

# Detection and Segmentation

# Core CV Tasks

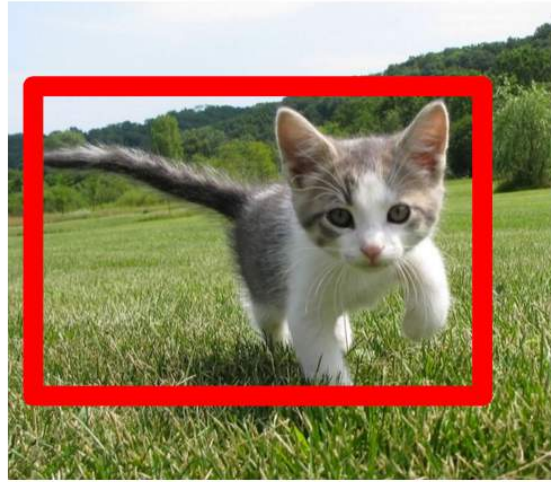
## Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

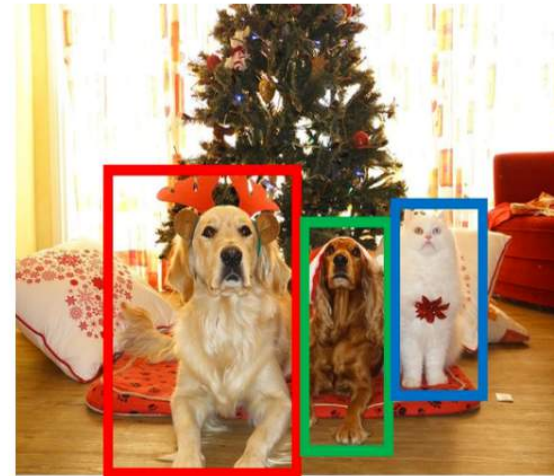
## Classification + Localization



CAT

Single Object

## Object Detection



DOG, DOG, CAT

Multiple Object

## Instance Segmentation



DOG, DOG, CAT

This image is CC0 public domain

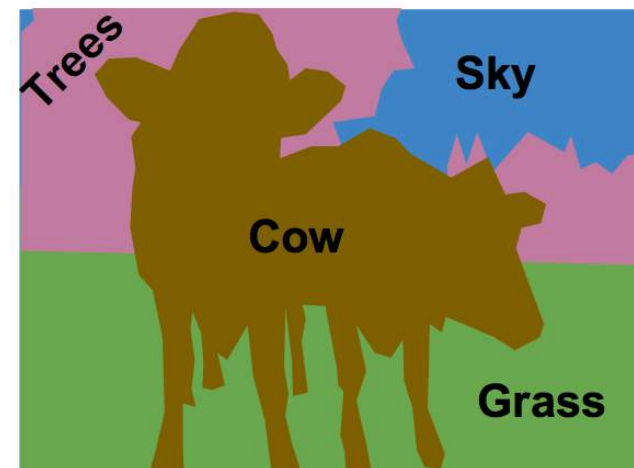
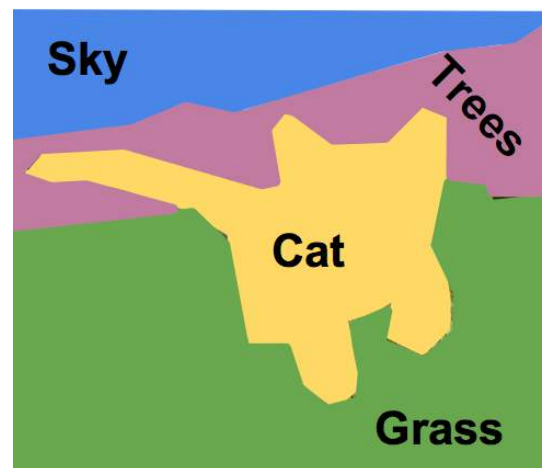
# Semantic Segmentation

## Задача:

Дать лейбл каждому пикселю на картинке.

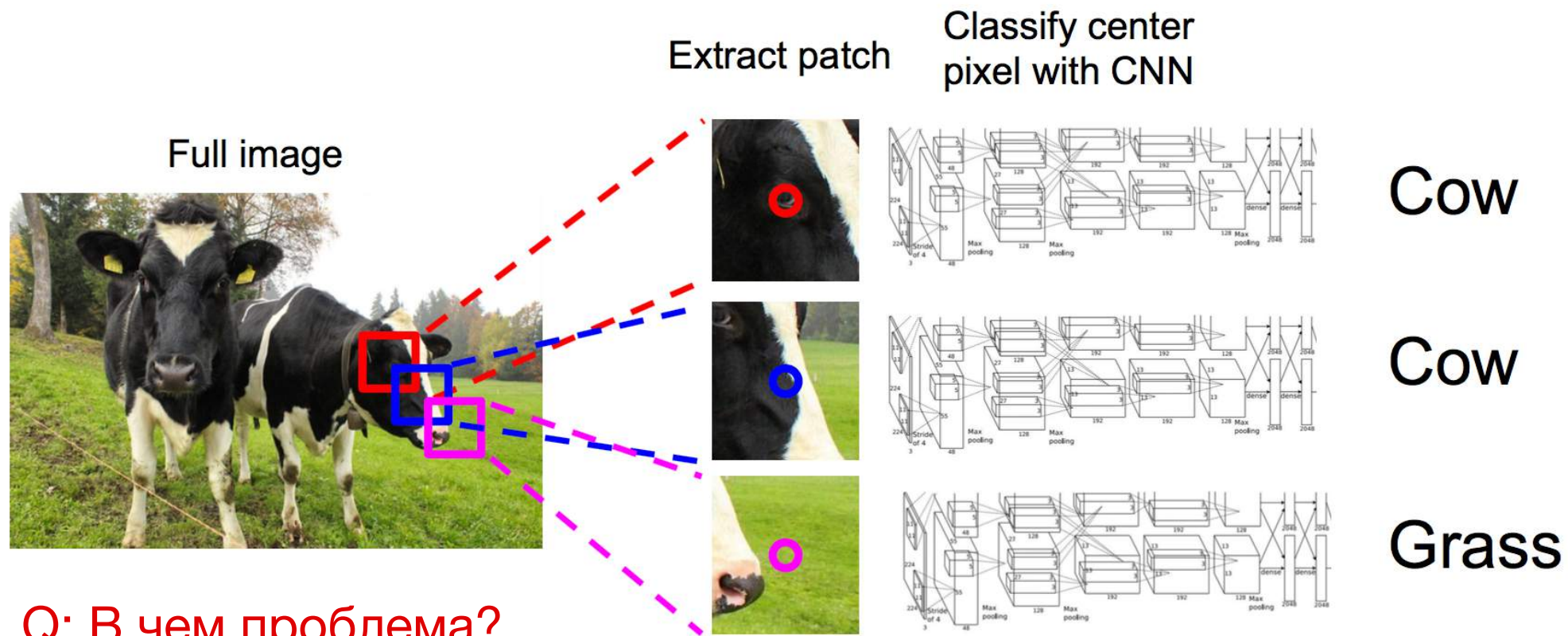
## Важная деталь:

Нас волнуют только пиксели, а не конкретные объекты.





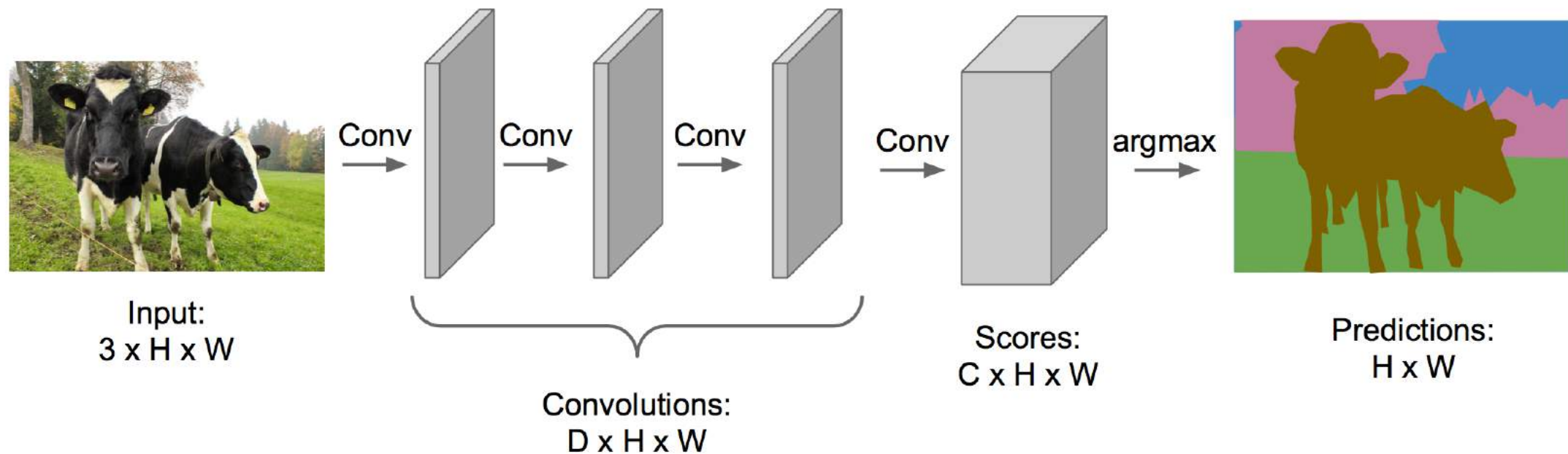
# Semantic Segmentation: Sliding Window



Очень дорого считать.

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

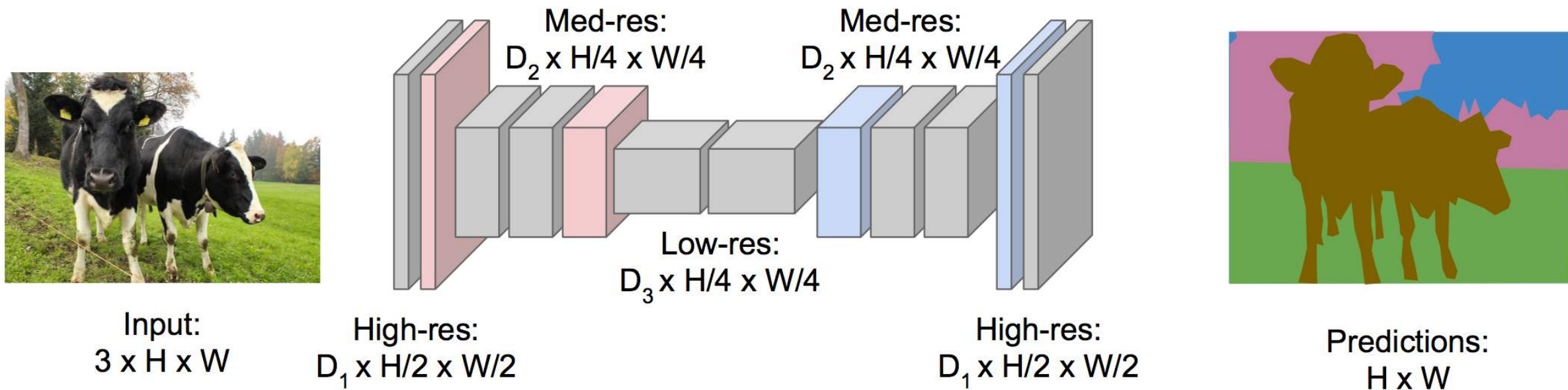
# Semantic Segmentation: Fully Convolutional Net



Q: В чем проблема?

Все еще дорого считать + нужно много памяти.

# Semantic Segmentation: Downsampling + Upsampling



Q: Как делать upsampling?

В основном двумя способами: unpooling и transpose convolution.

# 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

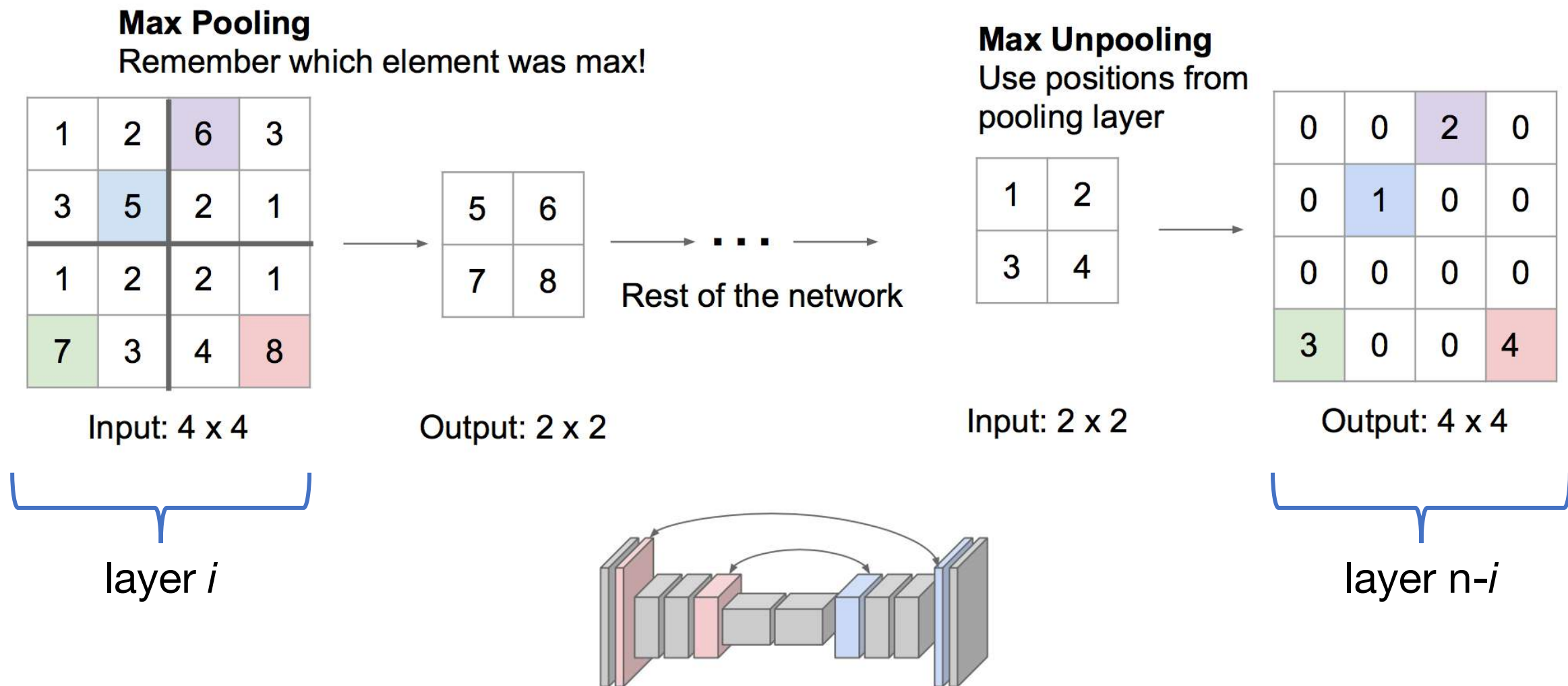
Input: 2 x 2

Output: 4 x 4

Q: Что плохо?

Это слишком просто, чтобы выучить сложные зависимости.

