# Object detection with R-CNN family

—

Computer vision

Comparison between image classification, object detection and instance segmentation.

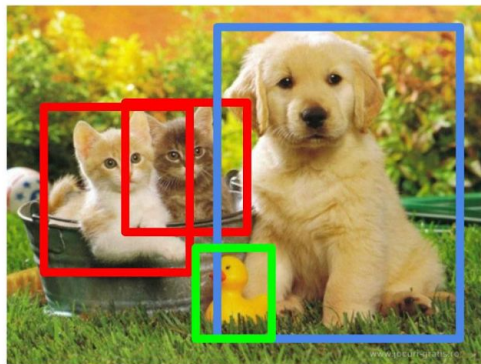**Classification** — CAT

**Classification + Localization** — CAT

**Object Detection** — CAT, DOG, DUCK

**Instance Segmentation** — CAT, DOG, DUCK

Single object
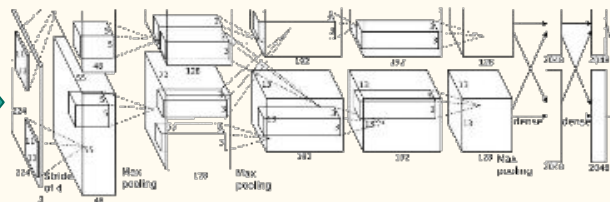
Multiple objects

# Classification



Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

**Vector:** 4096

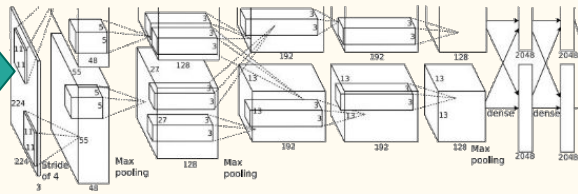**Fully-Connected:** 4096 to 1000

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

# Classification + Localization



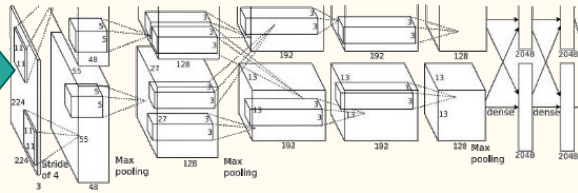**Fully Connected**: 4096 to 1000

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

**Vector:** 4096

**Fully Connected**: 4096 to 4

**Box Coordinates**
(x, y, w, h)

# Classification + Localization



**Correct label:** Cat

**Fully Connected:** 4096 to 1000

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

**Softmax Loss**

**Vector:** 4096

**Fully Connected:** 4096 to 4

**Box Coordinates** (x, y, w, h)

**L2 Loss**

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

# Classification + Localization



**Fully Connected:** 4096 to 1000

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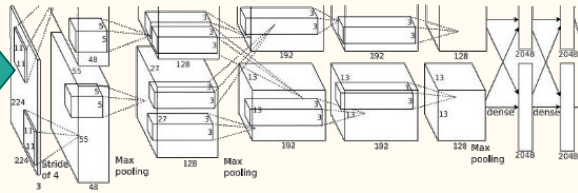...

**Correct label:** Cat

**Softmax Loss**

**Multitask Loss**

**Vector:** 4096

**Fully Connected:** 4096 to 4

**Box Coordinates** (x, y, w, h)

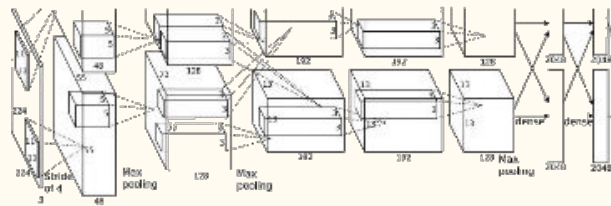**L2 Loss**

**Correct box:** (x', y', w', h')
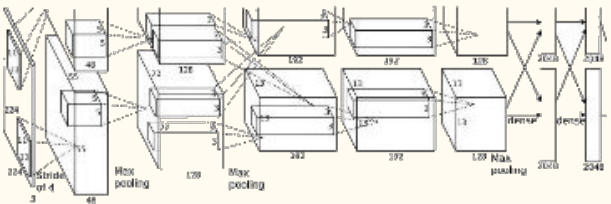
**+** → **Loss**

# Aside: 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
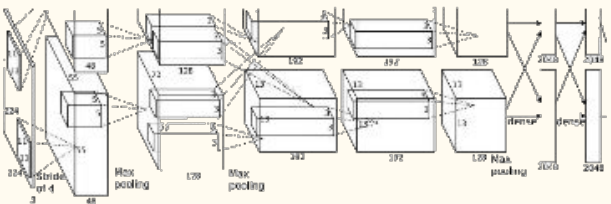Left/right hand
Neck
Head top

# Object detection as regression?
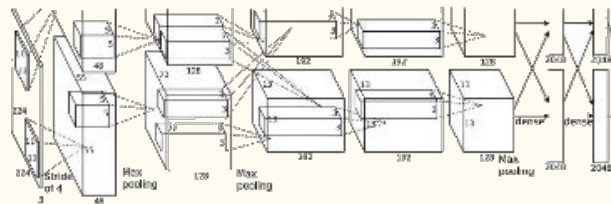


CAT: (x, y, w, h)

DOG: (x, y, w, h)
DUCK: (x, y, w, h)

DOG: (x, y, w, h)
CAT: (x, y, w, h)
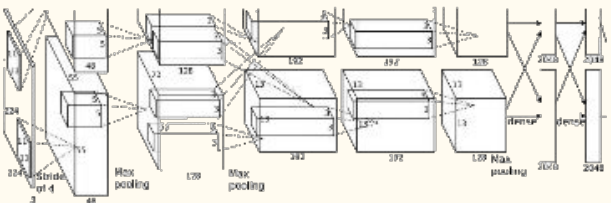DUCK: (x, y, w, h)
DUCK: (x, y, w, h)
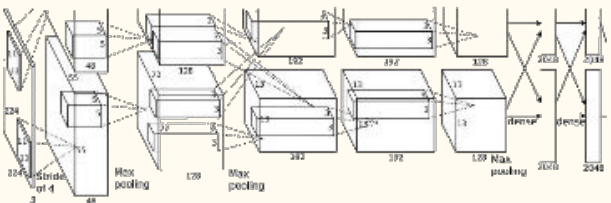...

# Object detection as regression?

CAT: (x, y, w, h)

**4 numbers**

DOG: (x, y, w, h)
DUCK: (x, y, w, h)

**8 numbers**

DOG: (x, y, w, h)
CAT: (x, y, w, h)
DUCK: (x, y, w, h)
DUCK: (x, y, w, h)
...

