer with 1x1 filter + sigmoid

Convolutional layer with 3x3 filters and subsampling by 2 + ReLu
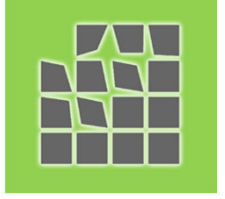
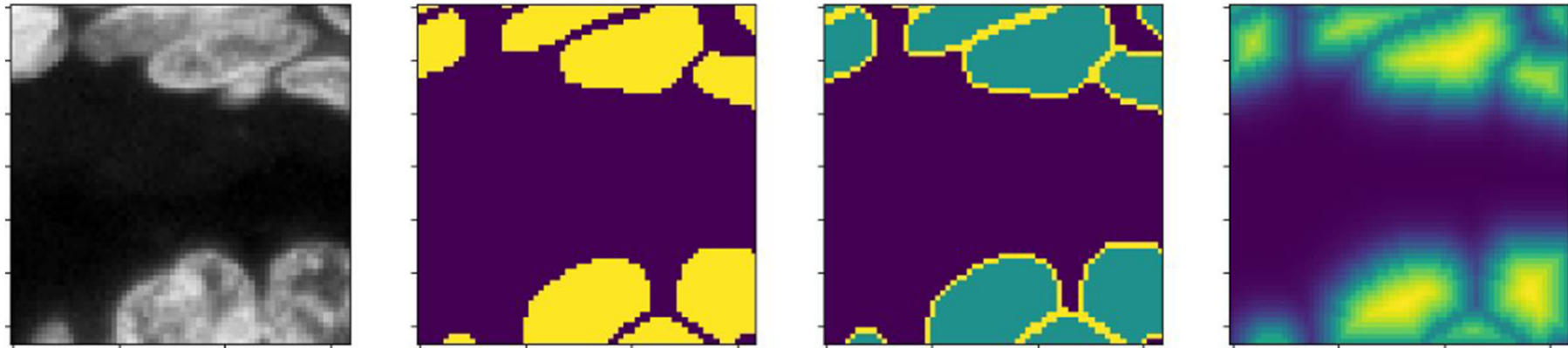Upsampling (2x)

Transfer and concatenation

Transfer and sum

Other architectures:  DeepLab (atrous conv.), PSPNet

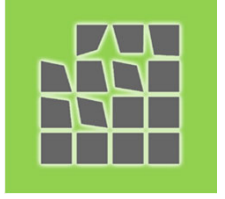# Bottom-Up IS

- Clustering in the (feature & spatial) domain
- Do not predict just object probability per pixel
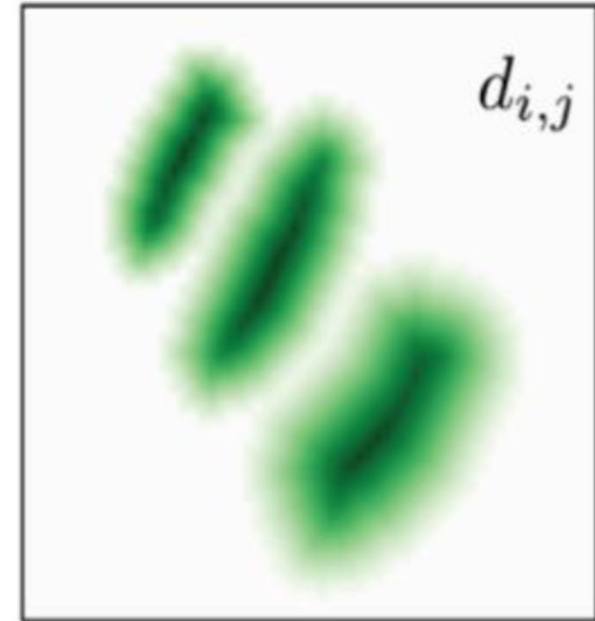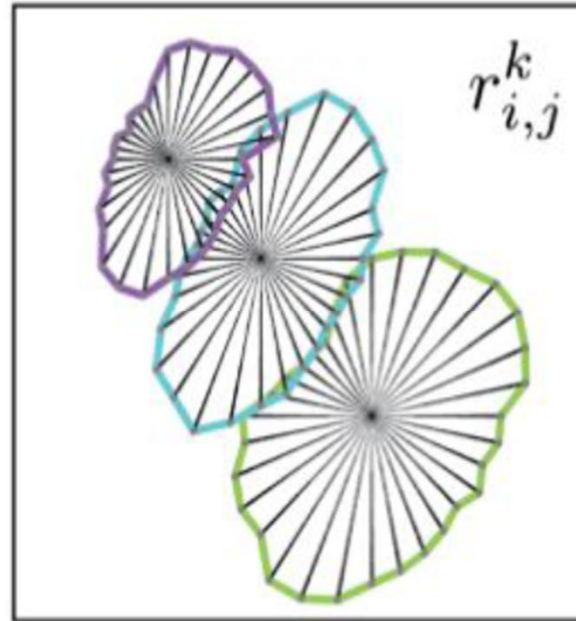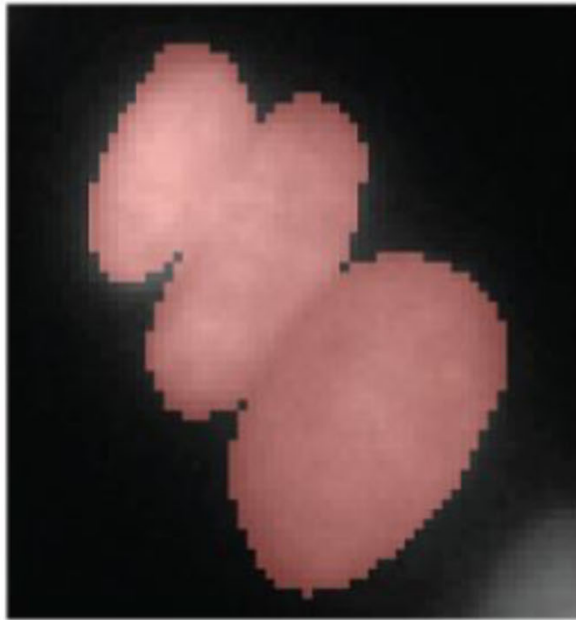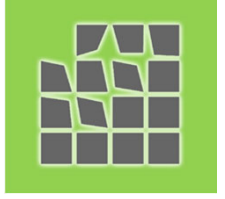- Predict distance from the border



