

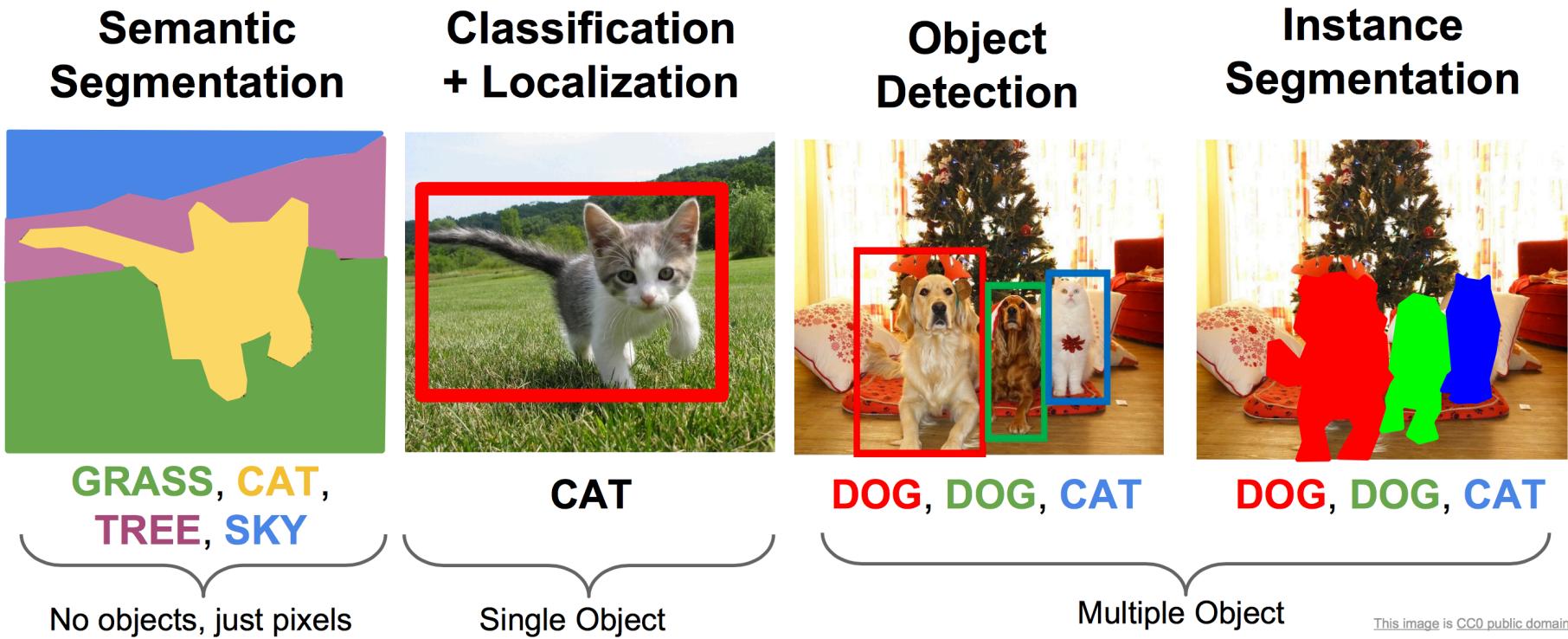
Сегментация и поиск объектов на изображениях

Белов Андрей
ФИВТ МФТИ, ABBYY

План лекции

- Обзор задач
- Датасеты
- Semantic Segmentation
- Object Localization
- Object Detection
- Instance Segmentation

Обзор задач



Датасеты

MS COCO

<http://cocodataset.org/#overview>

Instances:



Stuff:



Keypoints:



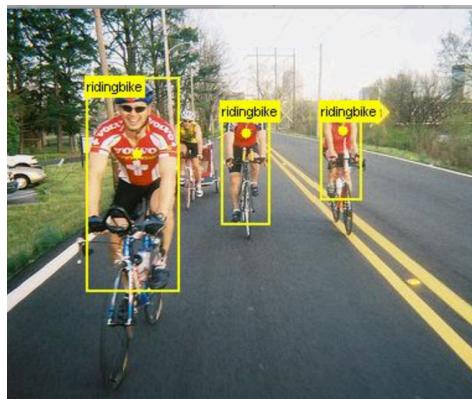
PASCAL VOC

<http://host.robots.ox.ac.uk/pascal/VOC/voc2012/index.html>

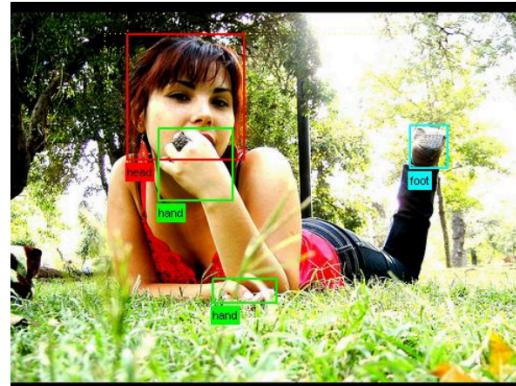
Image



Objects



Person Layout



Датасеты

COCO

<http://cocodataset.org/#overview>

Сегментация + Keypoints для людей

Pascal Voc

<http://host.robots.ox.ac.uk/pascal/VOC/>

ImageNet

<http://image-net.org/challenges/LSVRC/2017/index#comp>

Object localization + Objects Detection (в основном используется для классификации)

Places

<http://placeschallenge.csail.mit.edu>

CityScapes

<https://www.cityscapes-dataset.com>

CamVid

<http://mi.eng.cam.ac.uk/research/projects/VideoRec/CamVid/>

SpaceNet

<https://spacenetchallenge.github.io>

И еще тонны на Kaggle и других площадках

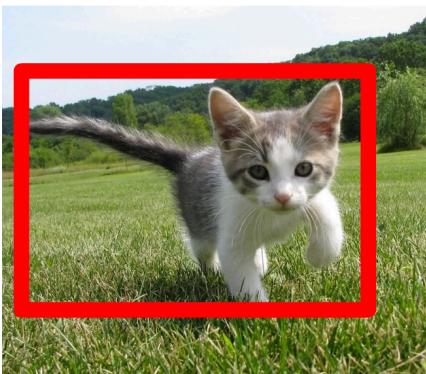
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



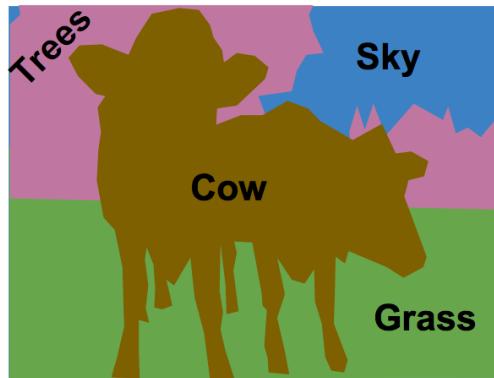
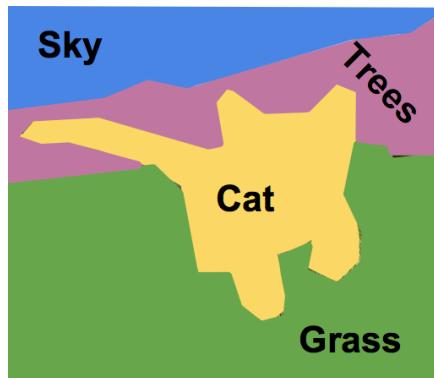
DOG, DOG, CAT

This image is CC0 public domain

Semantic Segmentation

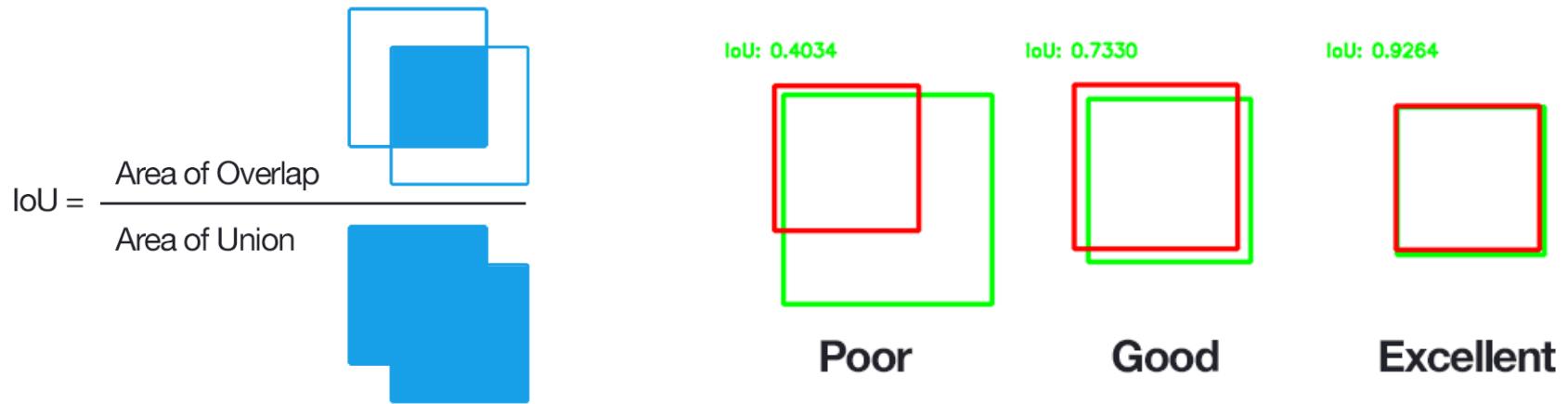
Размечаем каждый пиксель изображения

Не разделяем разные объекты одного класса



Основные метрики

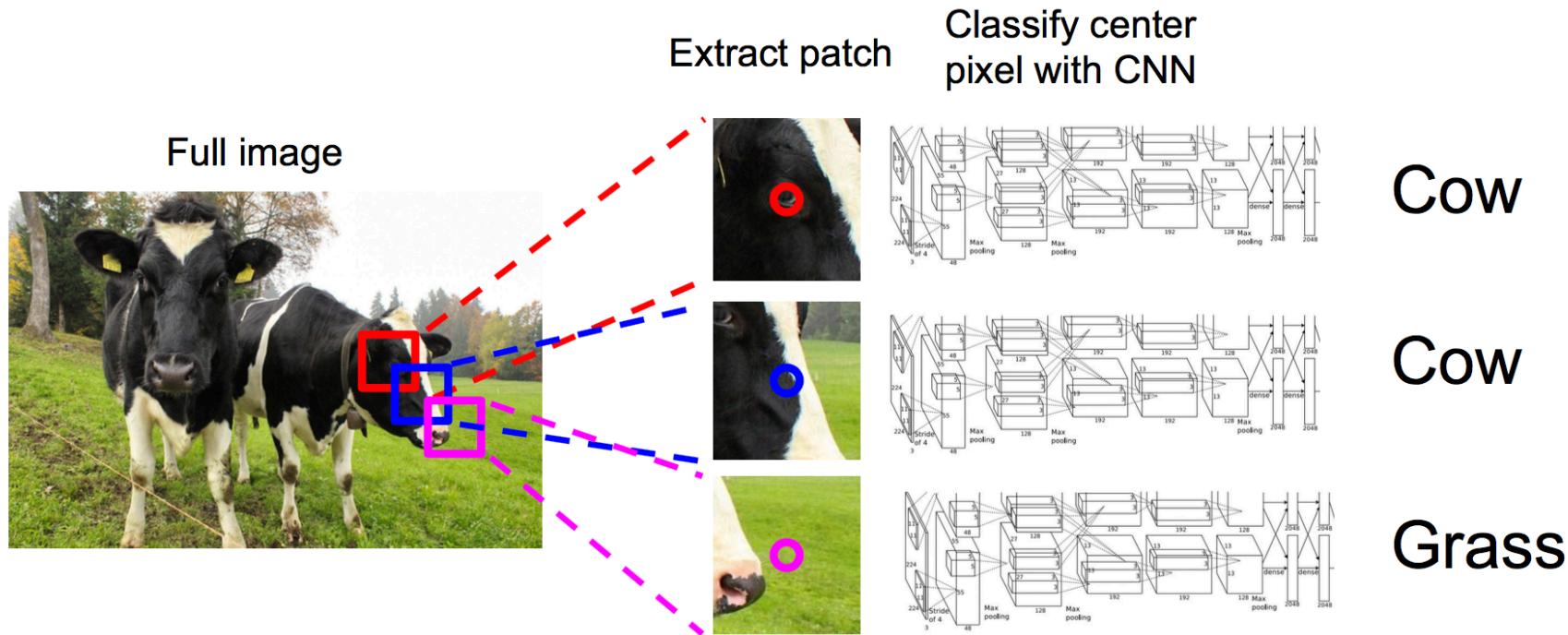
- Intersection over Union



- Dice coefficient

$$QS = \frac{2|X \cap Y|}{|X| + |Y|}$$

Сегментация: идея 1

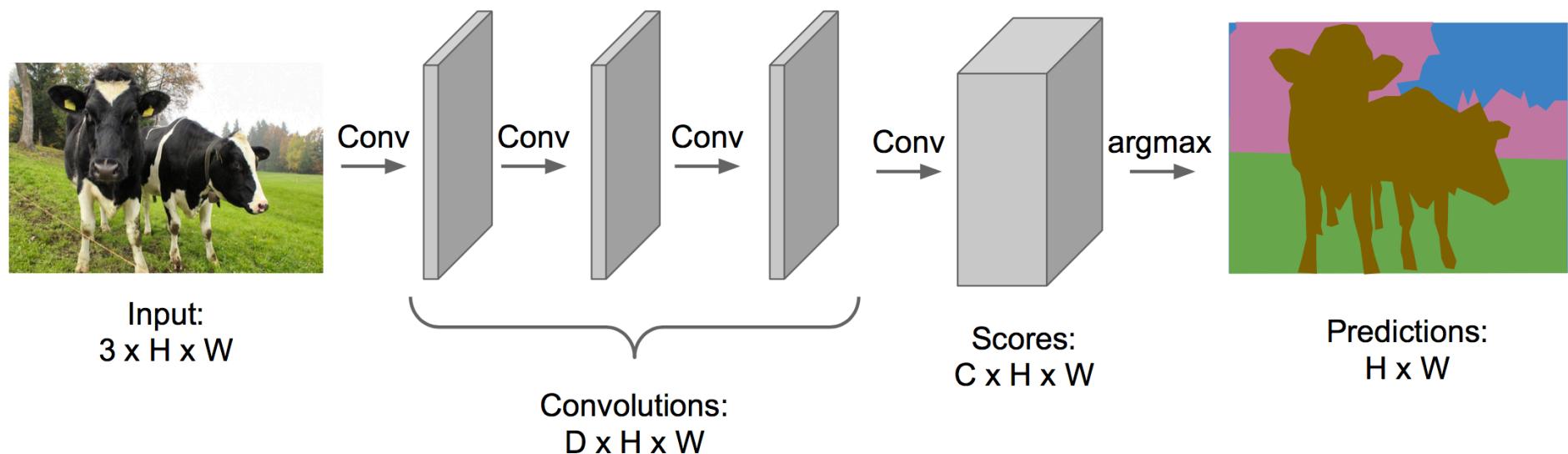


Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Неэффективно, т.к. признаки для пересекающихся фрагментов каждый раз вычисляются заново.

Сегментация: идея 2

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



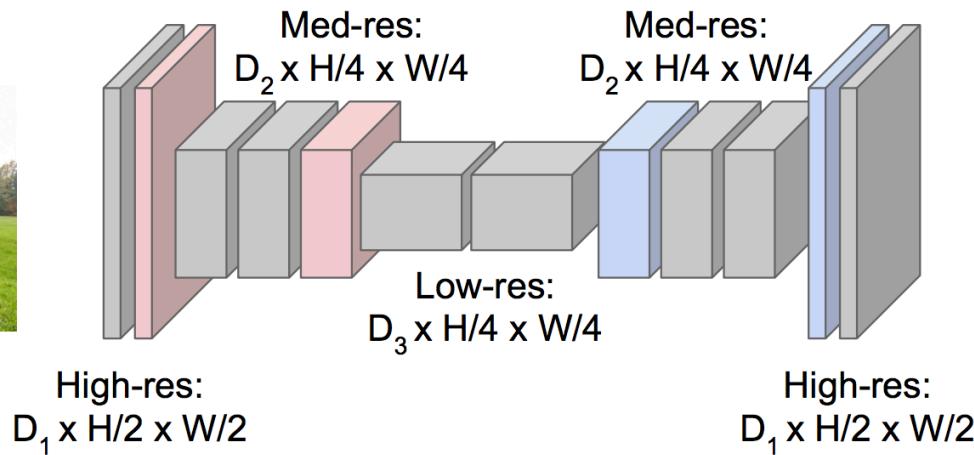
Очень много вычислений, если применять свертки на
полном разрешении изображения.

Сегментация: идея 3

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



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

Downsampling – pooling или strided convolutions
Но как делать upsampling?