**16+ numbers**

# Object detection as classification

## Sliding window

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

Cat: NO
Dog: NO
Background: YES

Cat: YES
Dog: NO
Background: NO

Cat: NO
Dog: YES
Background: NO
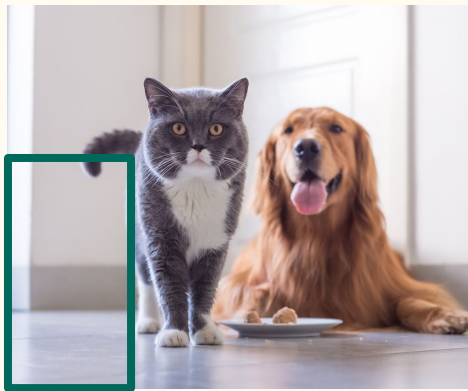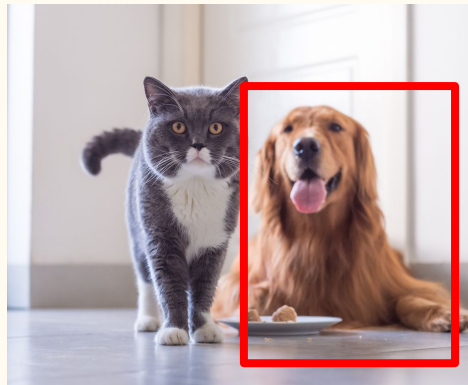
# Object detection as classification
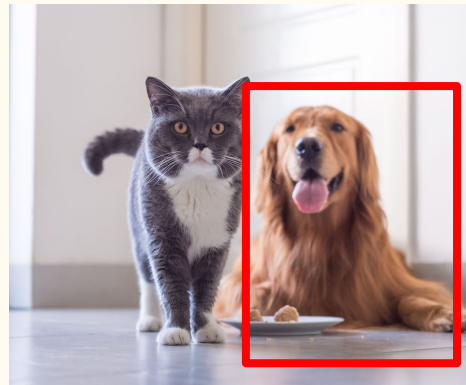
Sliding window

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

But is it good?

Cat: NO
Dog: NO
Background: YES

Cat: YES
Dog: NO
Background: NO
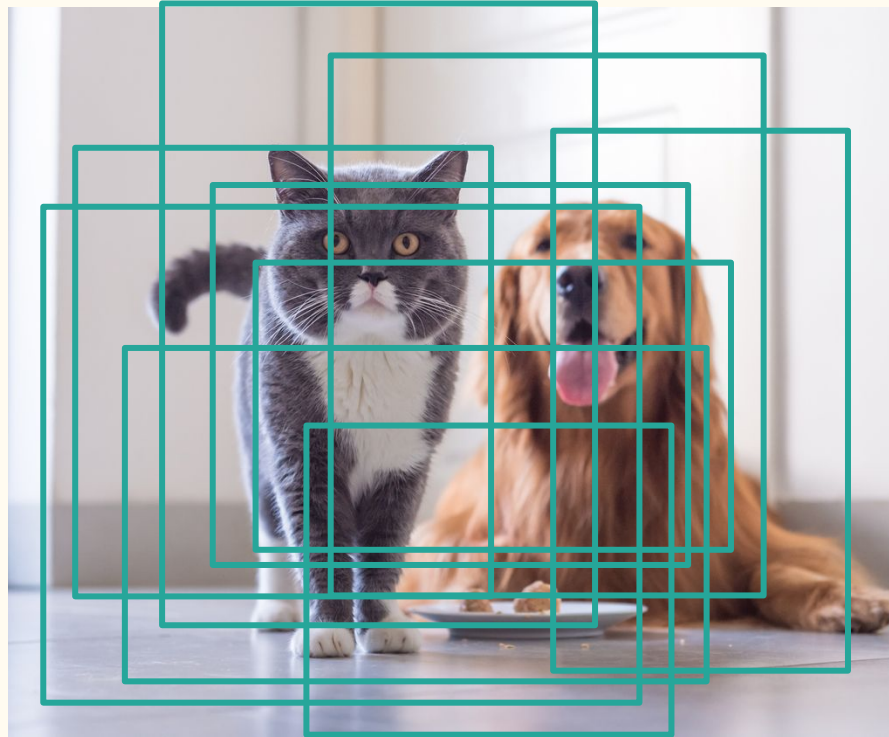
Cat: NO
Dog: YES
Background: NO

# Object detection as classification

Sliding window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background.
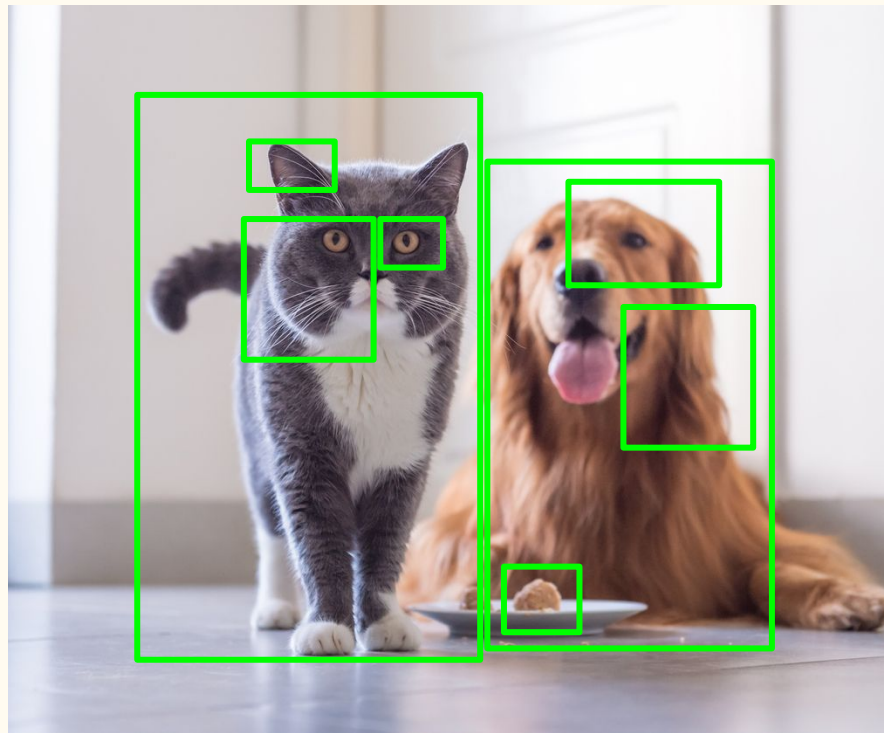
But is it good?