SDM, StarDist, Cellpose, Omnipose, DenoiSeg

SDM  2019

# StarDist

●Predict distance map and star-shape polygon



$r_{i,j}^k$

$d_{i,j}$
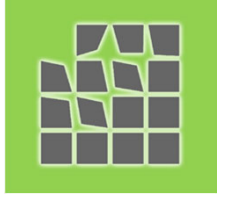
StarDist, 2018

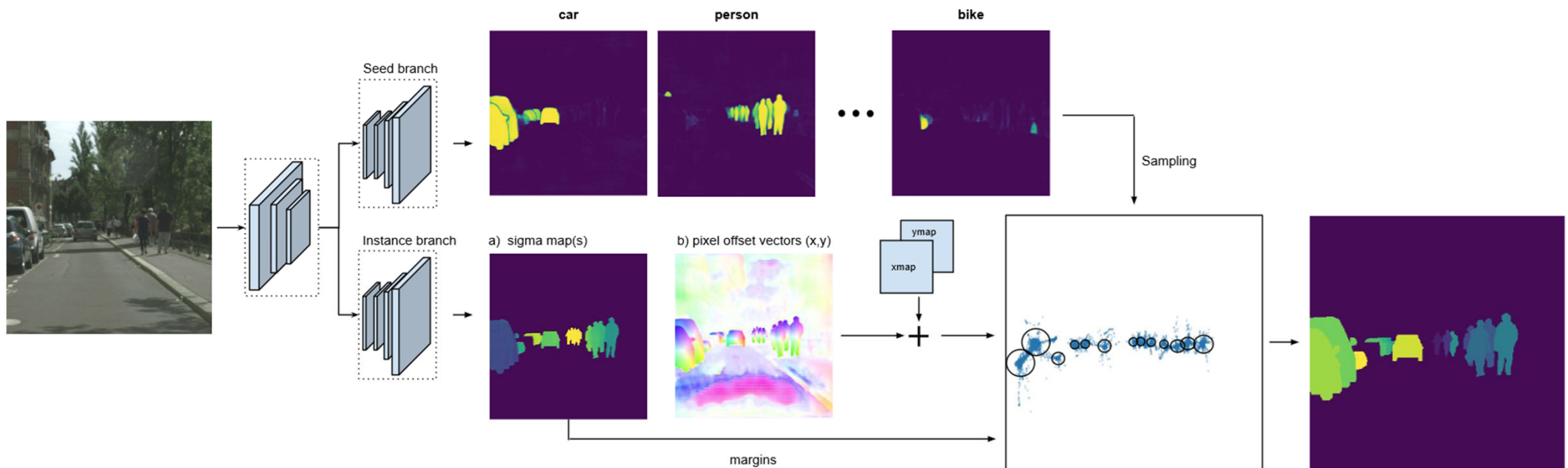# CellPose

• Predict binary mask and gradient flow
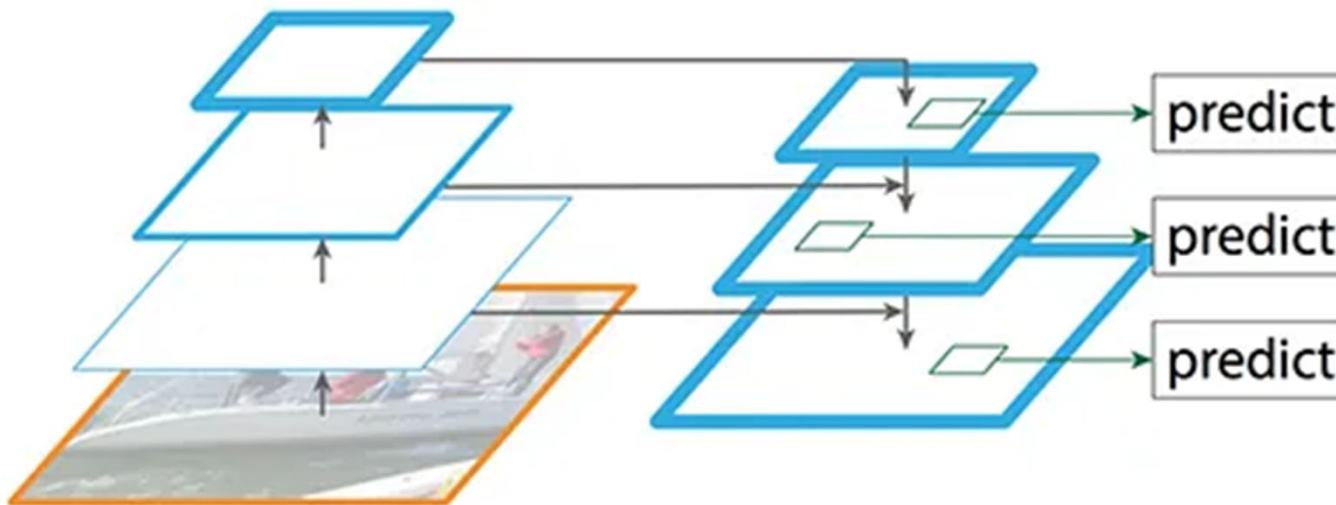


CellPose, 2020

# Pixel offset

- Predict distance, pixel offset from the center, sigma maps (cluster size)
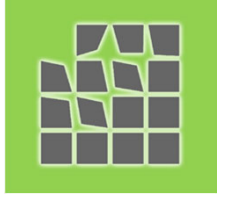
# Object Detection

- Extract features – Feature Pyramid Network

- (Select regions – Region Proposal Network)

- Classify each region (objectness, class, bounding box)
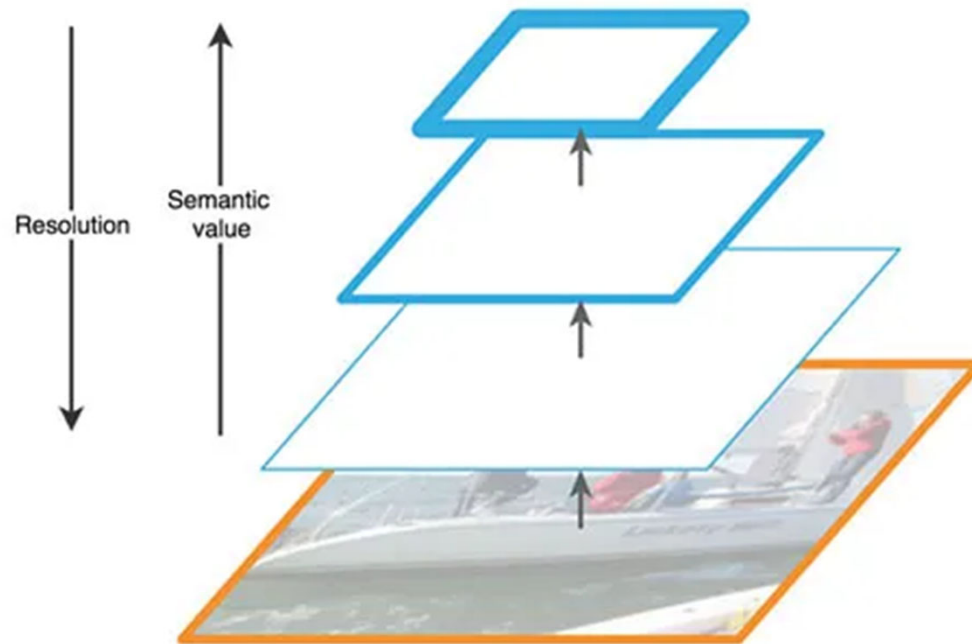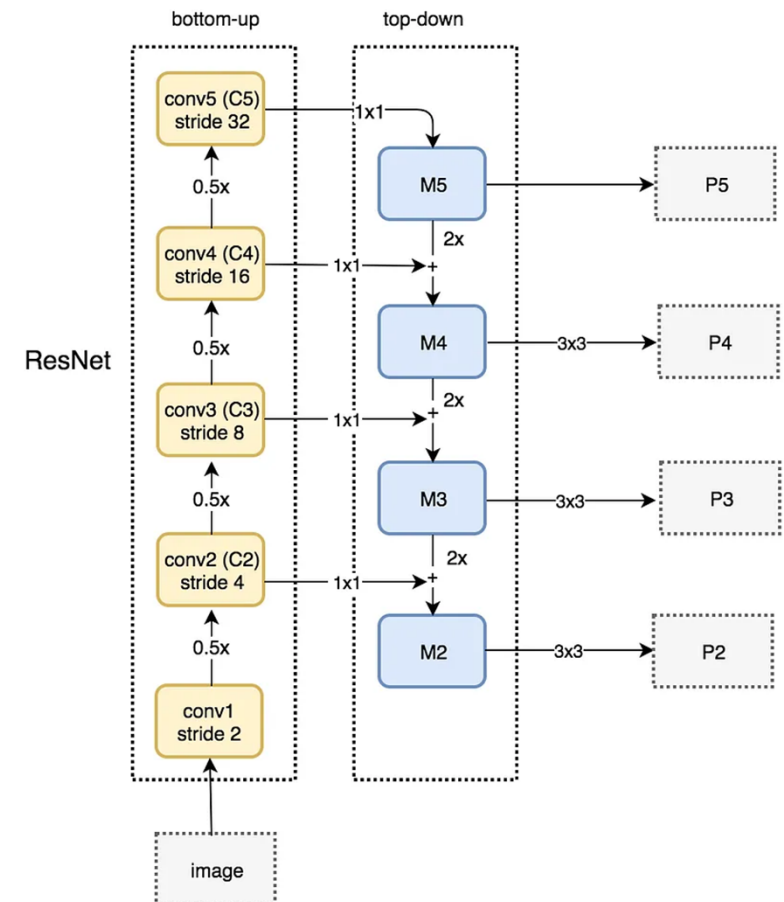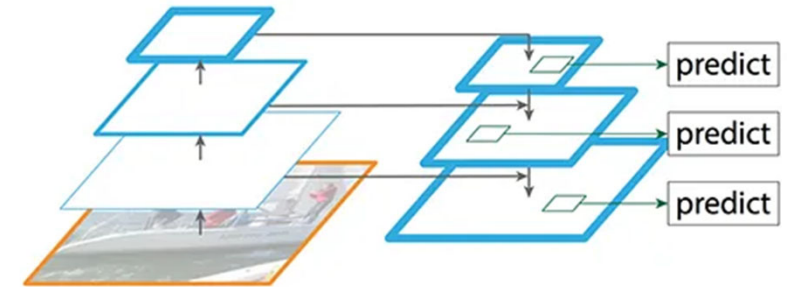
- Non-maximum suppression

