

Artificial Neural Networks

人工神经网络

权小军教授
中山大学计算机学院

quanjx3@mail.sysu.edu.cn

2023 年 5 月 4 日

Lecture 8 - Semantic Segmentation and Object Detection 语义分割和目标检测

(*Part of the slides are adapted from Stanford CS231N*)

Image Classification

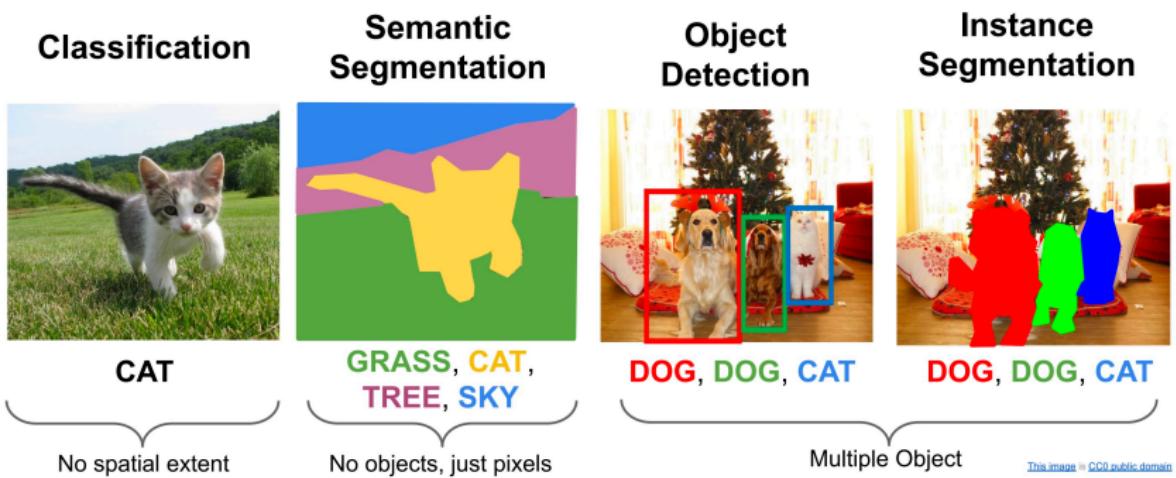
Classification: A core task in Computer Vision



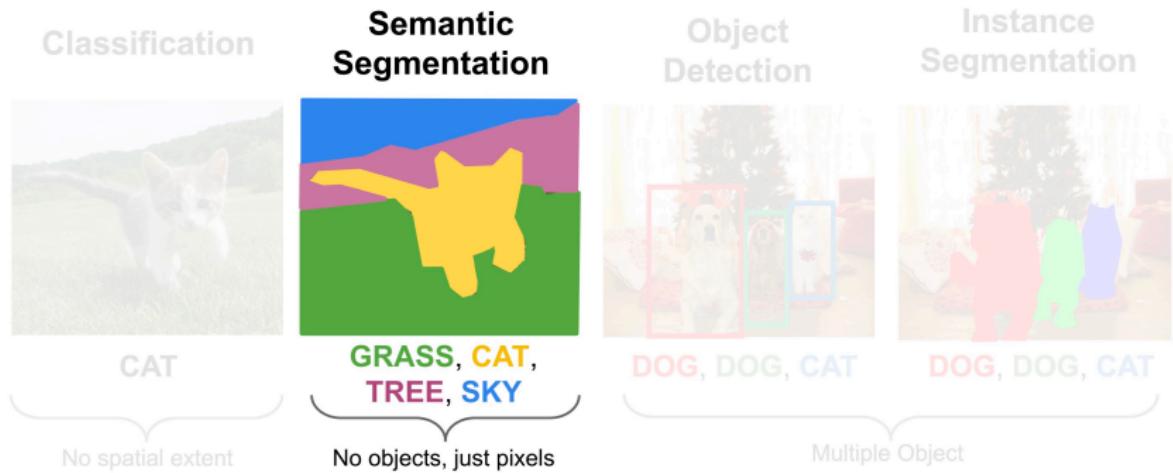
(assume given a set of possible labels)
{dog, cat, truck, plane, ...}



Computer Vision Tasks



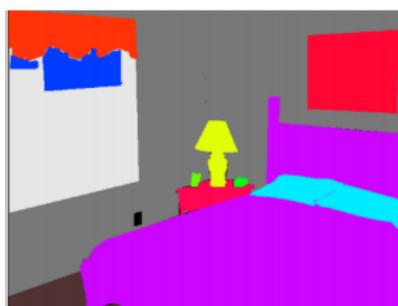
Semantic Segmentation



Lecture 8.1 Semantic Segmentation 语义分割

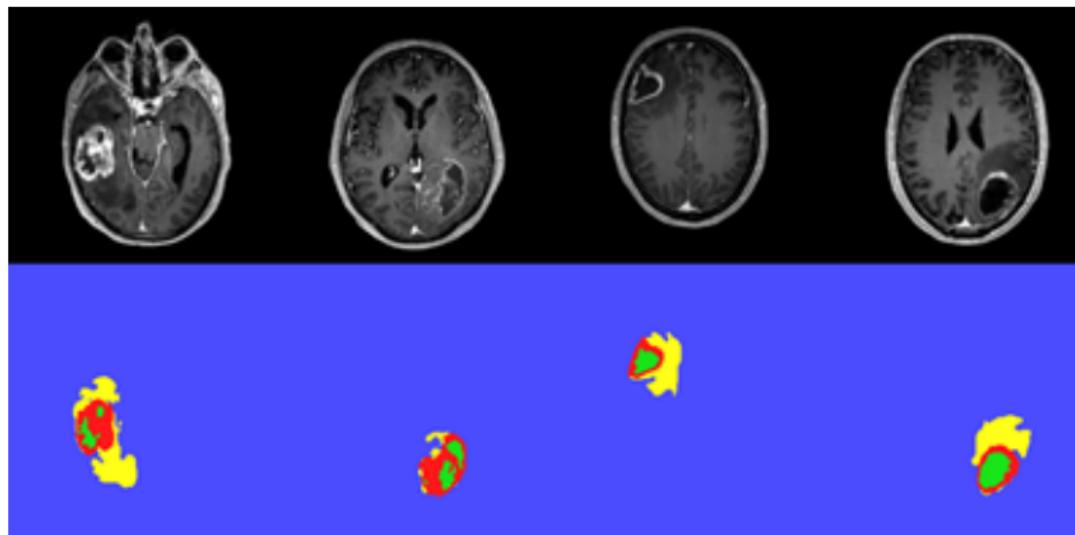
Semantic Segmentation

- Applications: self-driving, smart home



Semantic Segmentation

- More applications: intelligent medical analysis, etc.