NO! Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive

# Region Proposals / Selective Search

- Find "blobby" image regions that are likely to contain objects

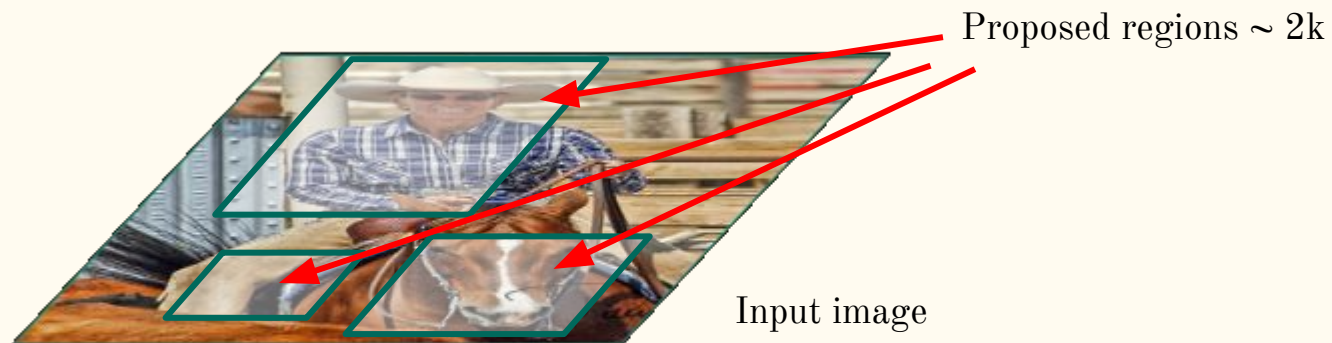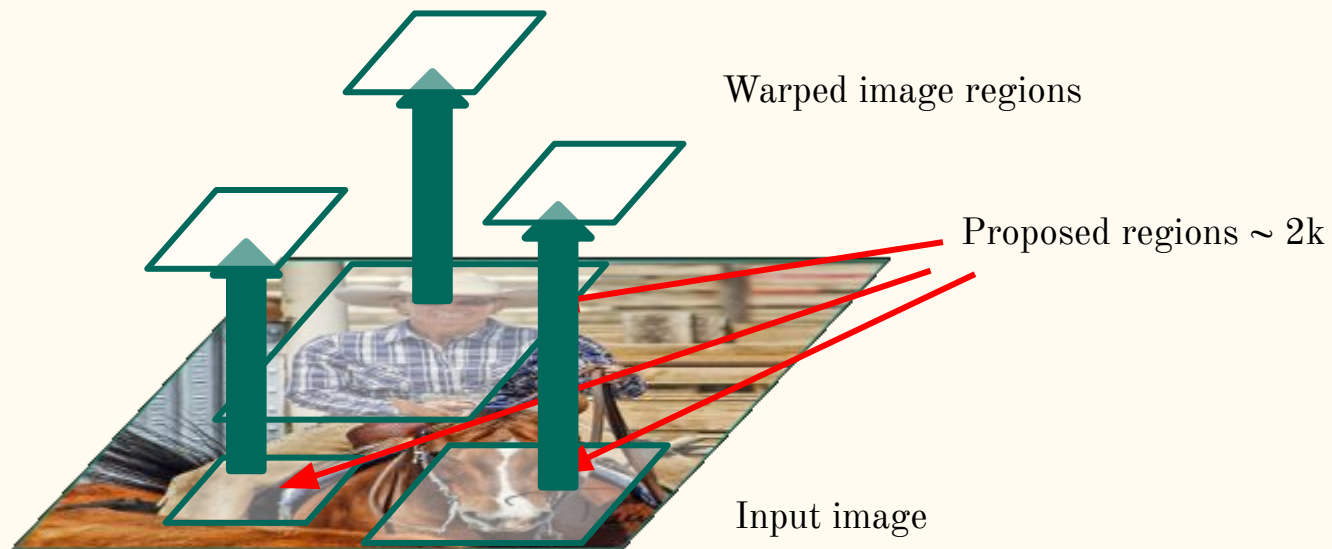- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU

# R-CNN
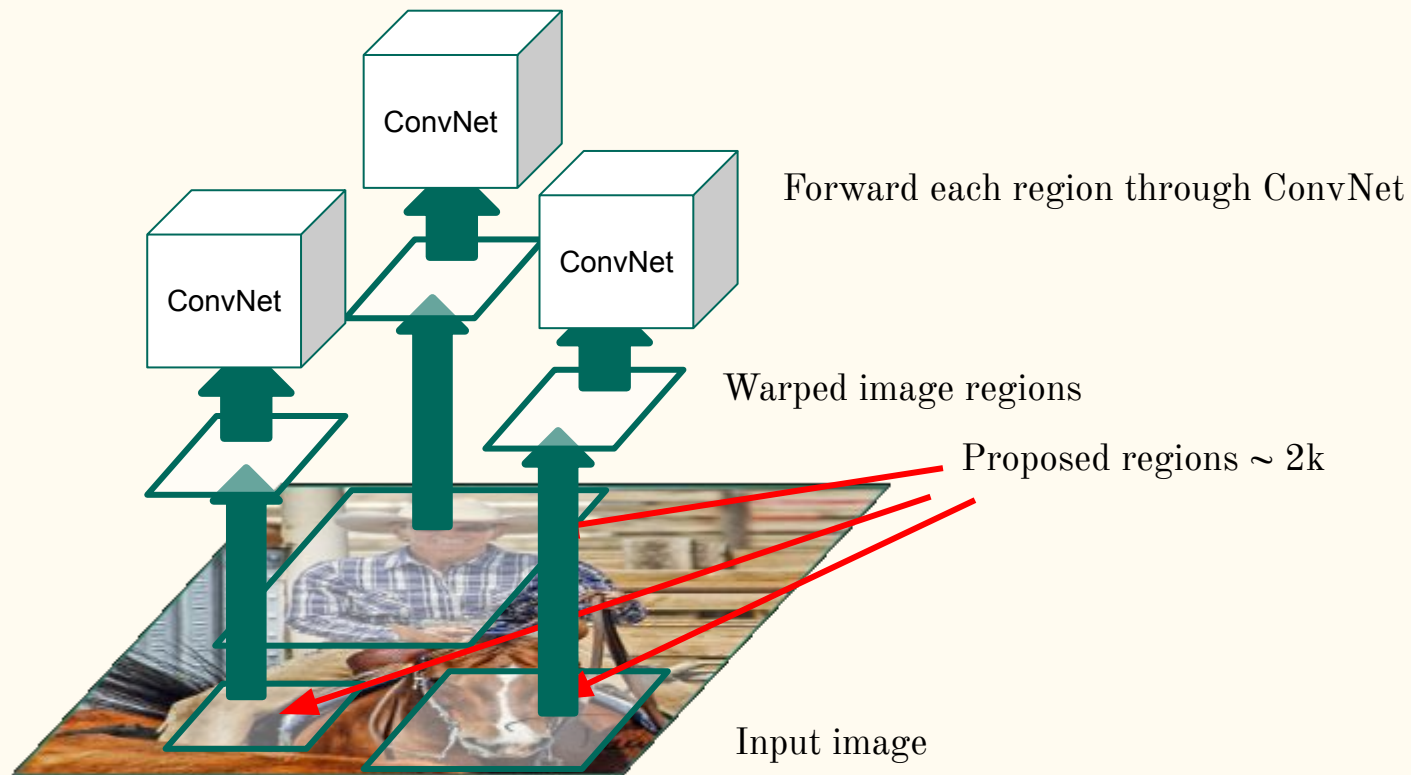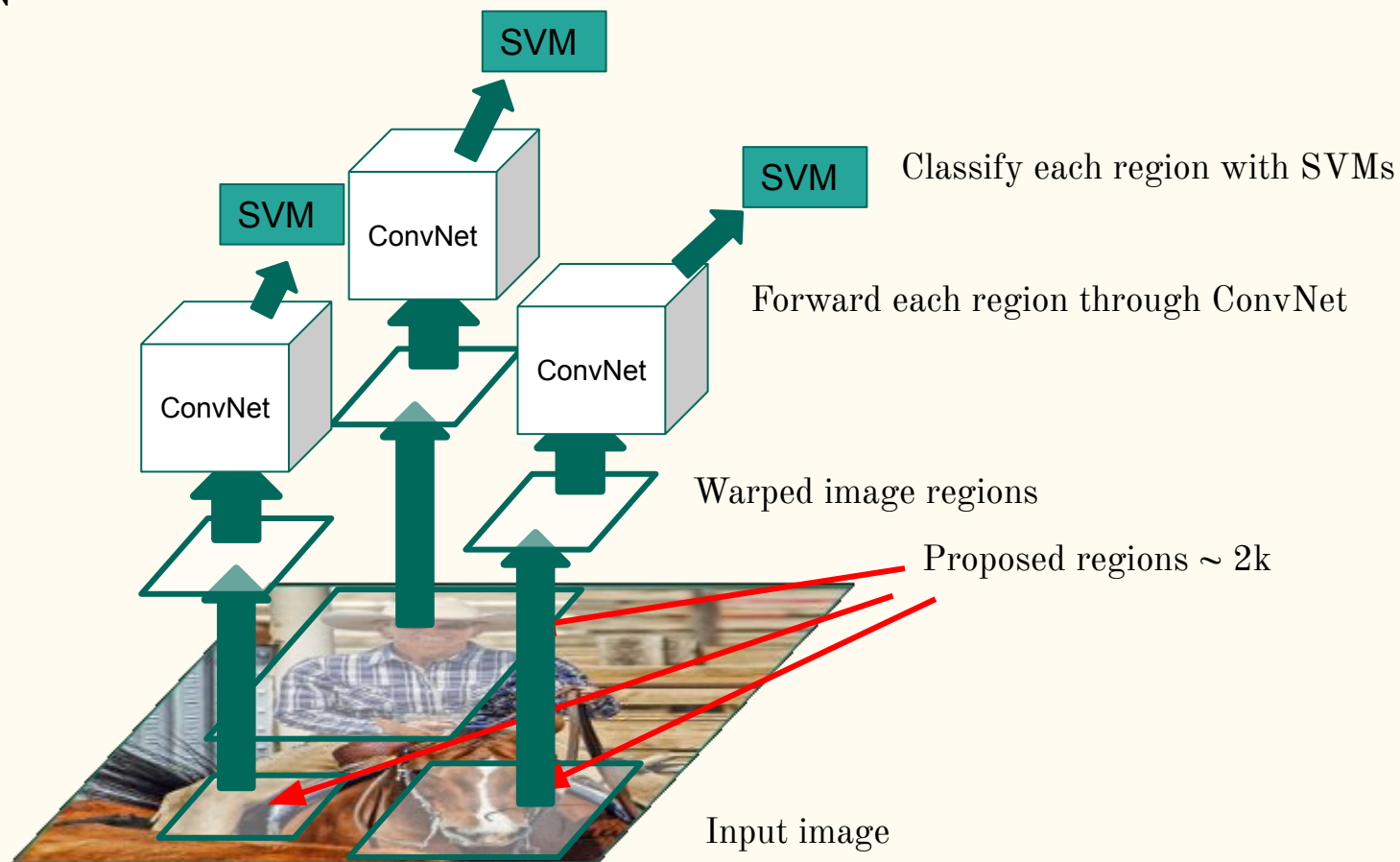
# R-CNN



Input image

# R-CNN



Proposed regions ~ 2k

Input image

# R-CNN



Warped image regions

Proposed regions ~ 2k

Input image

# R-CNN



Forward each region through ConvNet

Warped image regions

Proposed regions ~ 2k

Input image

# R-CNN



SVM

SVM

ConvNet

SVM          Classify each region with SVMs