# Feature Pyramid Network

- Similar to U-Net
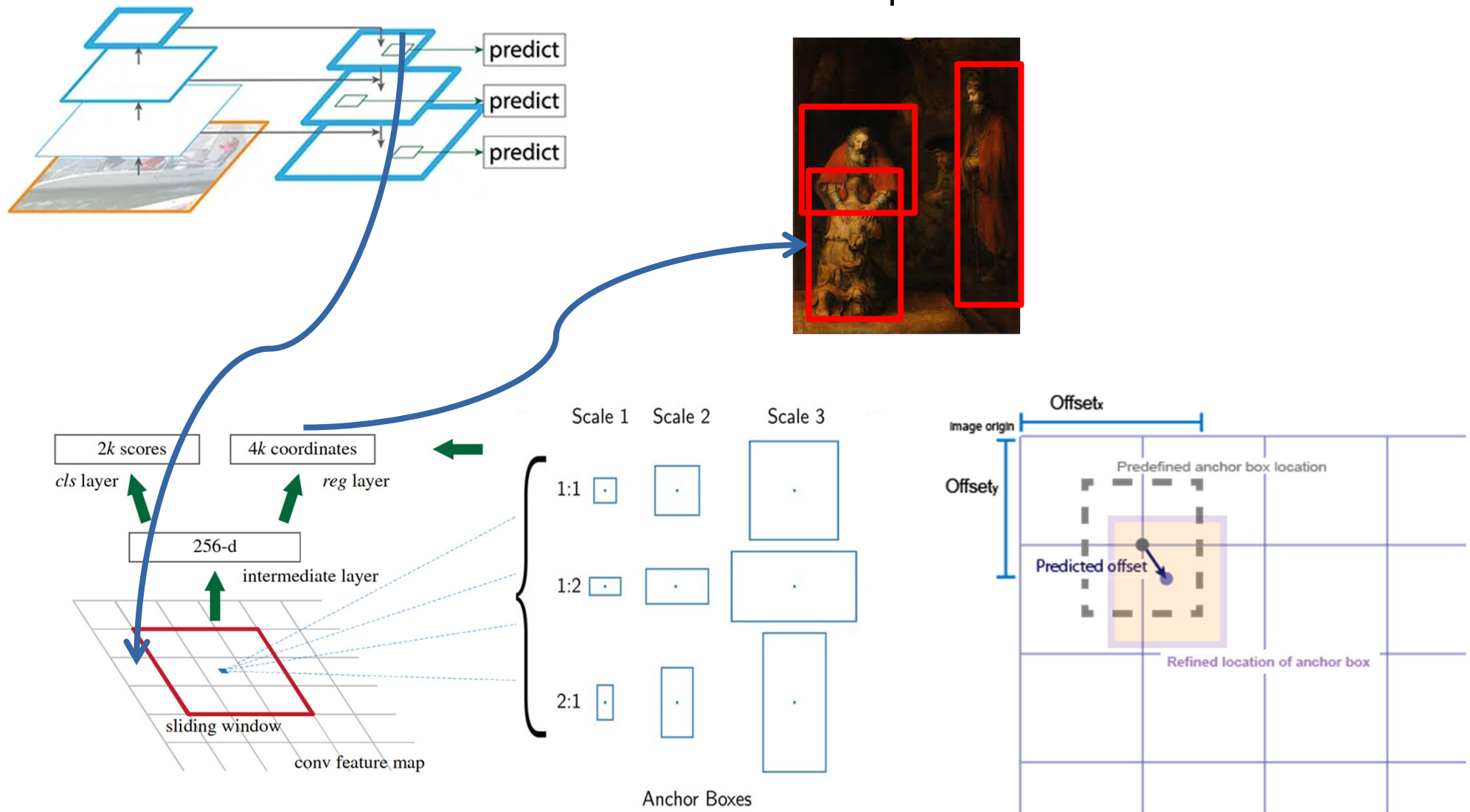- Asymmetric design
- Lateral connections use 1x1conv with addition.