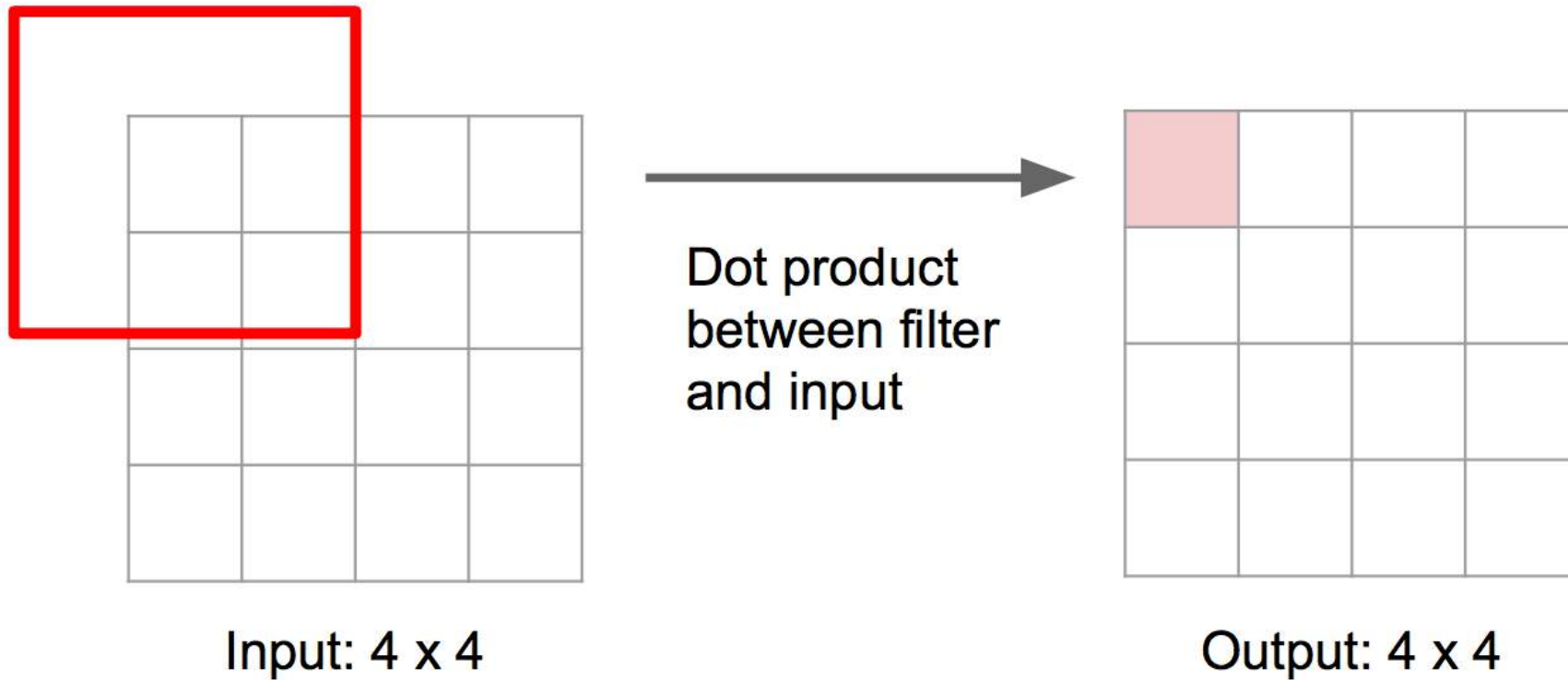
# Unpooling





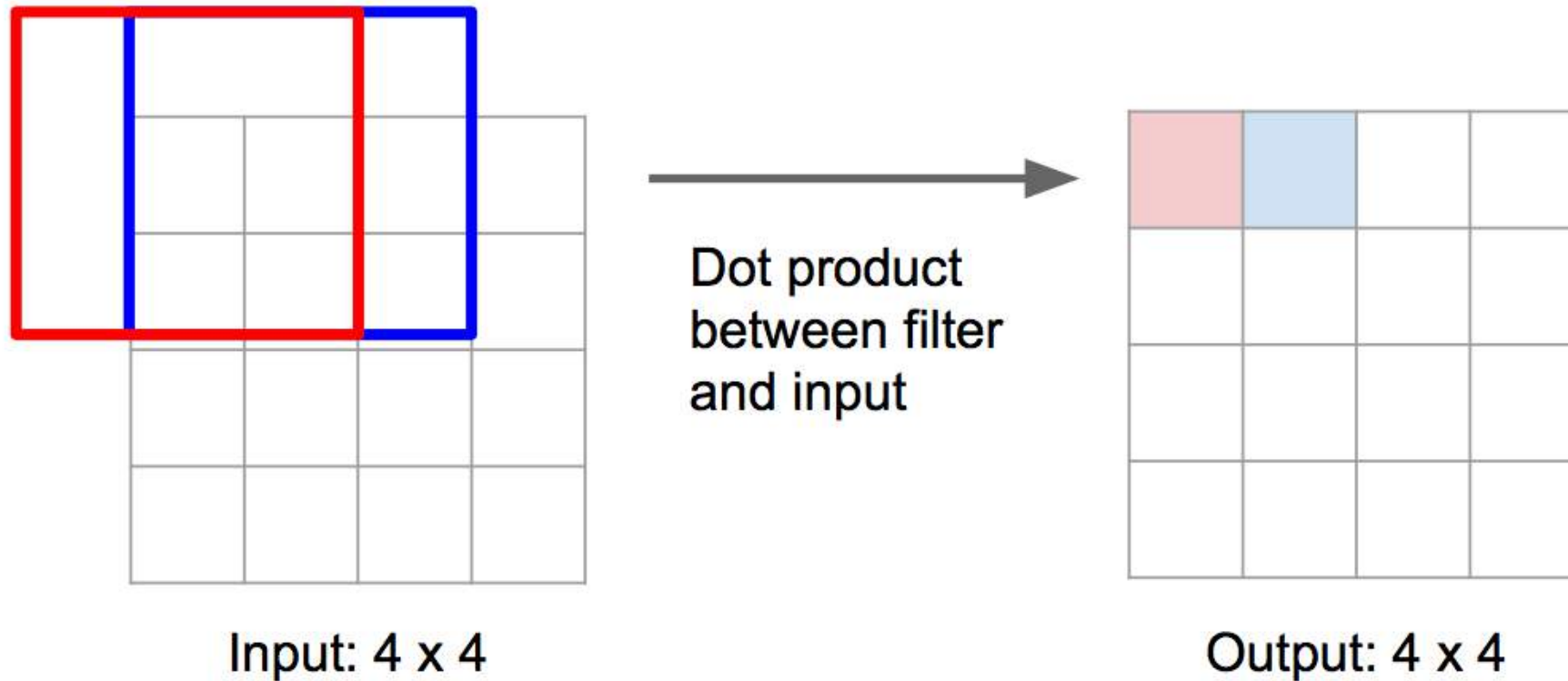
# Transpose Convolution

Обычная свертка filter = 3x3, stride = 1, pad = 1



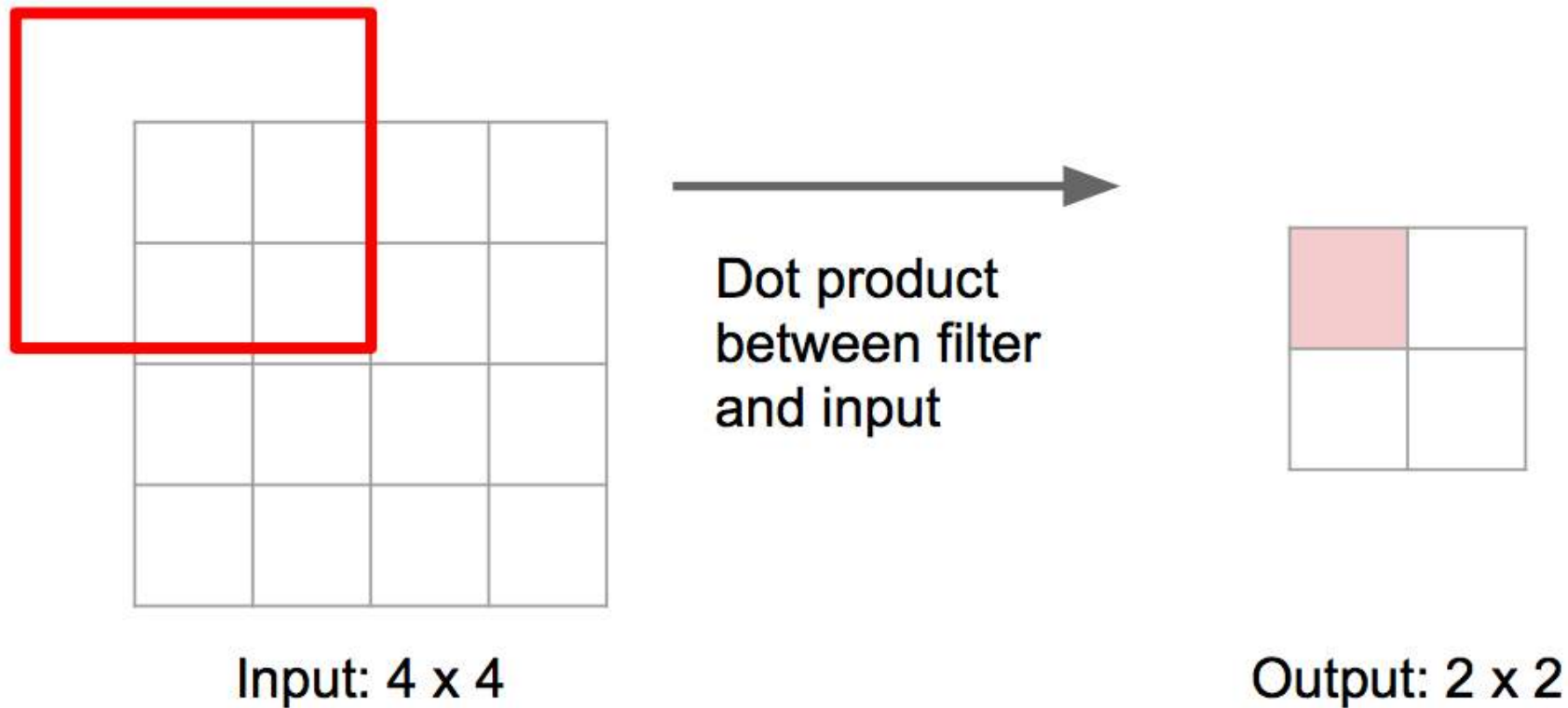
# Transpose Convolution

Обычная свертка filter = 3x3, stride = 1, pad = 1



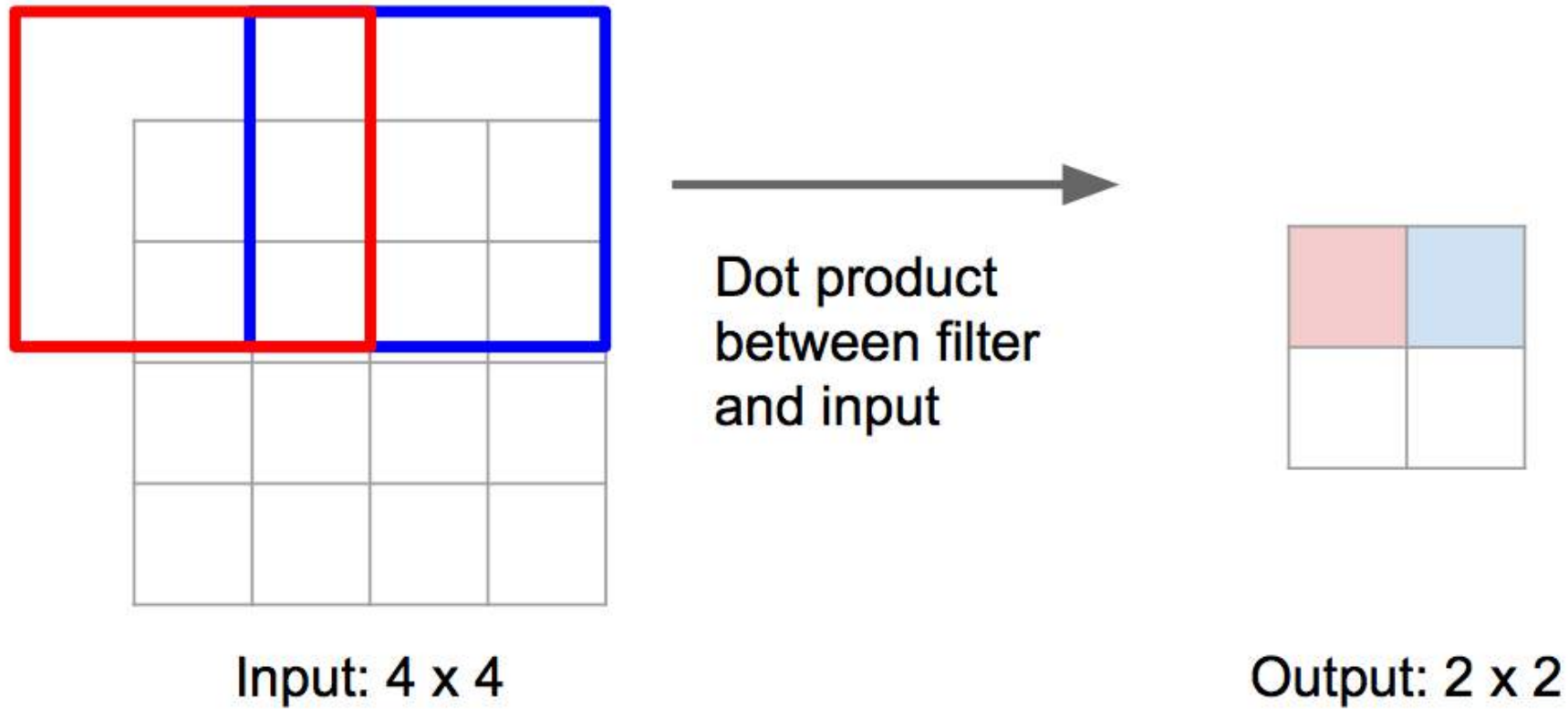
# Transpose Convolution

Обычная свертка filter = 3x3, stride = 2, pad = 1



# Transpose Convolution

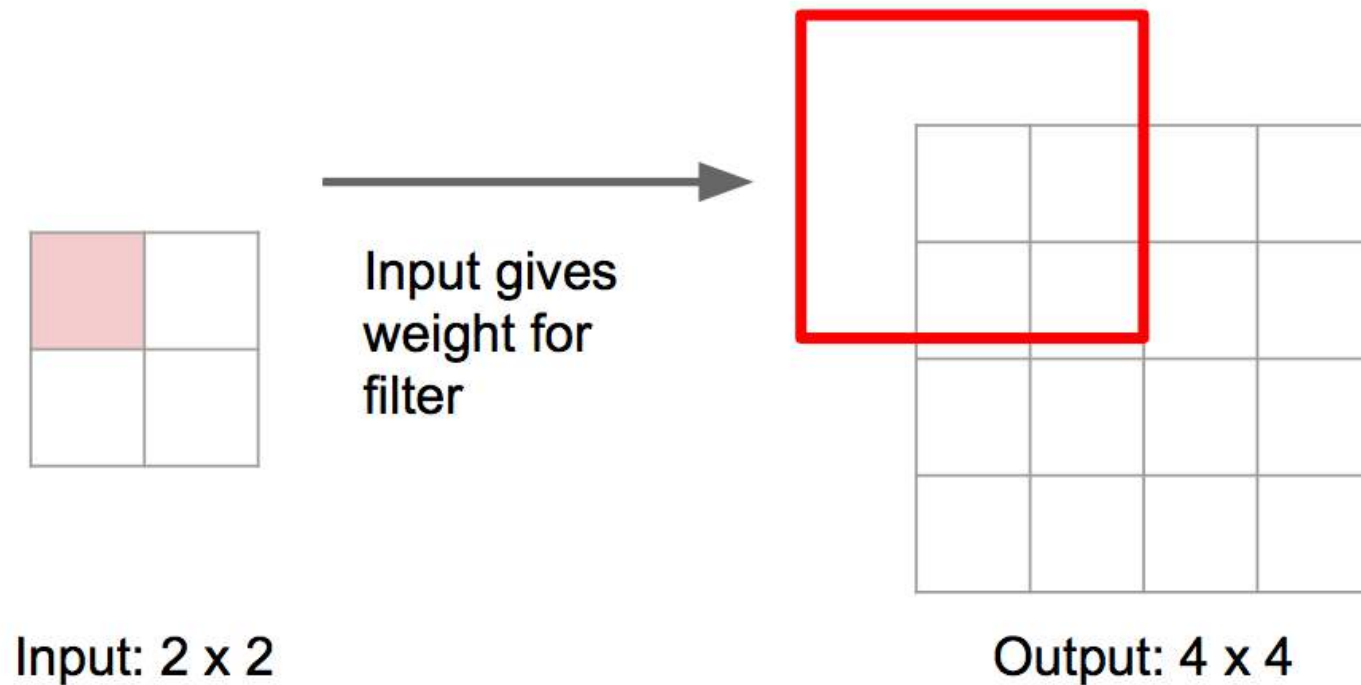
Обычная свертка filter = 3x3, stride = 2, pad = 1