ConvNet      Forward each region through ConvNet

ConvNet

Warped image regions

Proposed regions ~ 2k

Input image

# R-CNN



Linear regression for bounding box offsets

Classify each region with SVMs

Forward each region through ConvNet

Warped image regions

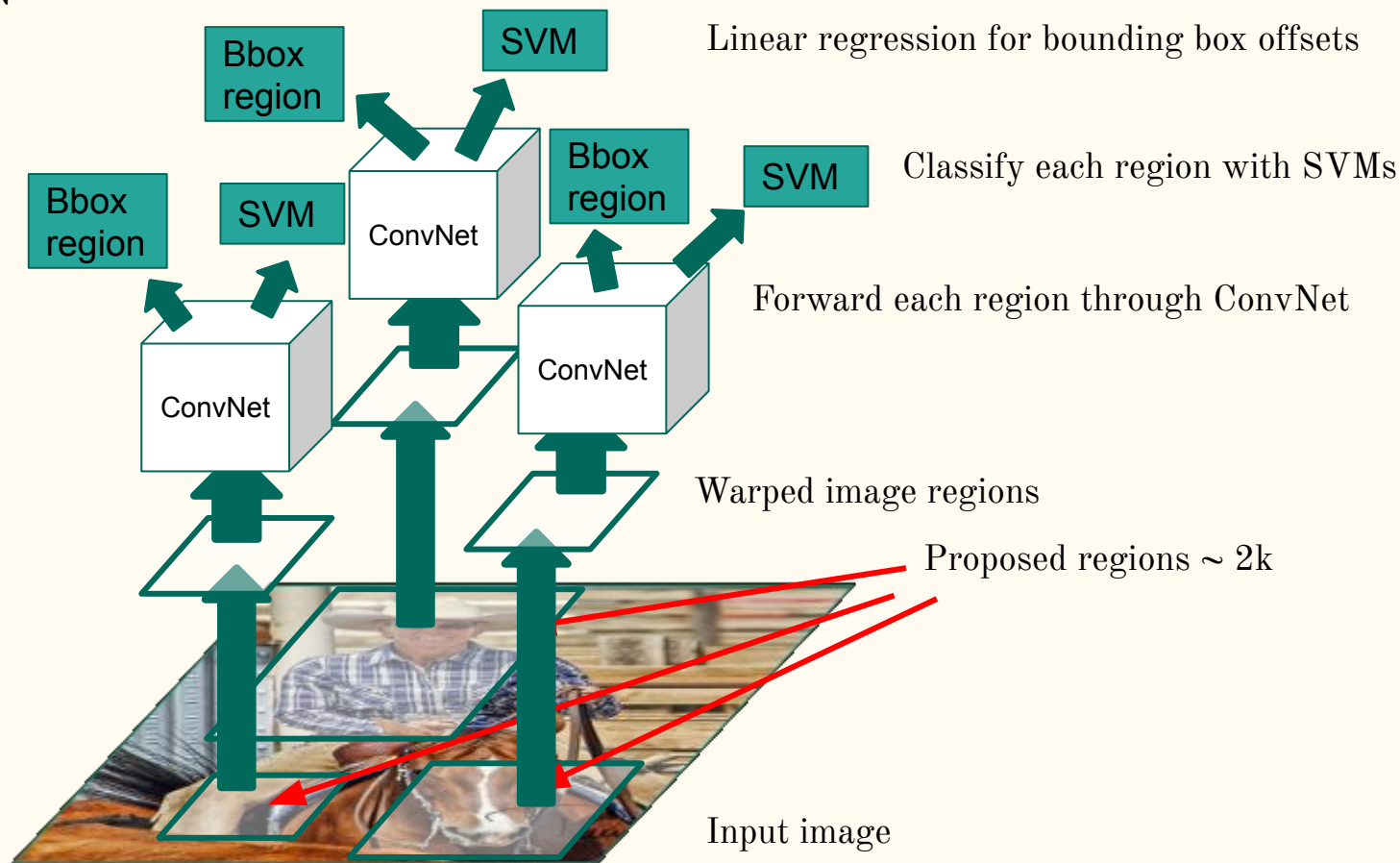Proposed regions ~ 2k
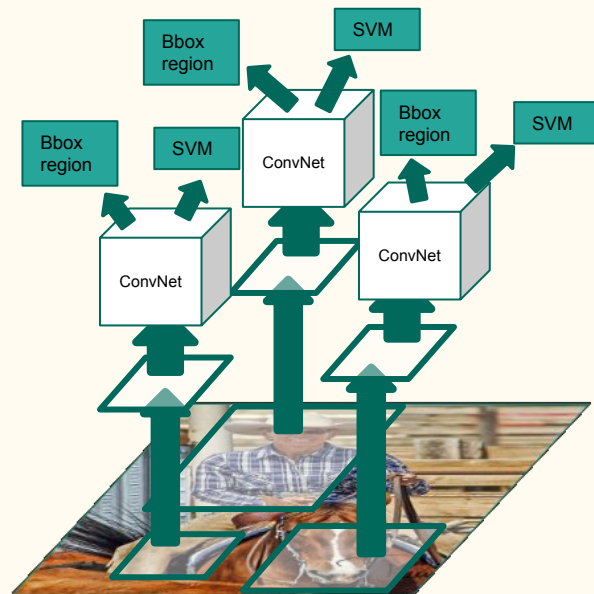
Input image