VGG , ResNet

# Region Proposal Network
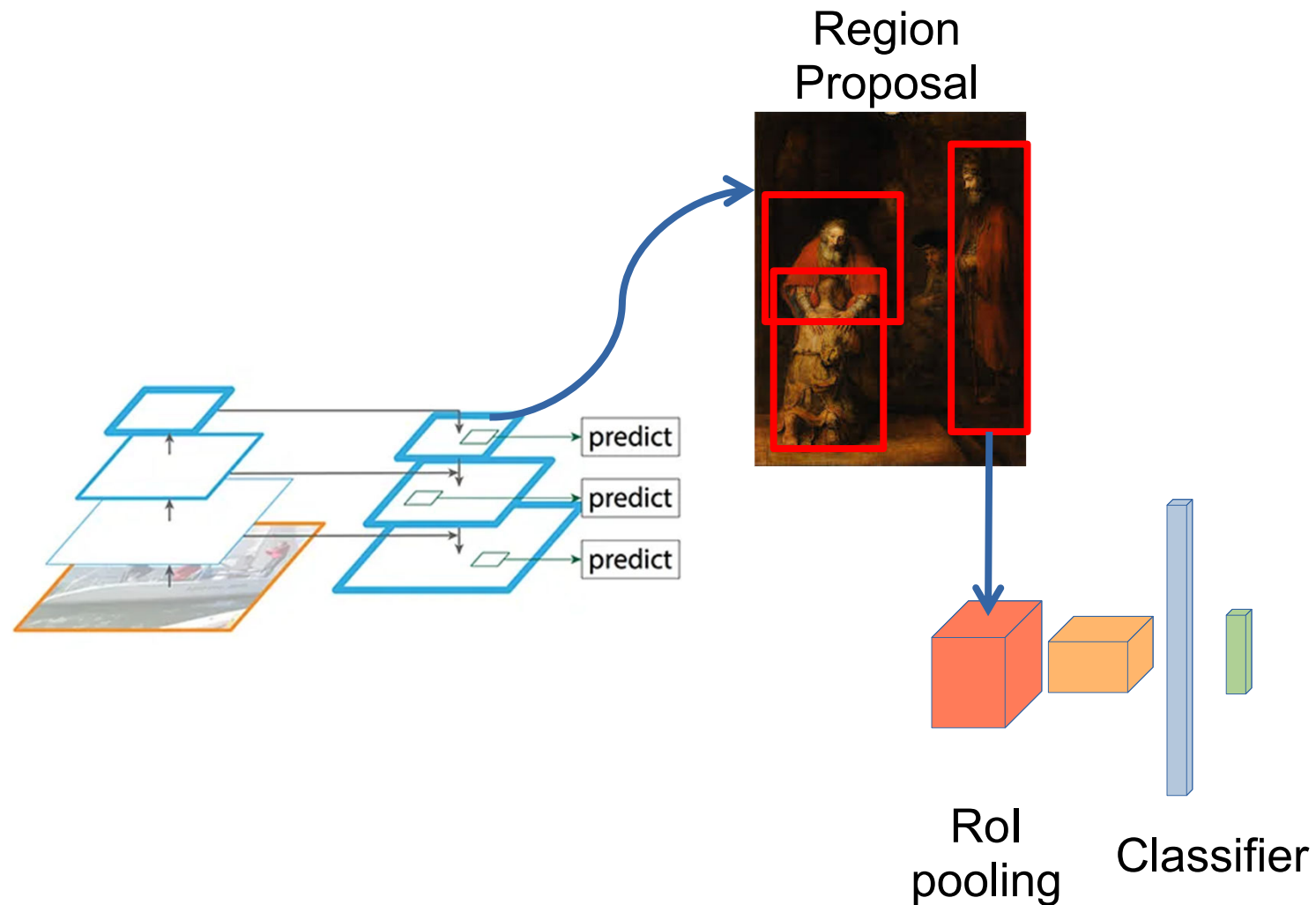


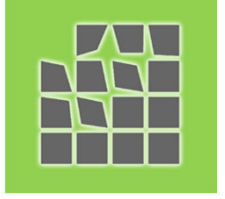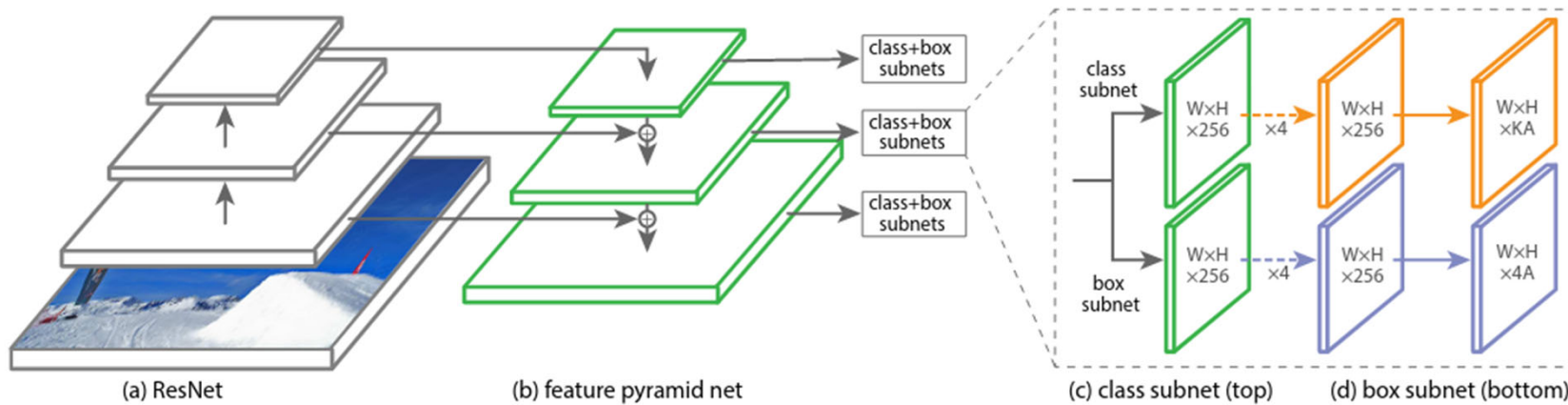Region Proposal

# Two-stage Detector



Region Proposal

RoI pooling

Classifier

predict

predict

predict

Girshick et al., R-CNN, 2015
Girshick, Fast R-CNN, 2015
Ren et al., Faster R-CNN, 2016

# One-stage Detector



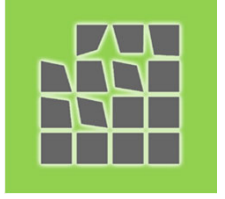(a) ResNet  (b) feature pyramid net  (c) class subnet (top)  (d) box subnet (bottom)

One-stage vs. two-stage:

faster but less accurate.