Find any common characteristics from these segmentation tasks?

Semantic Segmentation: The Problem

- ▶ Semantic segmentation (语义分割): classifying each pixel



GRASS, CAT,
TREE, SKY, ...

Paired training data: for each training image,
each pixel is labeled with a semantic category.



At test time, classify each pixel of a new image.

CNNs for Semantic Segmentation

How are CNNs applied to semantic segmentation?

Semantic Segmentation Idea: Sliding Window

Full image



Semantic Segmentation Idea: Sliding Window

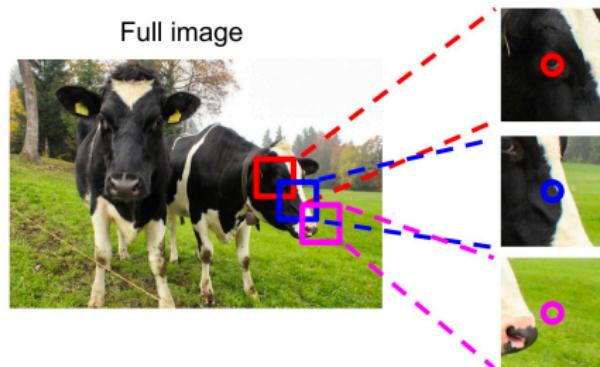
Full image



Impossible to classify without context

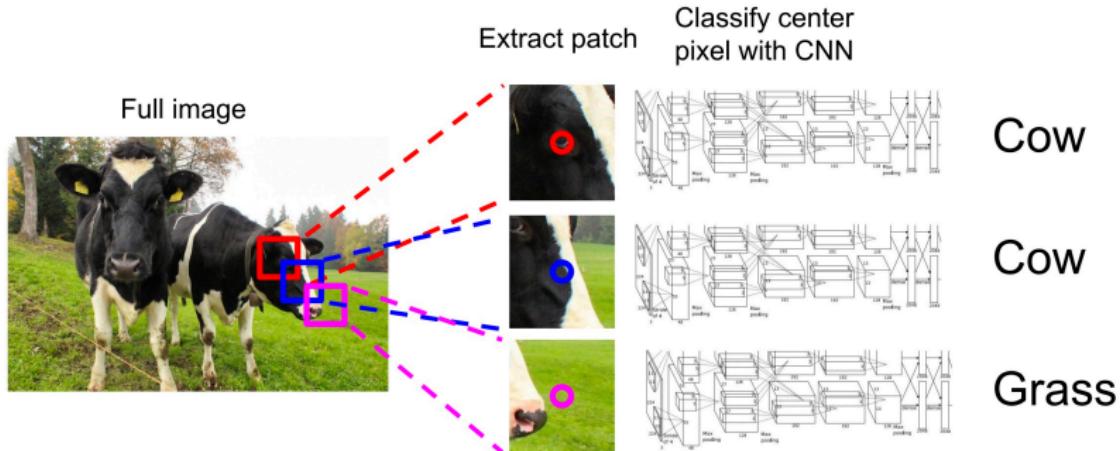
Q: how do we include context?

Semantic Segmentation Idea: Sliding Window

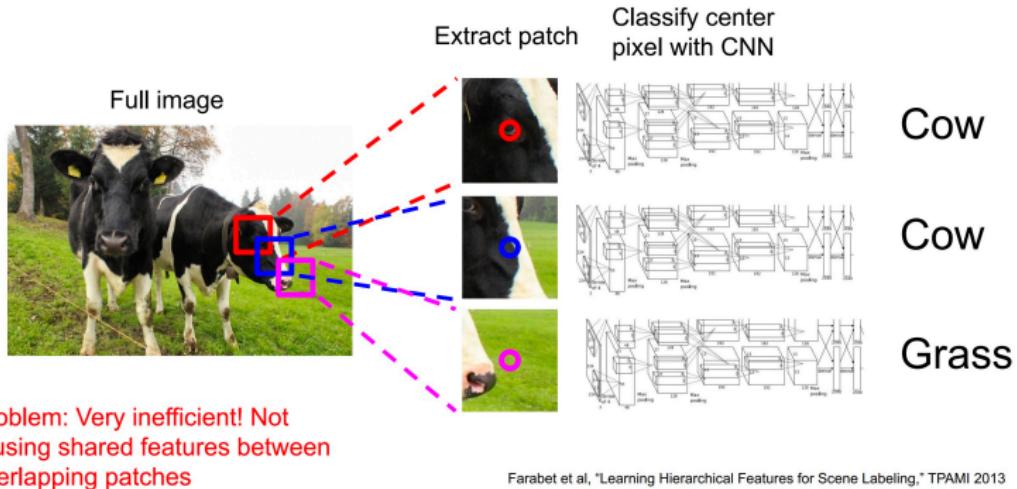


Q: how do we model this?