Bbox region

SVM

ConvNet

Input image

# R-CNN problems

- Ad hoc training objectives
  - Fine-tune network with softmax classifier (log loss)
  - Train post-hoc linear SVMs (hinge loss)
  - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
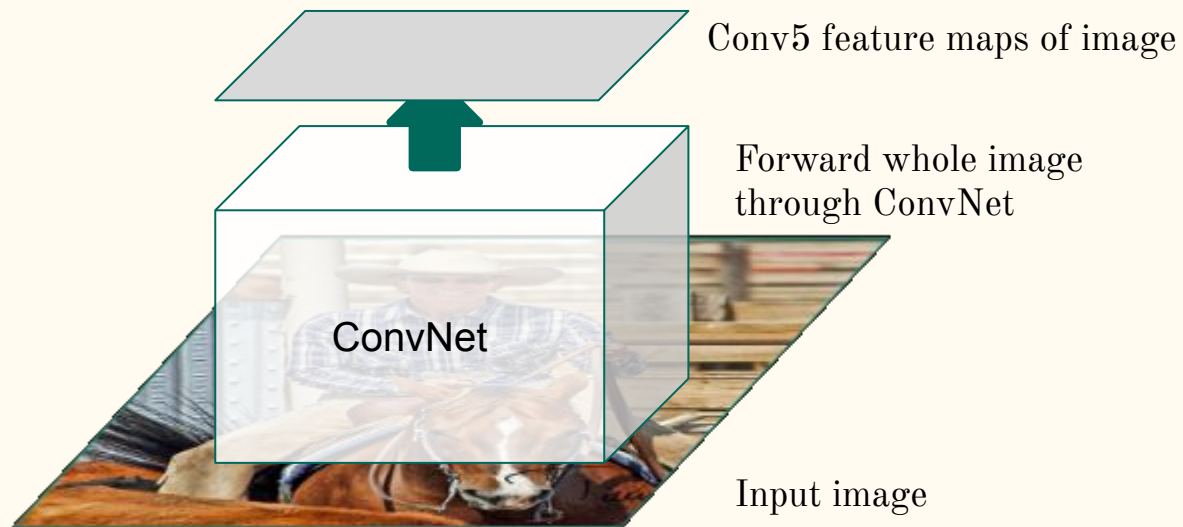  - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
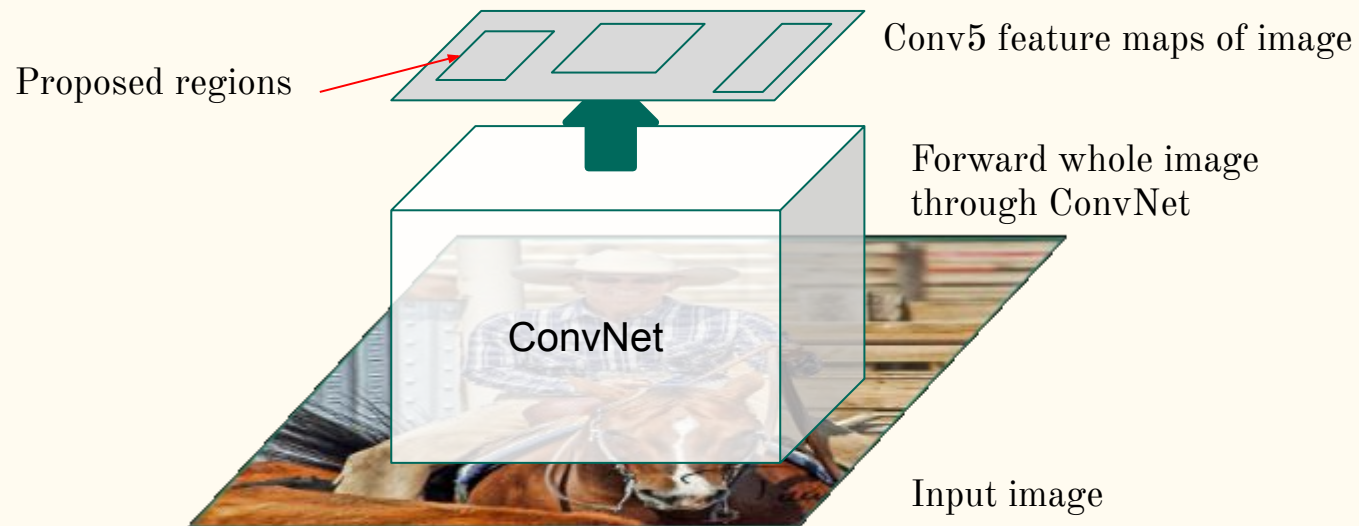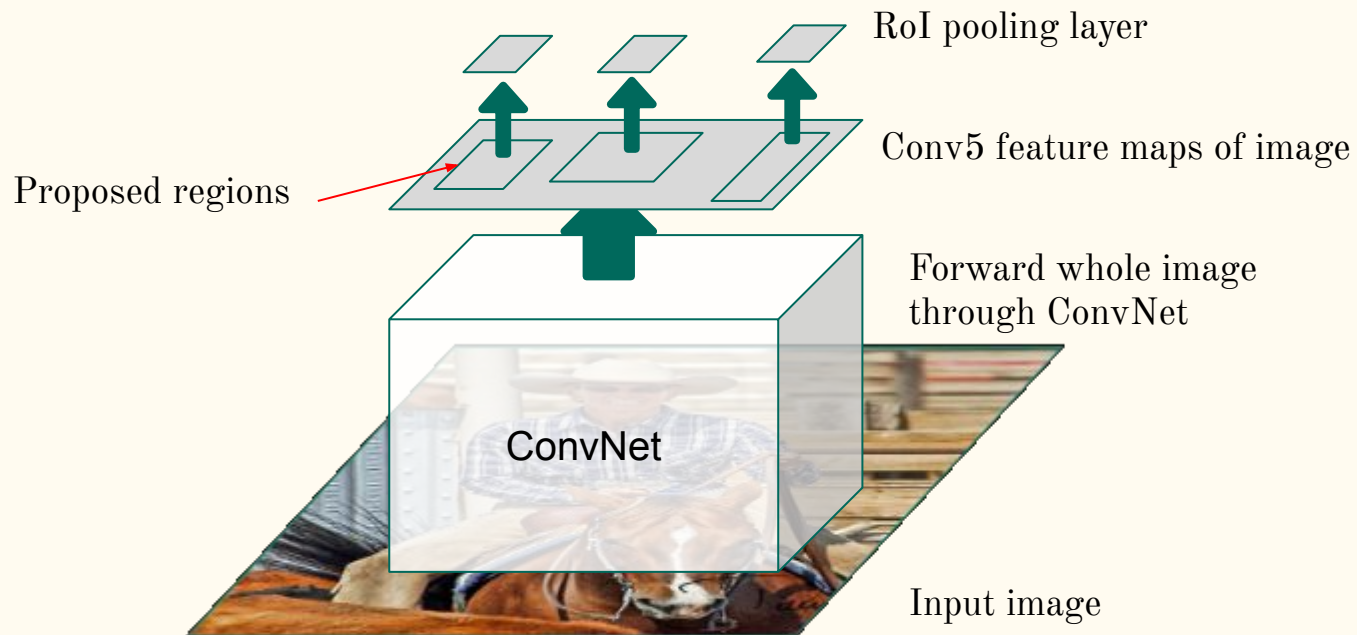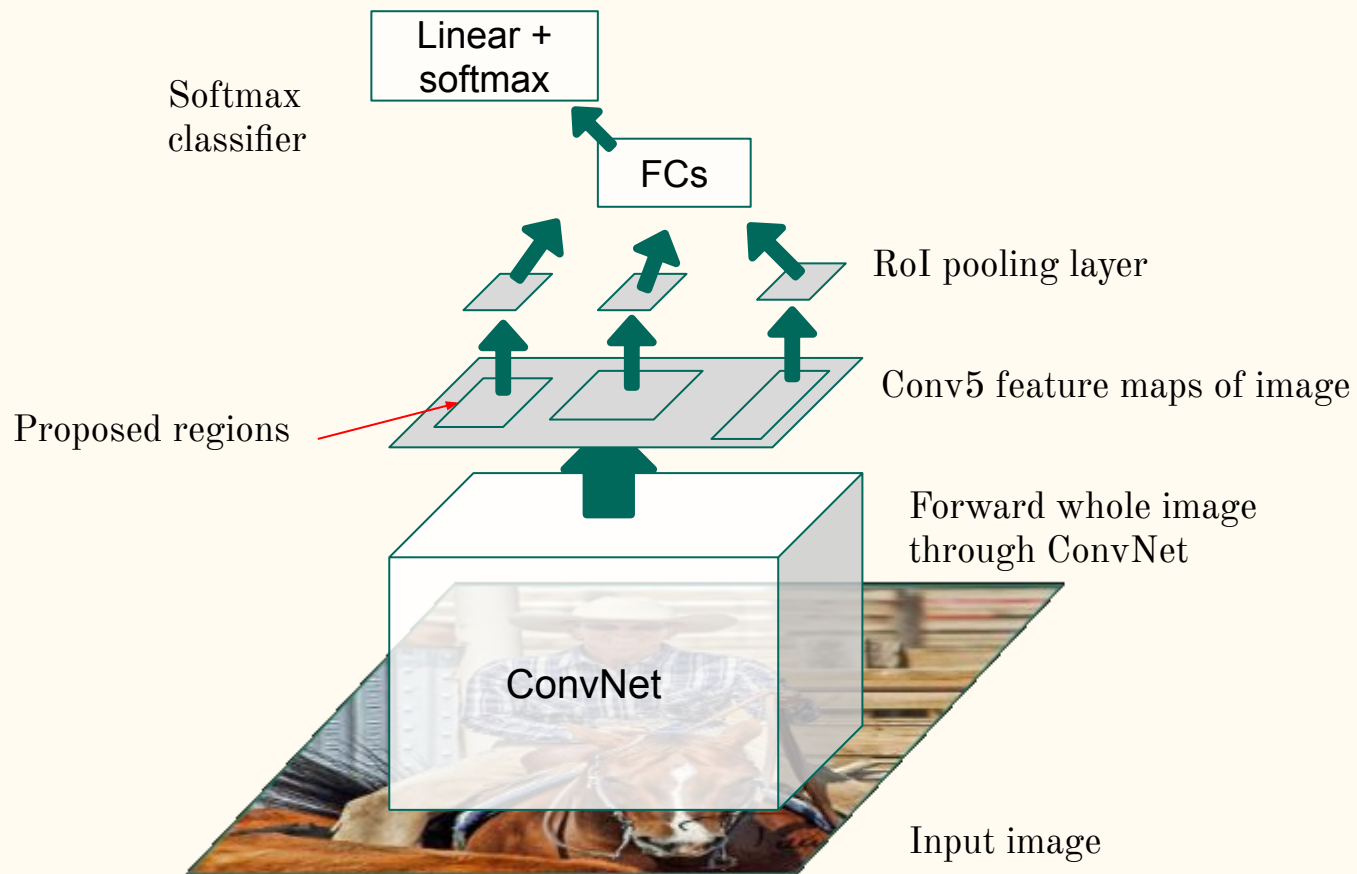
# Fast R-CNN



Input image

# Fast R-CNN



Conv5 feature maps of image

Forward whole image
through ConvNet

ConvNet

Input image

# Fast R-CNN



Conv5 feature maps of image

Proposed regions

Forward whole image
through ConvNet

ConvNet

Input image

# Fast R-CNN



RoI pooling layer

Conv5 feature maps of image

Proposed regions

Forward whole image
through ConvNet

ConvNet

Input image

# Fast R-CNN



Softmax classifier

Linear + softmax

FCs

RoI pooling layer

Proposed regions

Conv5 feature maps of image

ConvNet

Forward whole image through ConvNet

Input image

Fast R-CNN

# Fast R-CNN



Softmax classifier

Linear + softmax

Linear

Bounding-box regressors

FCs

RoI pooling layer

Proposed regions

Conv5 feature maps of image

ConvNet

Forward whole image through ConvNet

Input image

# R-CNN VS SPP-Net VS Fast R-CNN

**Training time (Hours)**

- R-CNN: 84
- SPP-Net: 25.5
- Fast R-CNN: 8.75

**Test time (seconds)**

Including Region propos... / Excluding Region Propo...