Non-Maximum Suppression

4A , A, KA

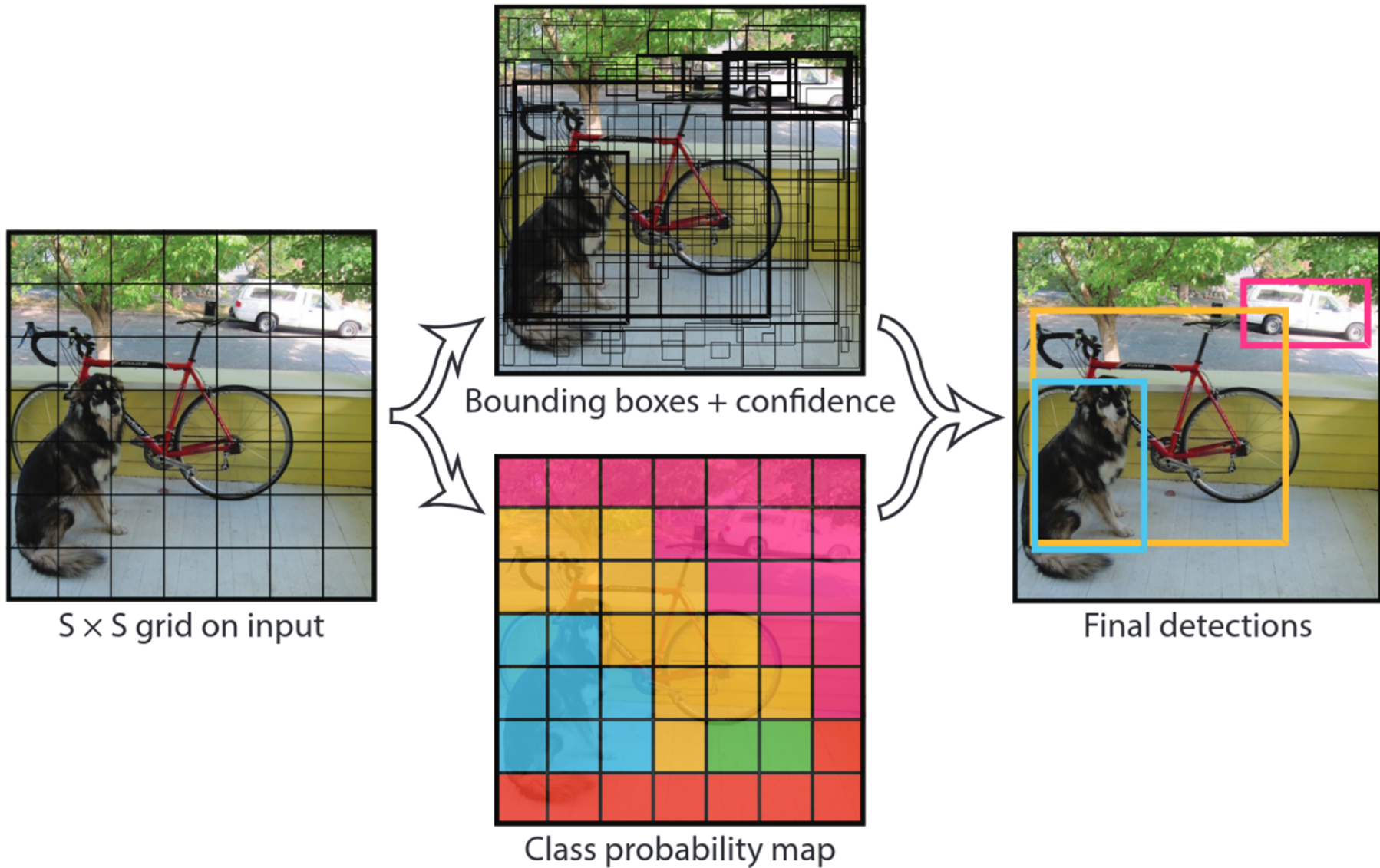- A regions (x,y,w,h) for every location
- P(classes + background) for every location
- Objectness (prediction of IoU)

Redmon et al., YOLO, CVPR 2016
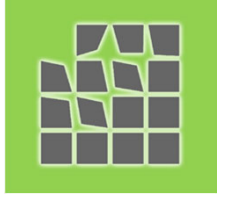Liu et al., SSD, ECCV 2017
Lin et al., RetinaNet, 2018

# YOLO

- You Only Live/Look Once



S × S grid on input

Bounding boxes + confidence

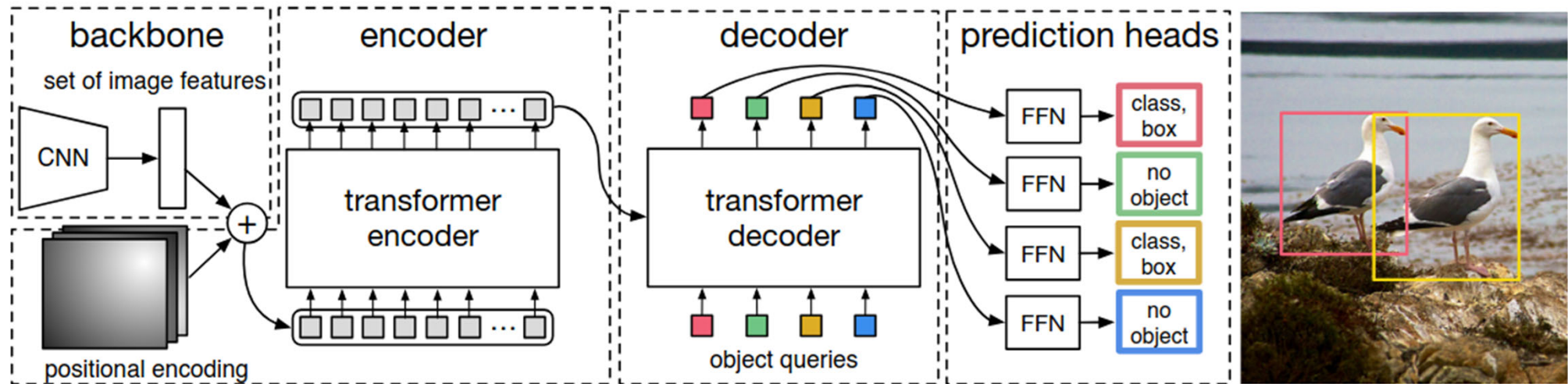Class probability map
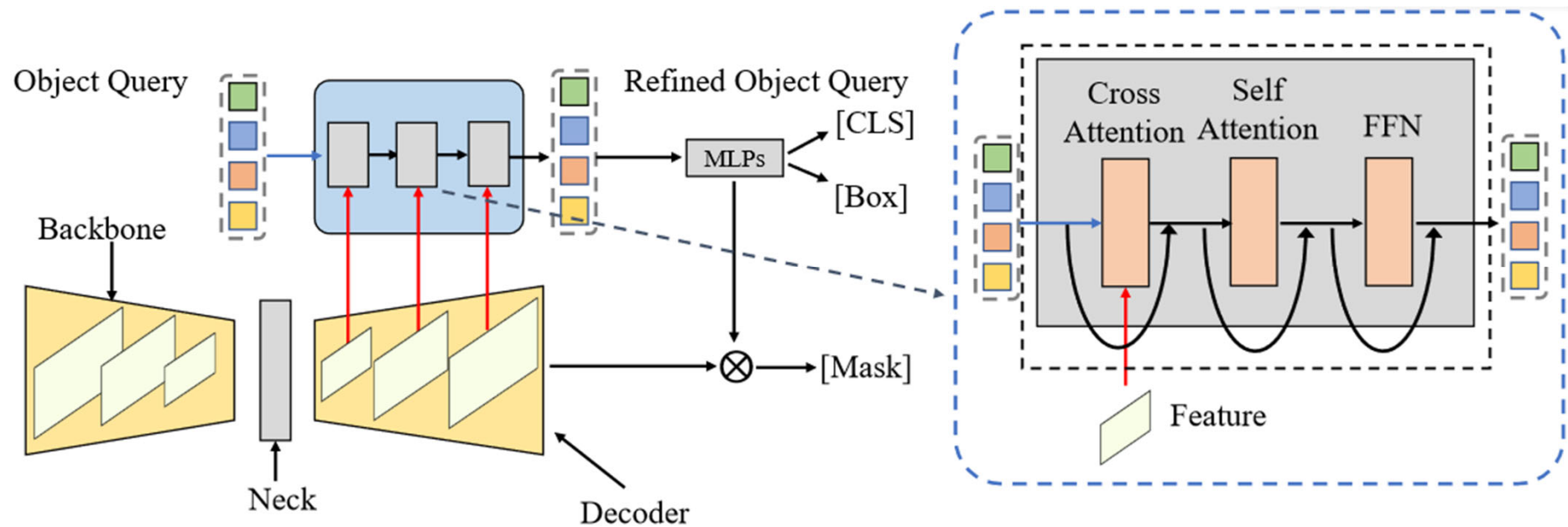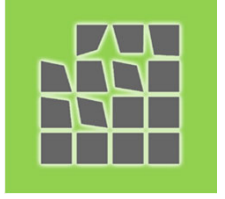
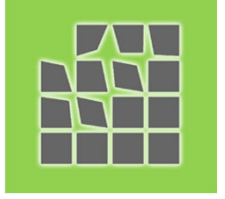Final detections