Unpooling

Nearest Neighbor

1	2
3	4



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

Input: 2 x 2

Output: 4 x 4

“Bed of Nails”

1	2
3	4



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

Input: 2 x 2

Output: 4 x 4

Max-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

Output: 2 x 2



Max Unpooling

Use positions from pooling layer

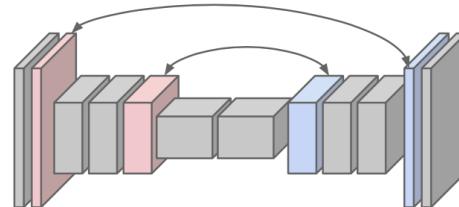
1	2
3	4

Input: 2 x 2

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

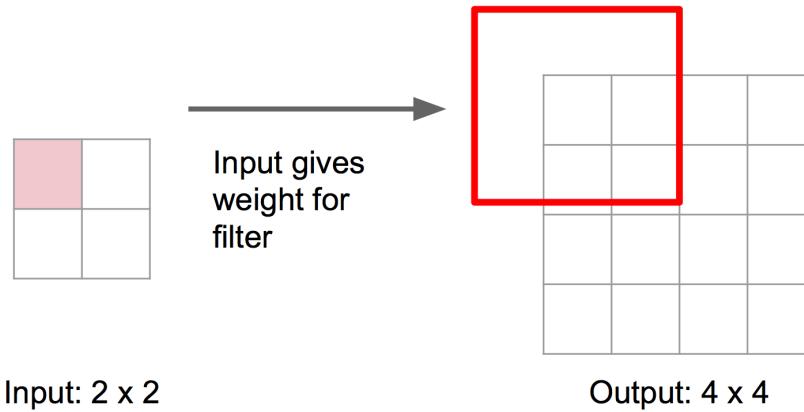
Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers

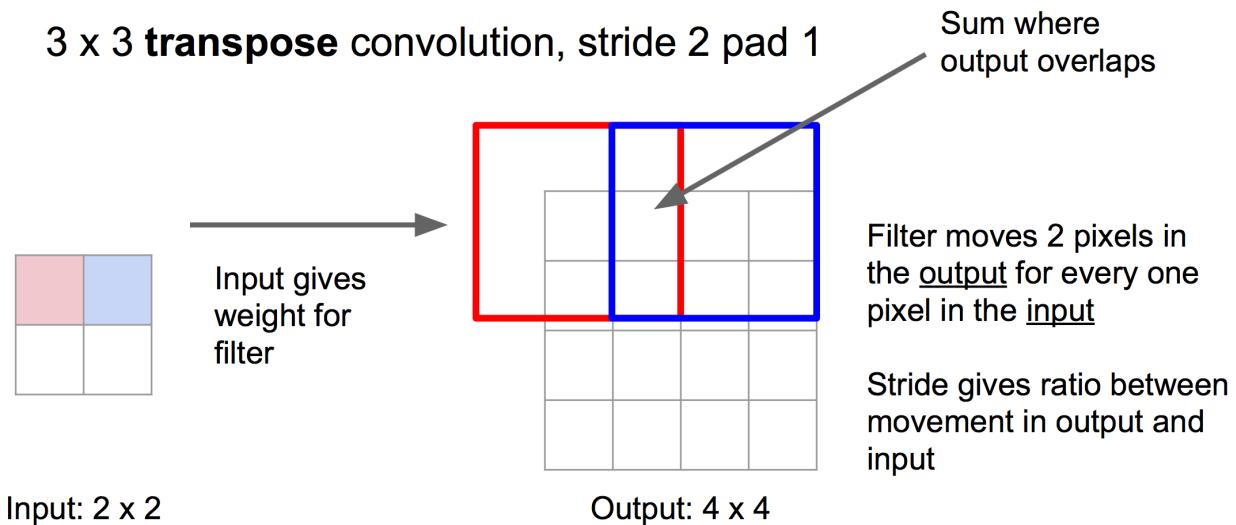


Transposed convolutions

3 x 3 **transpose** convolution, stride 2 pad 1

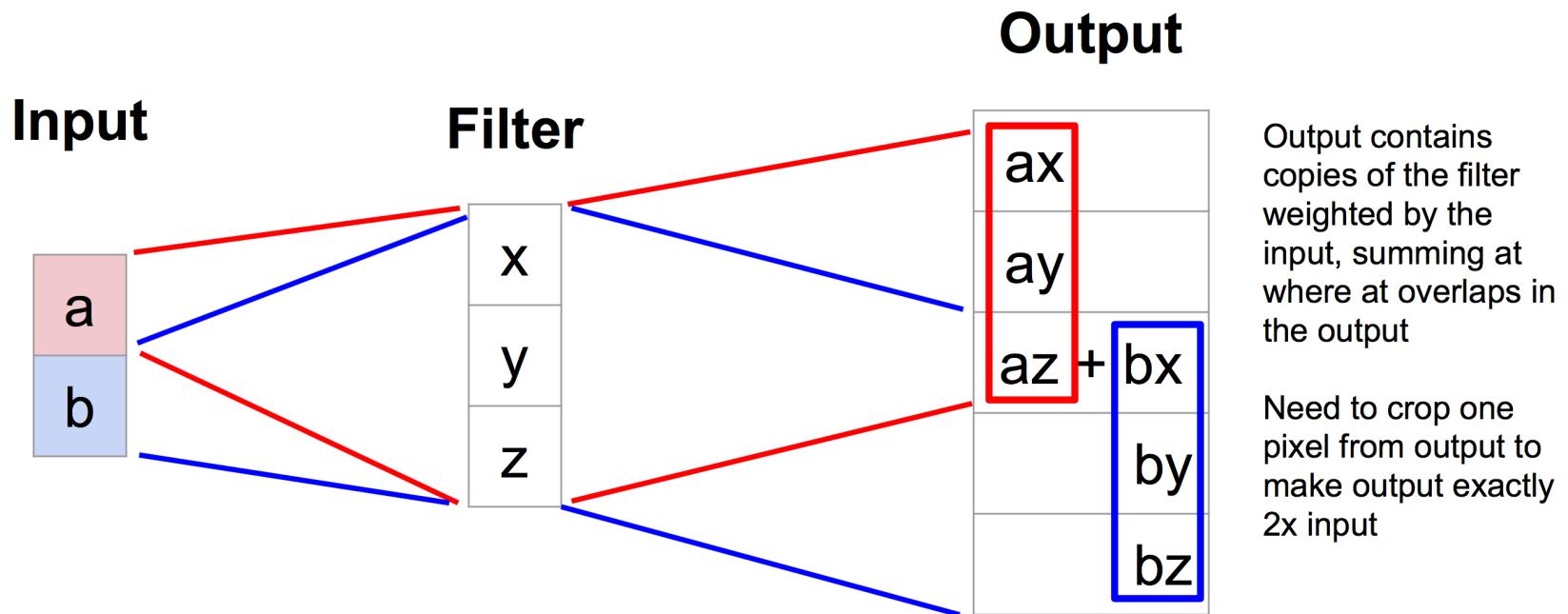


Transposed convolutions

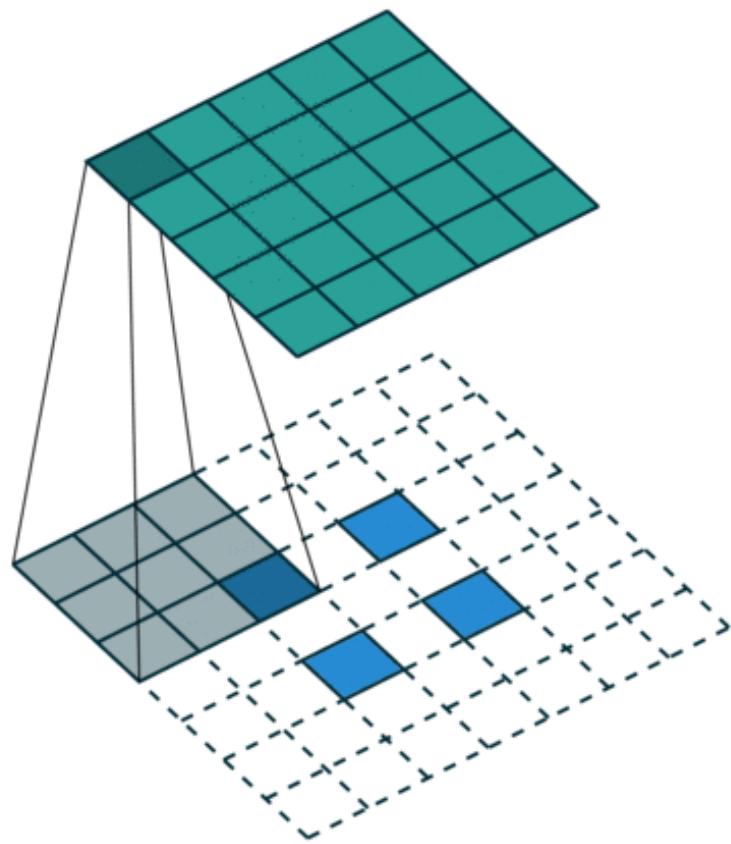


Transposed convolutions.

Одномерный пример



Суммирование может приводить к появлению решетки, поэтому иногда используют свертки kernel 2 stride 2 или kernel 4 stride 2



Архитектуры: U-Net

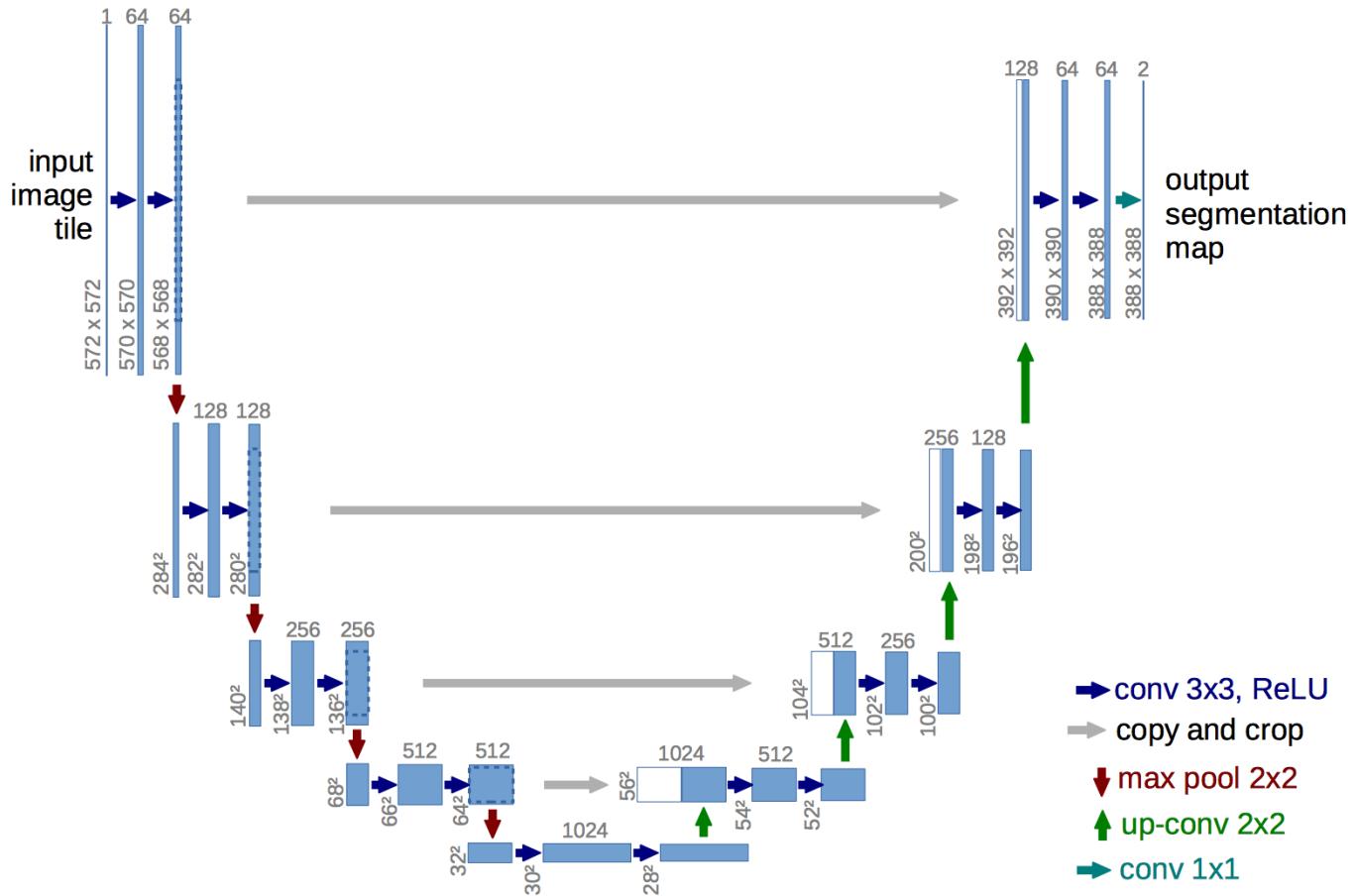


Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

FCN - Fully convolutional networks for semantic segmentation

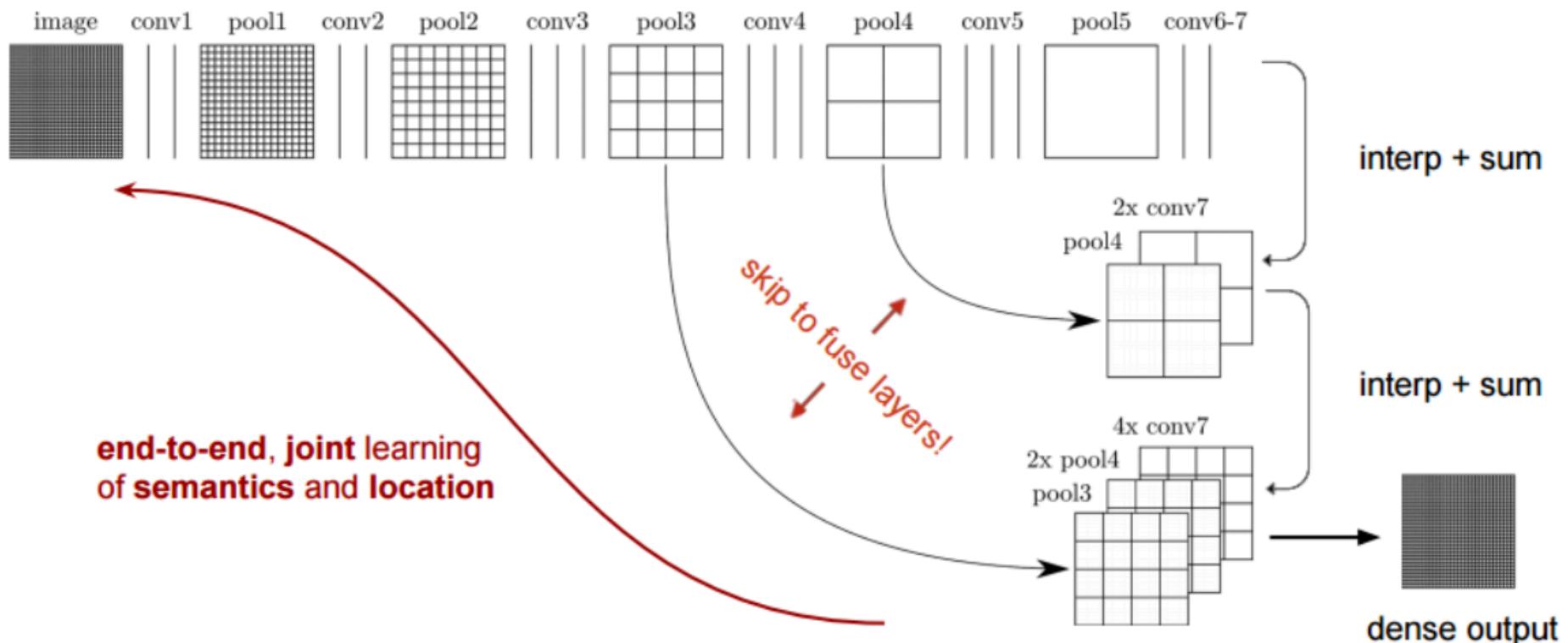
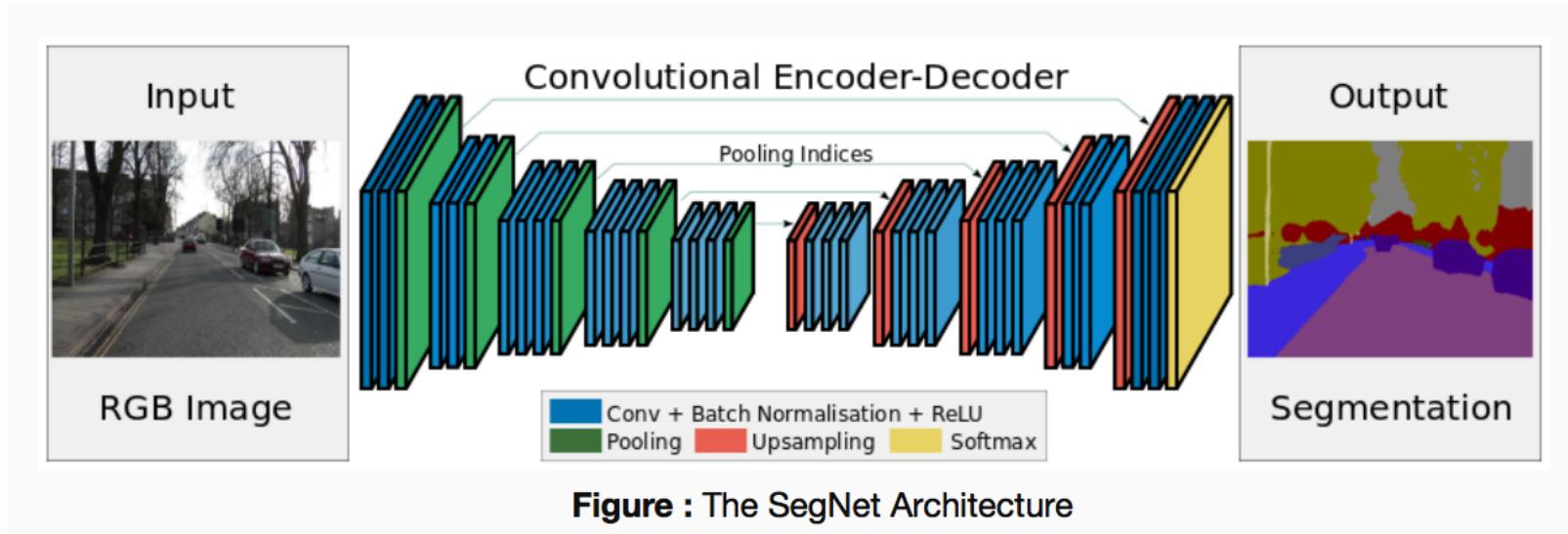
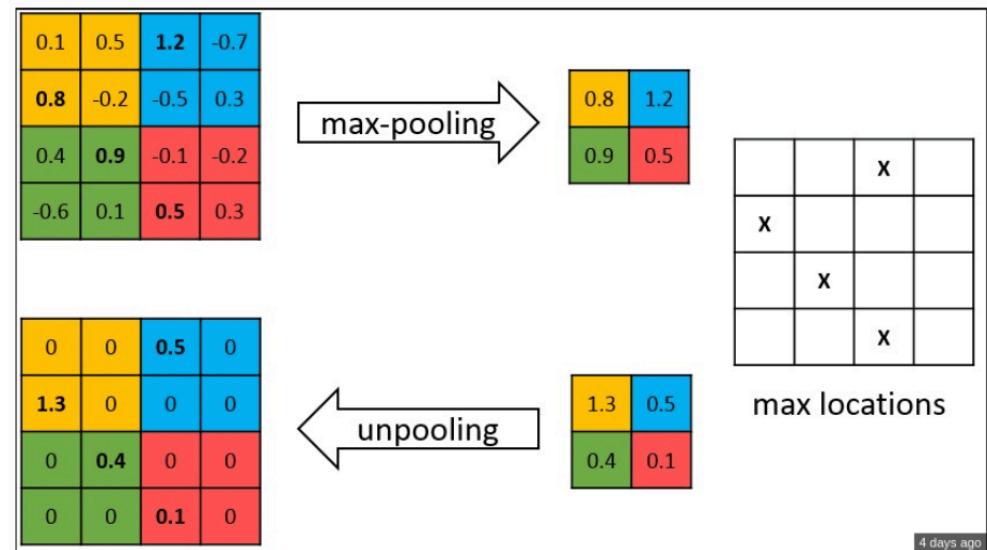