Semantic Segmentation Idea: Sliding Window

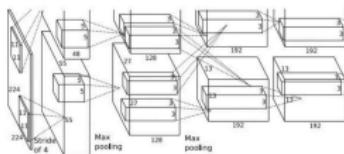


Semantic Segmentation Idea: Sliding Window



Semantic Segmentation Idea: Convolution

Full image

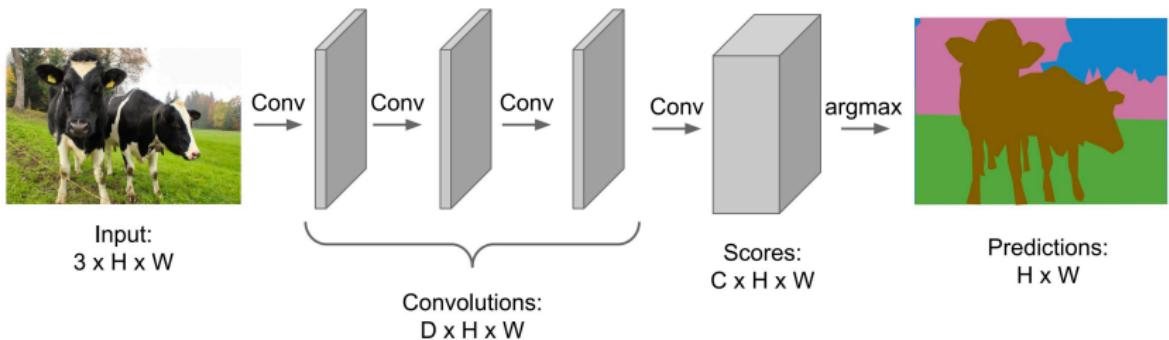


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

Problem: classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

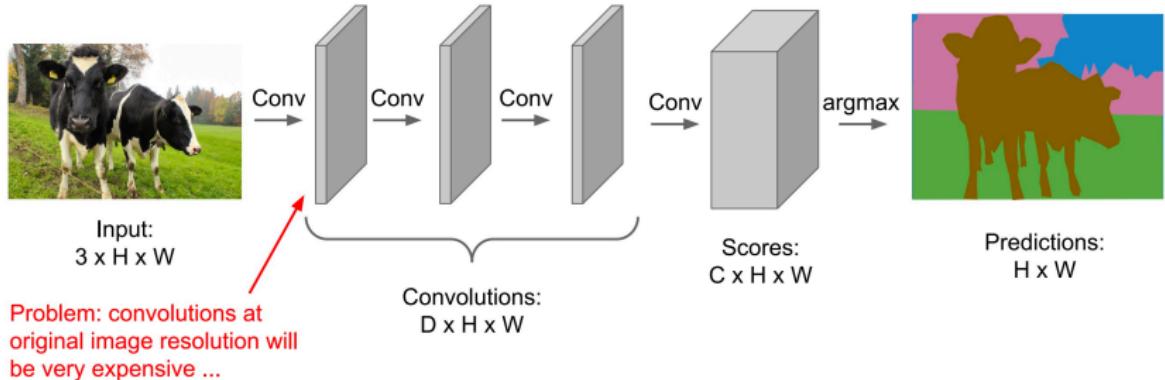
Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



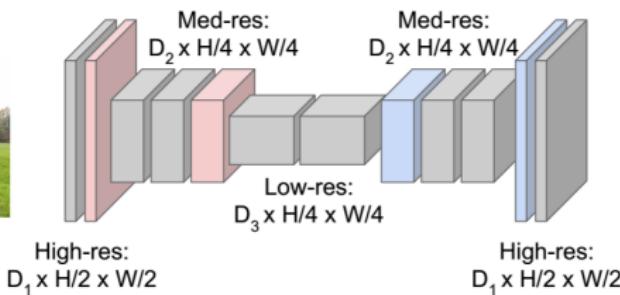
Fully Convolutional Network (FCN, 全卷积网络)

- ▶ Encoder-decoder framework: reduce size of feature maps, then increase to original size

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



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

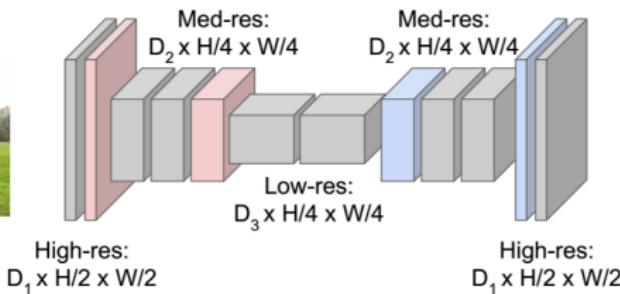
Fully Convolutional Network (FCN, 全卷积网络)

- ▶ Encoder-decoder framework: reduce size of feature maps, then increase to original size

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



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

- ▶ But how to realize “upsampling” (上采样)?

In-Network Upsampling: “Unpooling” (上池化)

Nearest Neighbor

1	2
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

Transposed convolution (反卷积): learnable upsampling

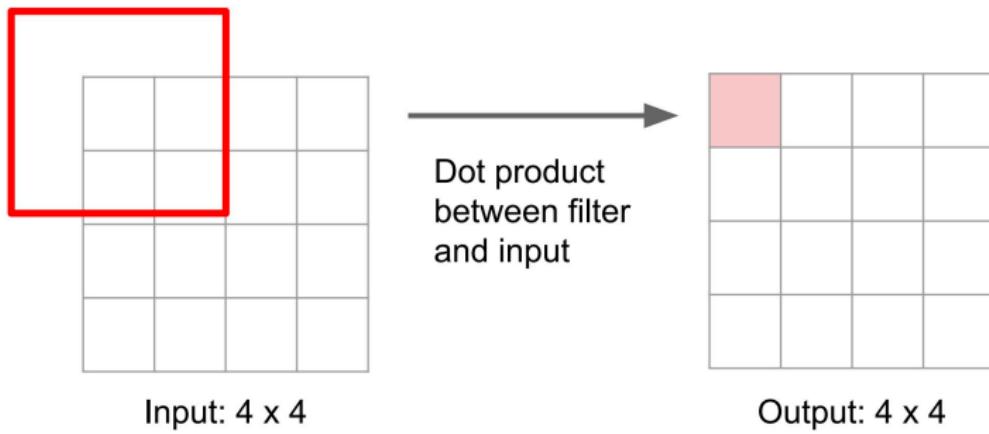
Recall: Normal 3×3 convolution, stride 1 pad 1

Input: 4×4

Output: 4×4

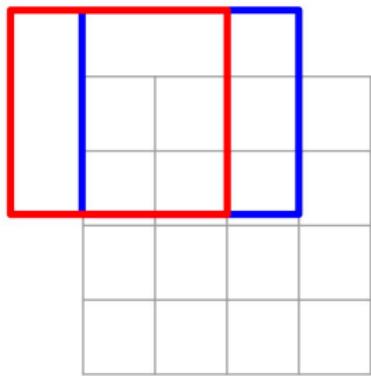
Transposed convolution (反卷积): learnable upsampling

Recall: Normal 3×3 convolution, stride 1 pad 1



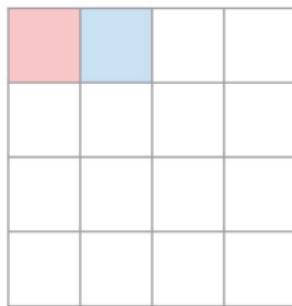
Transposed convolution (反卷积): learnable upsampling

Recall: Normal 3×3 convolution, stride 1 pad 1



Input: 4×4

Dot product
between filter
and input

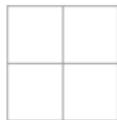


Output: 4×4

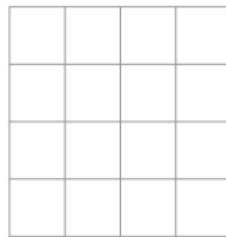
Transposed convolution (反卷积): learnable upsampling