Redmon et al., YOLO, CVPR 2016

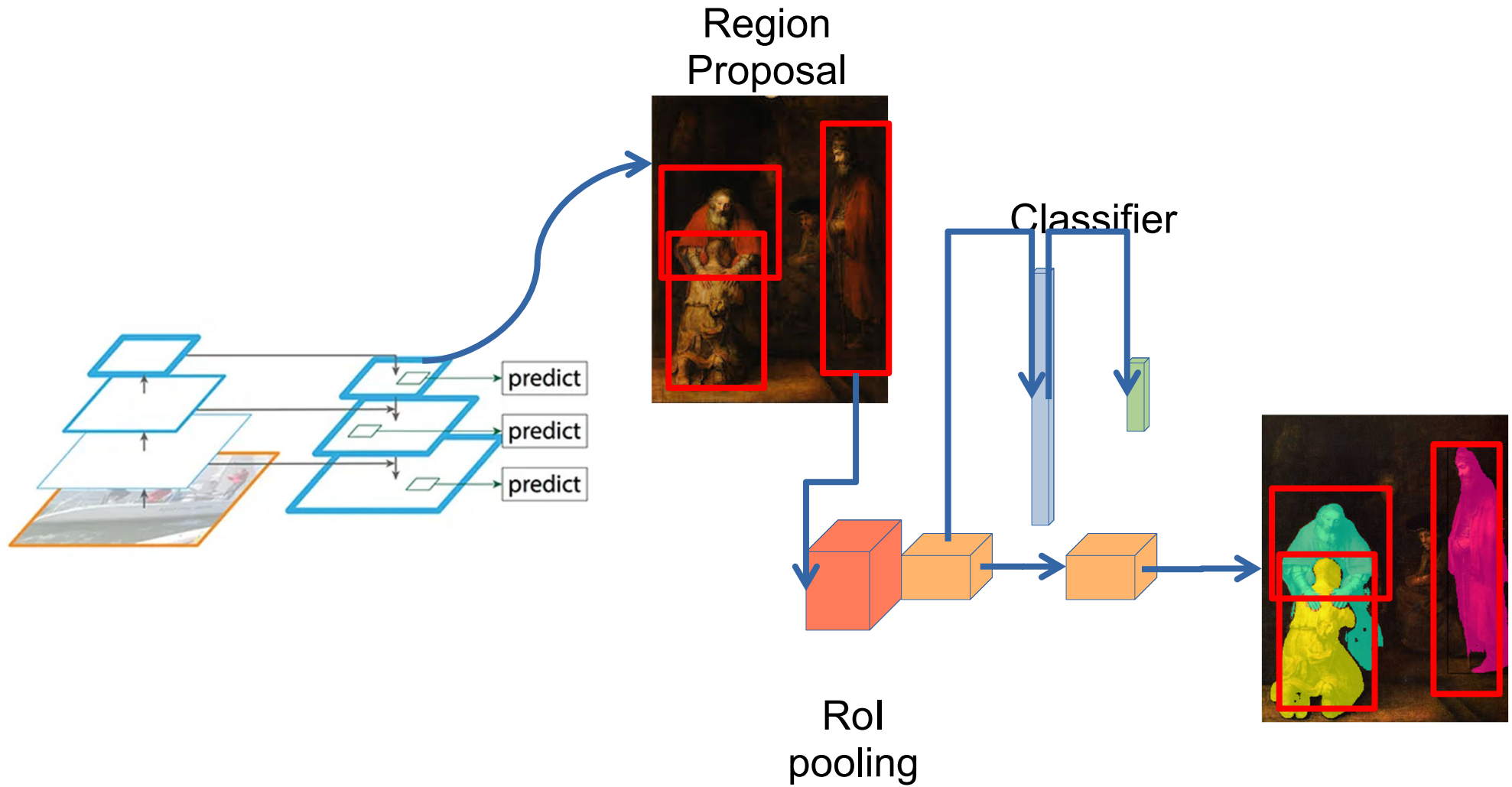# Detection with Transformers

- DETR, ECCV2020
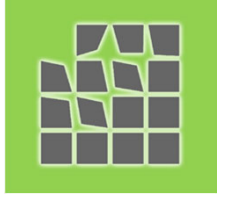
# Meta Architecture



ResNet, ViT, Swin, ConvNeXt

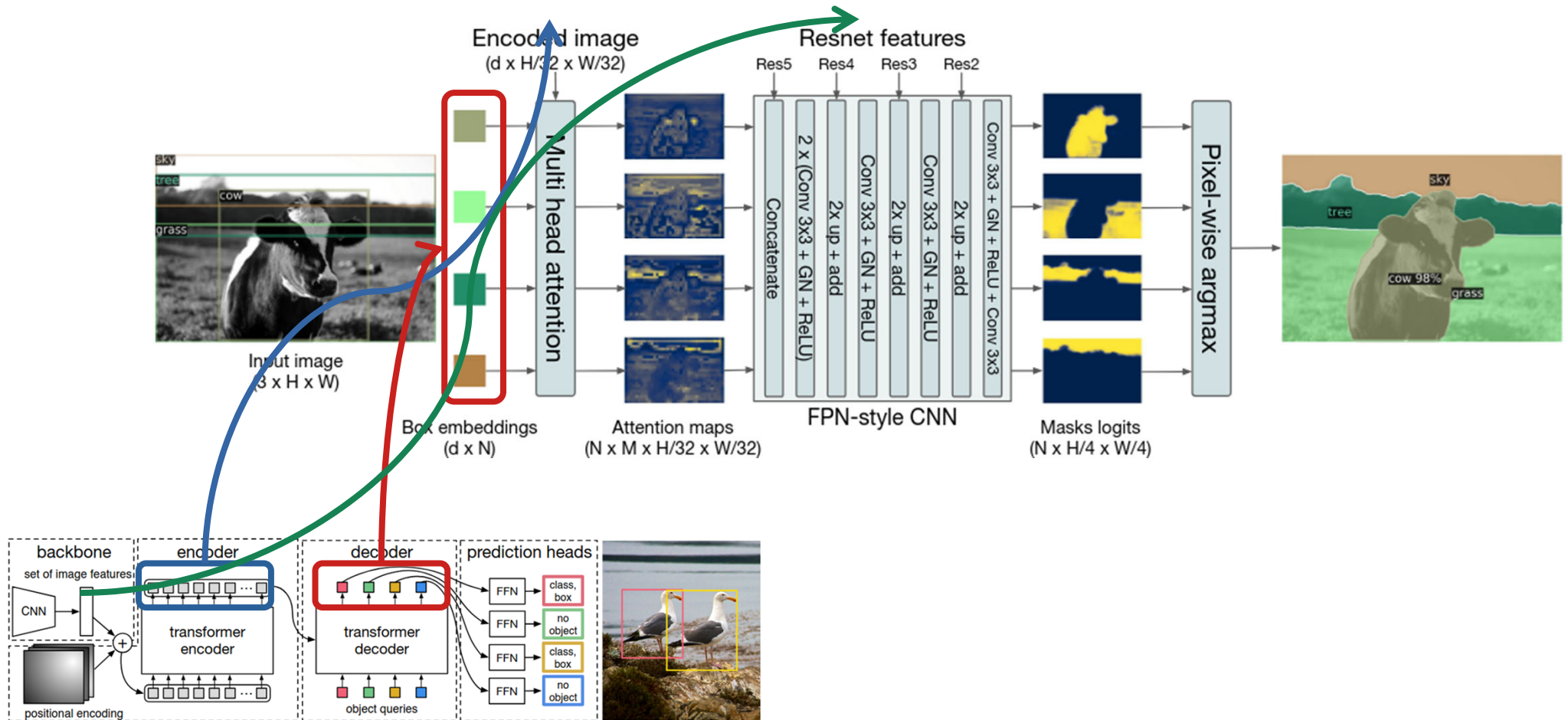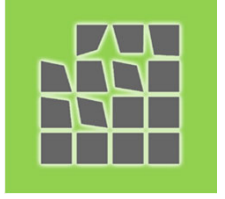- Mask R-CNN, 2018