# Transpose Convolution

---

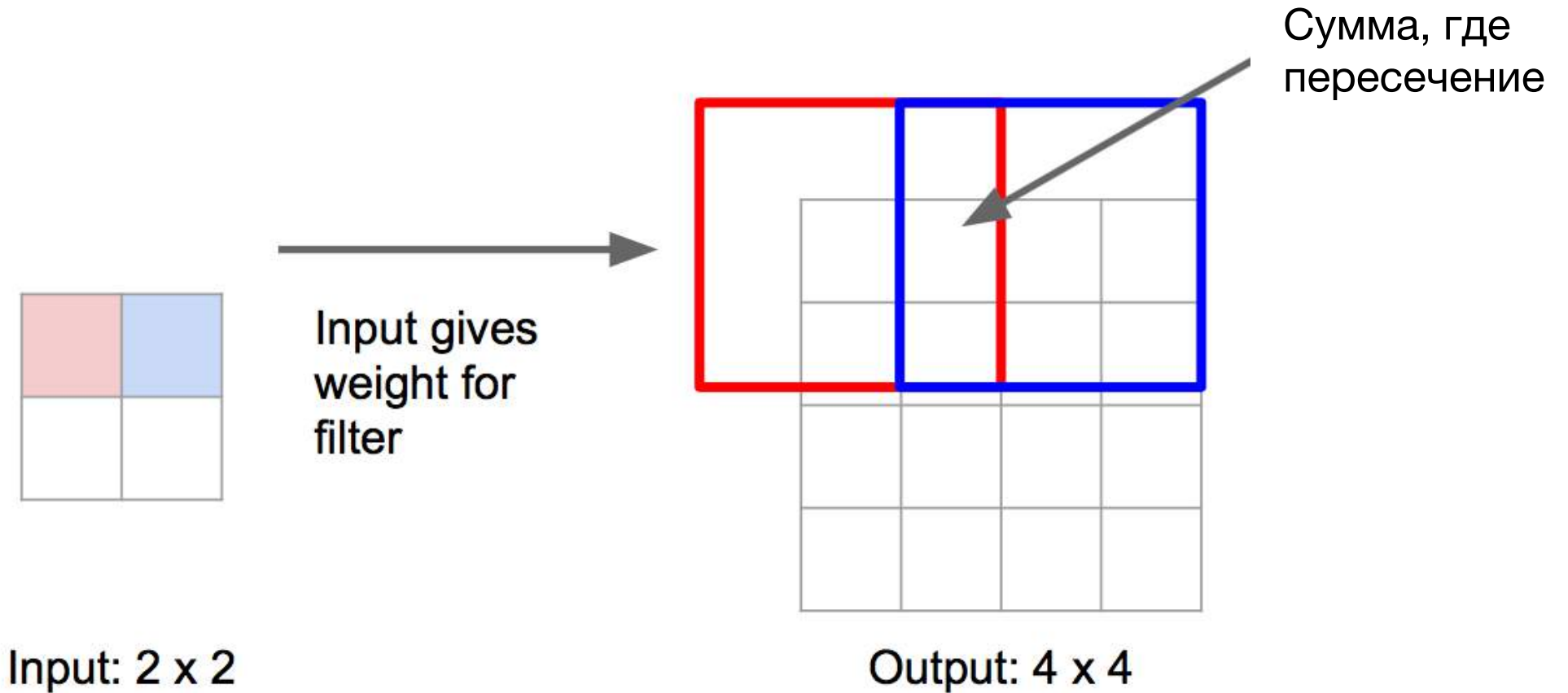
Transpose свертка filter = 3x3, stride = 2, pad = 1





# Transpose Convolution

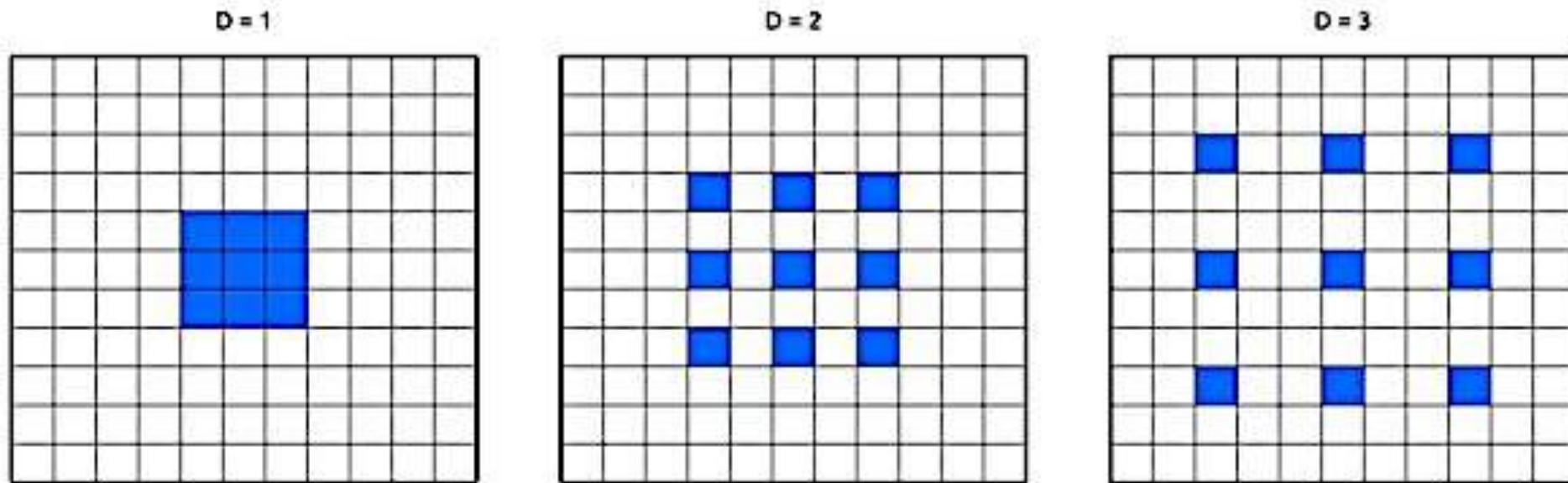
Transpose свертка filter = 3x3, stride = 2, pad = 1