Figure : The FCN-32s Architecture

SegNet

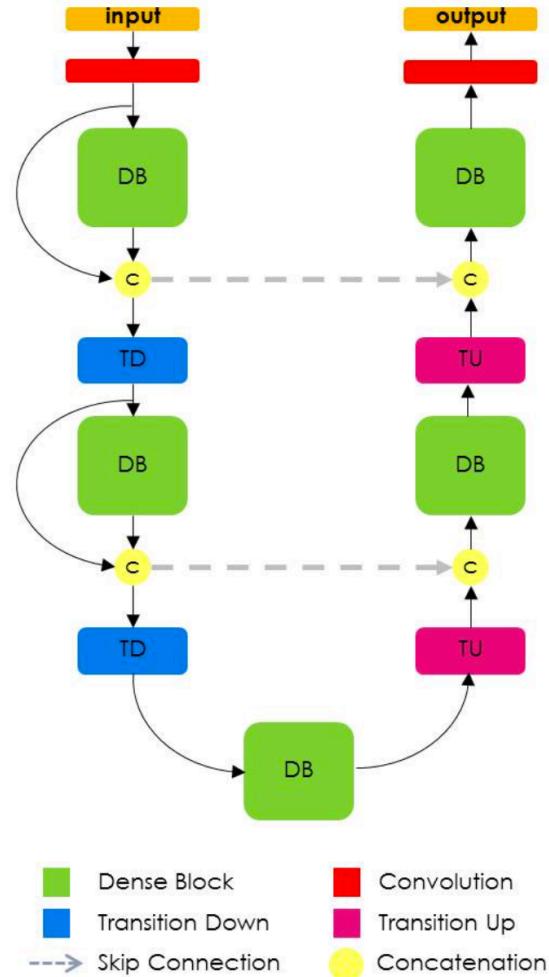
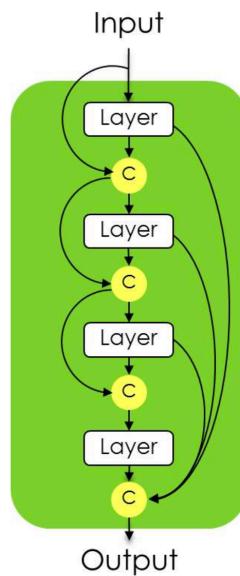


Для повышения
размерности используется
max-unpooling

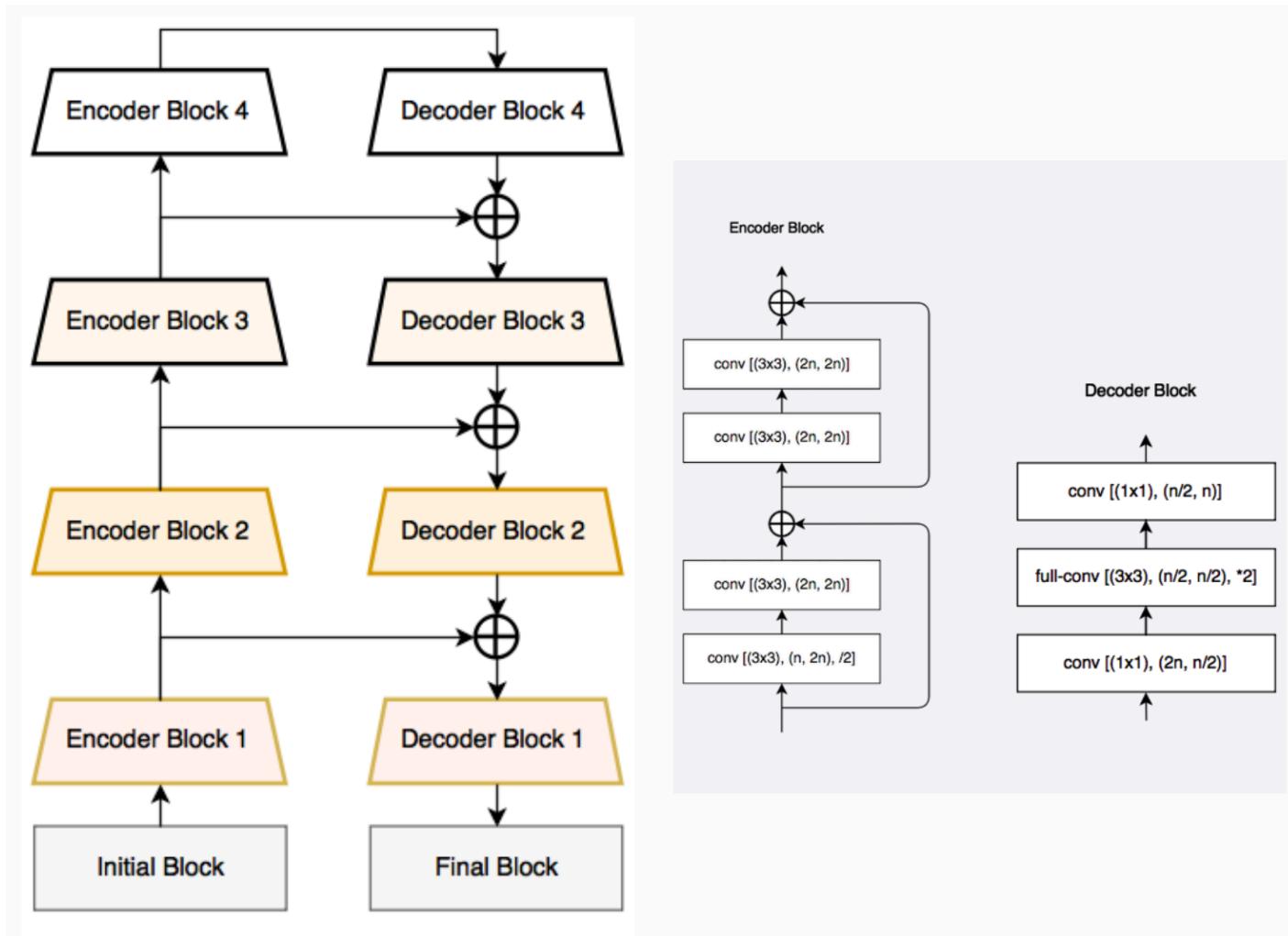


Fully Convolutional DenseNet

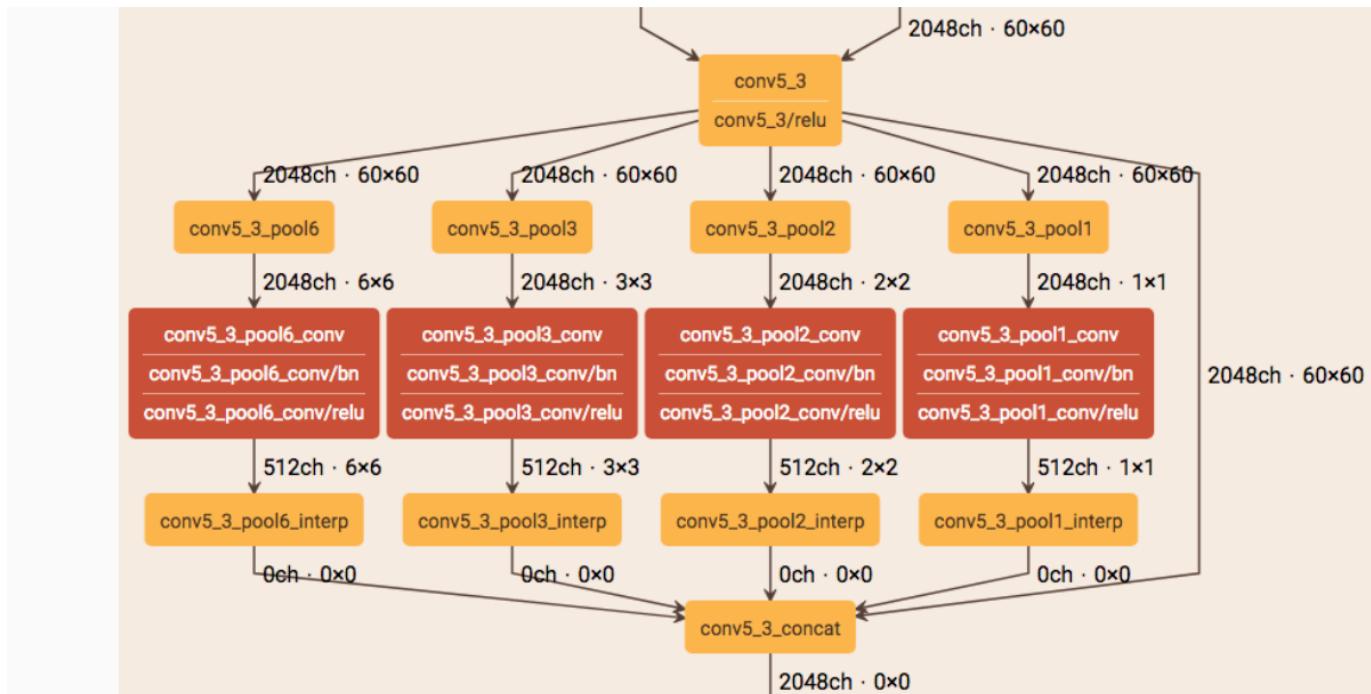
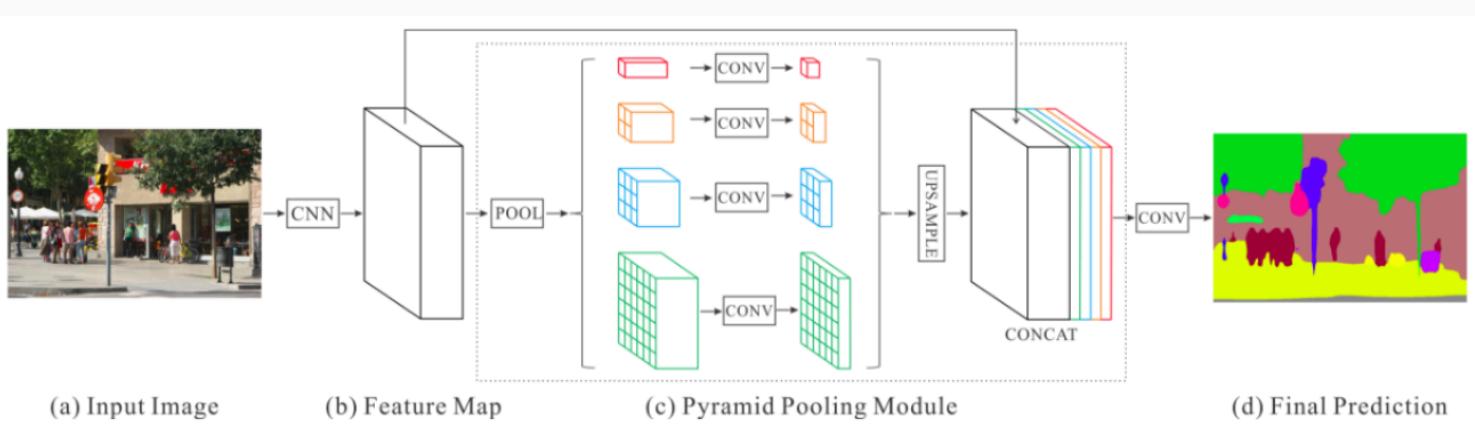
В качестве основы для декодера и энкодера используется dense-блок, как в архитектуре DenseNet



LinkNet (и ENet)



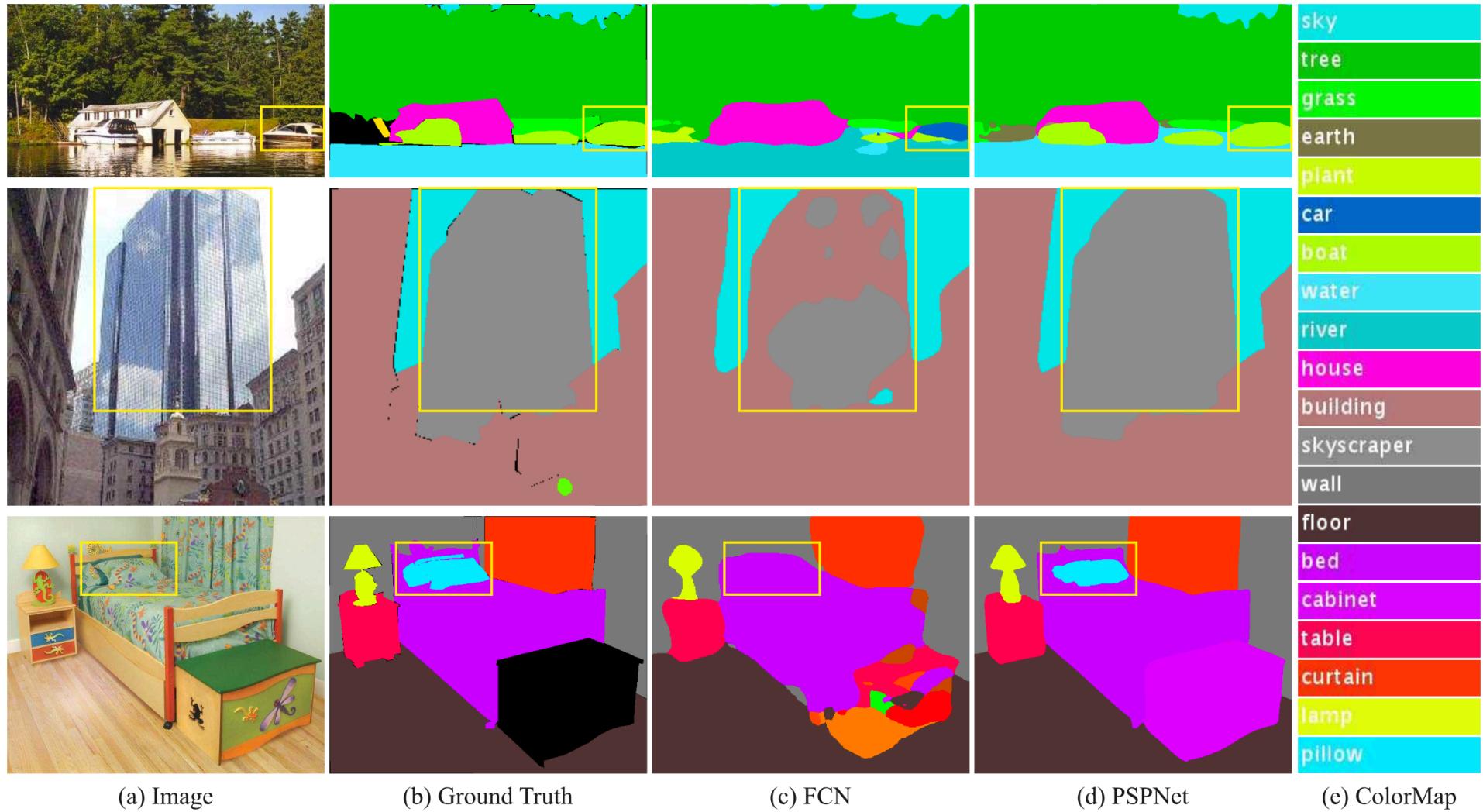
PSP-Net



Top : The PSPNet Architecture

Bottom : The Spatial Pyramid Pooling in visualized in detail using netscope

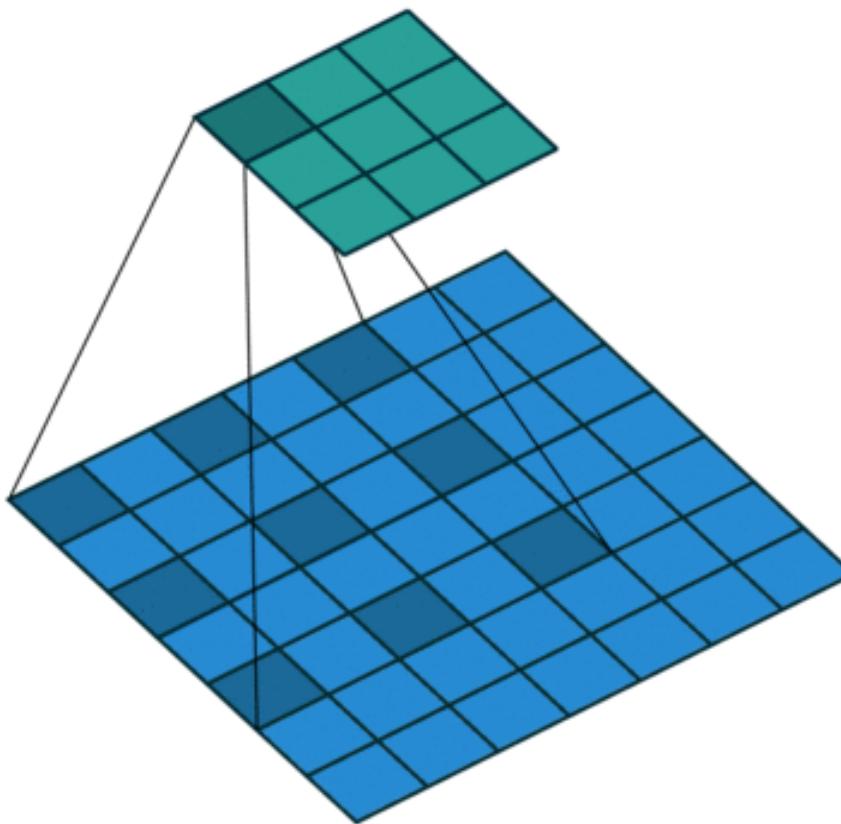
PSP-net. Результаты



ФУНКЦИИ ПОТЕРЬ

- crossentropy
- crossentropy – a * soft dice
- crossentropy – a * log(soft dice)
- фантазии на тему ($y_{true} - y_{pred}$)

Dilated (Atrous) convolutions



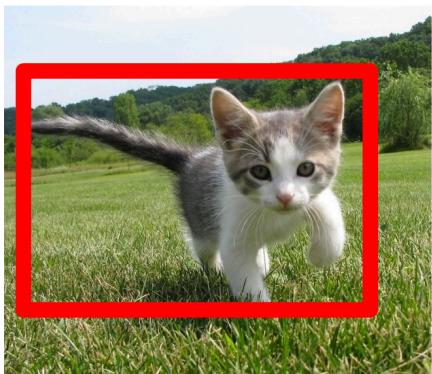
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



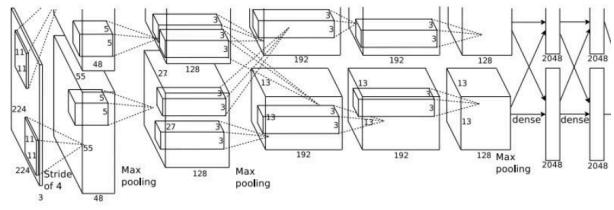
DOG, DOG, CAT

This image is CC0 public domain

Classification + Localization



This image is CC0 public domain



Treat localization as a
regression problem!

Fully
Connected:
4096 to 1000

Vector:
4096
Fully
Connected:
4096 to 4

Class Scores
Cat: 0.9
Dog: 0.05
Car: 0.01
...