Learnable Upsampling: Transposed Convolution

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



Input: 2 x 2

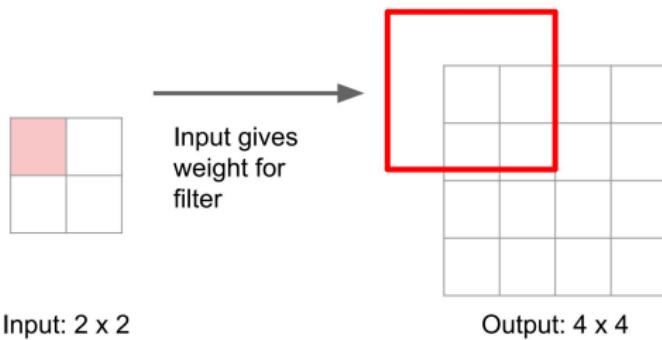


Output: 4 x 4

Transposed convolution (反卷积): learnable upsampling

Learnable Upsampling: Transposed Convolution

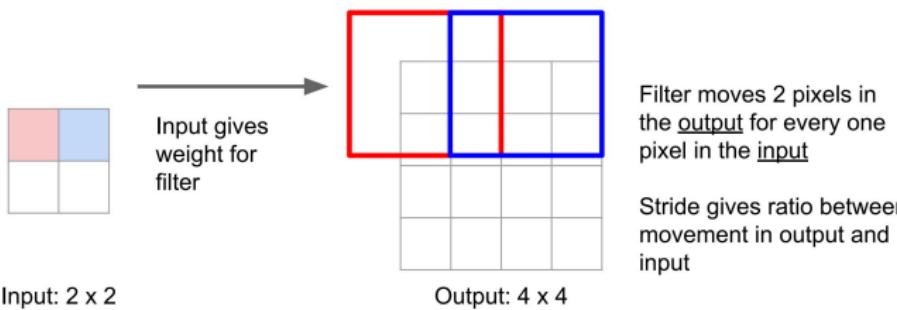
3 x 3 **transposed** convolution, stride 2 pad 1



Transposed convolution (反卷积): learnable upsampling

Learnable Upsampling: Transposed Convolution

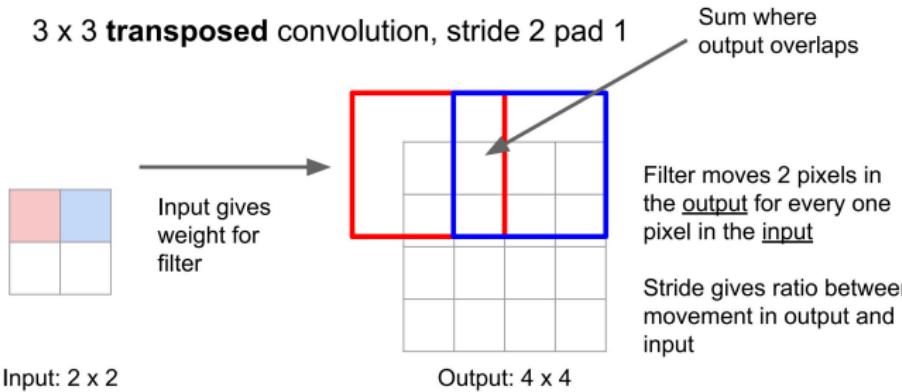
3 x 3 **transposed** convolution, stride 2 pad 1



Transposed convolution (反卷积): learnable upsampling

Learnable Upsampling: Transposed Convolution

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

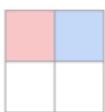


Transposed convolution (反卷积): learnable upsampling

Learnable Upsampling: Transposed Convolution

Q: Why is it called transposed convolution?

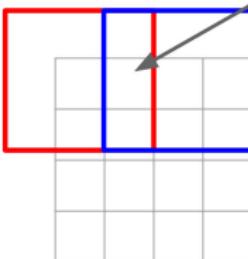
3 x 3 **transposed** convolution, stride 2 pad 1



Input: 2 x 2



Input gives
weight for
filter



Output: 4 x 4

Sum where
output overlaps

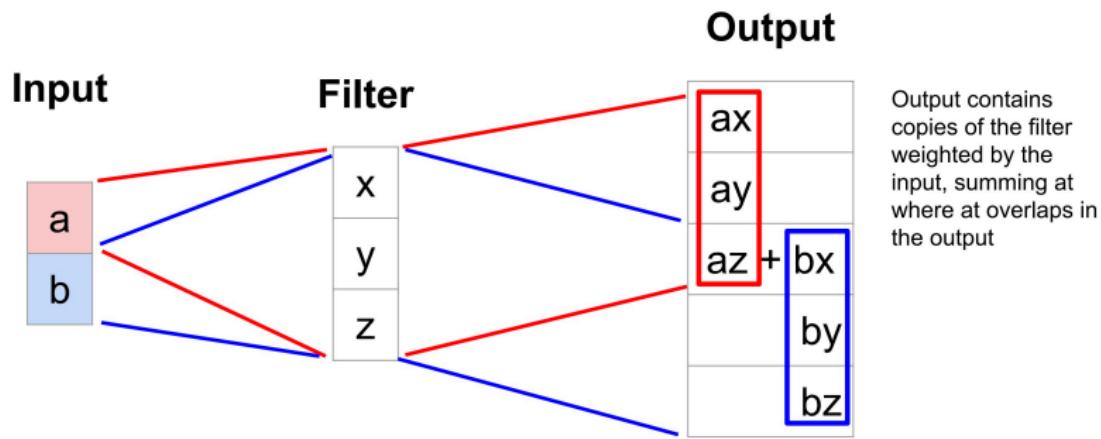
Filter moves 2 pixels in
the output for every one
pixel in the input

Stride gives ratio between
movement in output and
input

Padding: unlike in the regular convolution where padding is applied to input, it is applied to output in the transposed convolution.

Transposed convolution (反卷积): learnable upsampling

Learnable Upsampling: 1D Example



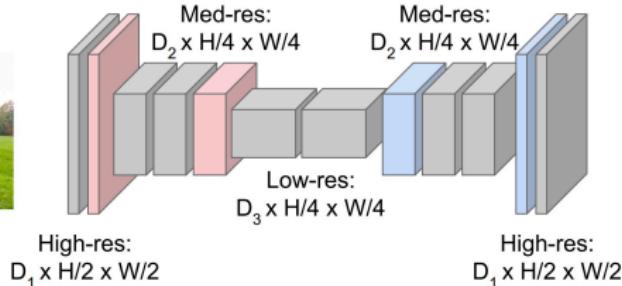
Fully Convolutional Network (FCN, 全卷积网络)