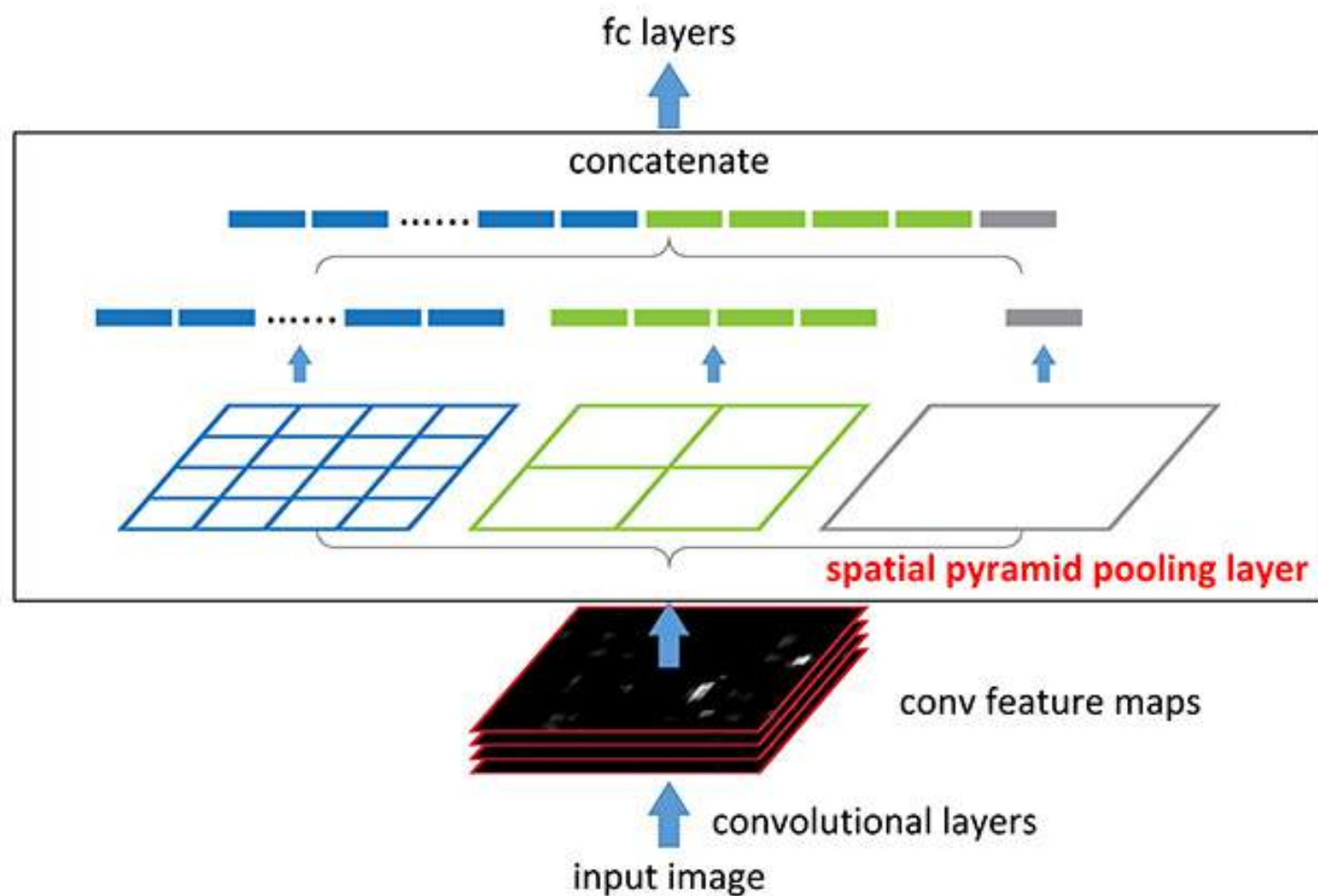
# Dilated Convolution

---

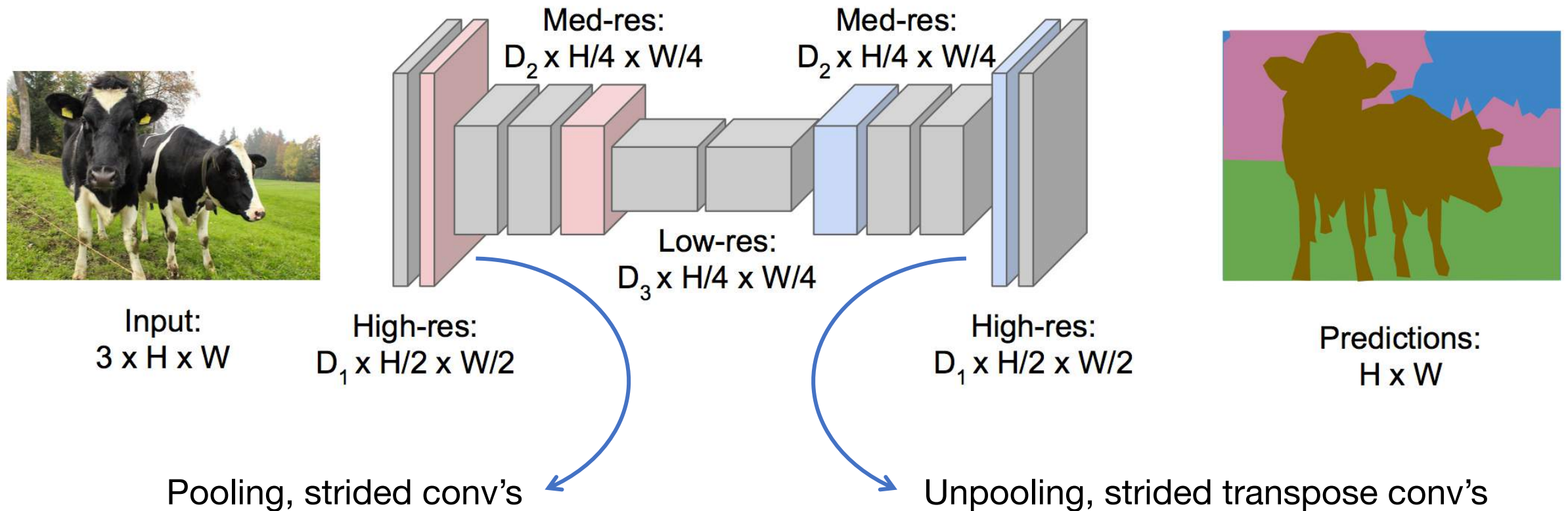
Dilated свертка filter = 3x3



# Pyramid Pooling

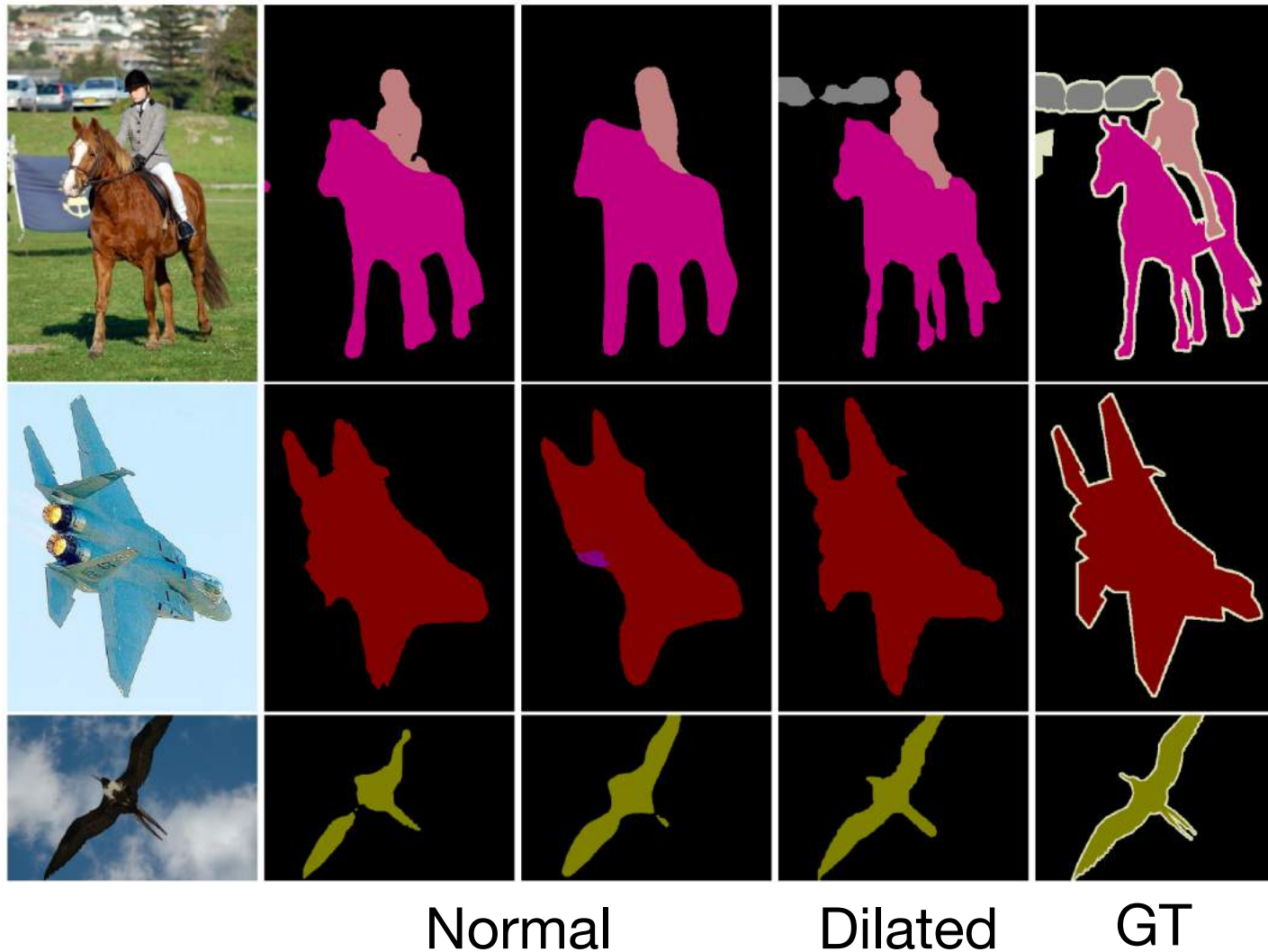


# Semantic Segmentation: General Architecture



# Semantic Segmentation: Outcome

---



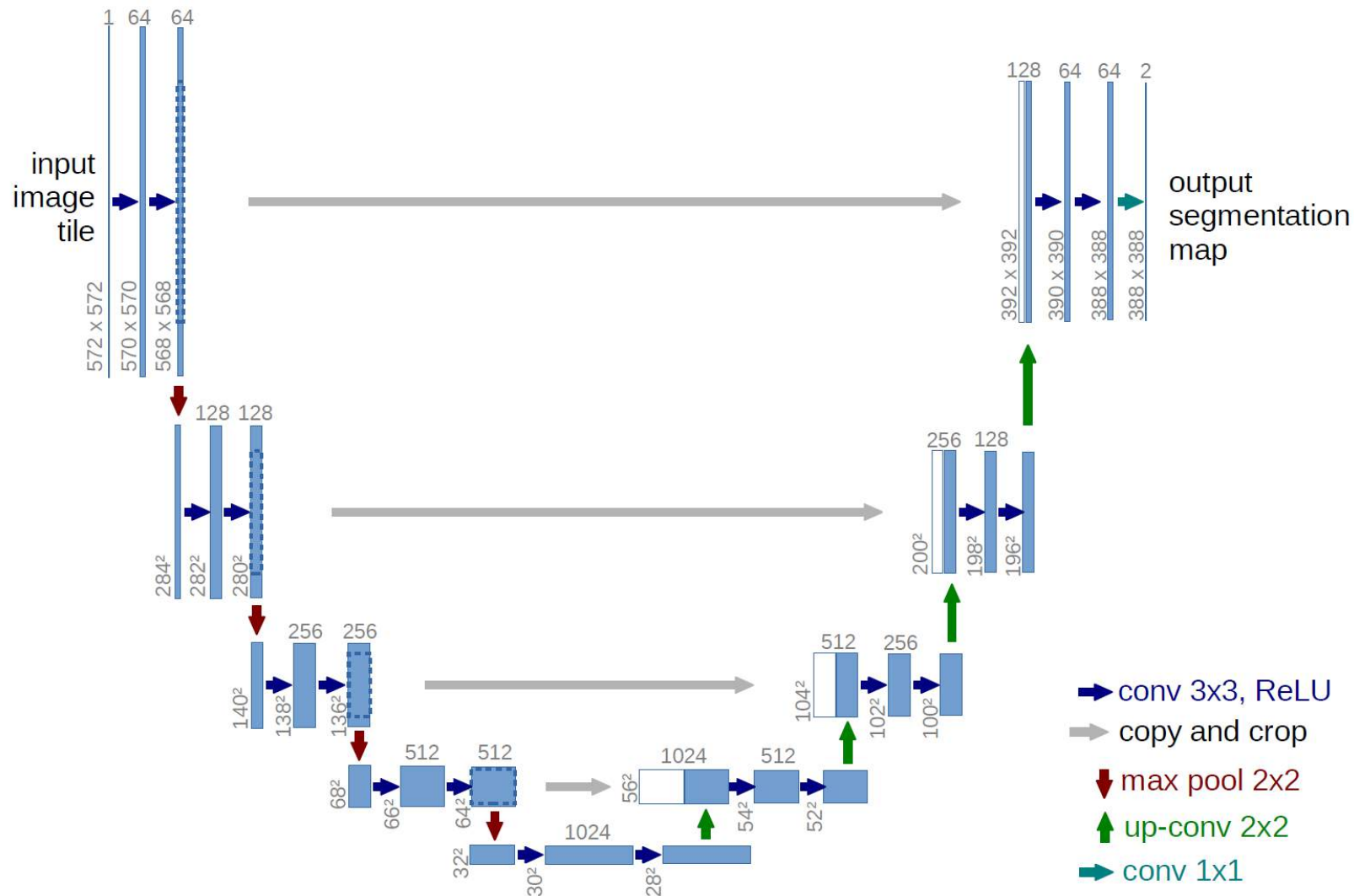


# Semantic Segmentation: Outcome

---



# Semantic Segmentation: U-Net (skip-connections)



# Semantic Segmentation: Basic Ideas

---

- Downsampling + Upsampling
- Transpose conv's вместо unpooling
- Dilated свертки вместо обычных
- Skip-connections
- Pyramid pooling

# Classification + Localization

---

## **Задача:**

Обвести объект рамкой и классифицировать его.

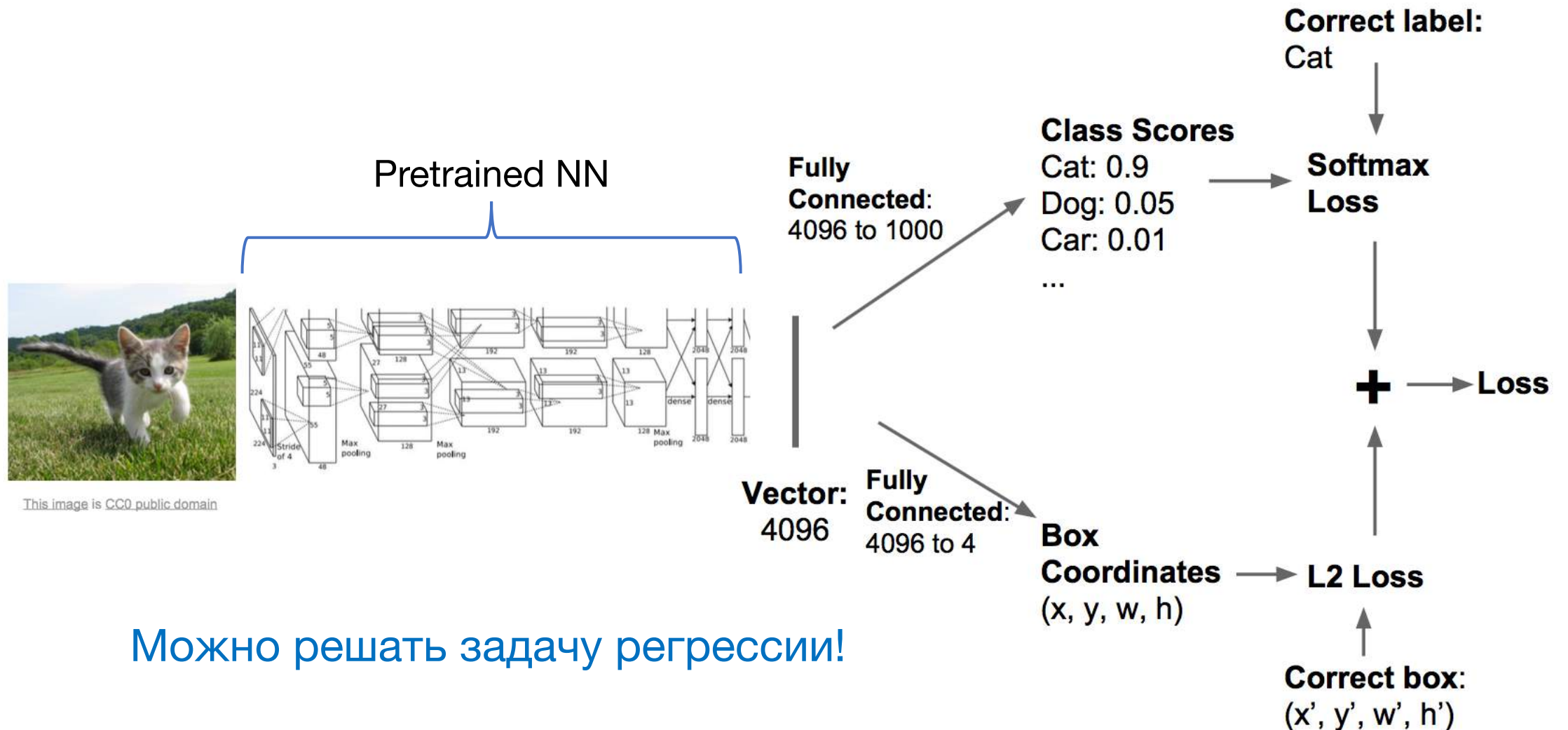
## **Важная деталь:**

Мы всегда имеем дело только с одним объектом на картинке.



**CAT**

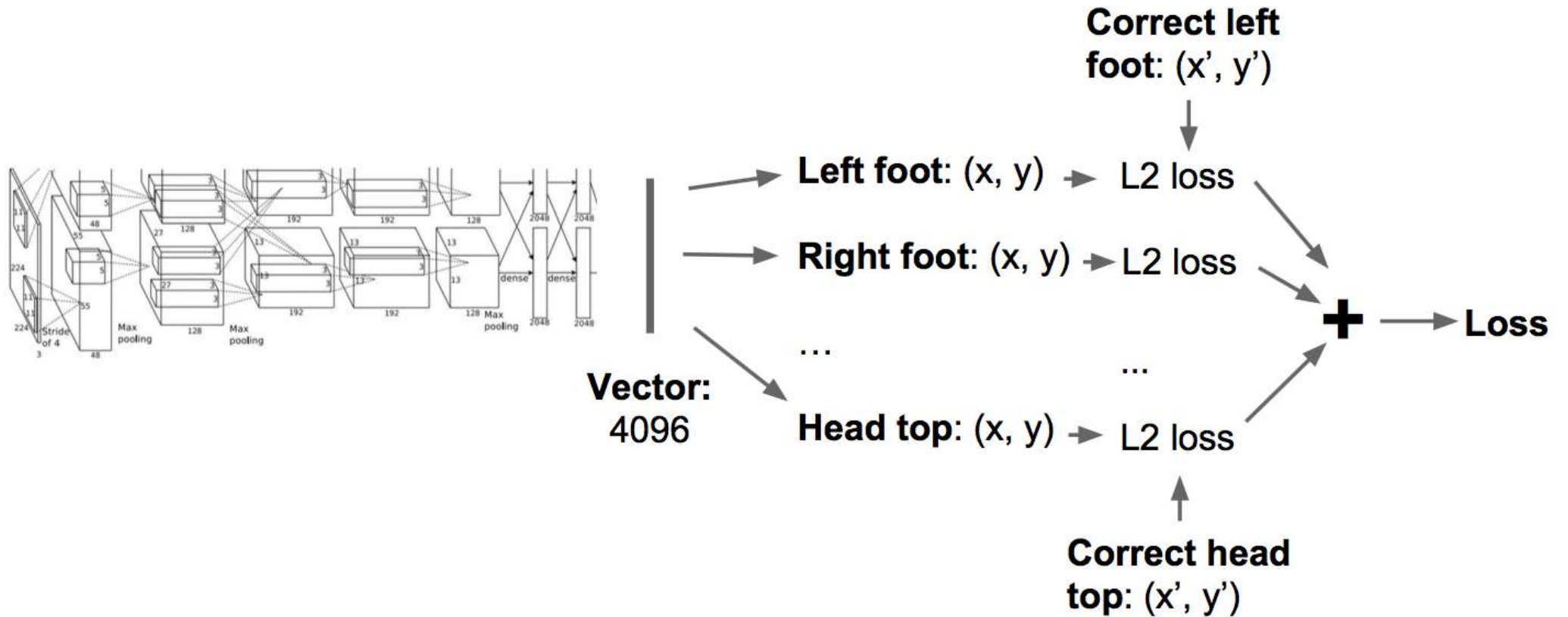
# Classification + Localization



Можно решать задачу регрессии!



# Human Pose Estimation



# Object Detection

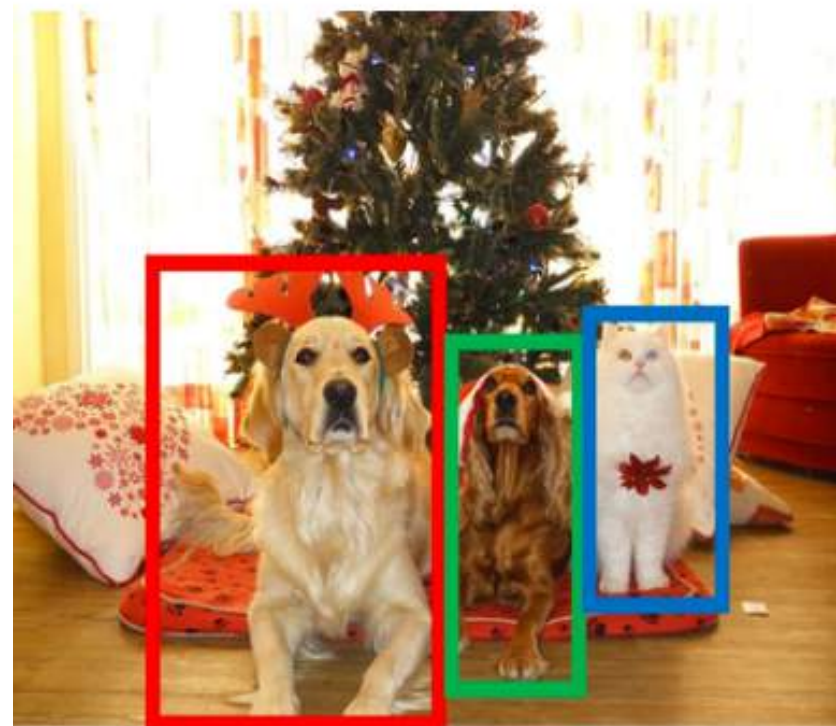
---

## Задача:

Обвести все объекты рамкой и классифицировать их.