Multitask Loss

Box
Coordinates
 (x, y, w, h)

Correct label:
Cat

Softmax
Loss

+

Correct box:
 (x', y', w', h')

L2 Loss

Human pose estimation



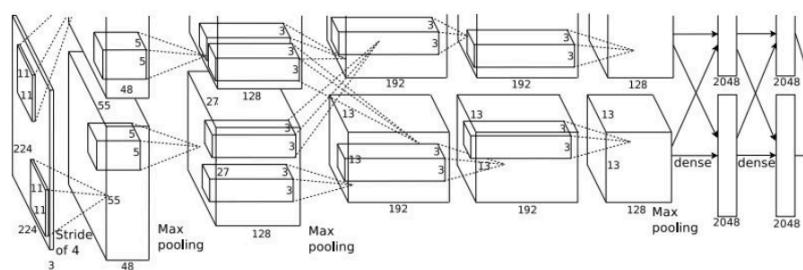
Represent pose as a set of 14 joint positions:

- Left / right foot
- Left / right knee
- Left / right hip
- Left / right shoulder
- Left / right elbow
- Left / right hand
- Neck
- Head top

This image is licensed under CC-BY 2.0.

Johnson and Everingham, "Clustered Pose and Nonlinear Appearance Models for Human Pose Estimation", BMVC 2010

Human pose estimation



→ **Left foot: (x, y)**
→ **Right foot: (x, y)**
...
Vector:
4096 → **Head top: (x, y)**

Toshev and Szegedy, "DeepPose: Human Pose Estimation via Deep Neural Networks", CVPR 2014

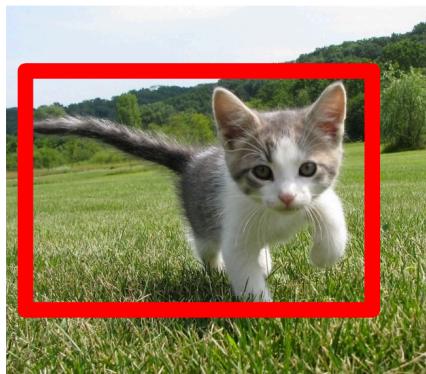
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

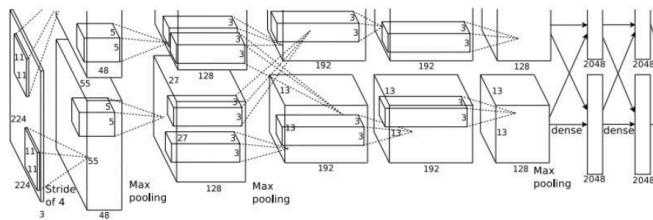
Instance Segmentation



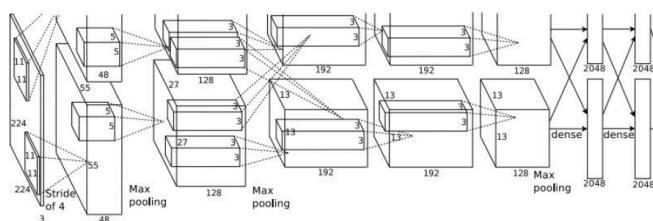
DOG, DOG, CAT

This image is CC0 public domain

Почему не подойдет подход с регрессией?



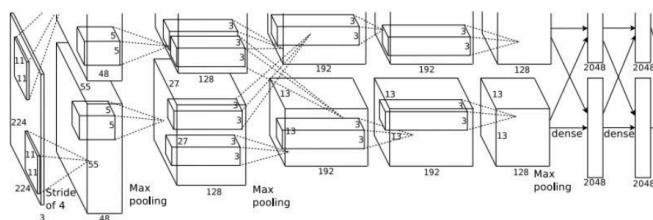
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



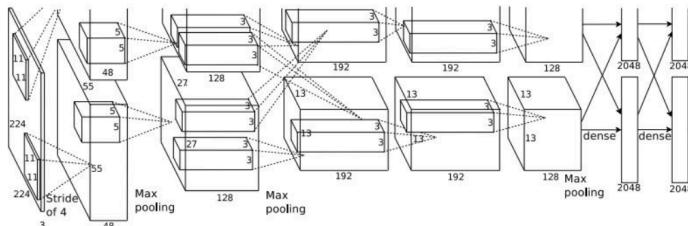
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

....

Object Detection as Classification: Sliding Window

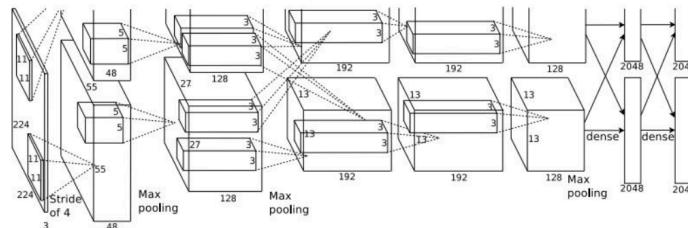
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

Object Detection as Classification: Sliding Window

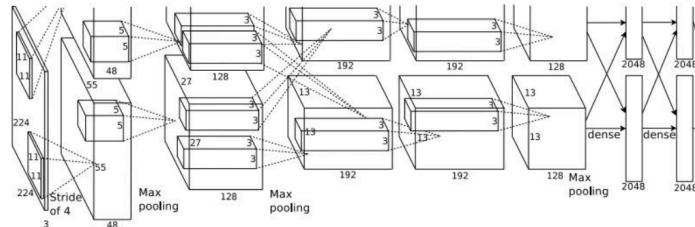
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection as Classification: Sliding Window

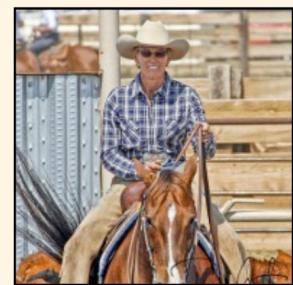
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



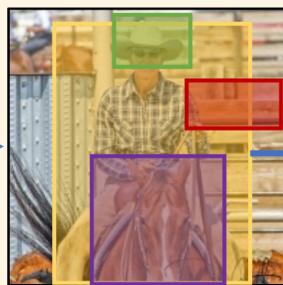
Dog? YES
Cat? NO
Background? NO

“Slow” R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...

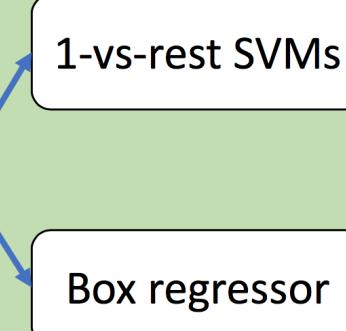


Crop &
warp

Per-region computation for each $r_i \in r(I)$



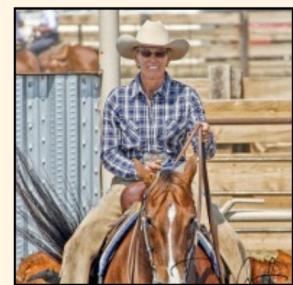
ConvNet(r_i)



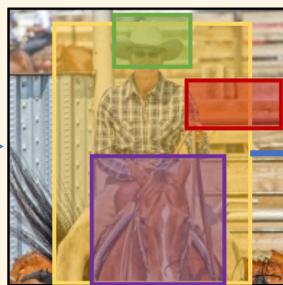
Very heavy *per-region* computation
E.g., 2000 full network evaluations

“Slow” R-CNN

Per-image computation



Selective search,
Edge Boxes,
MCG, ...



Crop &
warp

Per-region computation for each $r_i \in r(I)$



ConvNet(r_i)

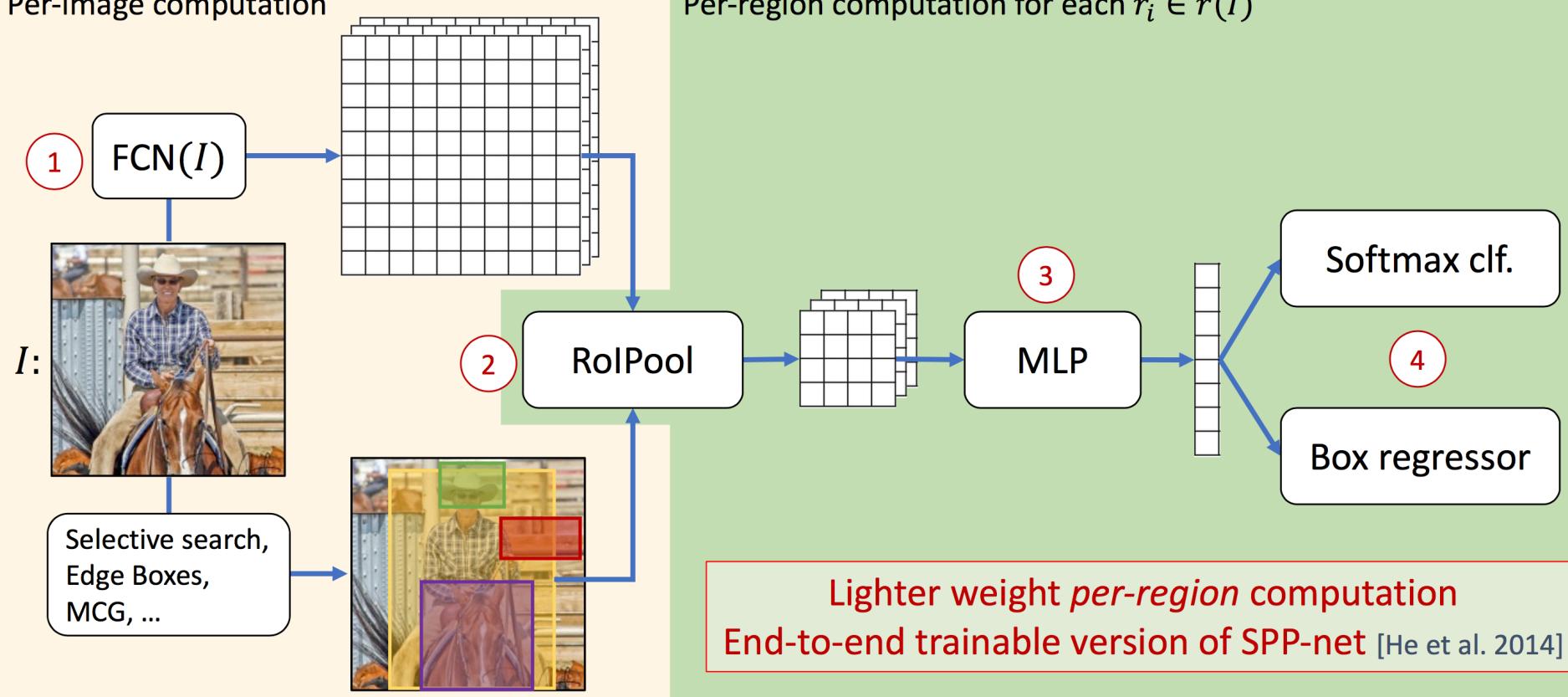
1-vs-rest SVMs

Box regressor

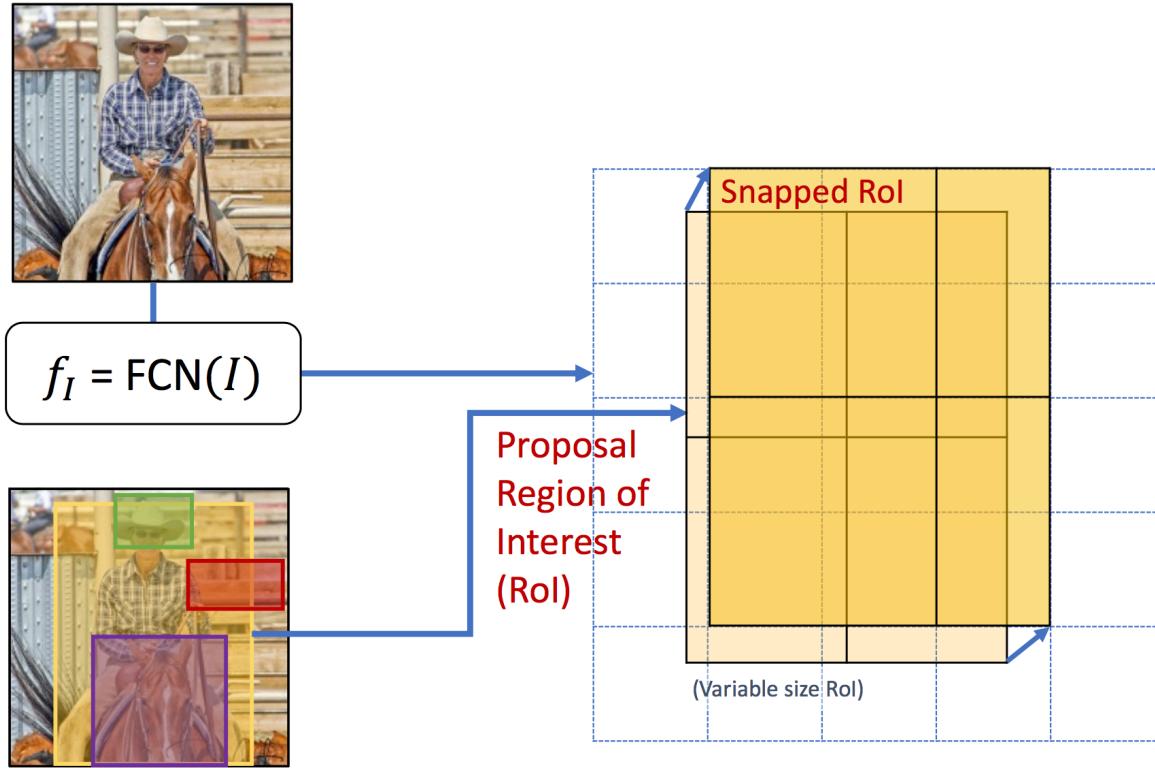
Very heavy *per-region* computation
E.g., 2000 full network evaluations

Generalized R-CNN → Fast R-CNN

Per-image computation



RoIPool (on each Proposal)



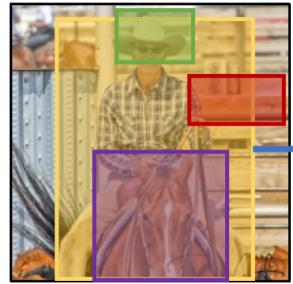
Key innovation in SPP-net
[He et al. 2014]

RoIPool (on each Proposal)

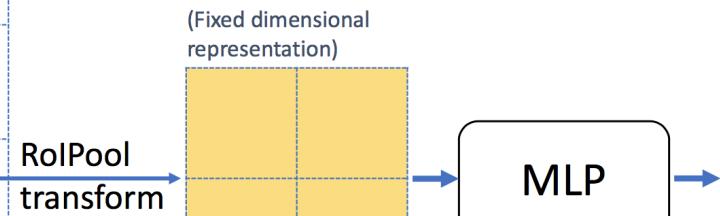
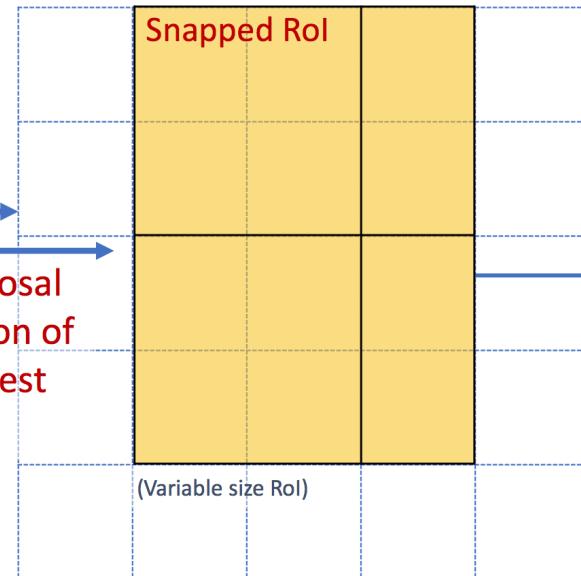
Transform arbitrary size proposal into a fixed-dimensional representation (e.g., 2x2)



$$f_I = \text{FCN}(I)$$



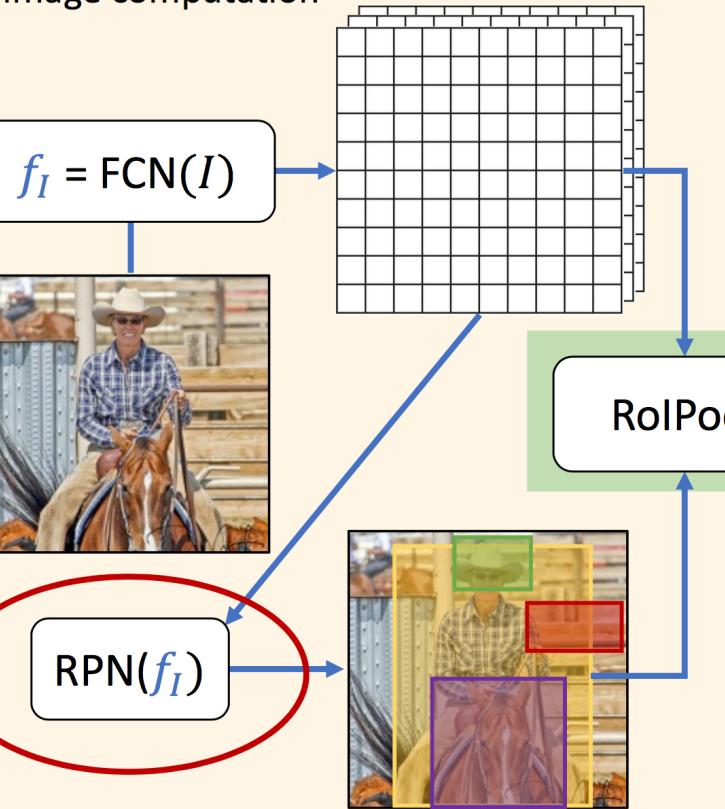
Proposal
Region of
Interest
(RoI)