- R-CNN: 49 / 47
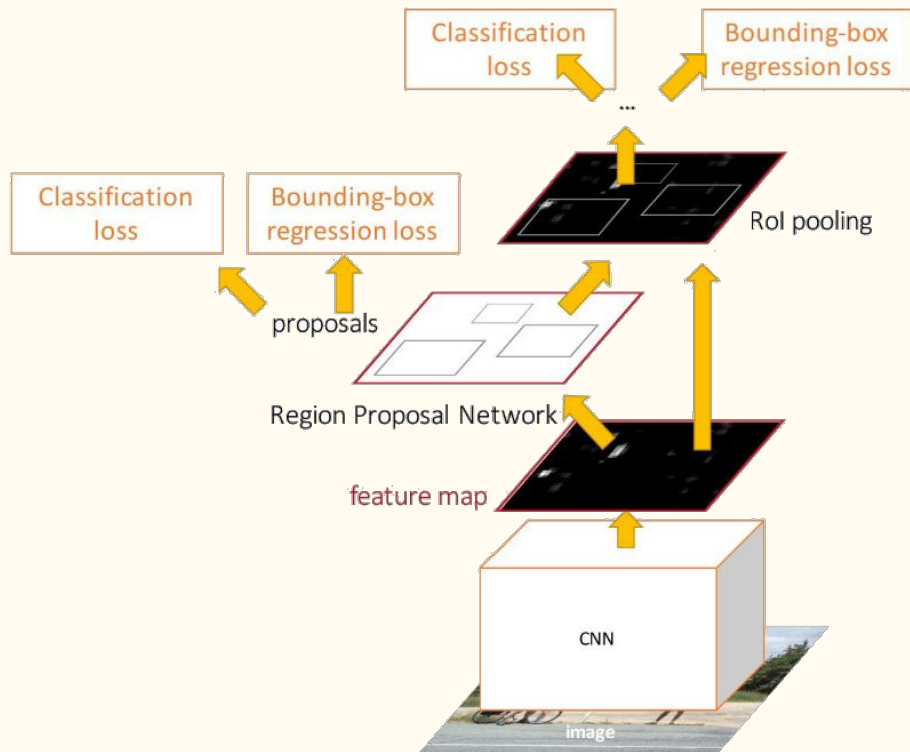- SPP-Net: 4.3 / 2.3
- Fast R-CNN: 2.3 / 0.32

# 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
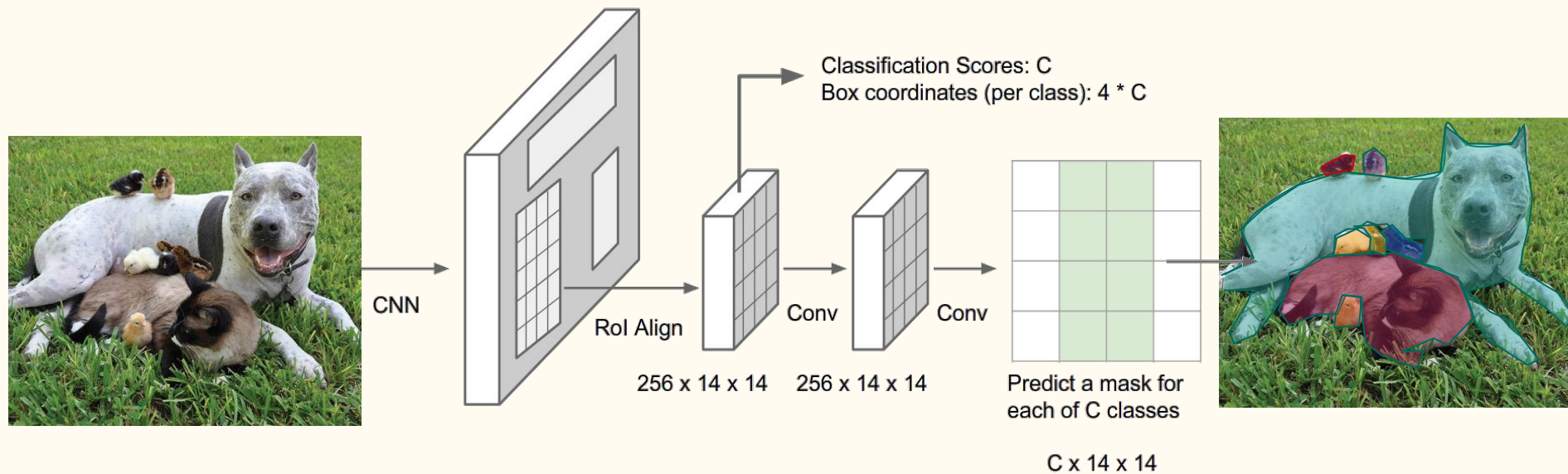3. Final classification score (object classes)
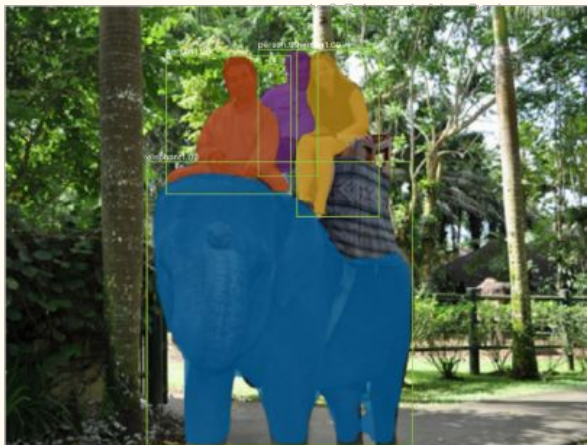4. Final box coordinates