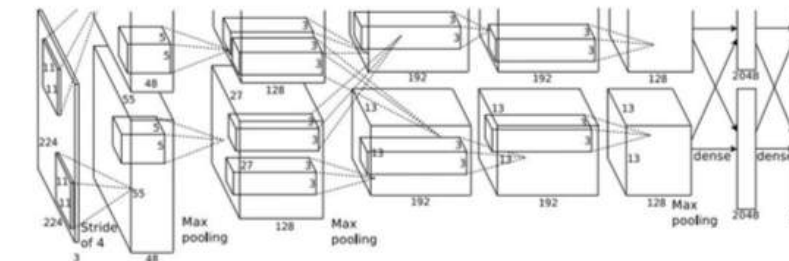
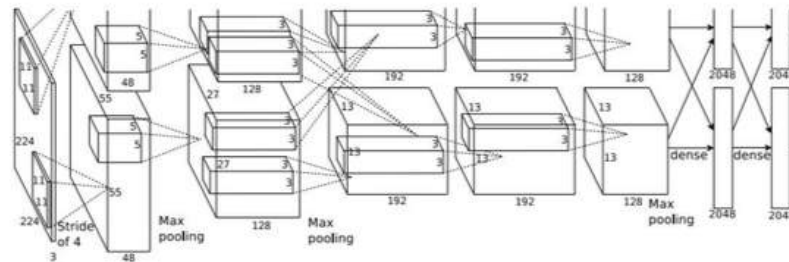
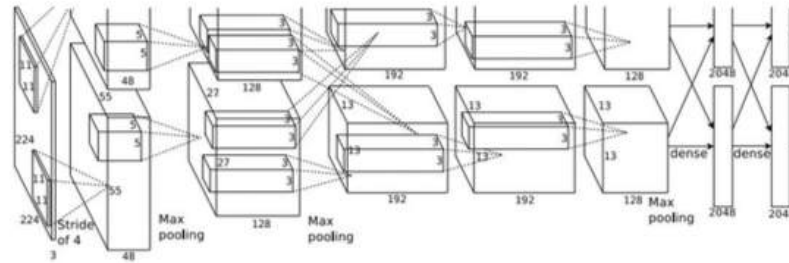
## Важная деталь:

Мы имеем дело с неизвестным количеством объектов.



**DOG, DOG, CAT**

# Object Detection



CAT: (x, y, w, h)

DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)

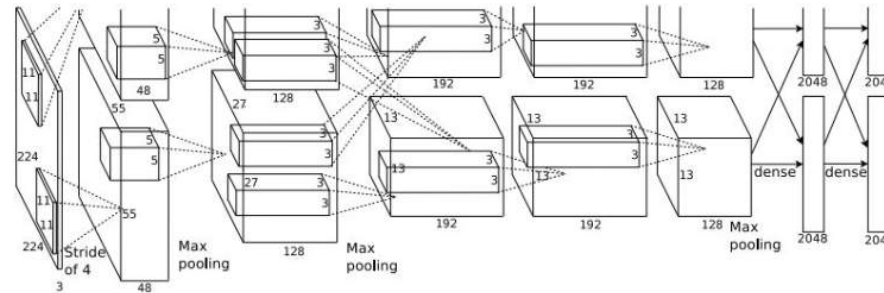
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

....

Разное  
количество  
выходов

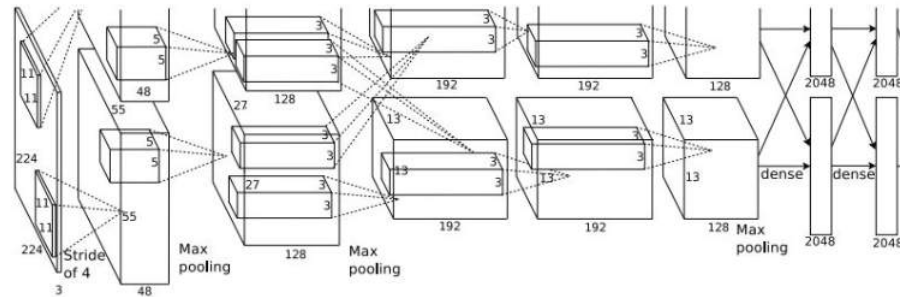
# Object Detection: Sliding Window



Dog? NO  
Cat? NO  
Background? YES



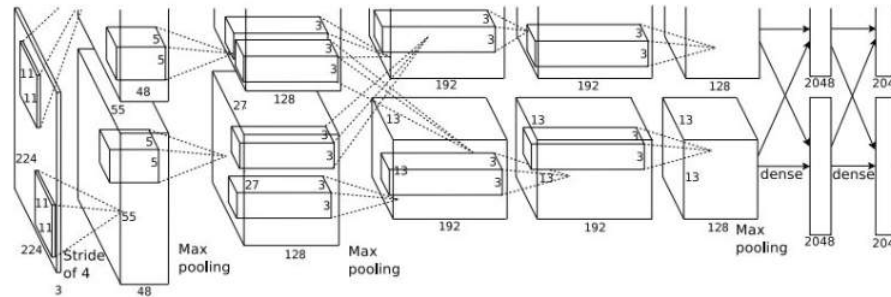
# Object Detection: Sliding Window



Dog? YES  
Cat? NO  
Background? NO

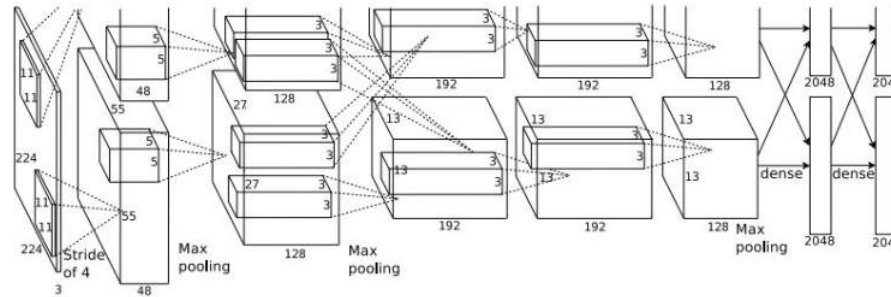


# Object Detection: Sliding Window



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Sliding Window



Dog? NO  
Cat? YES  
Background? NO

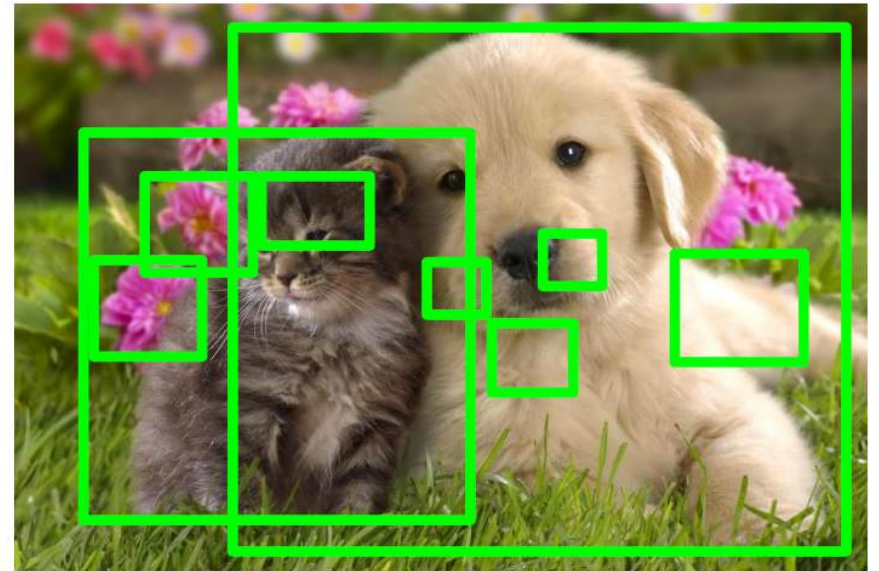
Q: В чем проблема?

Очень дорого считать.

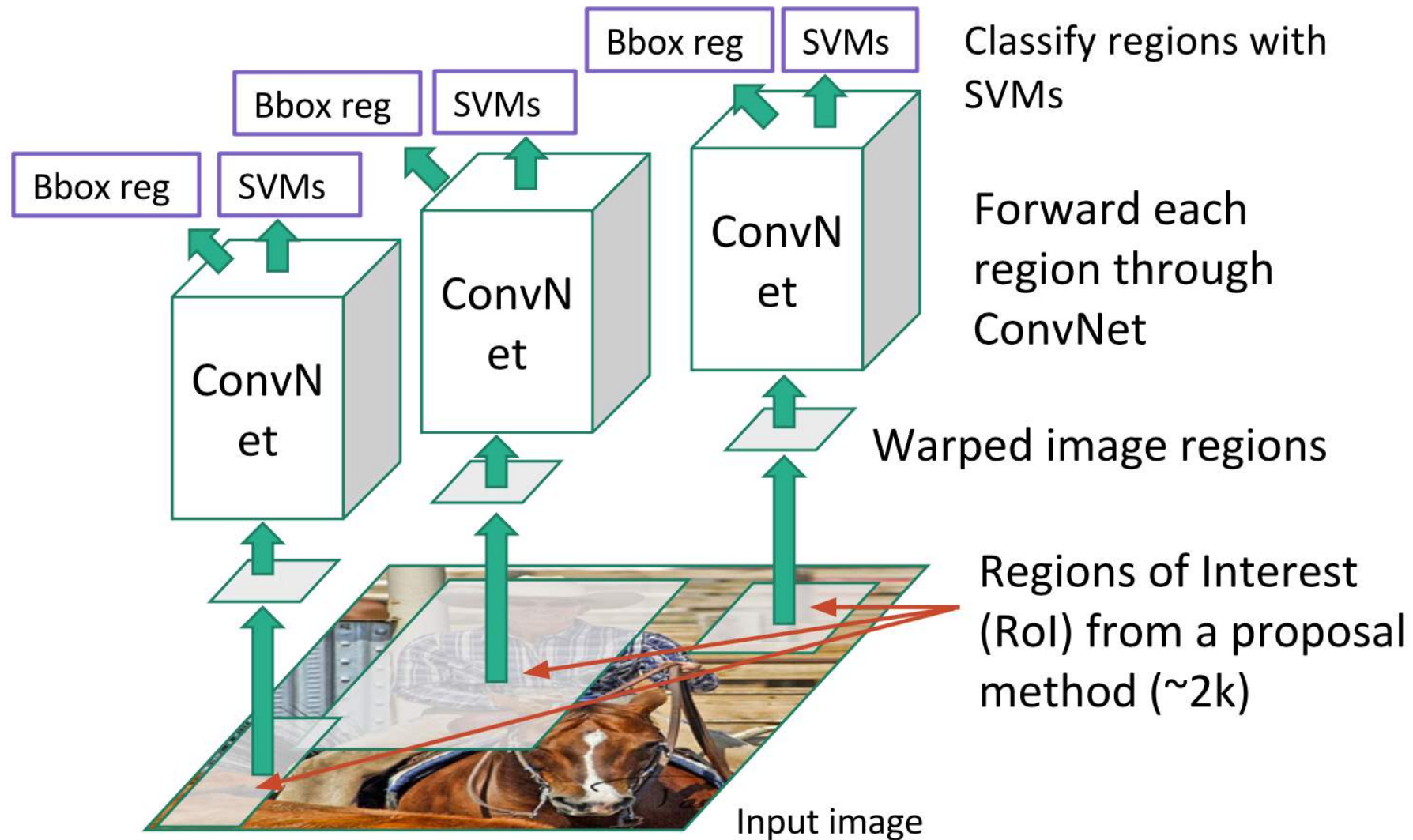
# Region Proposals

---

- Можно использовать методы, не связанные с DL. Например, Selective Search.
- Работает с нормальной скоростью, примерно 2000 регионов за пару секунд.



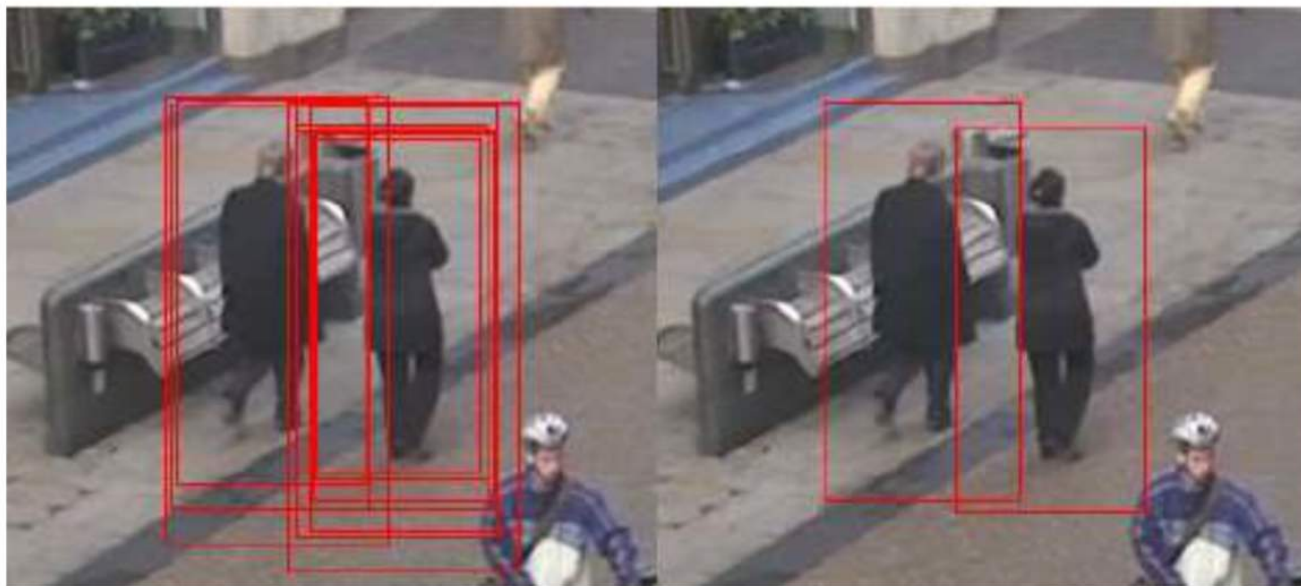
# Object Detection: R-CNN





# Object Detection: Non-Maximum Suppression

---

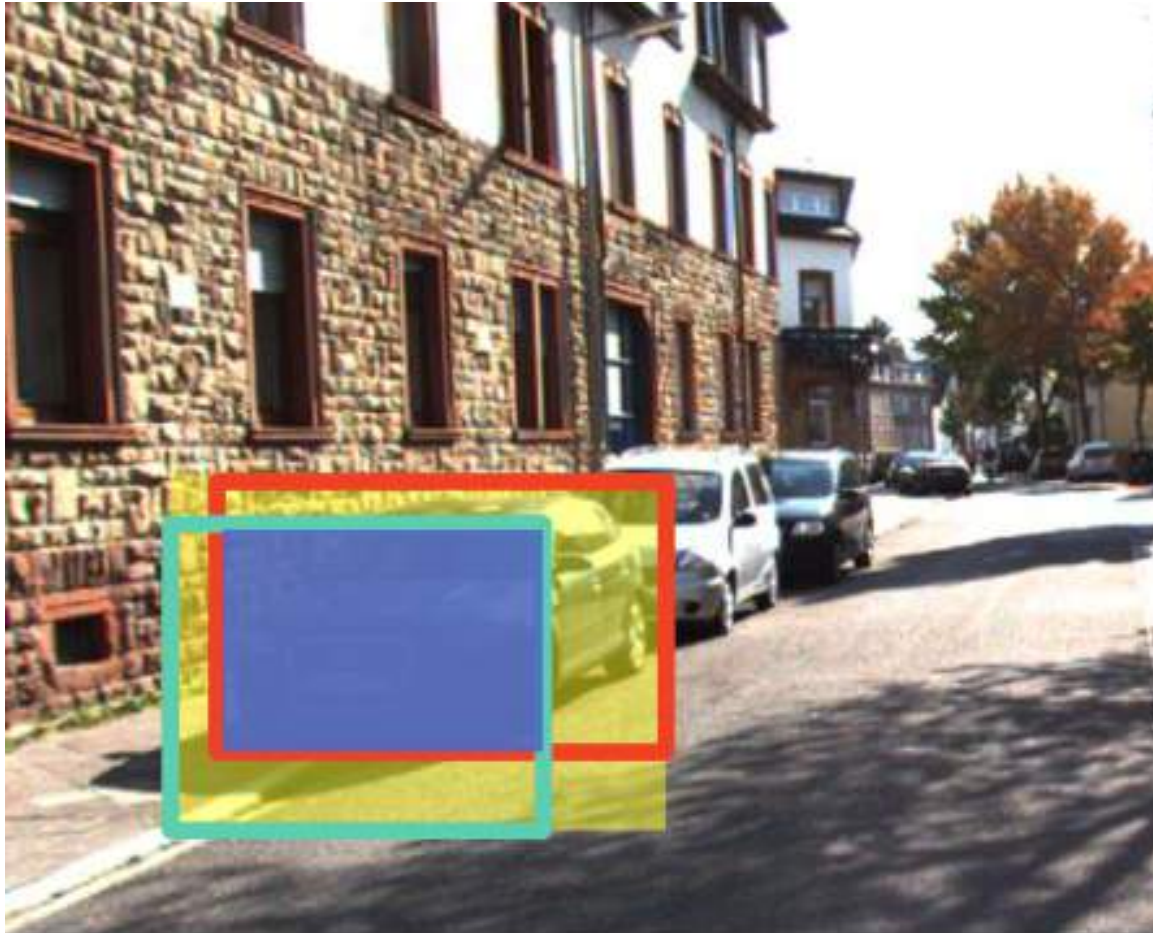


От детектора получается много пересекающихся bbox, покрывающие одни и те же объекты:

1. Фильтруем по убыванию вероятностей
2. Идём по порядку и удаляем все прочие bbox, пересекающие с текущим на  $>50\%$

# Object Detection: IoU

---



intersection



union

$$\text{IoU} = \text{intersection} / \text{Union}$$

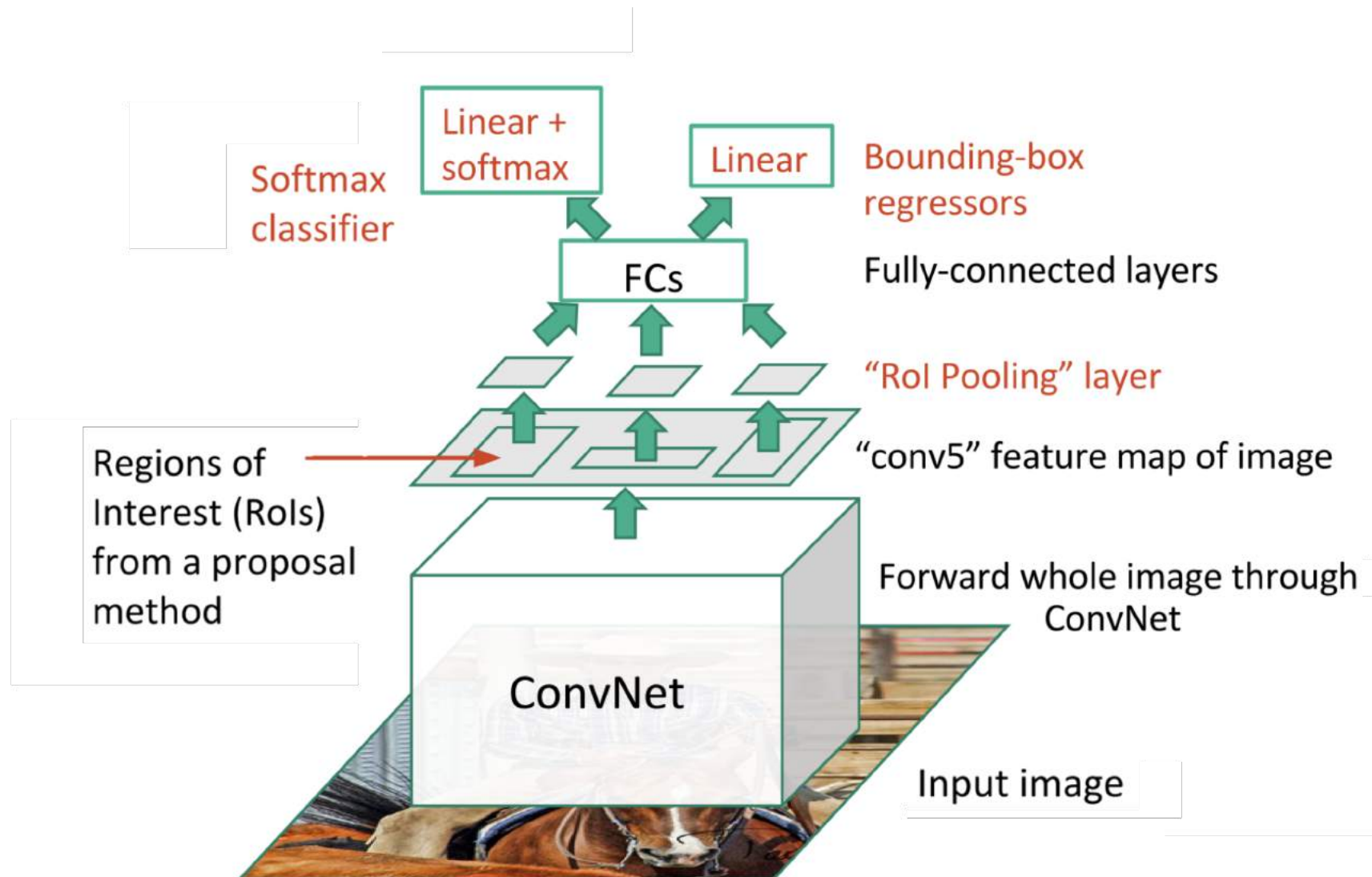


# Object Detection: R-CNN Problems

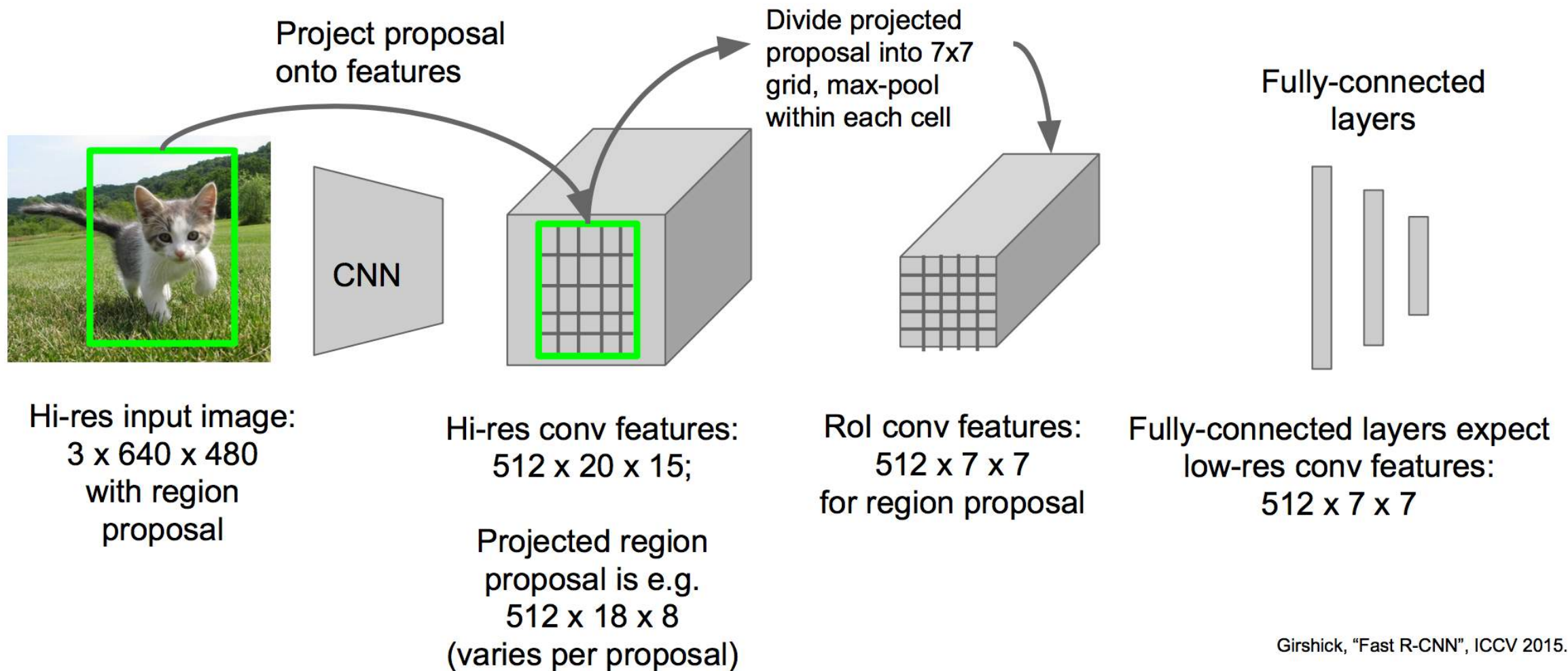
---

- Очень много разных моделей в одном фреймворке.
  - CNN, которую мы доучиваем в процессе.
  - SVM, который мы тоже доучиваем.
  - BBox Regression, опять же учим.
- Долго учится и требует много памяти.
- Долго работает на inference стадии (47s на картинку).

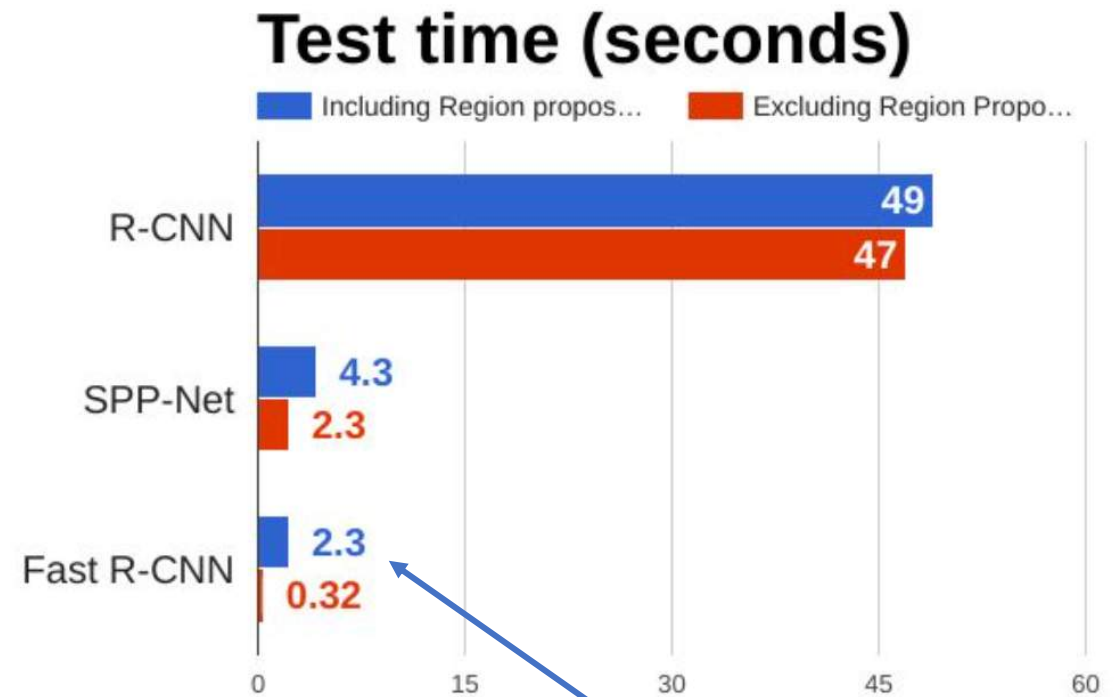
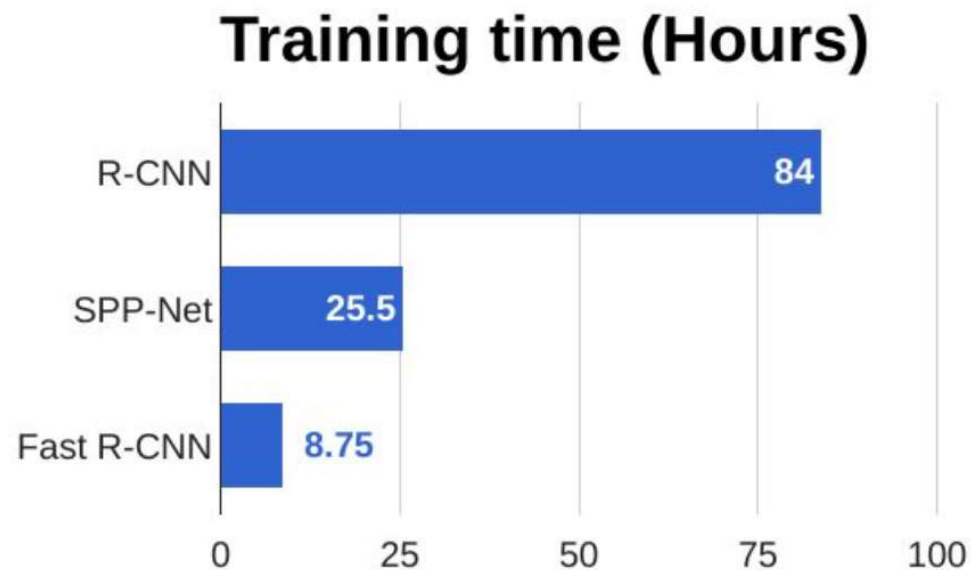
# Object Detection: Fast R-CNN