Region Proposal

Classifier

RoI pooling
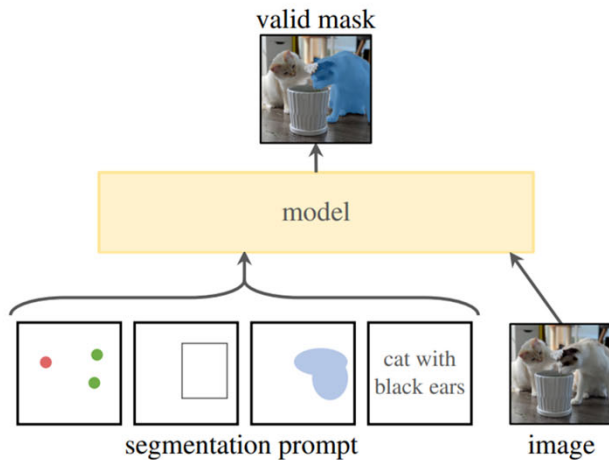
# Top-Down IS

## ●DETR - Panoptic

# Panoptic Segmentation

- Universal segmentation

- OneFormer
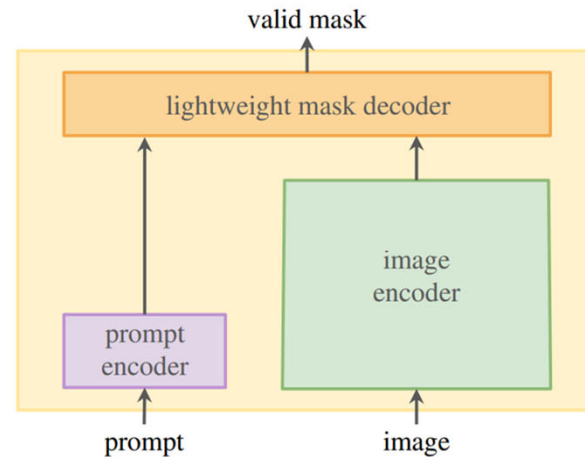
- Masked-attention Mask Transformer for Universal Image Segmentation

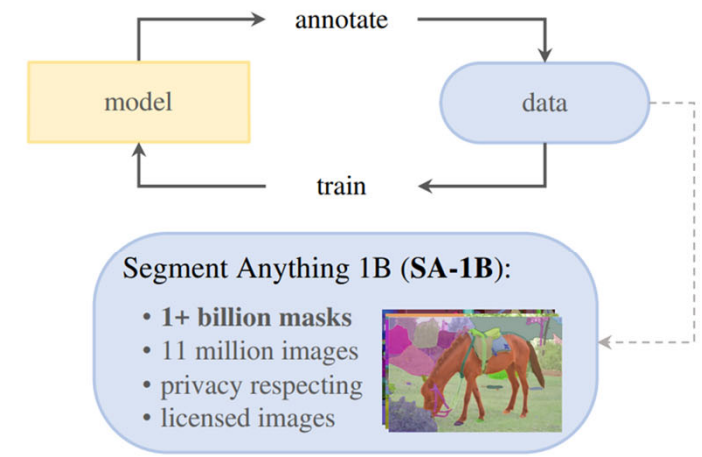# Segment Anything Model



(a) **Task**: promptable segmentation

(b) **Model**: Segment Anything Model (**SAM**)

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

Segment Anything 1B (**SA-1B**):
- **1+ billion masks**
- 11 million images
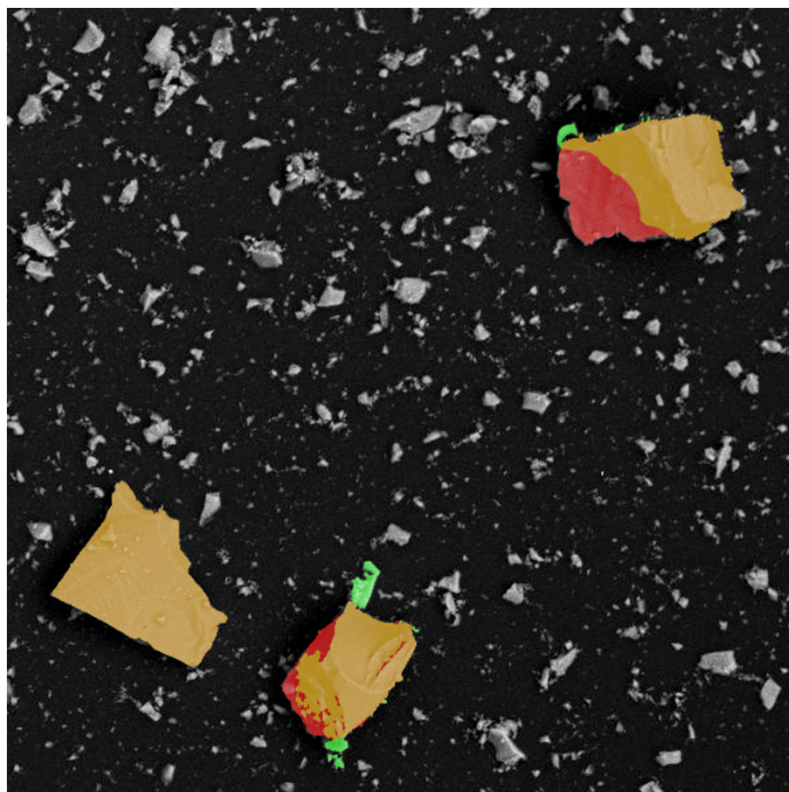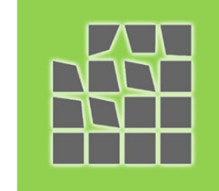- privacy respecting
- licensed images

CVAT.AI (annotation platform)