Downsampling:
Pooling, strided convolution



Input:
 $3 \times H \times W$

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



Upsampling:
Unpooling or strided transposed convolution



Predictions:
 $H \times W$

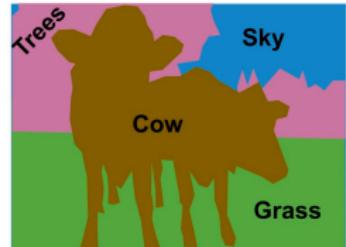
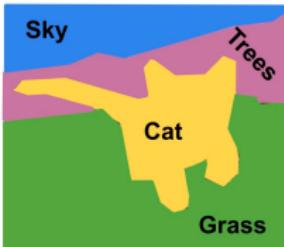
Lecture 8.2 Object Detection 目标检测

Semantic Segmentation (语义分割)

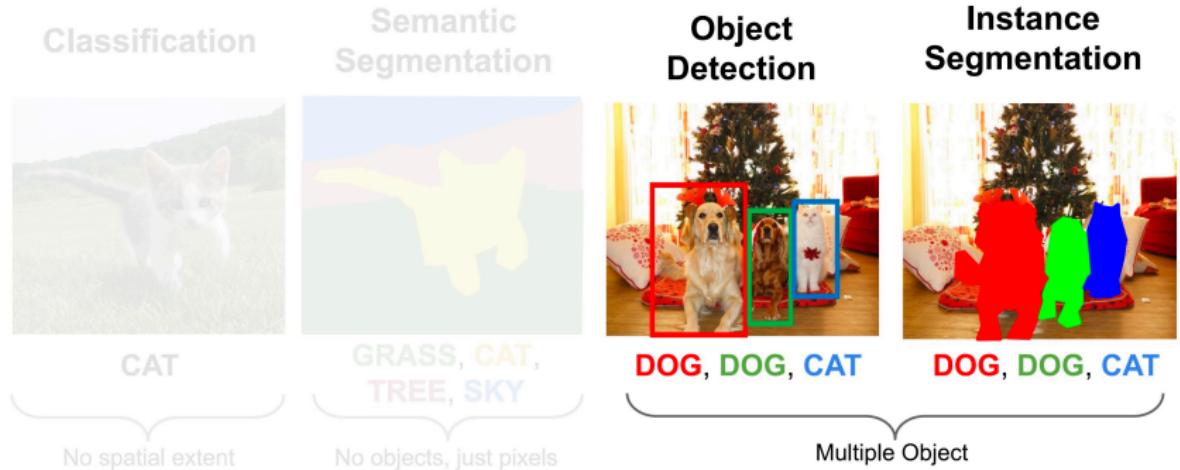
Semantic Segmentation

Label each pixel in the image with a category label

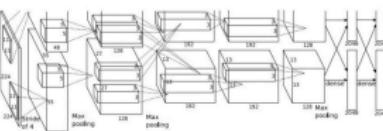
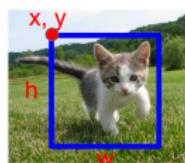
Don't differentiate instances, only care about pixels



Object Detection



Object Detection: Single Object (Classification + Localization)



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

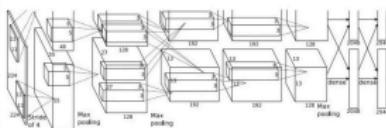
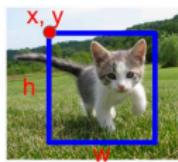
Fully Connected:
4096 to 1000

Vector:
4096

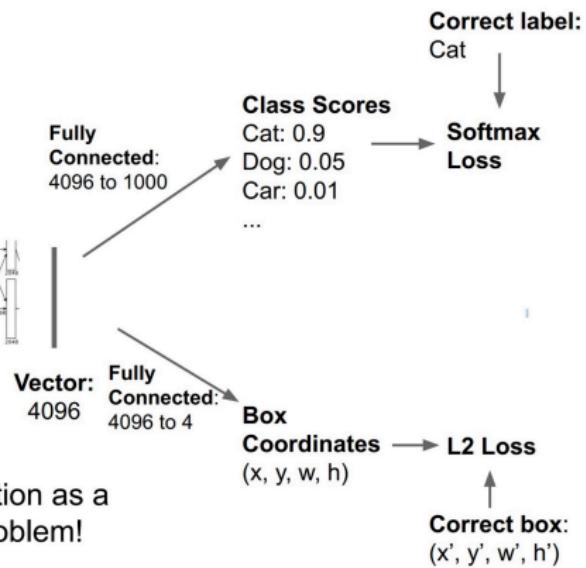
Fully Connected:
4096 to 4

Box Coordinates
(x , y , w , h)

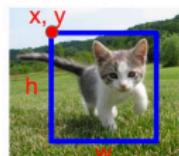
Object Detection: Single Object (Classification + Localization)



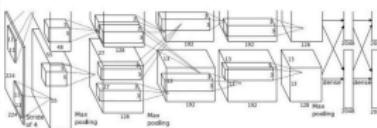
Treat localization as a regression problem!



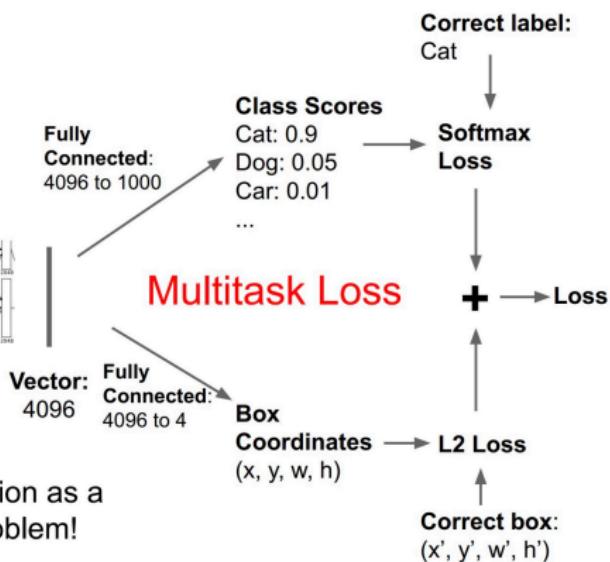
Object Detection: Single Object (Classification + Localization)



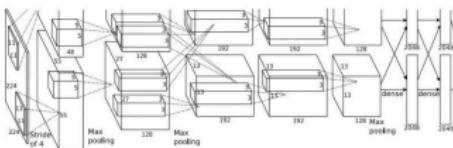
This image © CC0 public domain



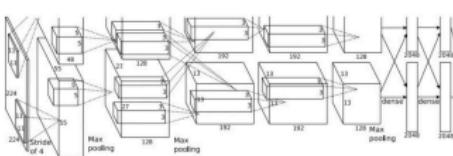
Treat localization as a regression problem!