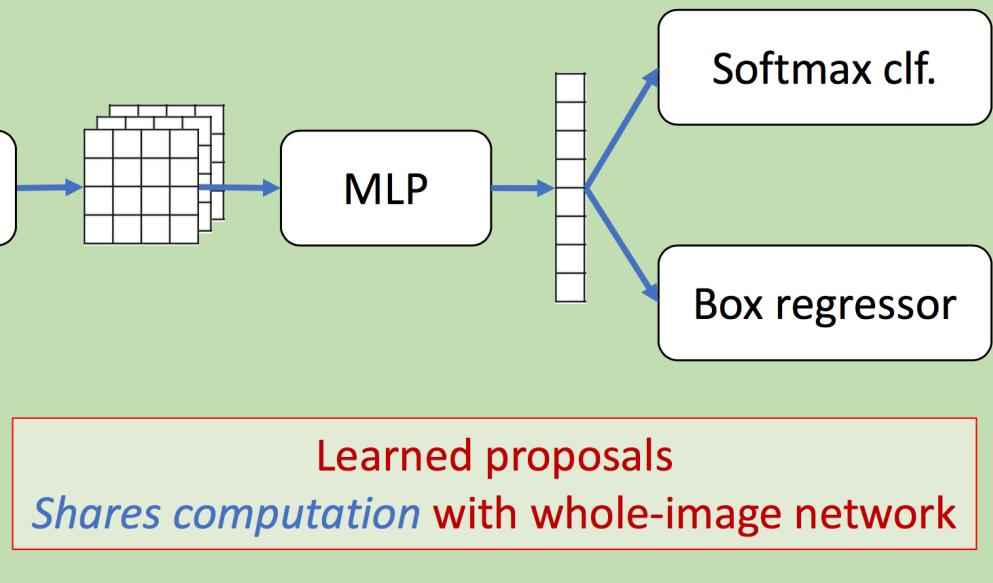
Feature value
is **max** over input
cells

Faster R-CNN

Per-image computation



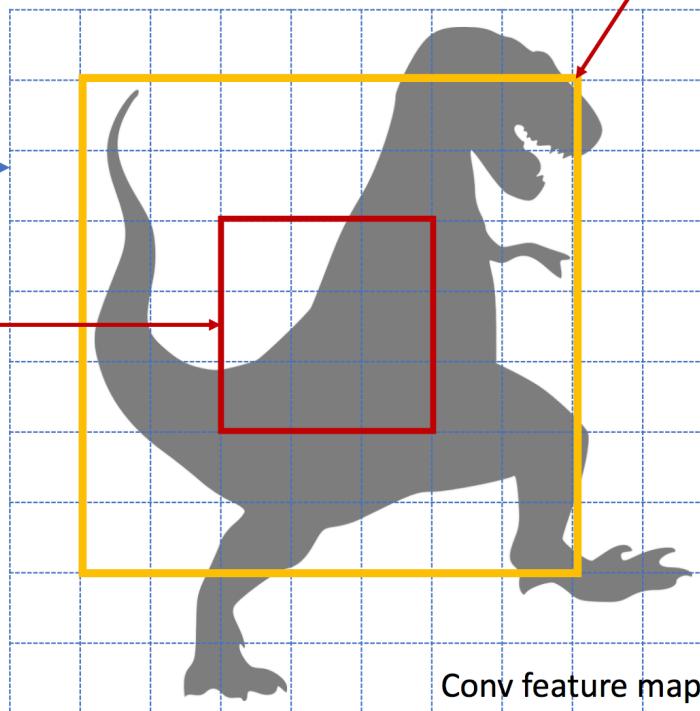
Per-region computation for each $r_i \in r(I)$



RPN: Anchor Box



$$f_I = \text{FCN}(I)$$



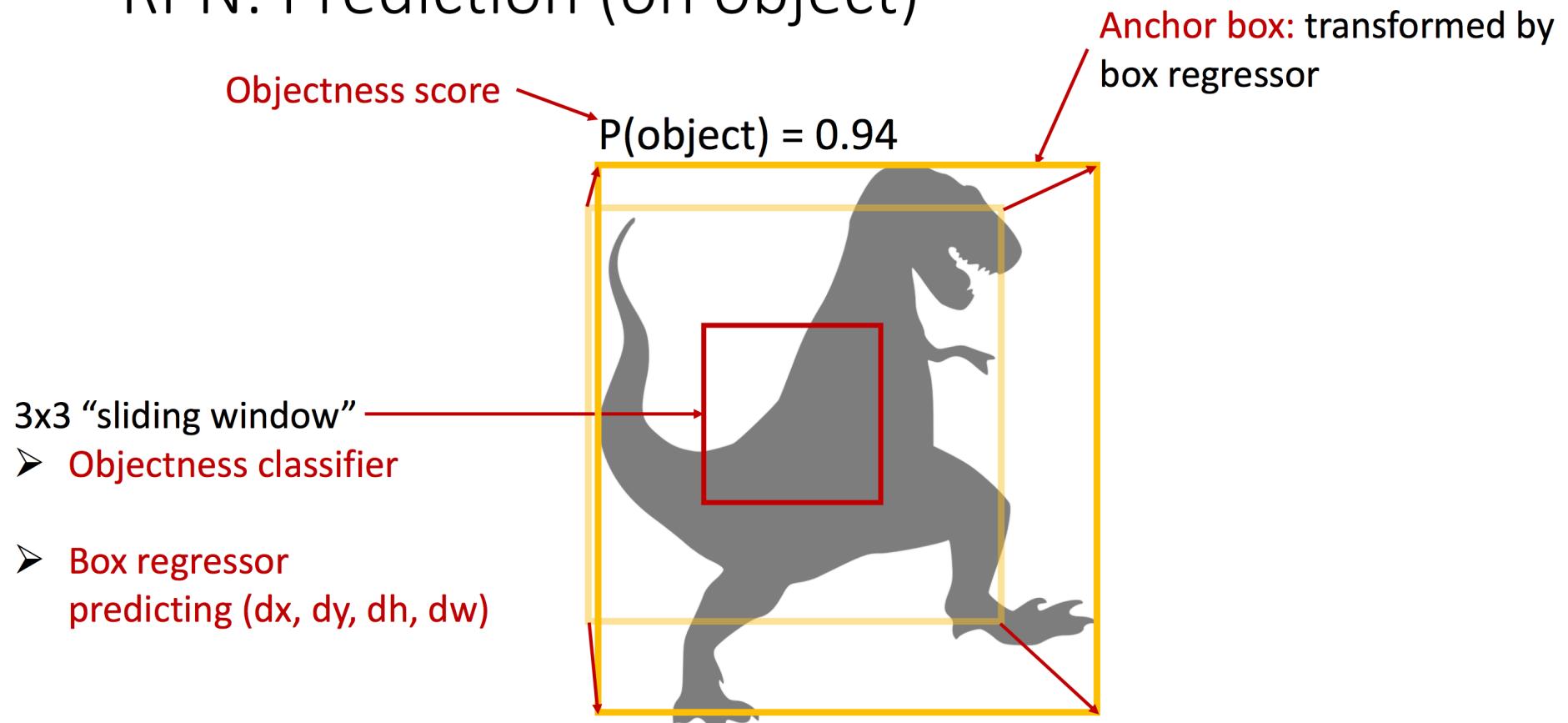
Anchor box: predictions are
w.r.t. this box, *not the 3x3
sliding window*

3x3 “sliding window”
➤ Objectness classifier

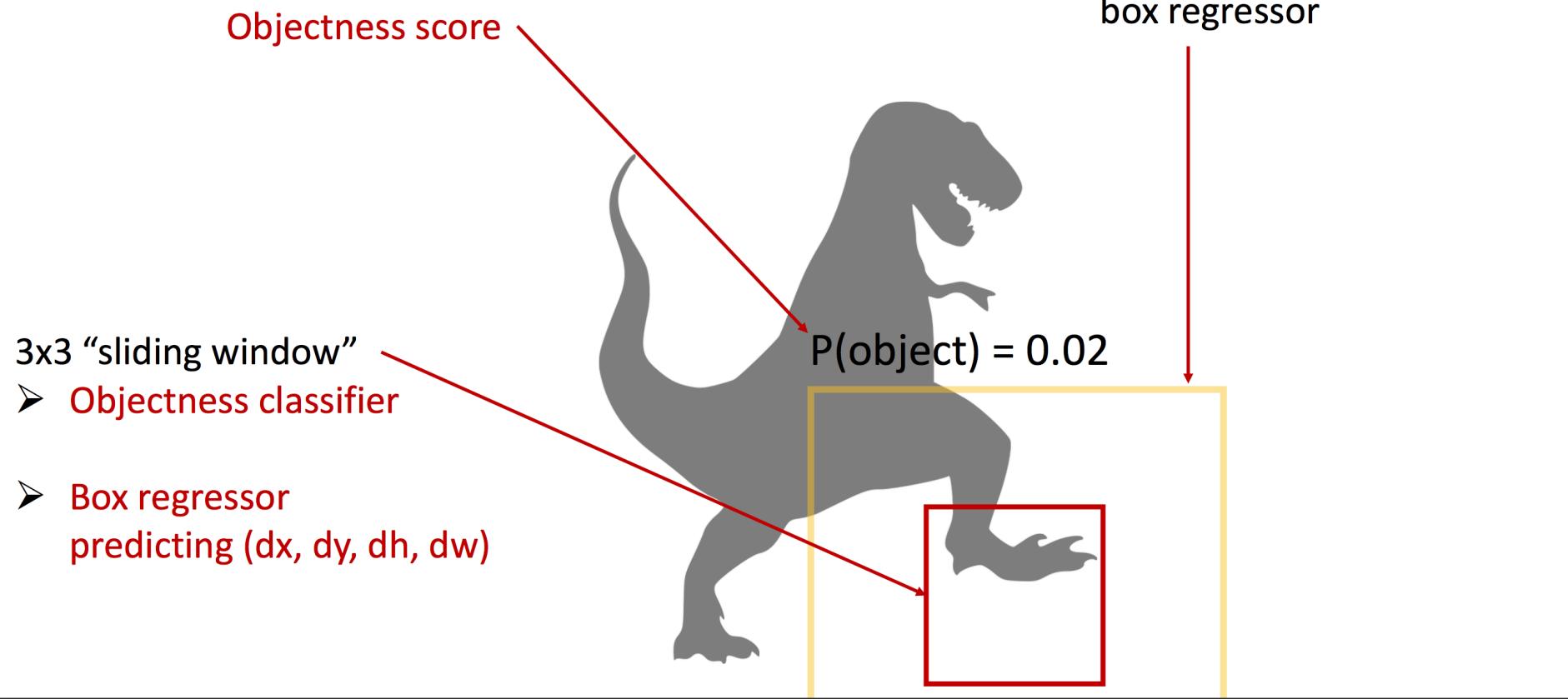
➤ Box regressor
predicting (dx, dy, dh, dw)

Conv feature map

RPN: Prediction (on object)



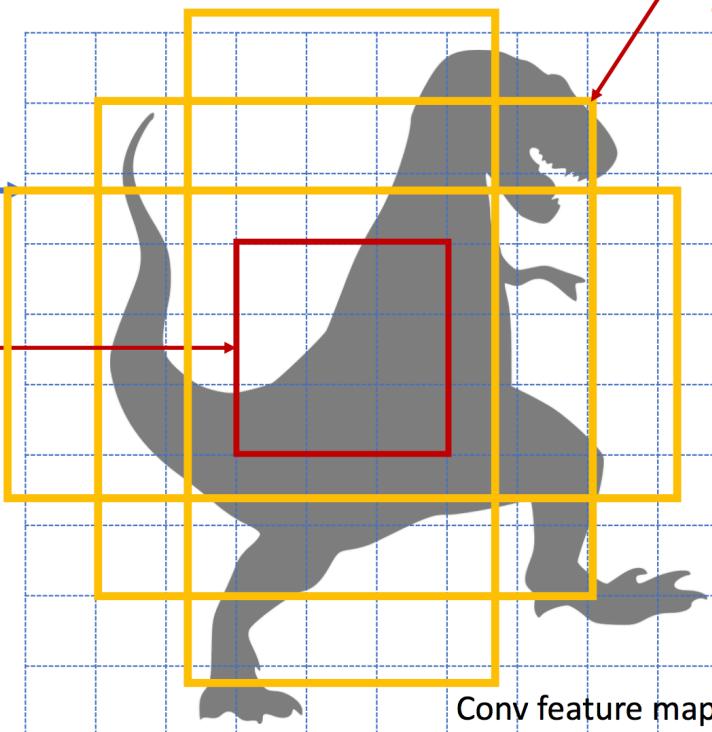
RPN: Prediction (off object)



RPN: Multiple Anchors



$$f_I = \text{FCN}(I)$$

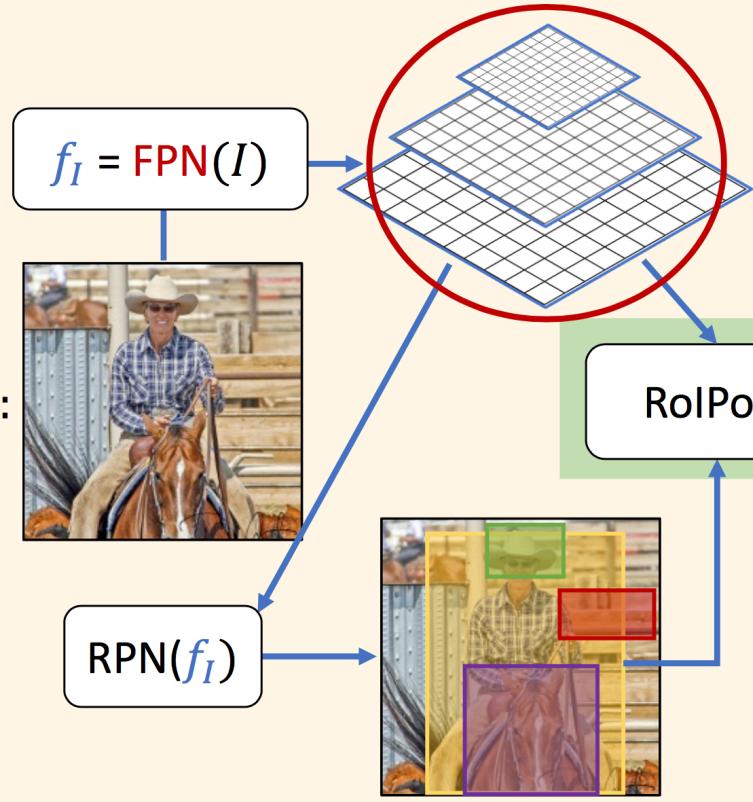


Anchor boxes: K anchors per location with different scales and aspect ratios

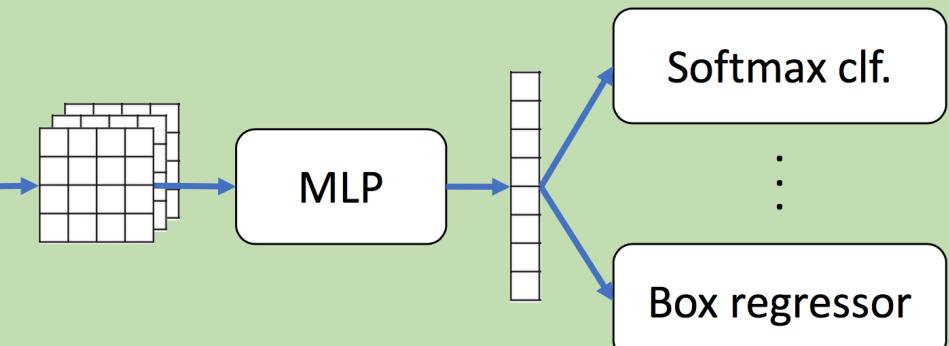
- 3x3 “sliding window”
 - **K objectness classifiers**
 - **K box regressors**

Faster R-CNN with a Feature Pyramid Network

Per-image computation



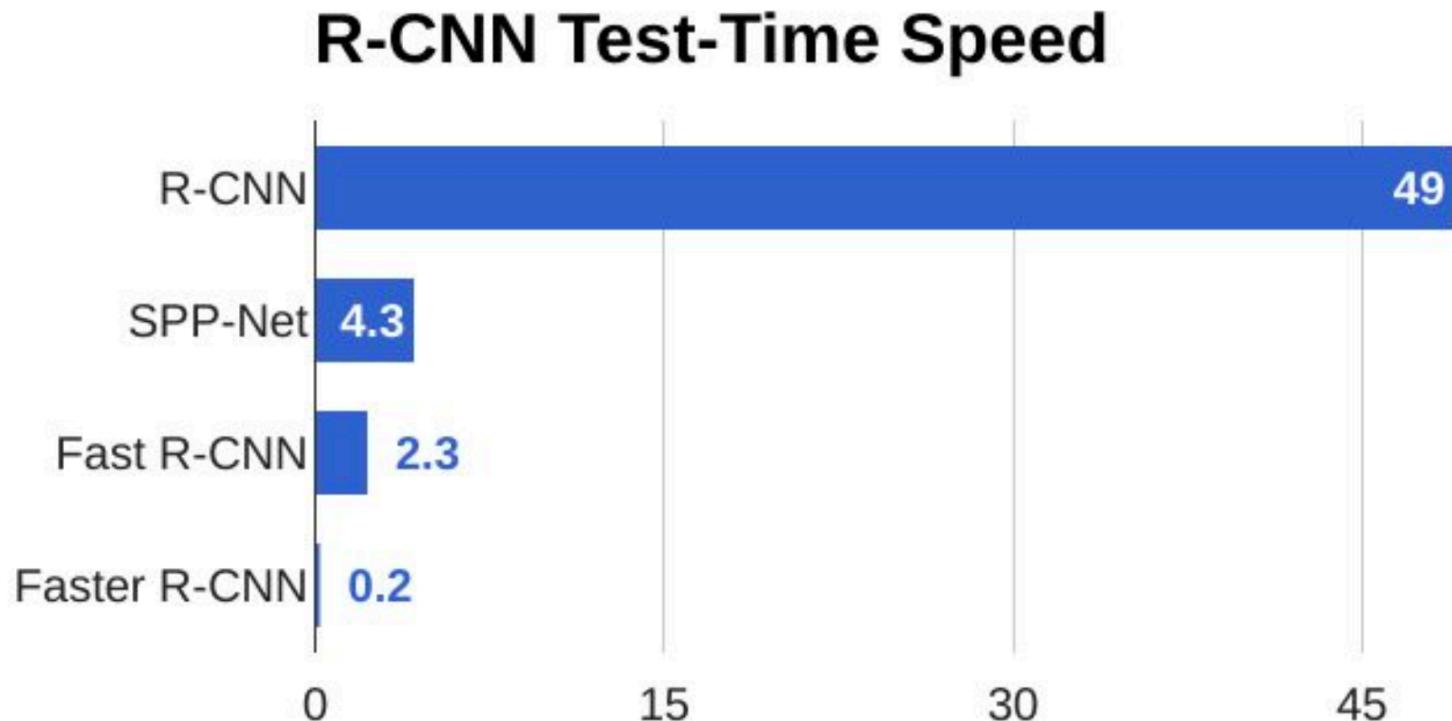
Per-region computation for each $r_i \in r(I)$



The whole-image feature representation can be improved by making it *multi-scale*

Faster R-CNN:

Make CNN do proposals!