# Object Detection: RoI Pooling

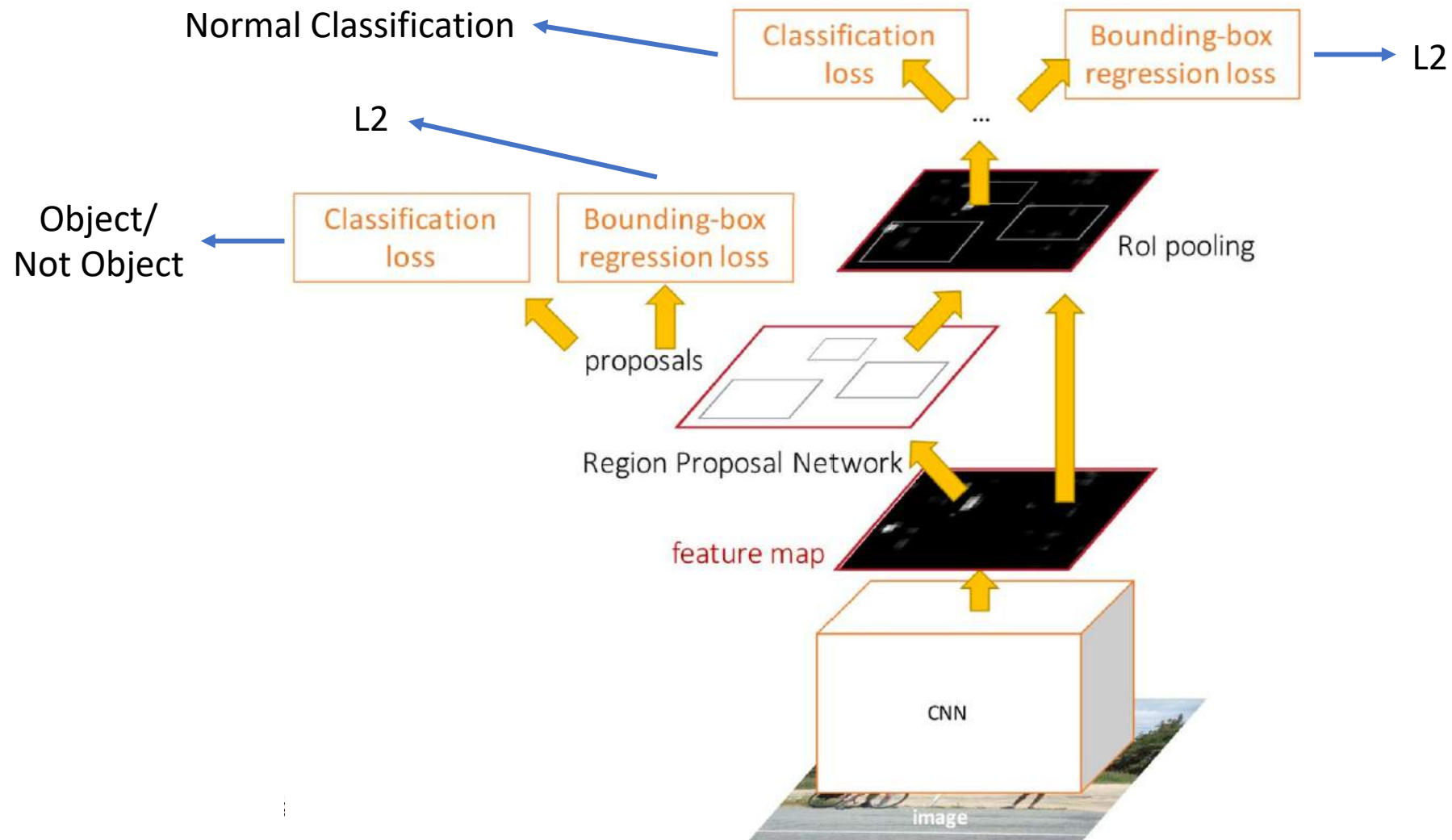


# Object Detection: Runtimes

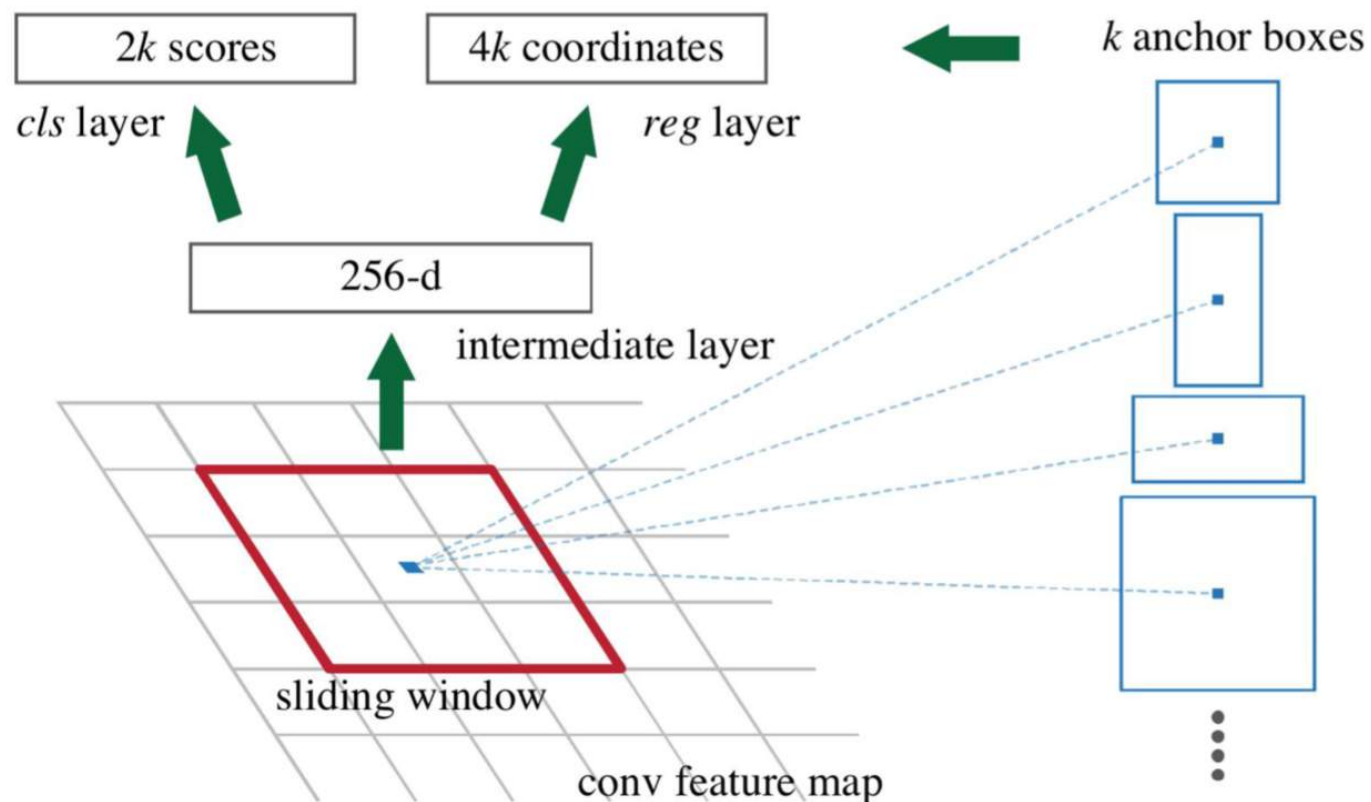


Медленно из-за  
алгоритма  
выделения RoI.

# Object Detection: Faster R-CNN



# Object Detection: Region Proposal Network

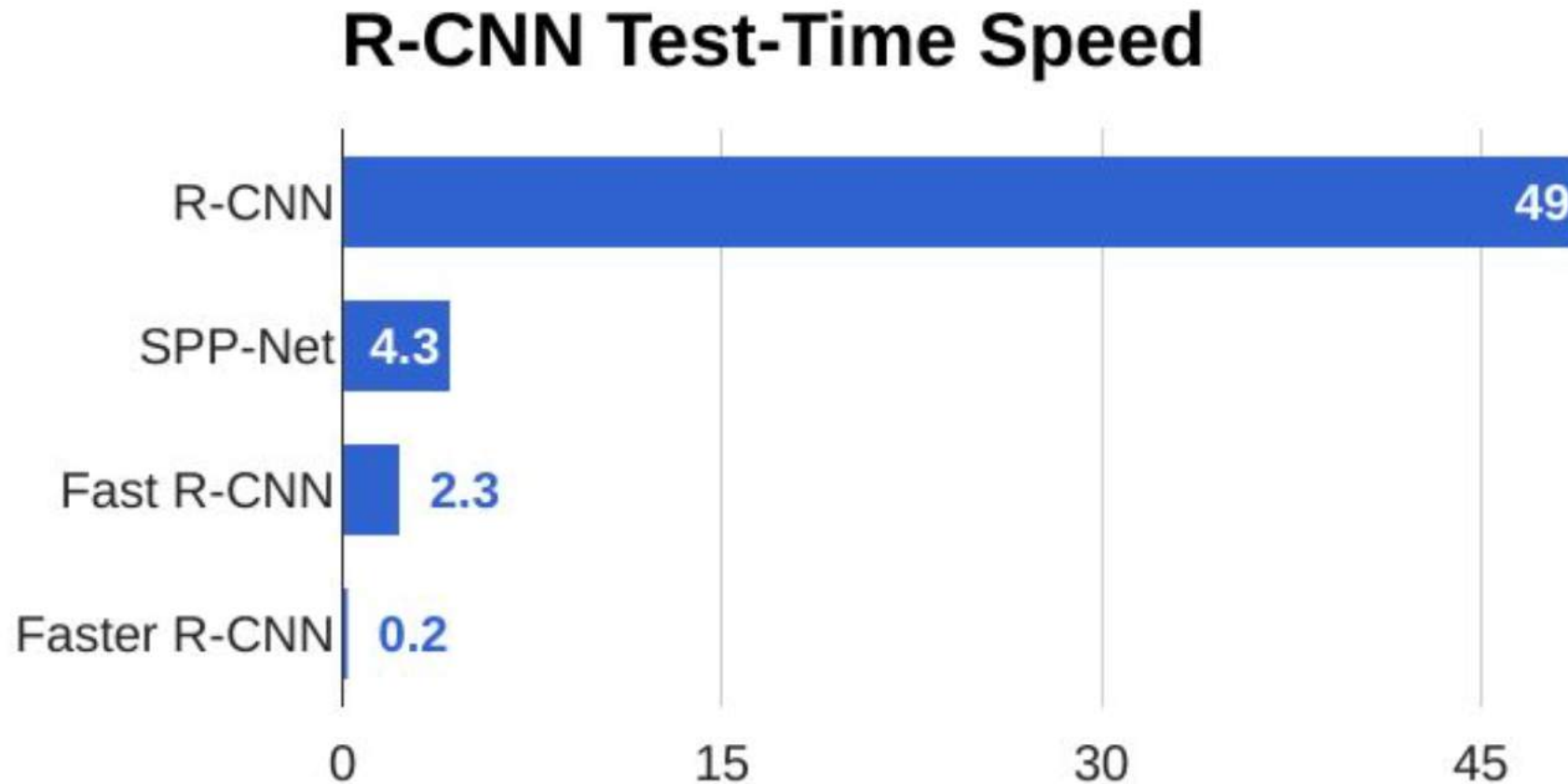


1. Небольшим скользящим окном двигаемся по признакам-активациям
2. Оцениваем шансы на объект + корректируем рамки
3. Окно разных размеров и соотношений
4. Из 17000 оставим 300 proposal

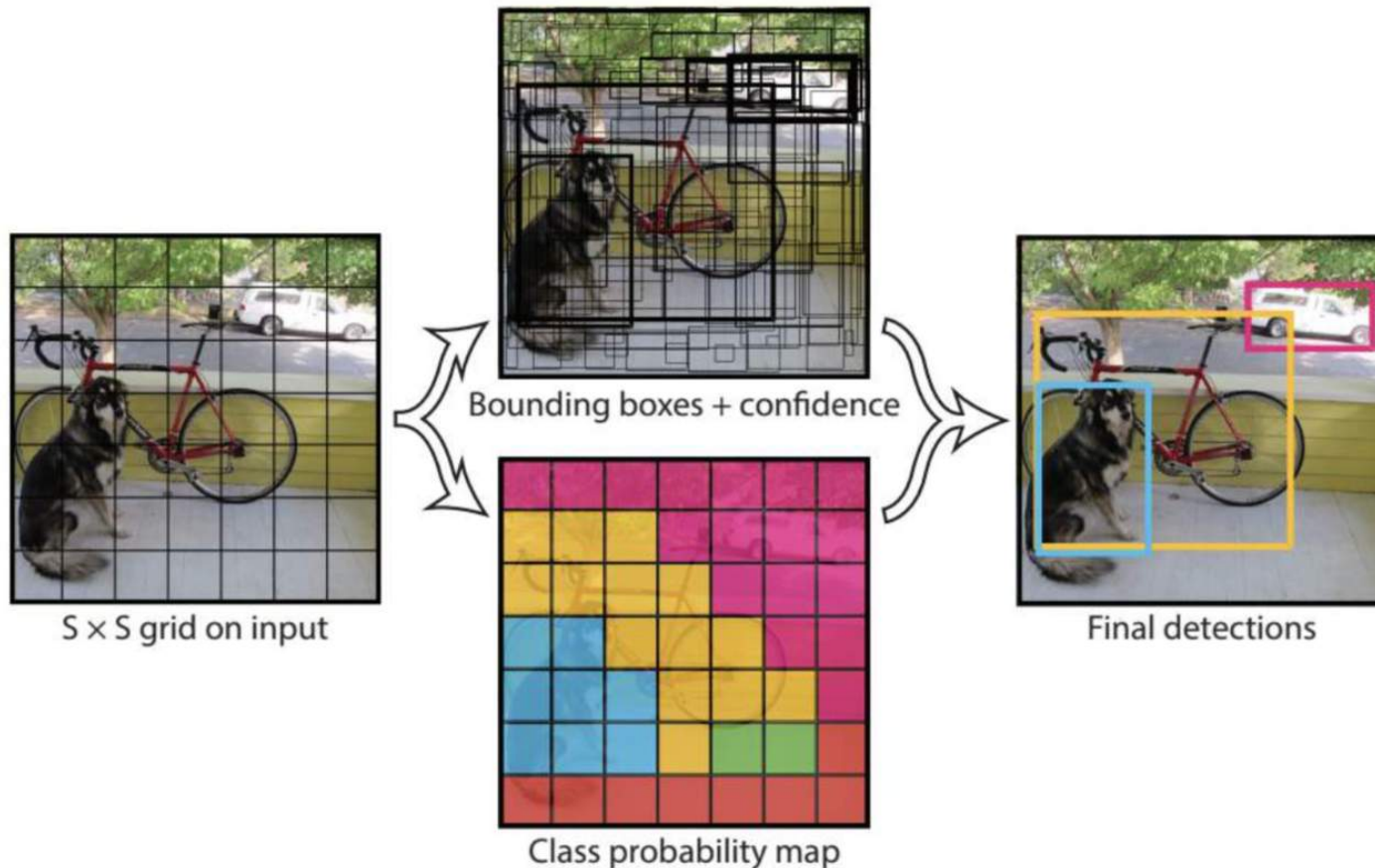


# Object Detection: Runtimes

---



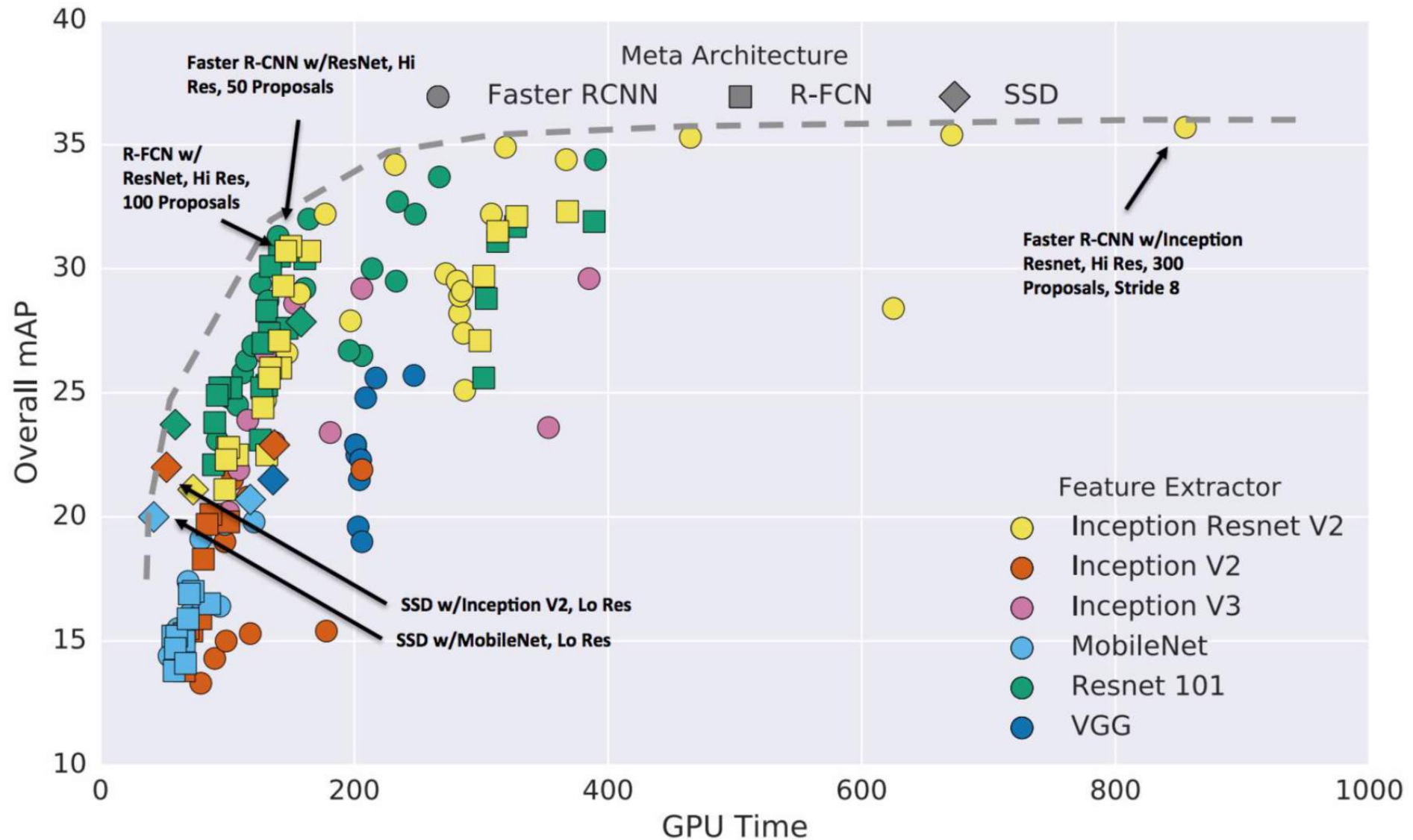
# Object Detection: YOLO



- Делим изображение на блоки
- Для каждого блока предсказывается распределение классов + координаты  
В штук bbox с уверенностями
- Фильтруем и находим основные