Object Detection: Multiple Objects

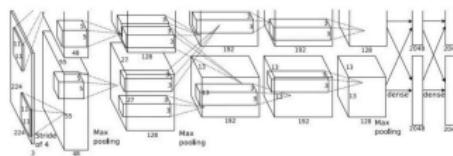


CAT: (x, y, w, h)



DOG: (x, y, w, h)

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



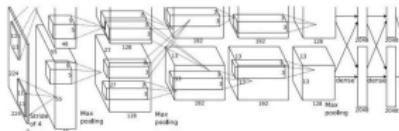
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

....

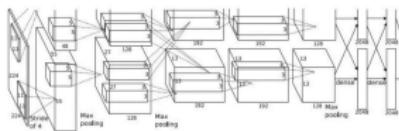
Object Detection: Multiple Objects

Each image needs a different number of outputs!



CAT: (x, y, w, h)

4 numbers

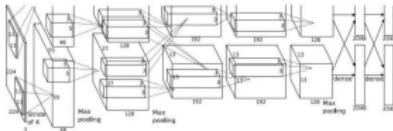


DOG: (x, y, w, h)

12 numbers

DOG: (x, y, w, h)

CAT: (x, y, w, h)



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

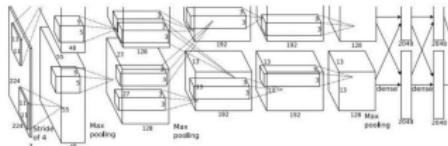
DUCK: (x, y, w, h) numbers!

....

Object Detection: Multiple Objects



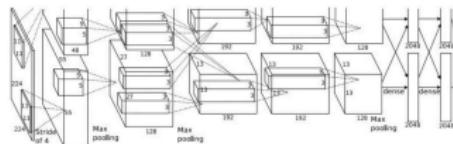
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

Object Detection: Multiple Objects

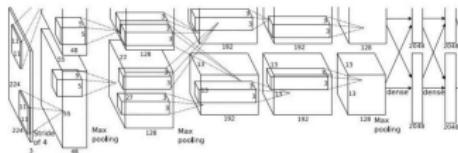
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection: Multiple Objects

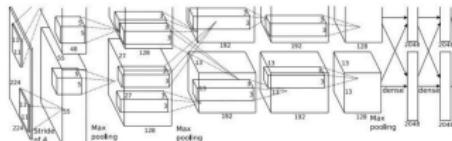
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection: Multiple Objects

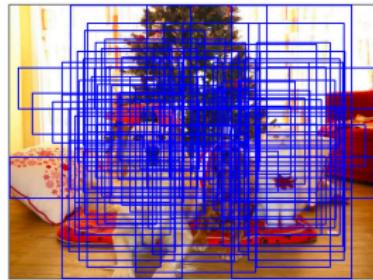
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



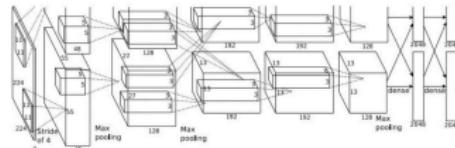
Dog? NO
Cat? YES
Background? NO

Q: What's the problem with this approach?

Object Detection: Multiple Objects



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

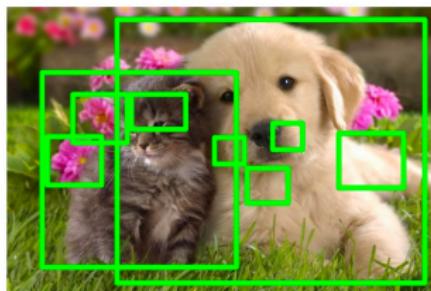


Dog? NO
Cat? YES
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

Region Proposals: Selective Search 候选区域: 选择搜索

- 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



Input image

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

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

R-CNN

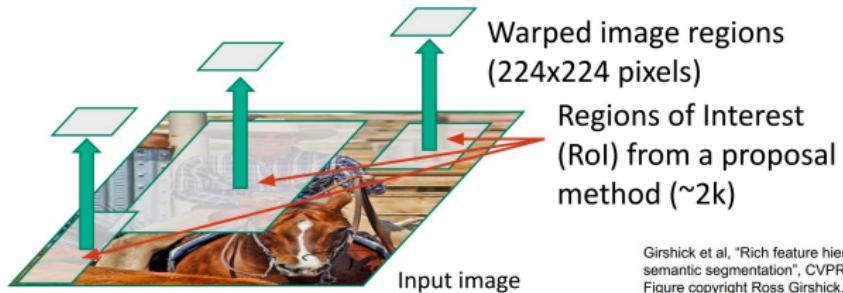


Input image

Regions of Interest
(RoI) from a proposal
method (~2k)

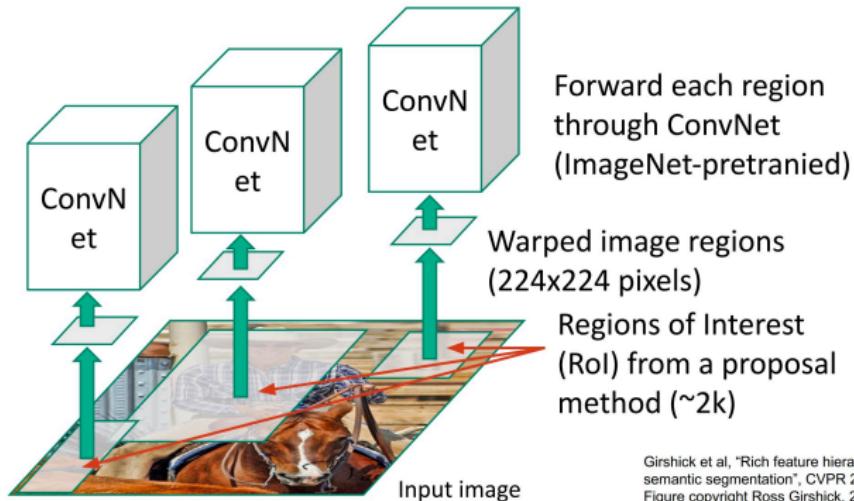
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN



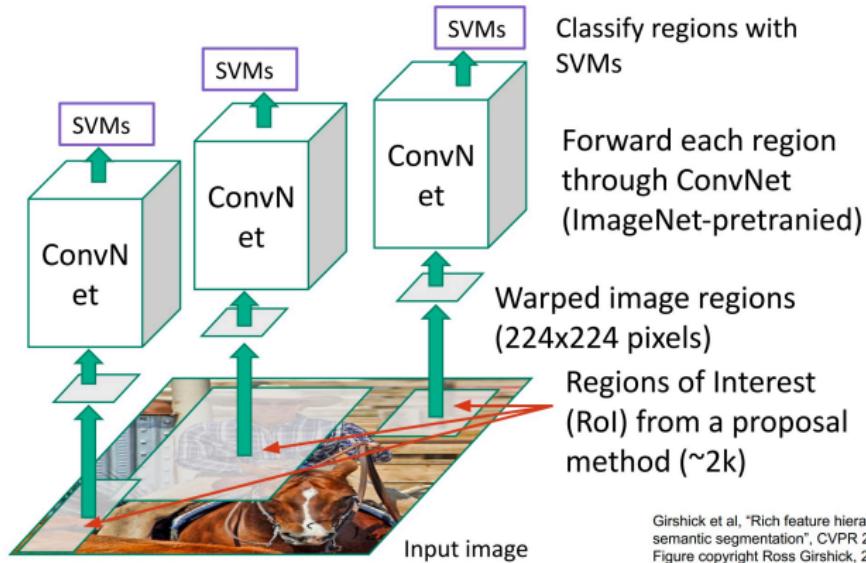
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN



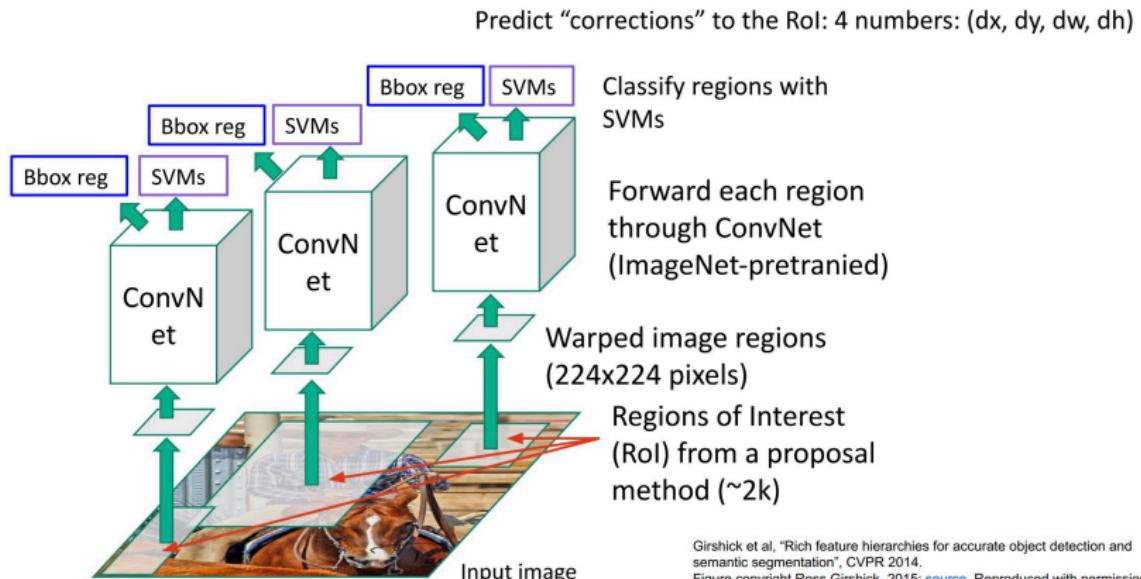
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015: [source](#). Reproduced with permission.

R-CNN



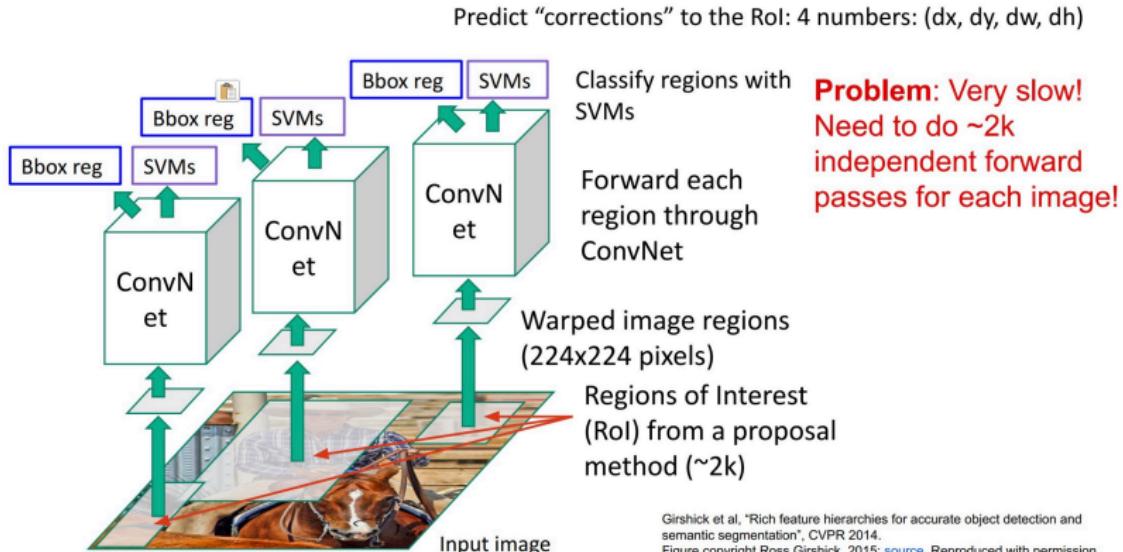
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

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.

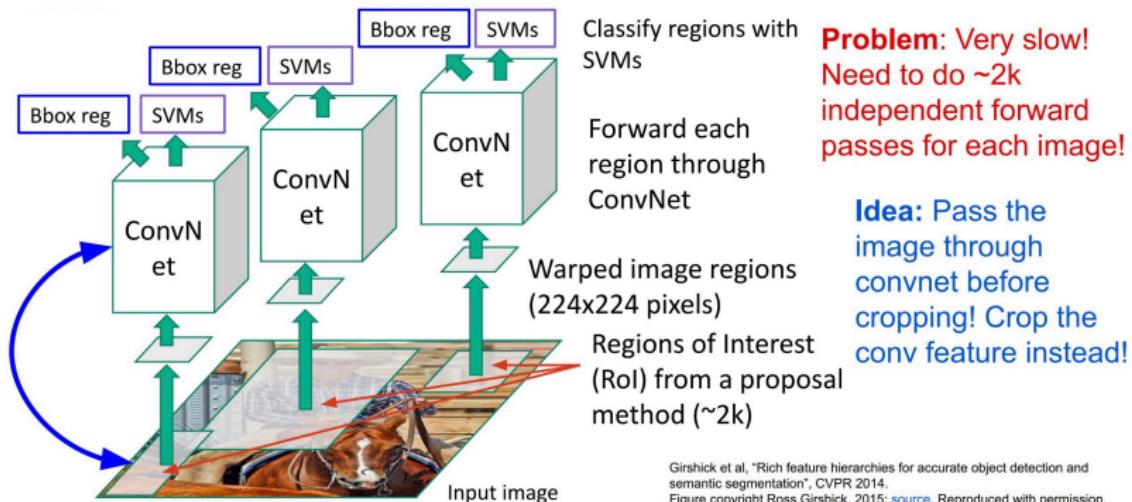
R-CNN



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

“Slow” R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)