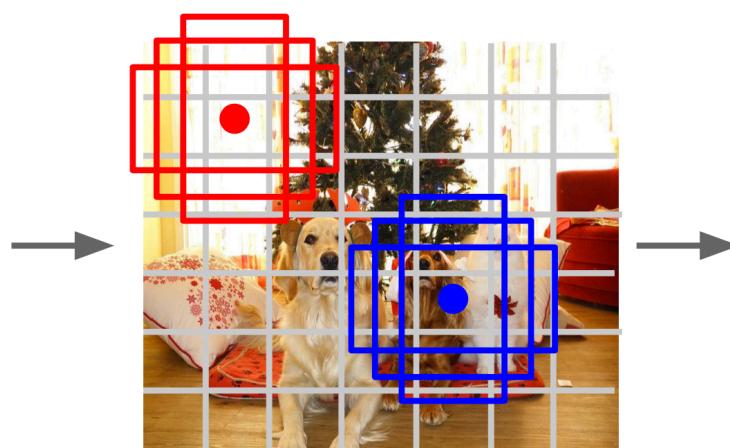


Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network!



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 $B = 3$

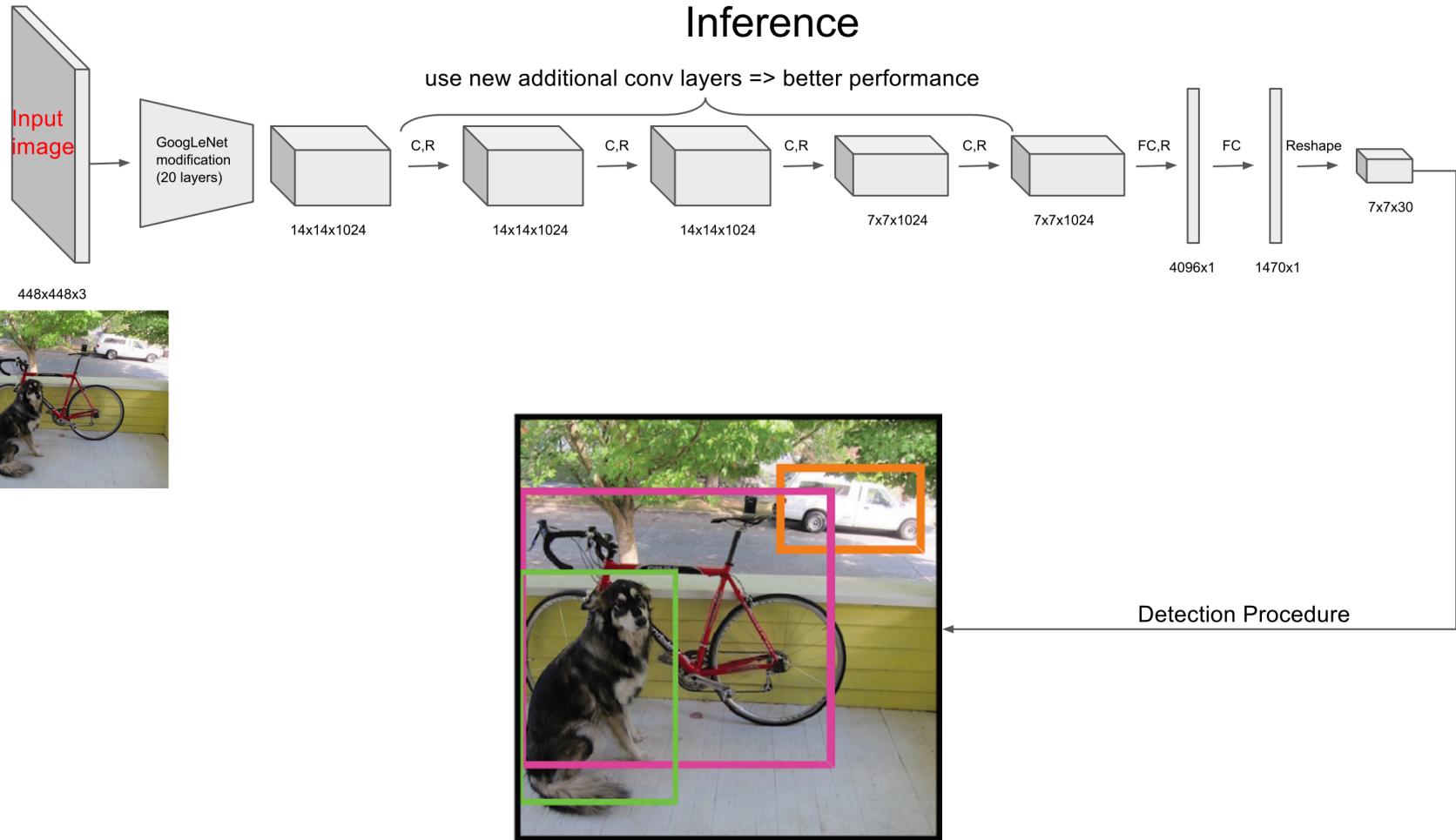
Within each grid cell:

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

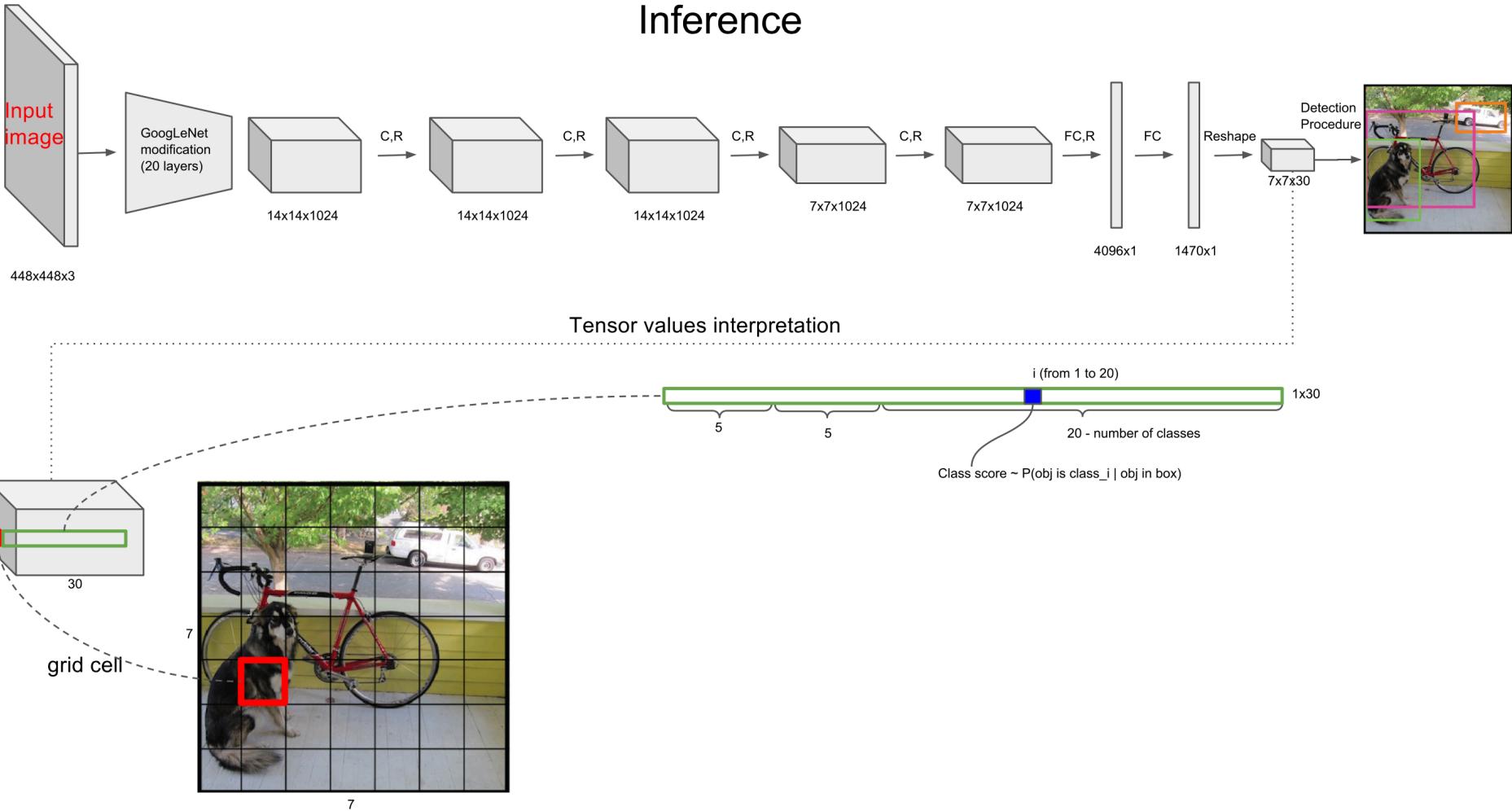
Output:
 $7 \times 7 \times (5 * B + C)$

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

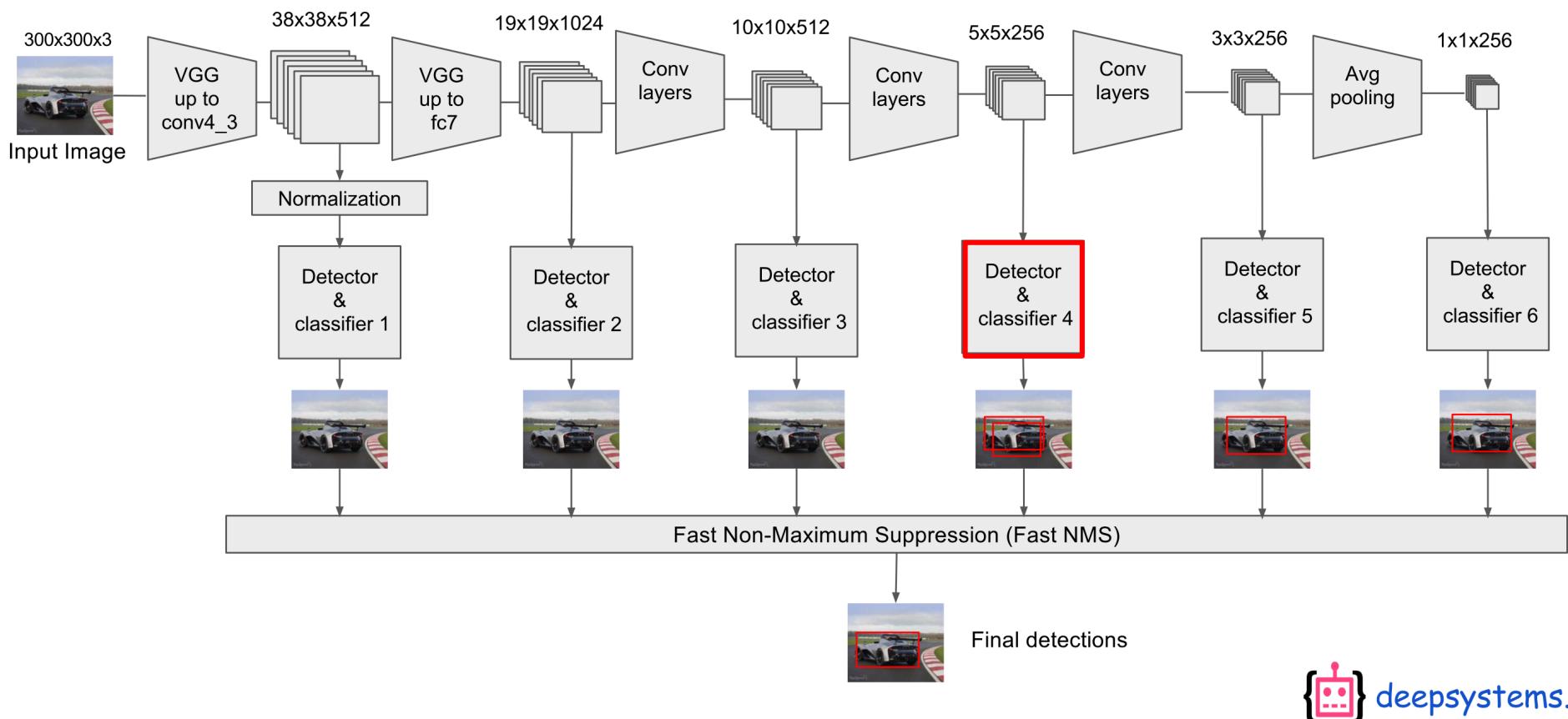
Yolo



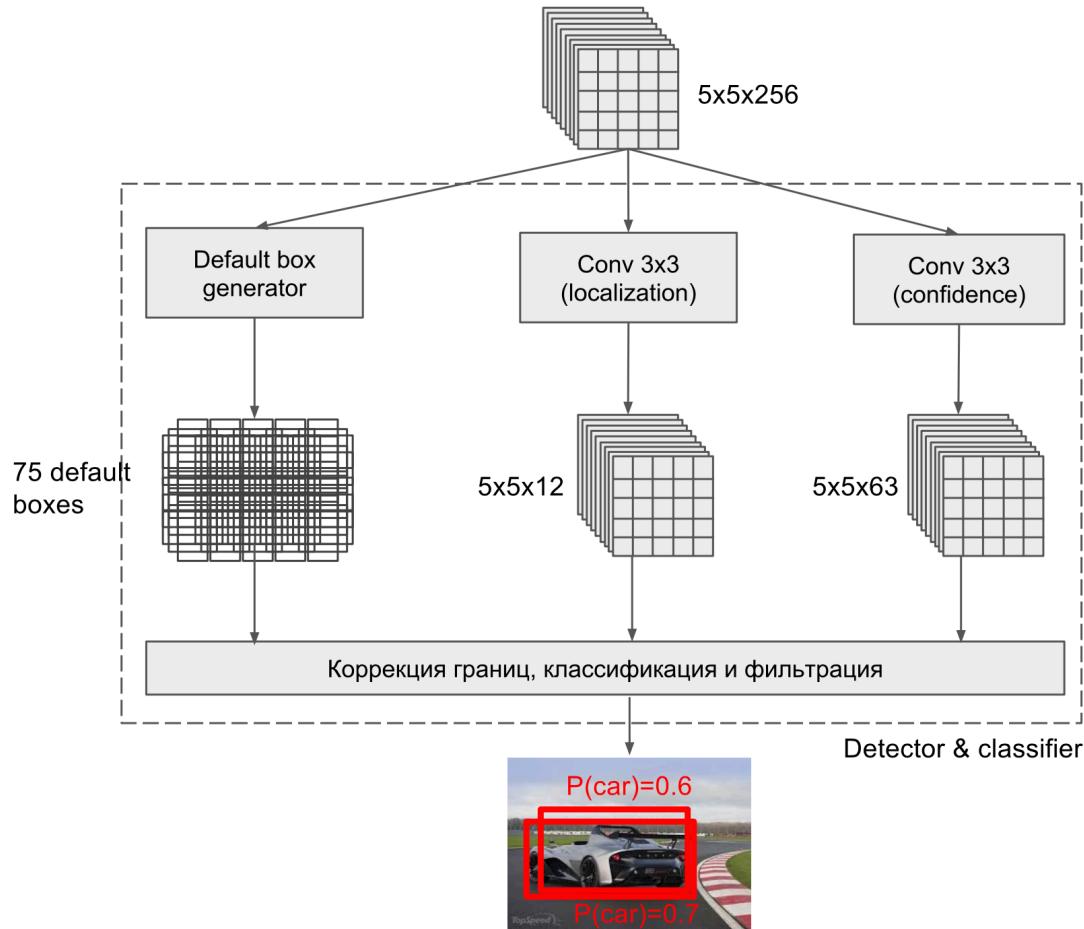
Yolo



Архитектура SSD 300



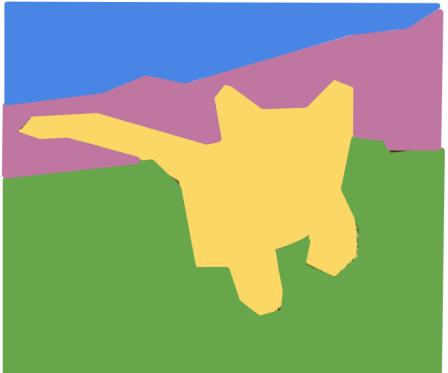
Архитектура SSD 300. Detector & classifier



Пусть заданы следующие параметры:

- Размер исходного изображения (300×300)
- Размерность feature maps ($5 \times 5 \times 256$)
- #default boxes = 3

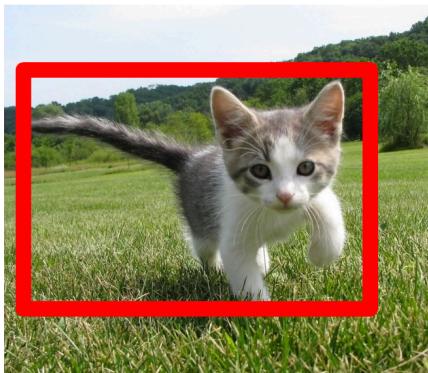
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

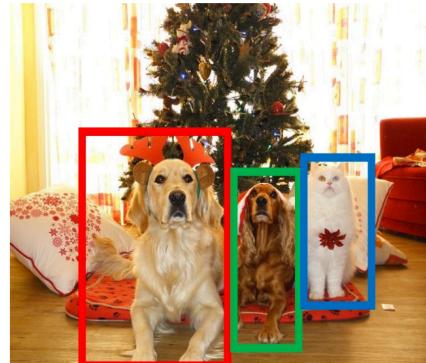
Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