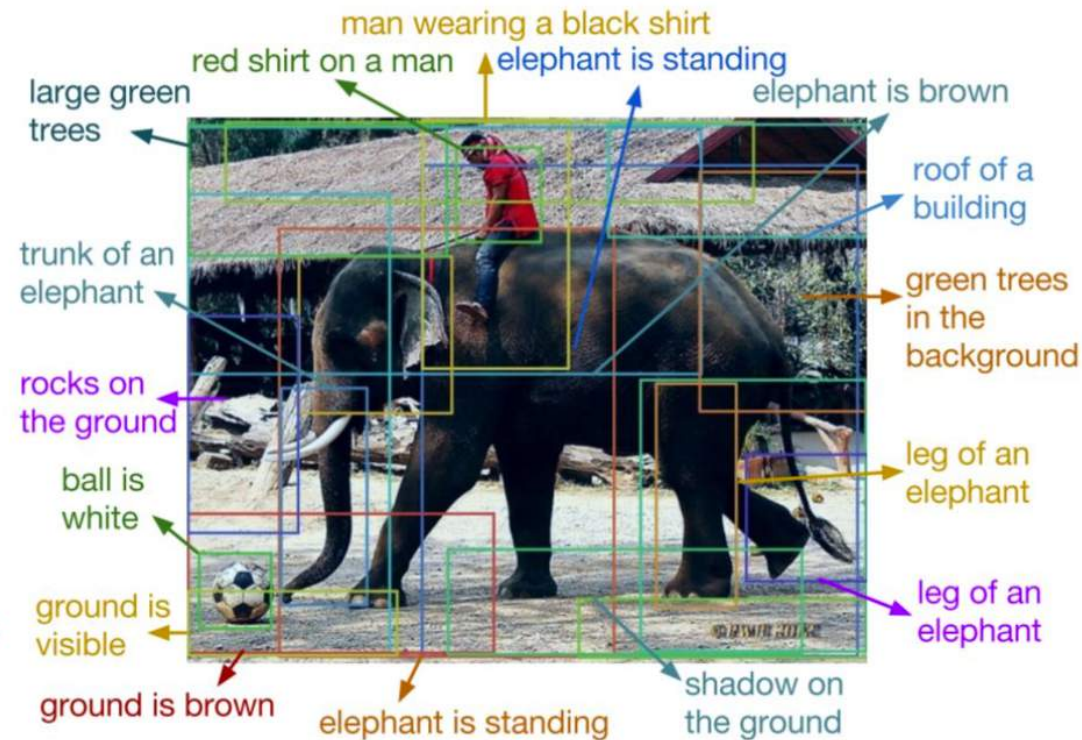
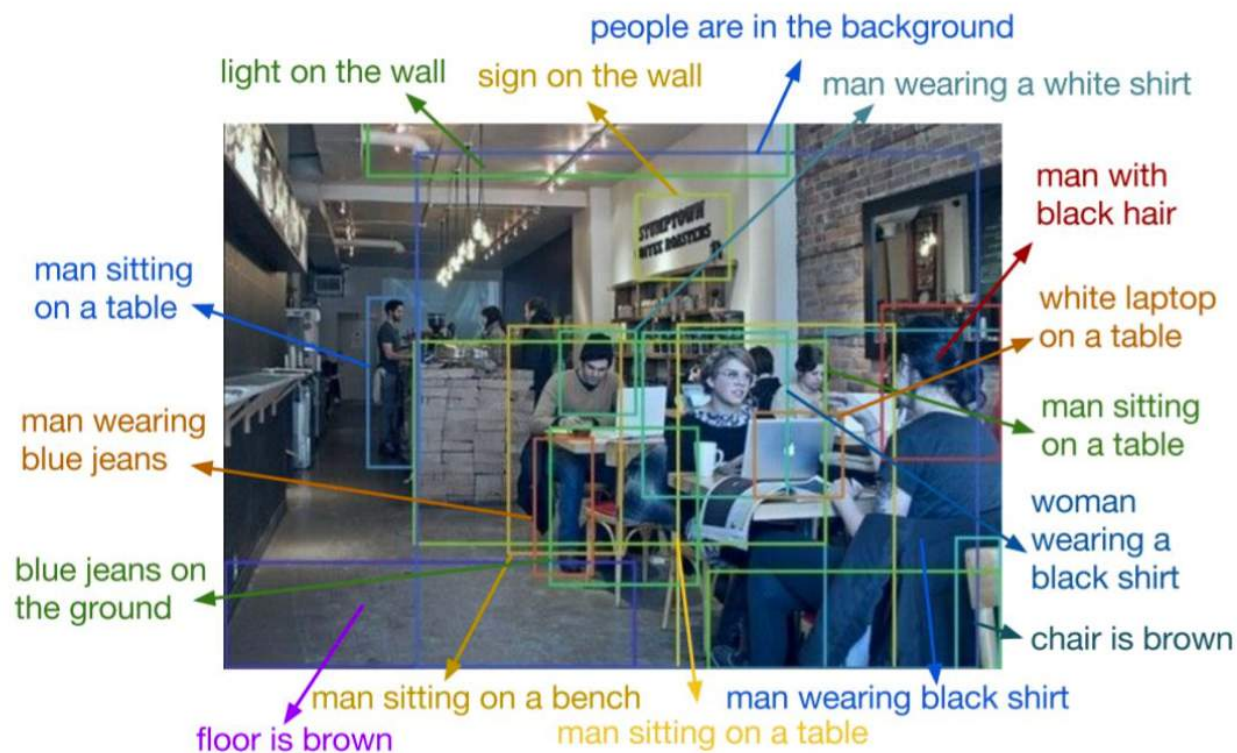
Главное преимущество в том, что это быстрее, чем Real-Time.  
Хотя и не так точно.

# Object Detection: Architectures





# Dense Captioning



# Instance Segmentation

---

## Задача:

Дать лейбл каждому пикселю на картинке. При этом отличая объекты одного класса друг от друга.

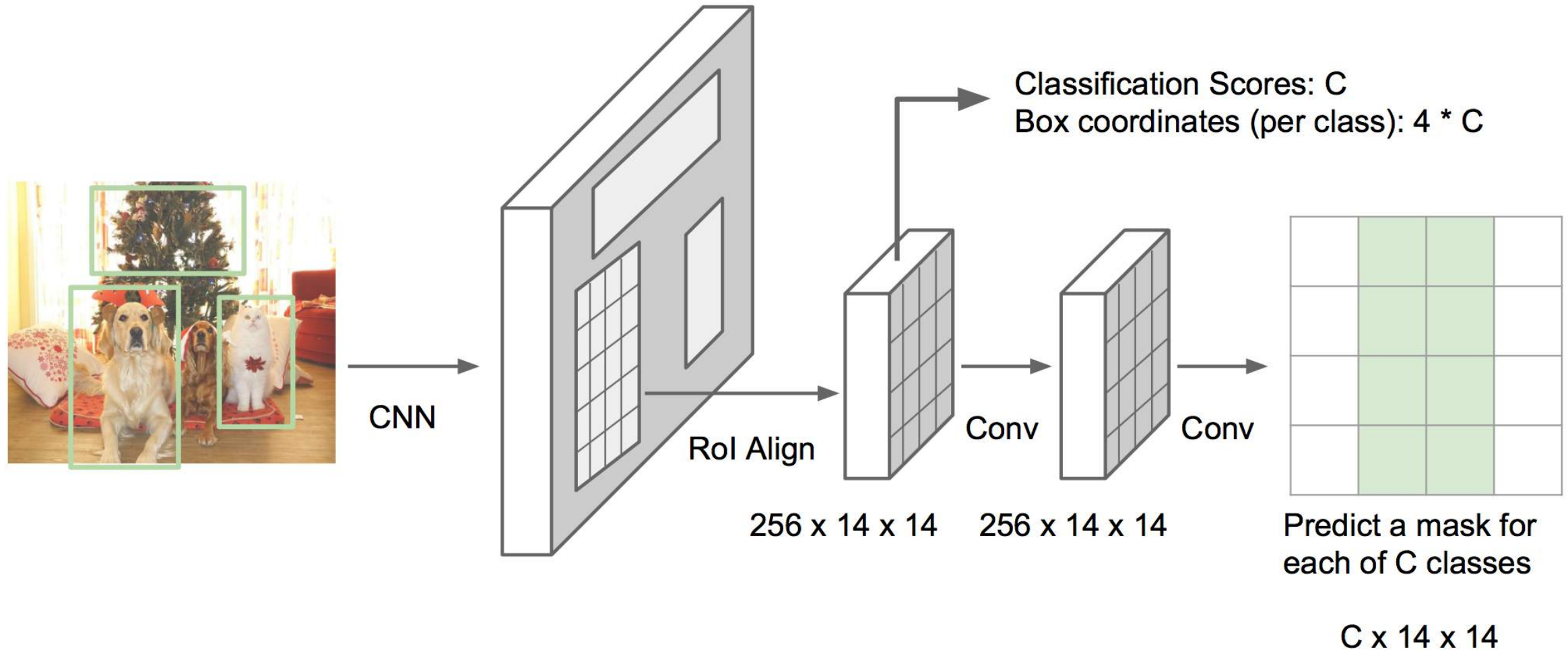
## Важная деталь:

Мы имеем дело с неизвестным количеством объектов.



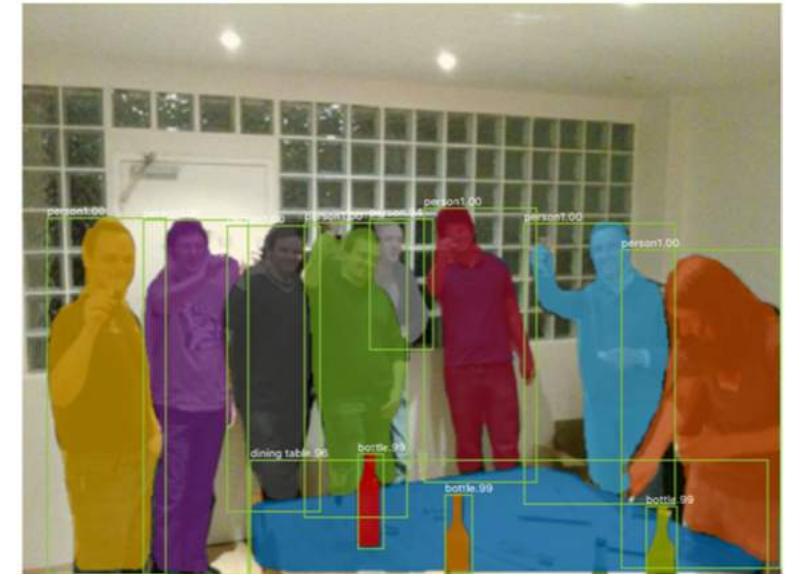
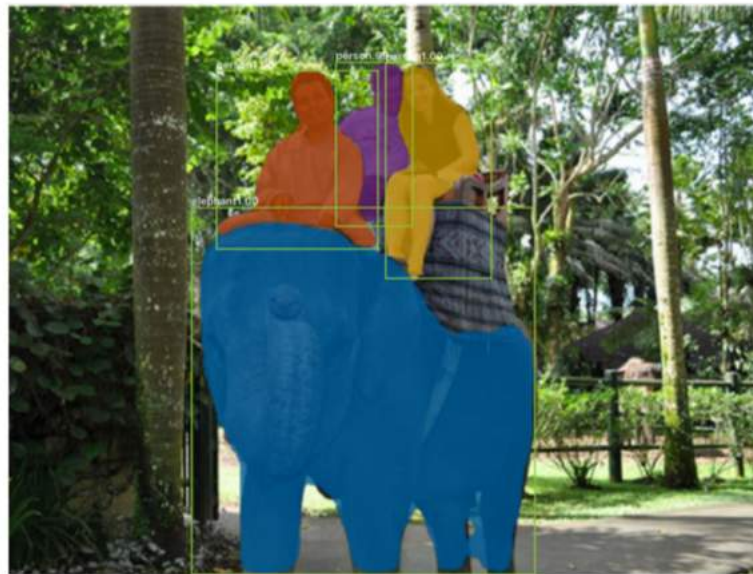
**DOG, DOG, CAT**

# Instance Segmentation: Mask R-CNN





# Instance Segmentation: Mask R-CNN



# Instance Segmentation: Mask R-CNN





# Recap

## Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

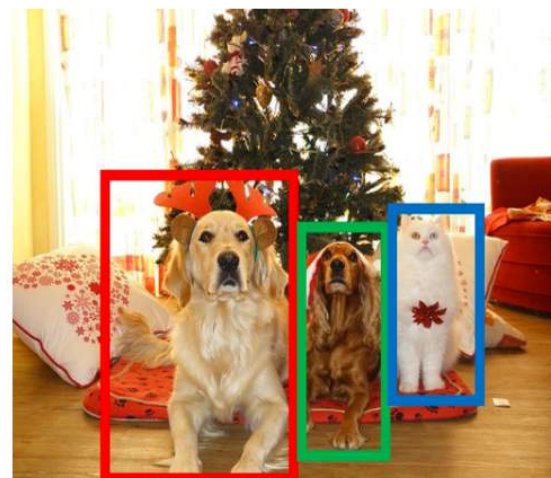
## Classification + Localization



CAT

Single Object

## Object Detection



DOG, DOG, CAT

Multiple Object

## Instance Segmentation



DOG, DOG, CAT

This image is CC0 public domain

# Useful Materials

---

## Лекции

- cs231n 2017 Lecture 11 video:  
<https://www.youtube.com/watch?v=nDPWywWRIRo&list=PL3FW7Lu3i5JvHM8ljYj-zLfQRF3EO8sYv>
- cs231n 2017 Lecture 11 Slides: [http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture11.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture11.pdf)

## Статьи

- Mask-RCNN: <https://arxiv.org/pdf/1703.06870.pdf>
- U-Net: <https://arxiv.org/pdf/1505.04597.pdf> 和 %5bTiramisu%5d(<https://arxiv.org/abs/1611.09326.pdf>)
- Faster R-CNN: <http://papers.nips.cc/paper/5638-faster-r-cnn-towards-real-time-object-detection-with-region-proposal-networks.pdf>
- YOLO-9000:  
[http://openaccess.thecvf.com/content\\_cvpr\\_2017/papers/Redmon\\_YOLO9000\\_Better\\_Faster\\_CVPR\\_2017\\_paper.pdf](http://openaccess.thecvf.com/content_cvpr_2017/papers/Redmon_YOLO9000_Better_Faster_CVPR_2017_paper.pdf)