SAM 2023

# Evaluation Metric

- Semantic: mIoU

- Instance: mAP

- Panoptic: Panoptic Quality

# Intersection over Union
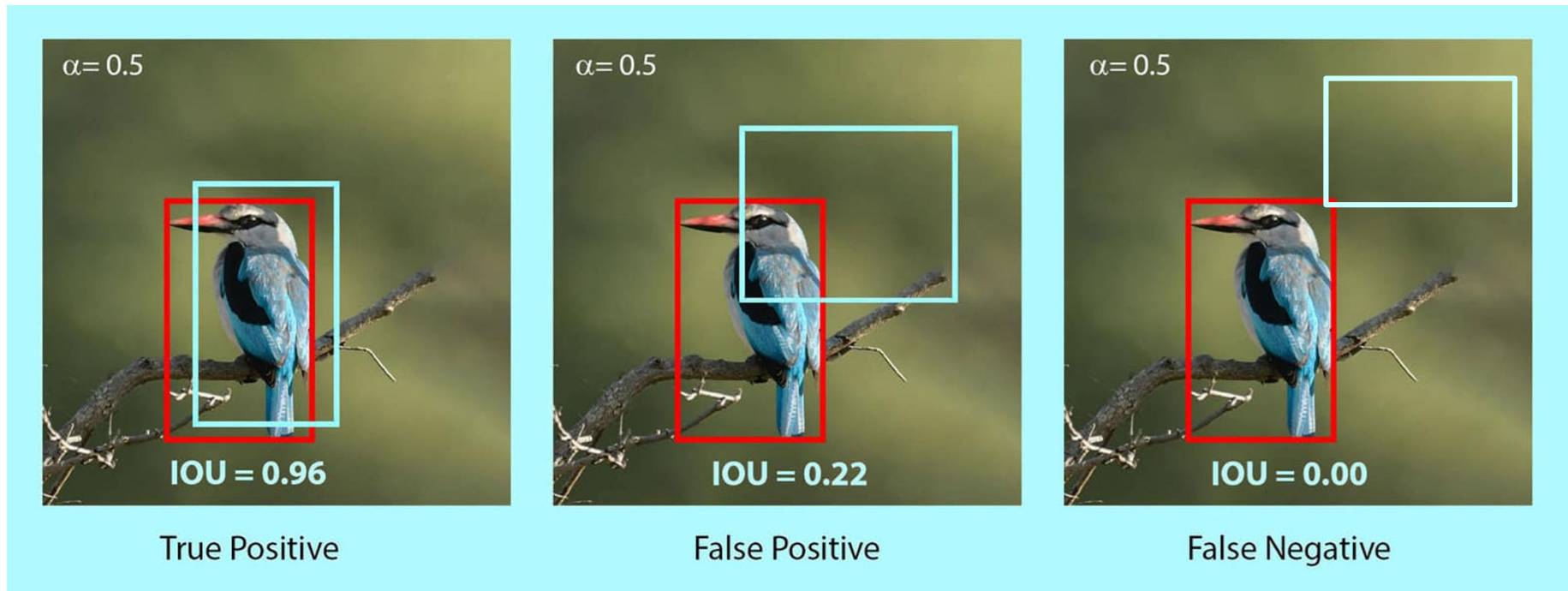


$$\mathrm{IoU} = \frac{|\mathrm{G} \cap \mathrm{P}|}{|\mathrm{G}| + |\mathrm{P}| - |\mathrm{G} \cap \mathrm{P}|}$$

- Mean IoU:   $$\mathrm{mIoU} = \sum_c \mathrm{IoU}_c$$

# AP - Average Precision

• Threshold: $IoU > \alpha$



$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = 2\frac{Precision \cdot Recall}{Precision + Recall}$$

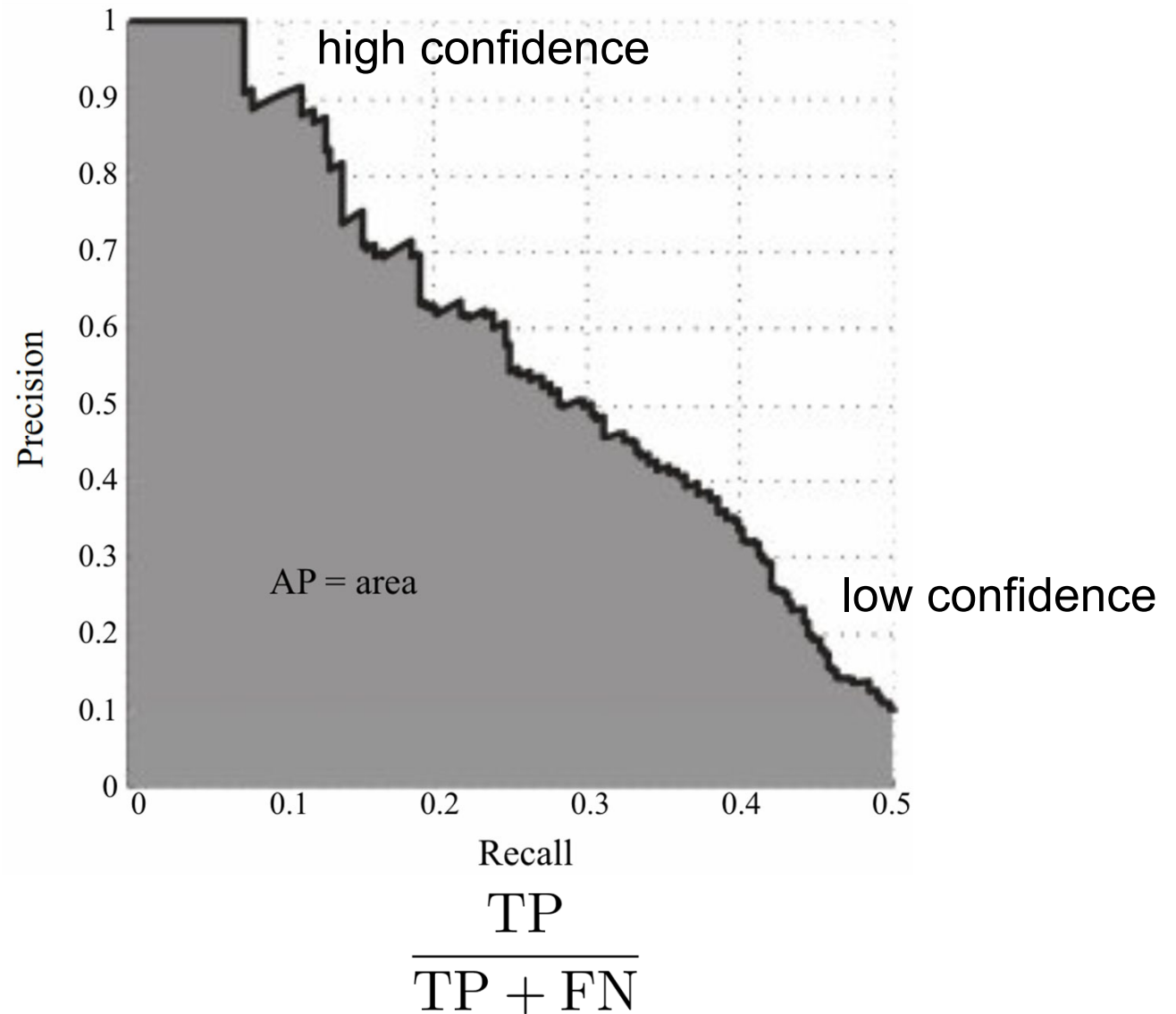https://www.ridgerun.ai/post/mean-average-precision-map-and-other-object-detection-metrics
https://jss367.github.io/what-do-these-different-ap-values-mean.html

42

# mAP

- For the given threshold of IoU:

$$\frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{mAP} = \sum_{c=1}^{C} \text{AP}_c$$



high confidence

AP = area

low confidence

Precision

Recall

$$\frac{\text{TP}}{\text{TP} + \text{FN}}$$

# Thank You

# Panoptic Quality

- PQ per class:

$$\text{PQ}_c = \frac{\sum_{(P,G)\in\text{TP}} \text{IoU}(P,G)}{|\text{TP}| + \frac{1}{2}|\text{FP}| + \frac{1}{2}|\text{FN}|}$$

$$\text{PQ} = \sum_c \text{PQ}_c$$