# Faster R-CNN
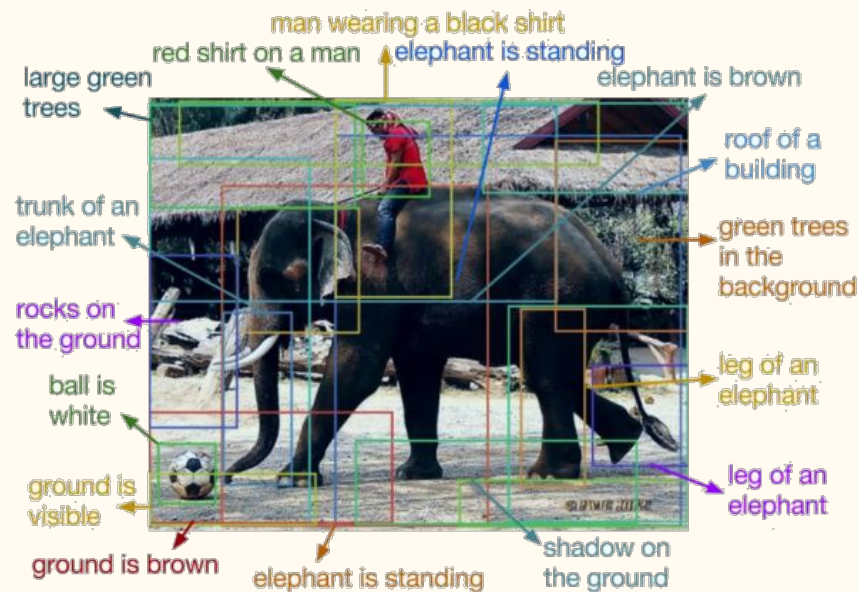
Make CNN do proposals!

~ 7 FPS

# Mask R-CNN



Classification Scores: C
Box coordinates (per class): 4 * C

CNN

RoI Align

256 x 14 x 14

Conv

256 x 14 x 14

Conv

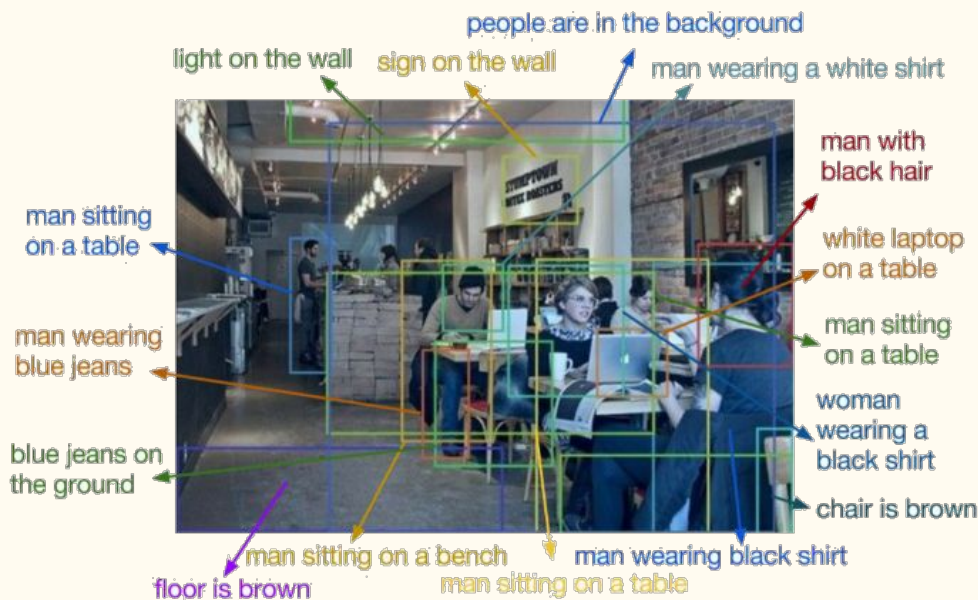Predict a mask for
each of C classes
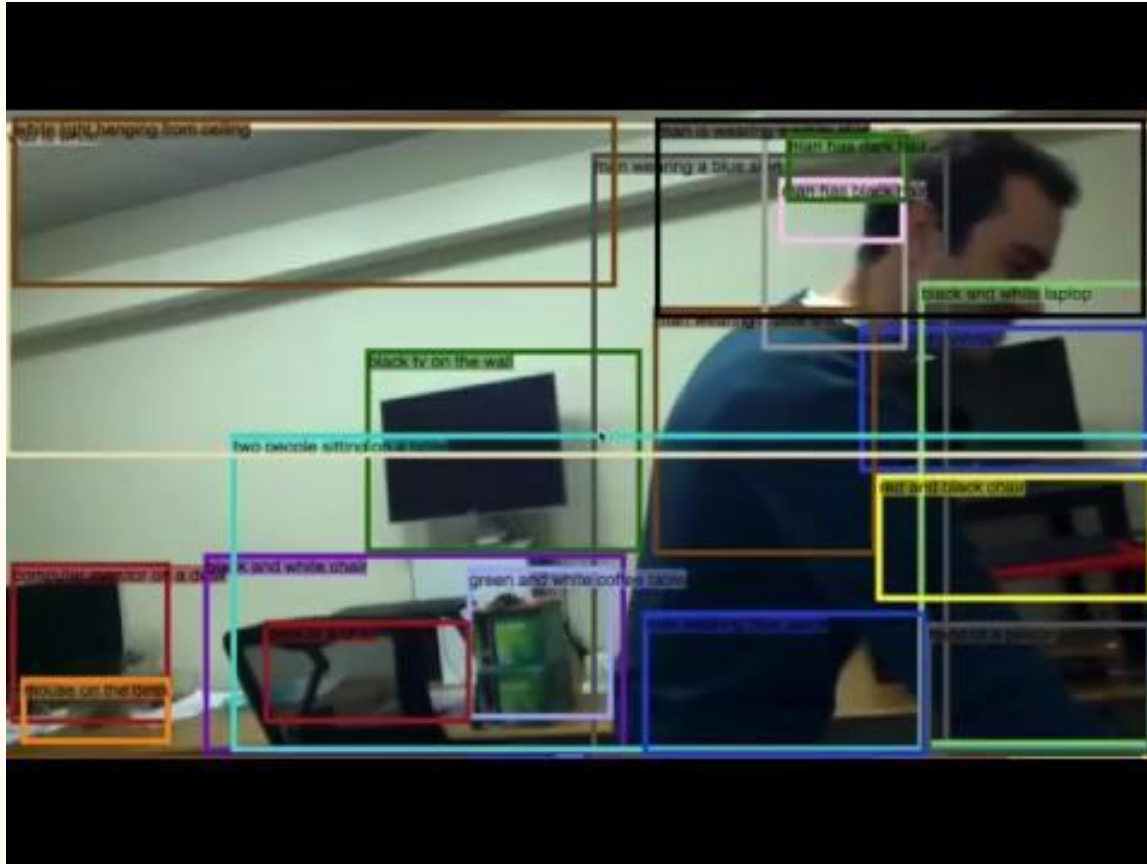
C x 14 x 14

# Mask R-CNN: Very good results

# Aside: Object Detection + Captioning= Dense Captioning

# Aside: Object Detection + Captioning= Dense Captioning

# T.HANKS!