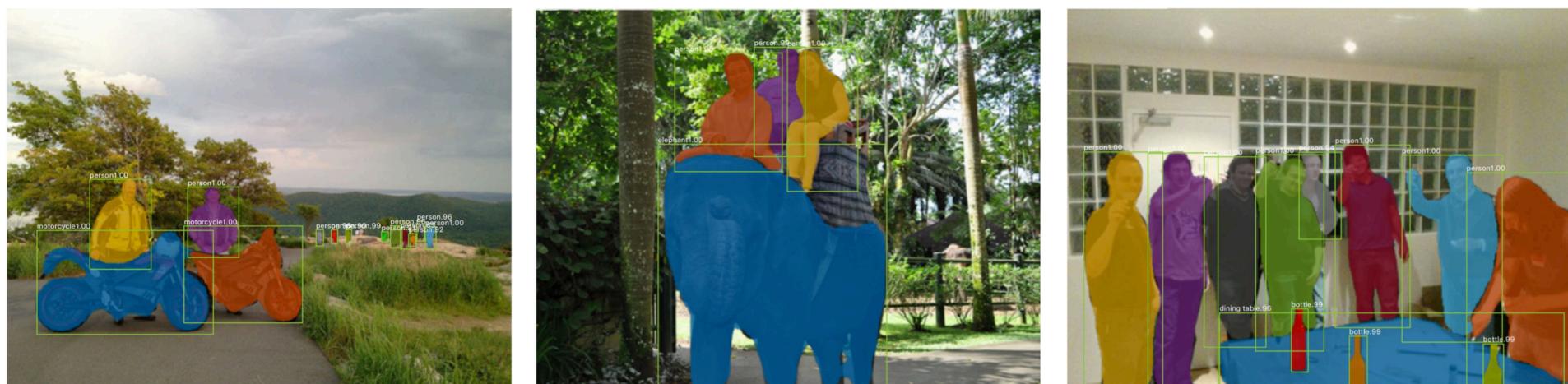
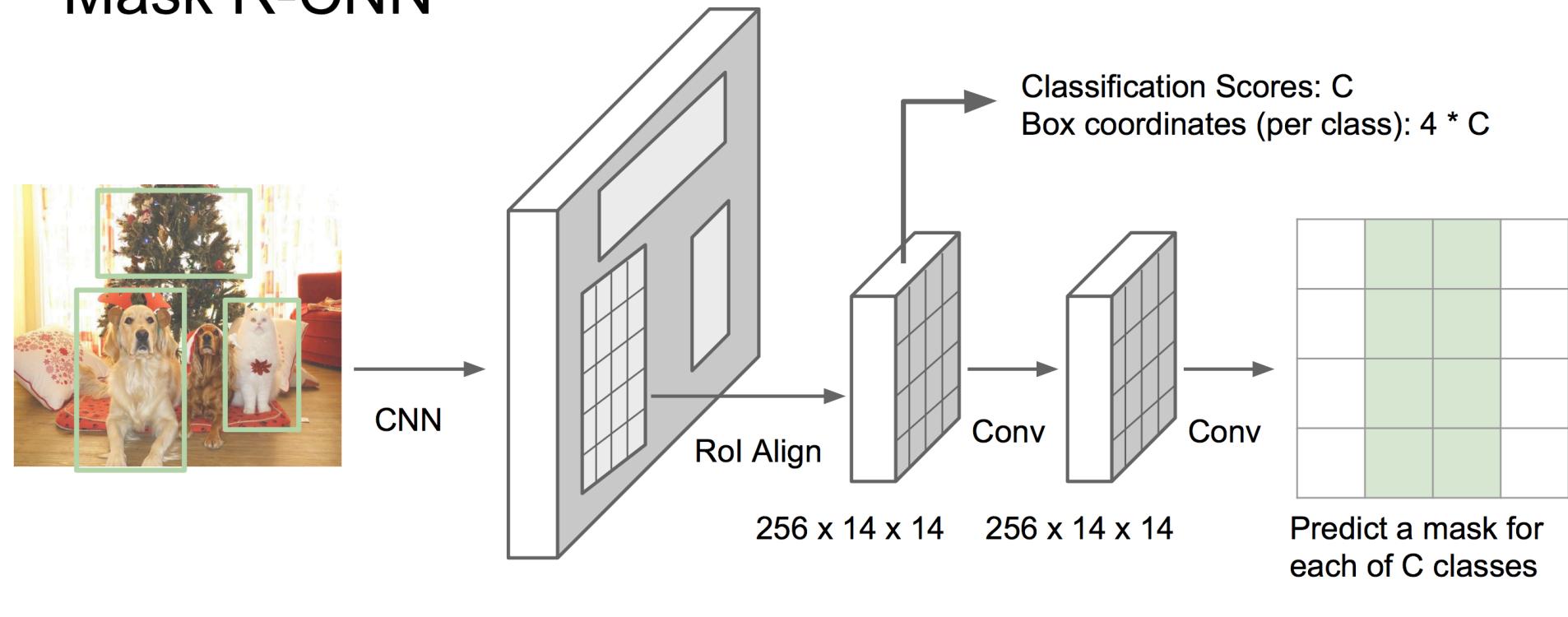
Instance Segmentation



DOG, DOG, CAT

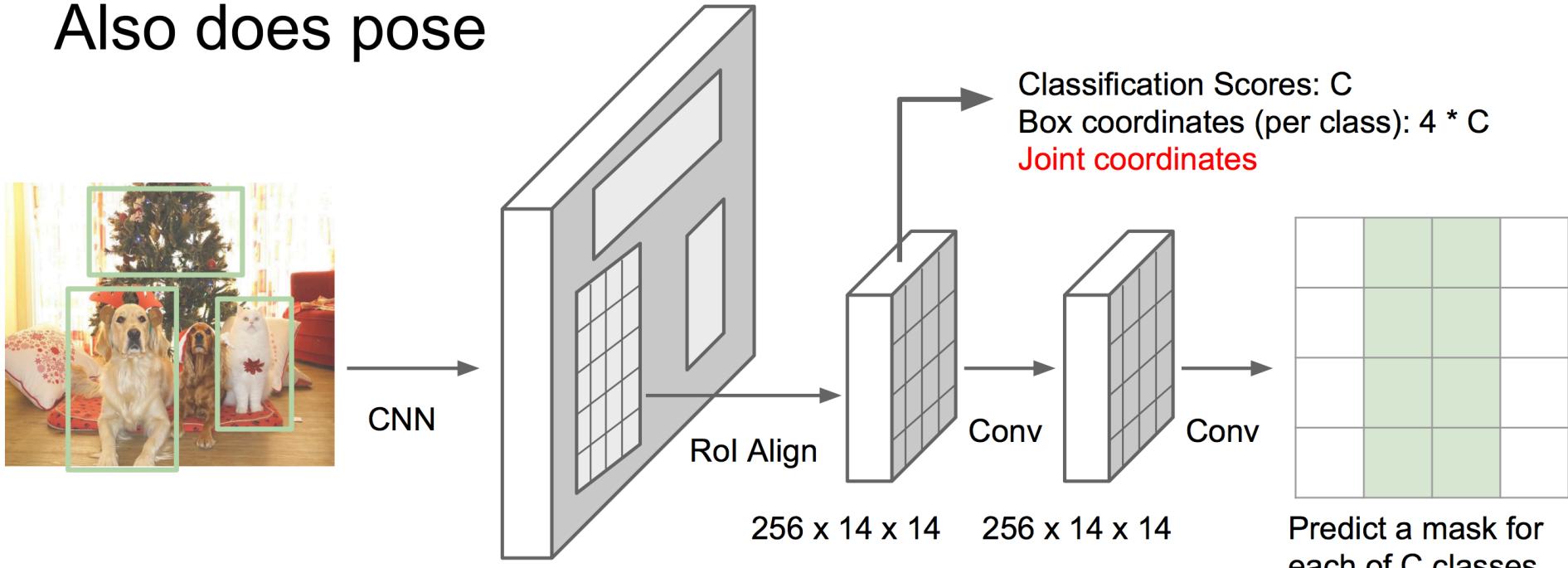
This image is CC0 public domain

Mask R-CNN



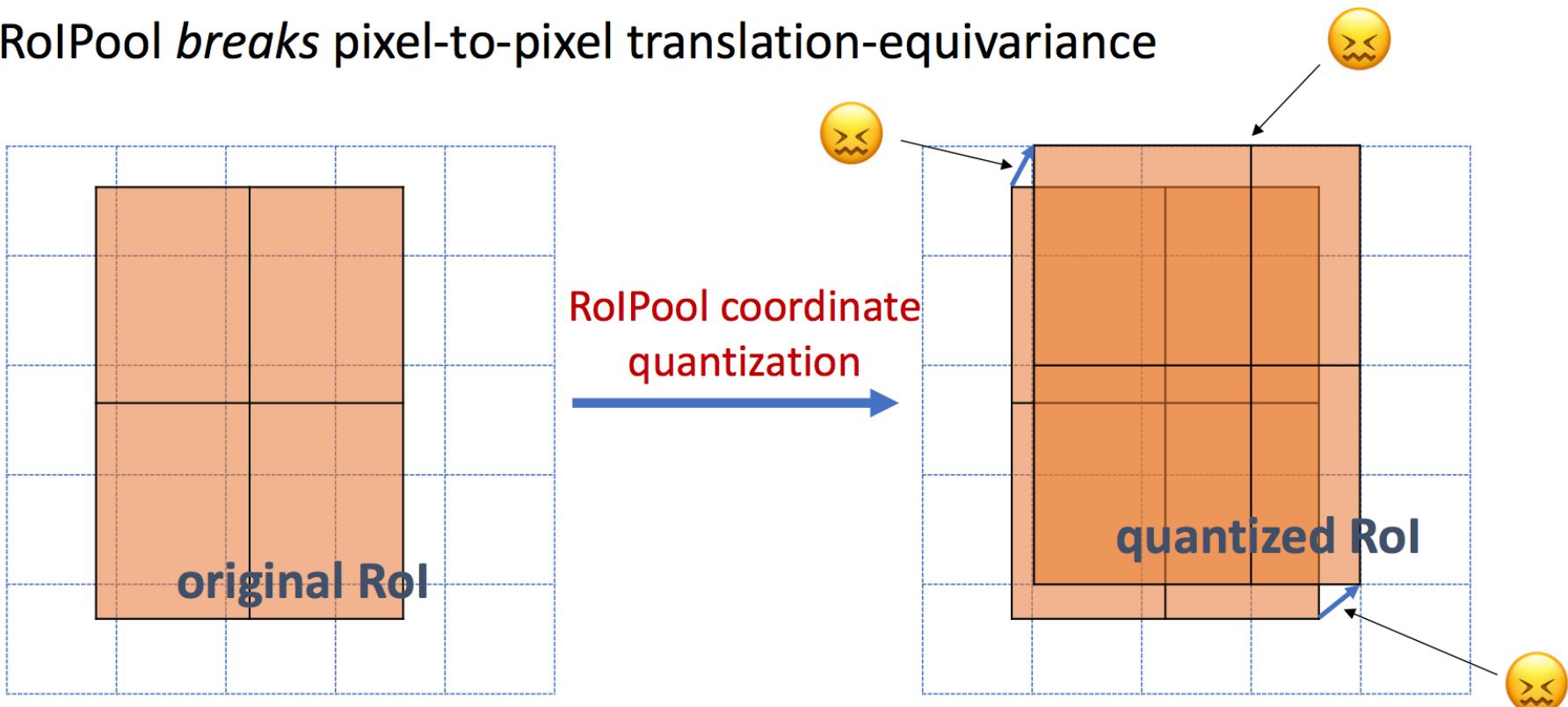
Mask R-CNN

Also does pose



RoIAlign vs. RoIPool

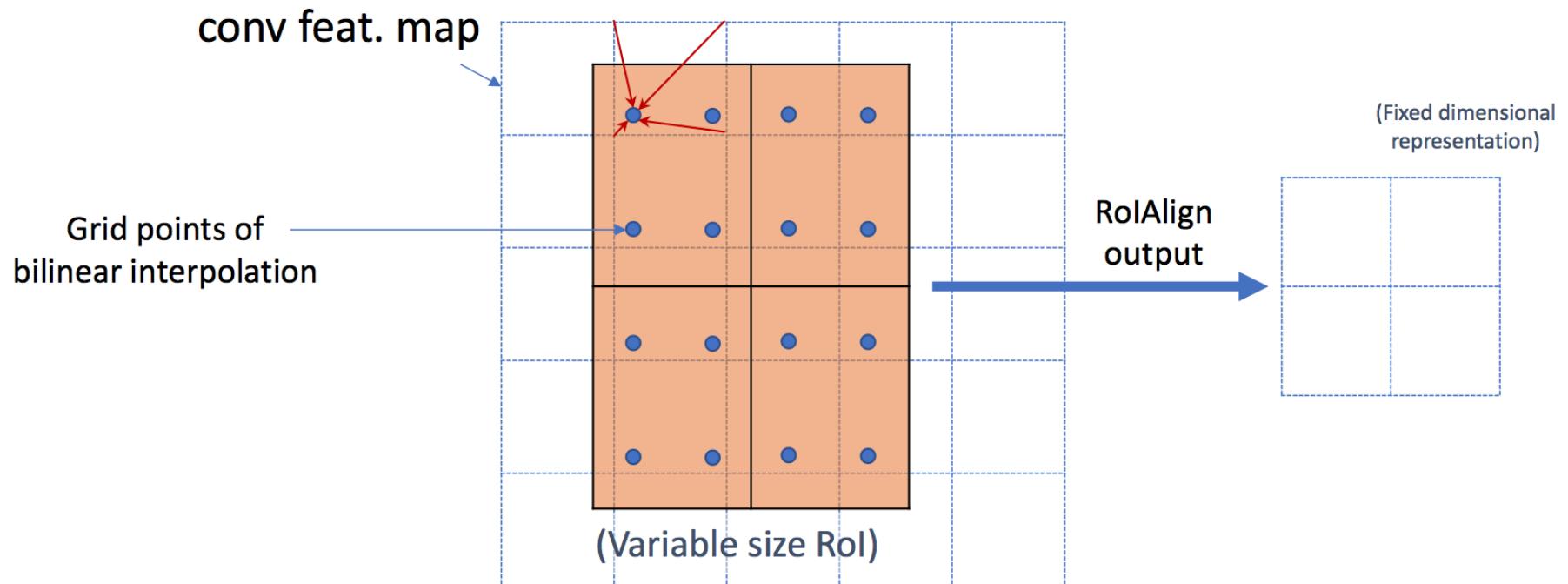
- RoIPool *breaks* pixel-to-pixel translation-equivariance

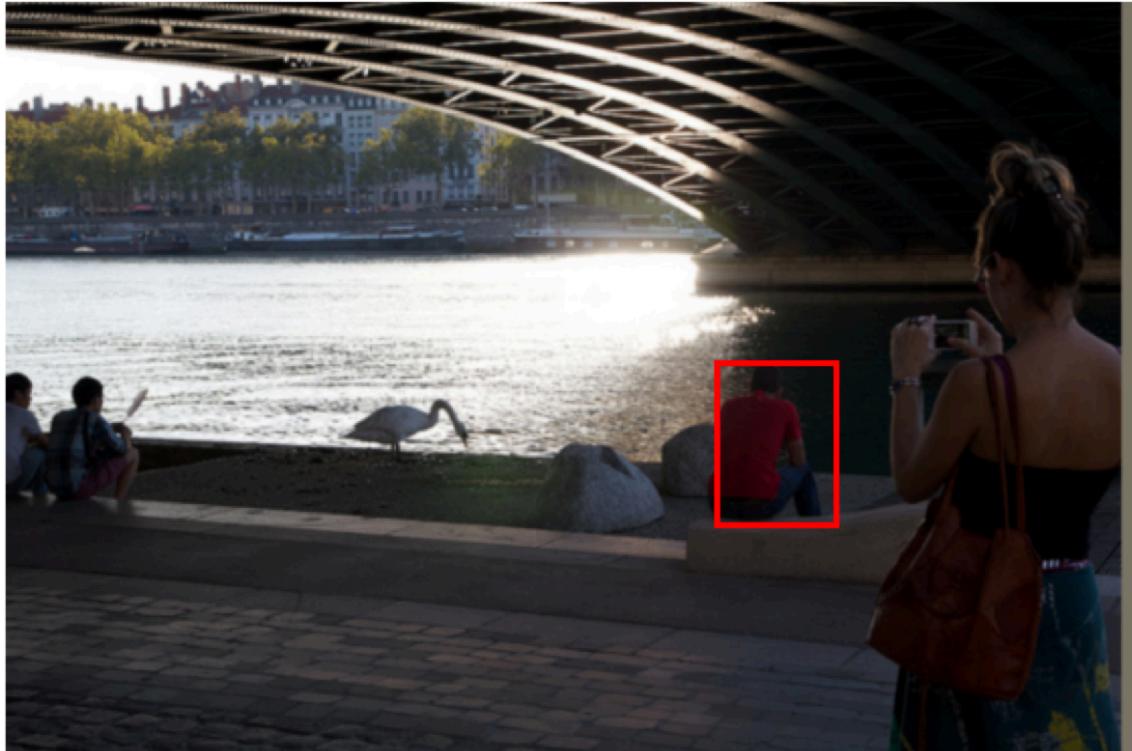


RoIAlign

FAQs: how to sample grid points within a cell?

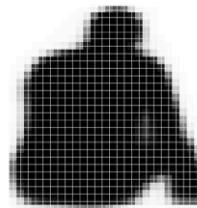
- 4 regular points in 2x2 sub-cells
- other implementation could work





Validation image with box detection shown in red

28x28 soft prediction from Mask R-CNN
(enlarged)



Soft prediction resampled to image coordinates
(bilinear and bicubic interpolation work equally well)

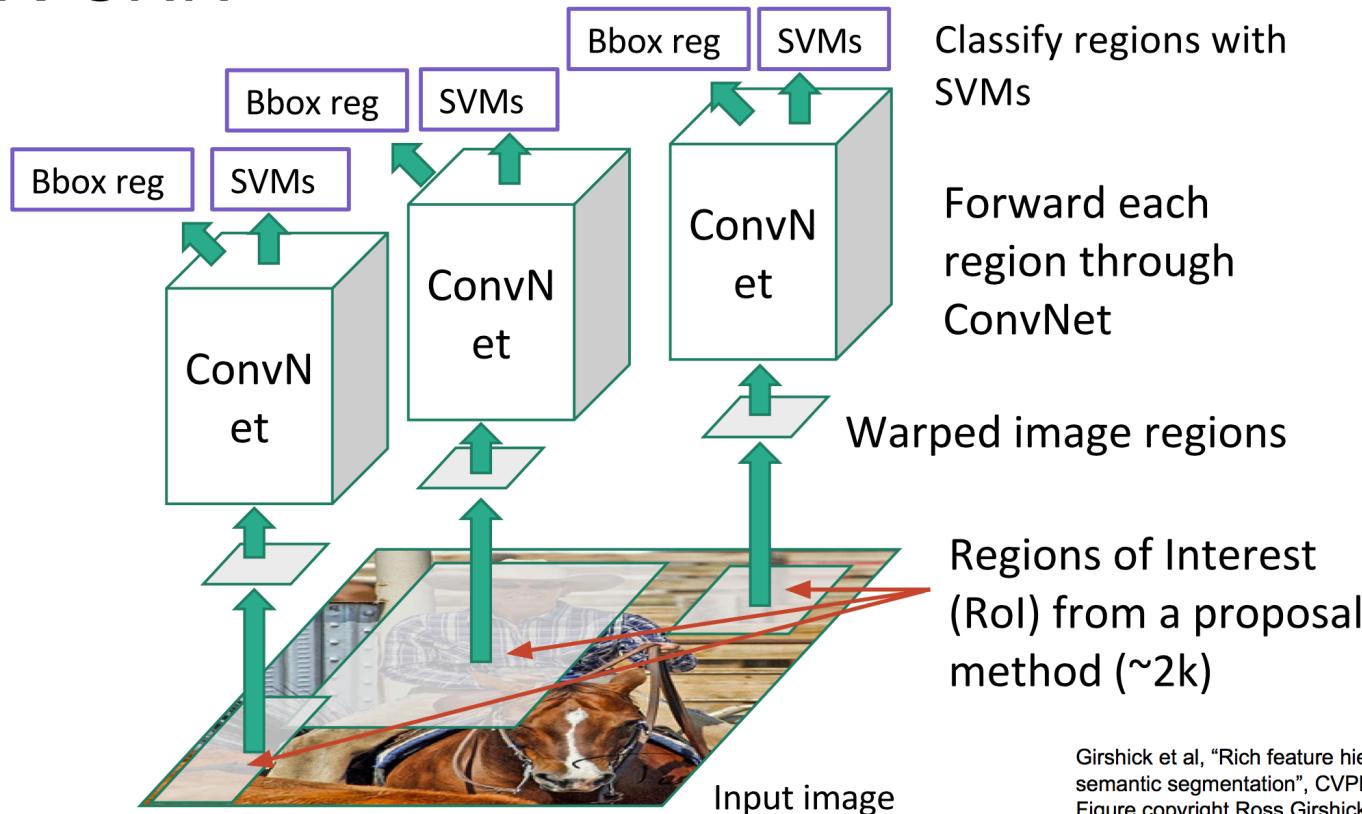


Final prediction (threshold at 0.5)