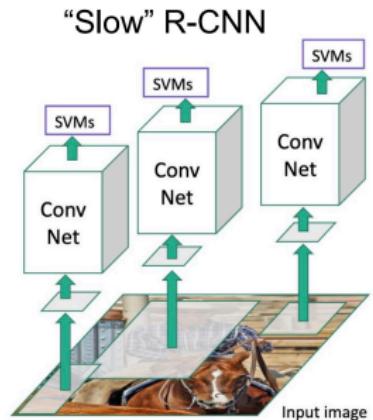


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

Fast R-CNN

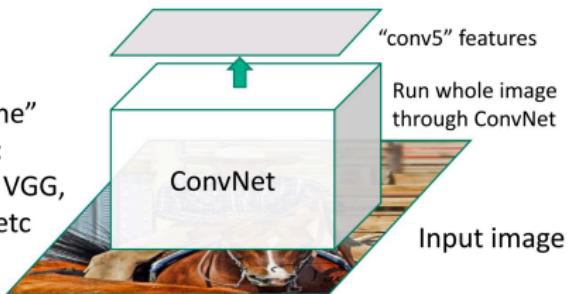


Input image

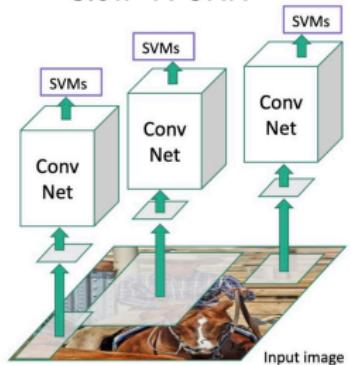


Fast R-CNN

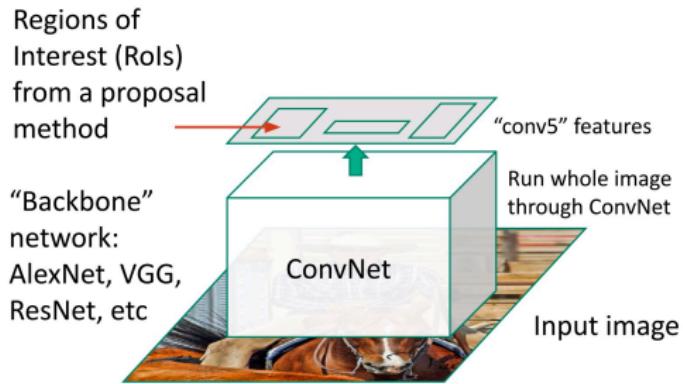
“Backbone”
network:
AlexNet, VGG,
ResNet, etc



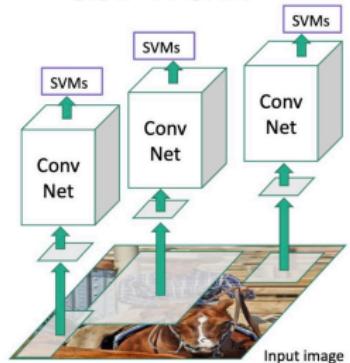
“Slow” R-CNN



Fast R-CNN



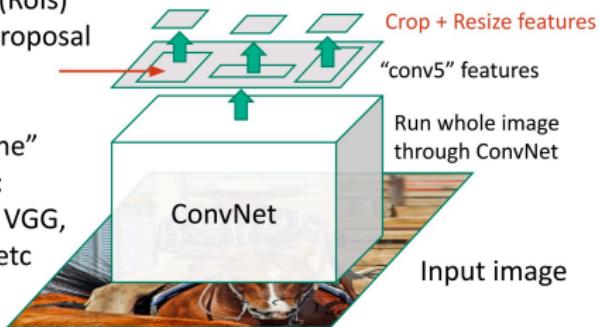
“Slow” R-CNN



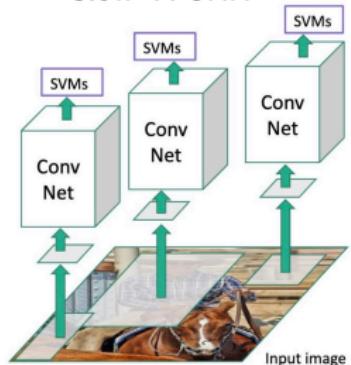
Fast R-CNN

Regions of Interest (Rois)
from a proposal
method

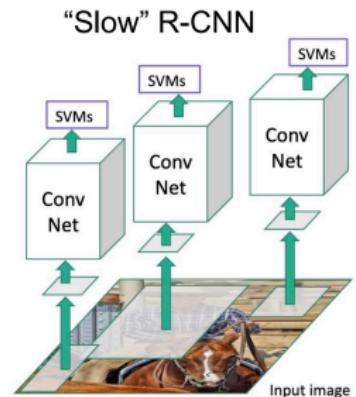
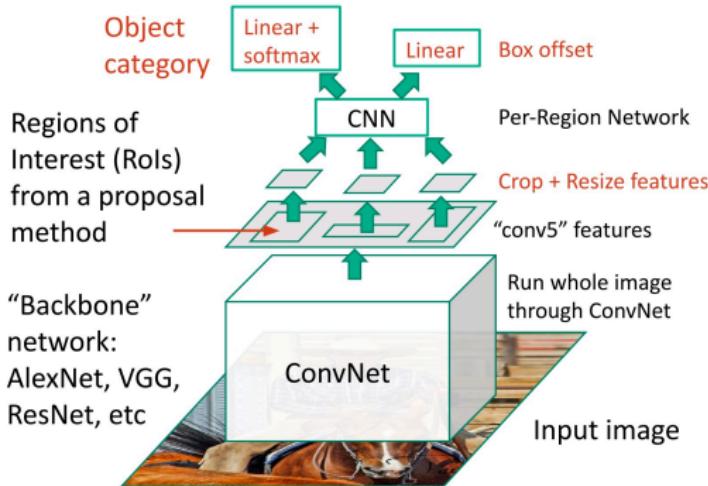
“Backbone”
network:
AlexNet, VGG,
ResNet, etc



“Slow” R-CNN



Fast R-CNN



Thank you!