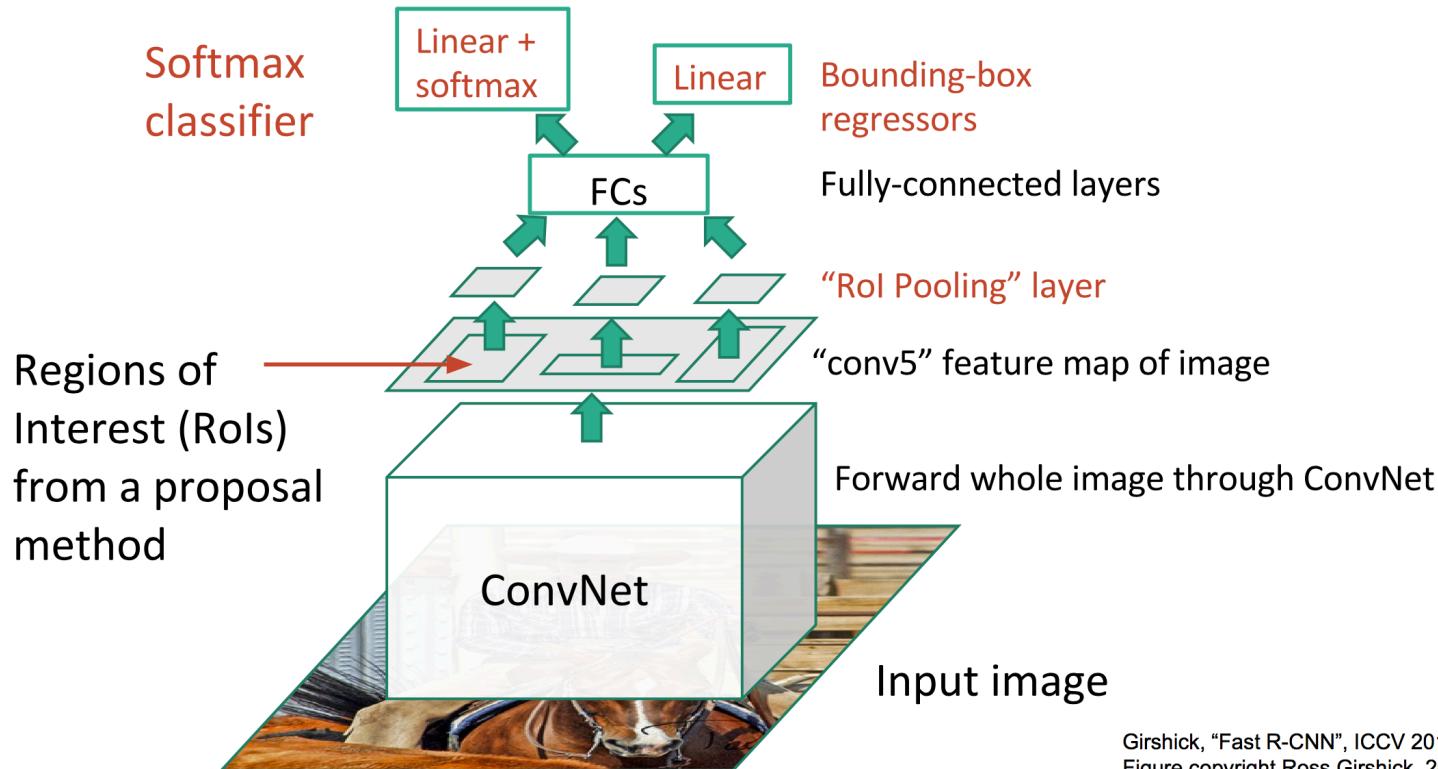
Спасибо за внимание

R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

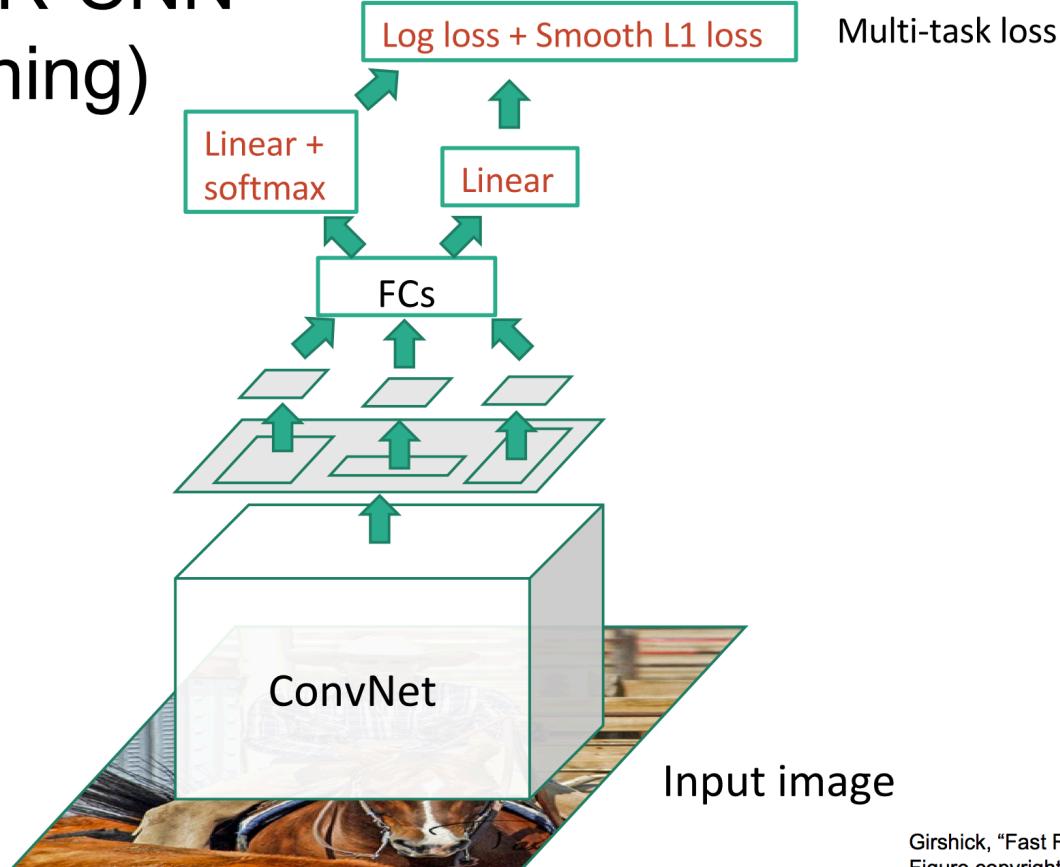
Fast R-CNN



Girshick, "Fast R-CNN", ICCV 2015.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN (Training)



Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

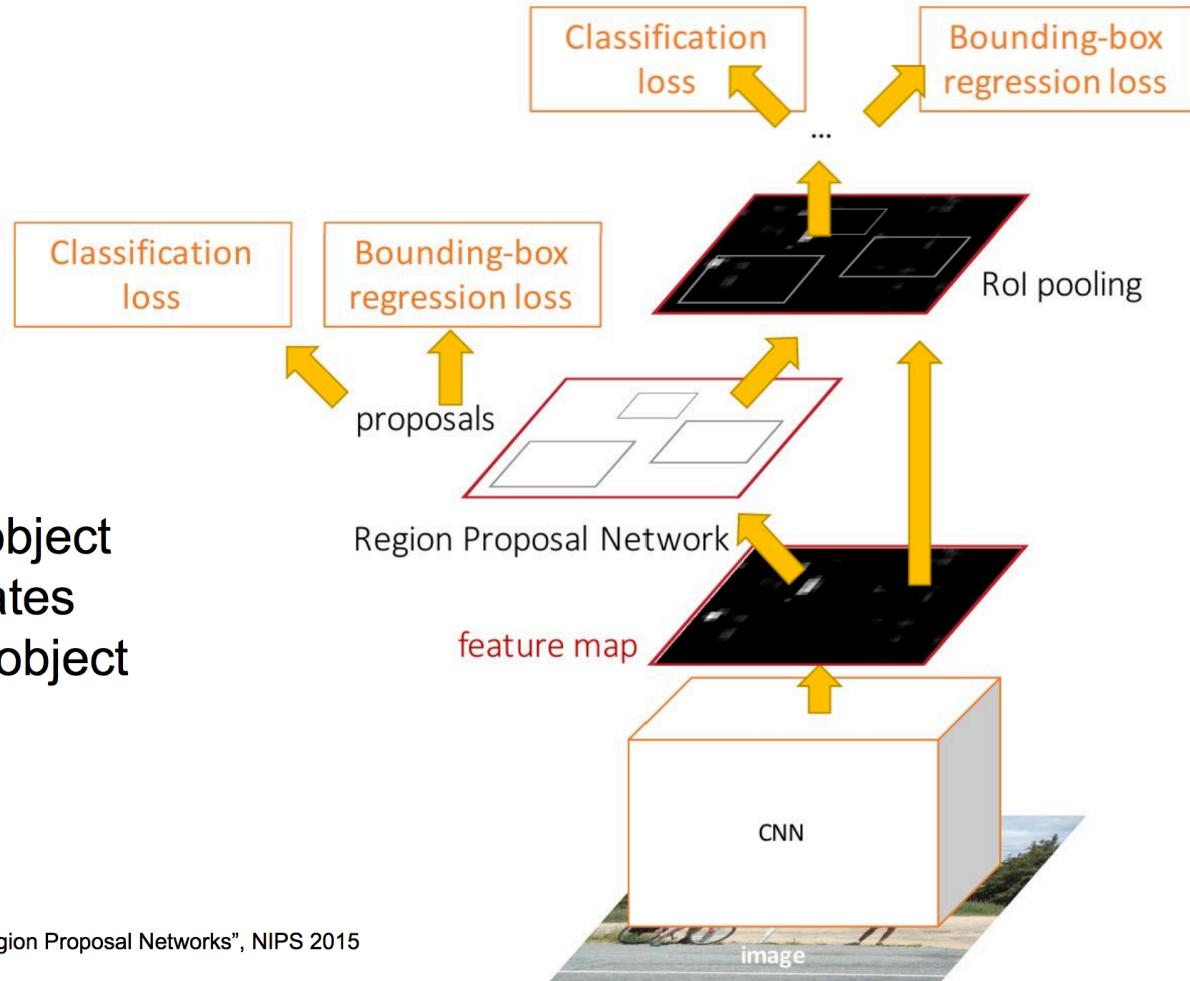
Faster R-CNN:

Make CNN do proposals!

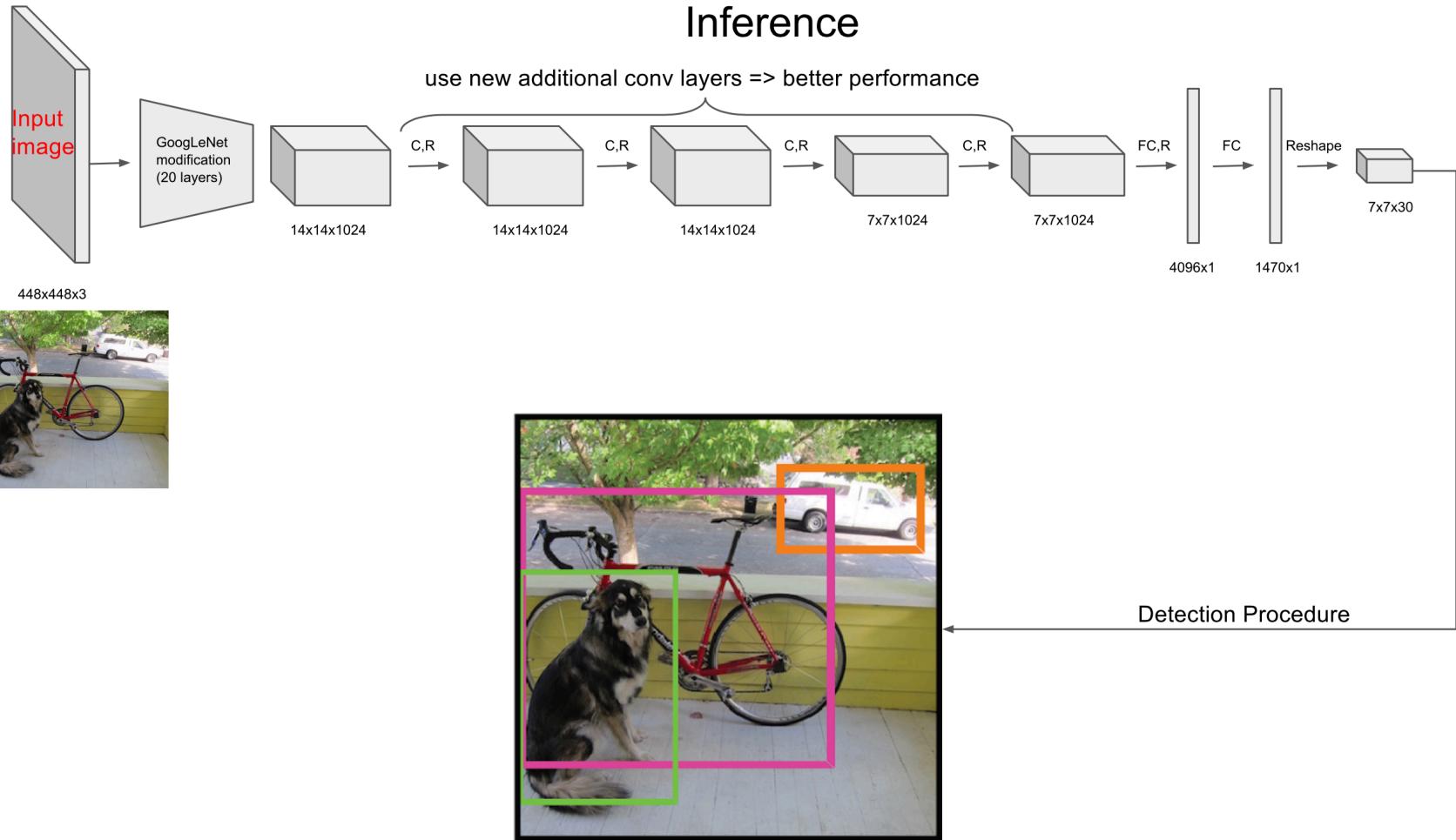
Insert **Region Proposal Network (RPN)** to predict proposals from features

Jointly train with 4 losses:

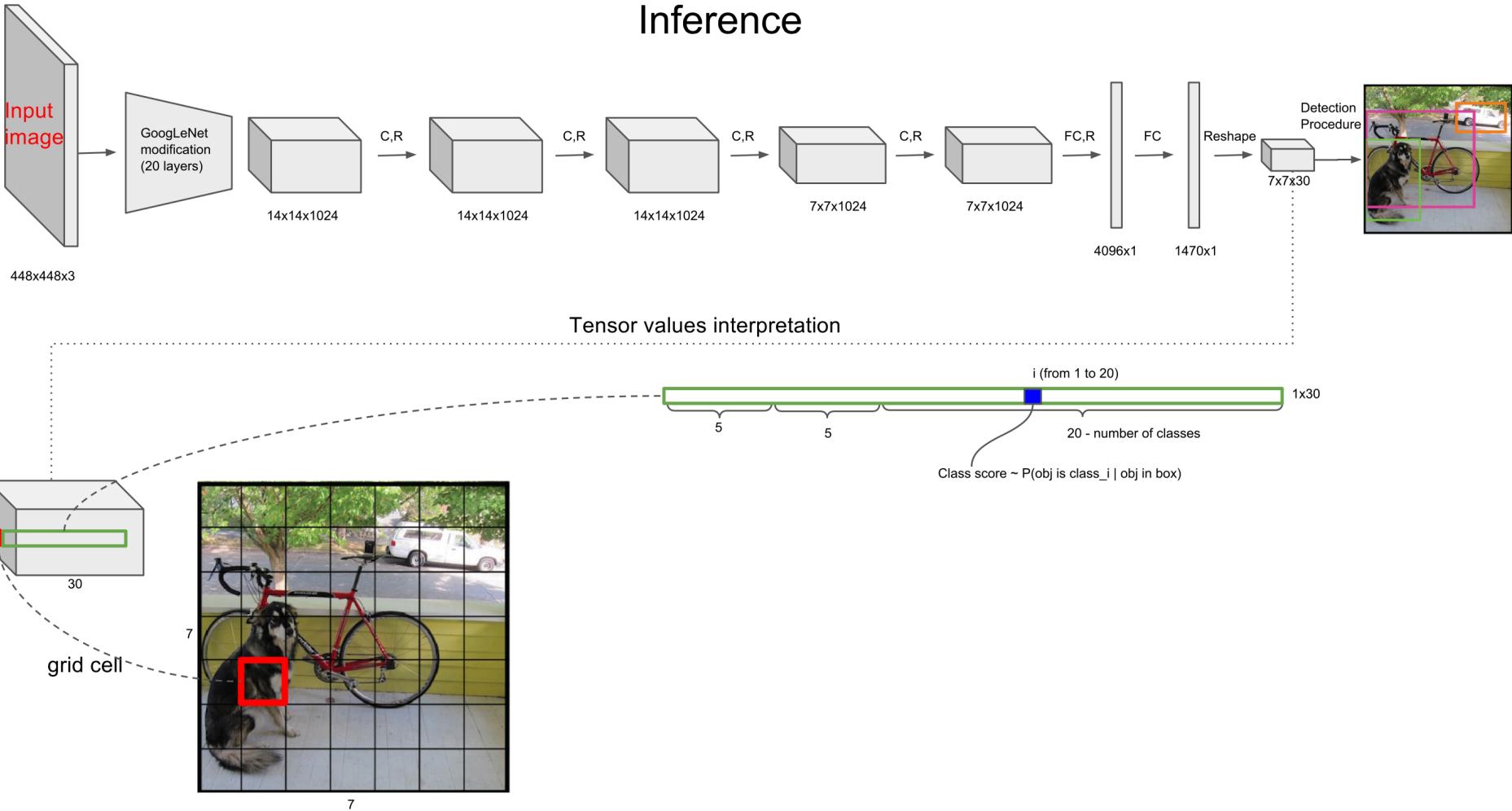
1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



Yolo



Yolo



mAP = mean Average Precision