

Lecture 11: Object Detection and Image Segmentation

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 1

May 9, 2023

Computer Vision Tasks

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

This image is CC0 public domain

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 3

May 9, 2023

Semantic Segmentation

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,
TREE, SKY

Object Detection



DOG, DOG, CAT

Instance Segmentation



DOG, DOG, CAT

Semantic Segmentation: The Problem



GRASS, CAT,
TREE, SKY, ...



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

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

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 4

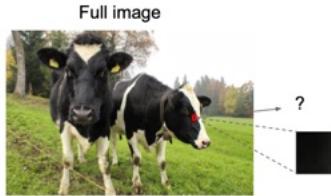
May 9, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 5

May 9, 2023

Semantic Segmentation Idea: Sliding Window



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 6

May 9, 2023

Semantic Segmentation Idea: Sliding Window



Impossible to classify without context

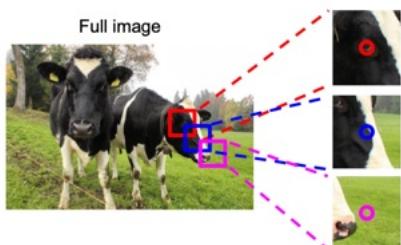
Q: how do we include context?

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 7

May 9, 2023

Semantic Segmentation Idea: Sliding Window



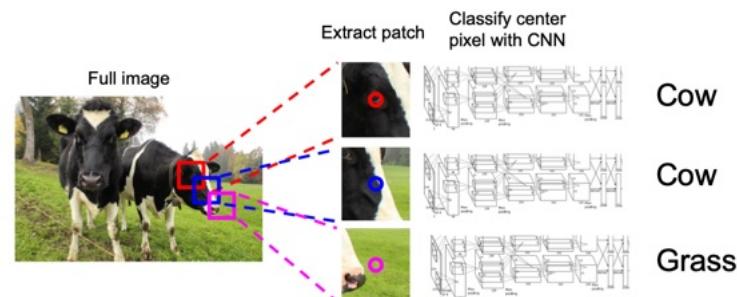
Q: how do we model this?

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 8

May 9, 2023

Semantic Segmentation Idea: Sliding Window



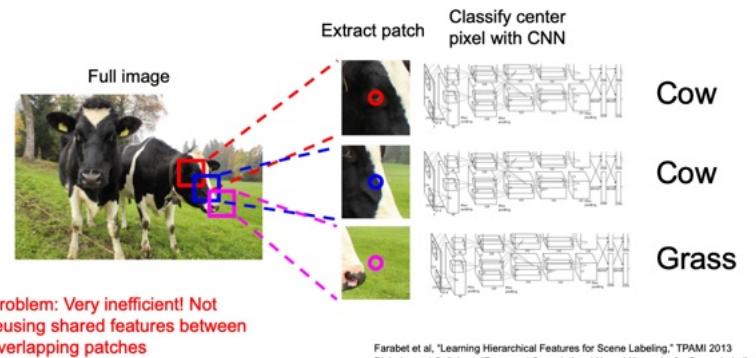
Farabet et al., "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 9

May 9, 2023

Semantic Segmentation Idea: Sliding Window

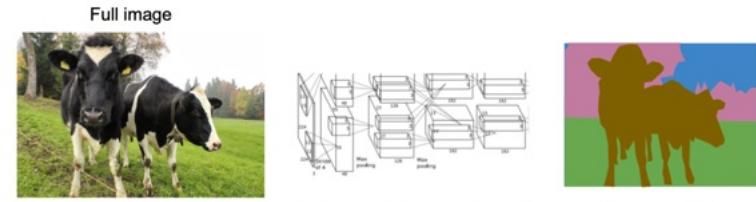


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 10

May 9, 2023

Semantic Segmentation Idea: Convolution



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

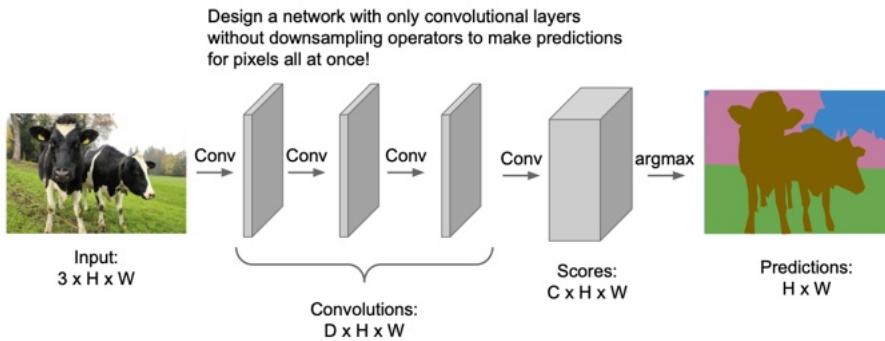
Problem: classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 11

May 9, 2023

Semantic Segmentation Idea: Fully Convolutional

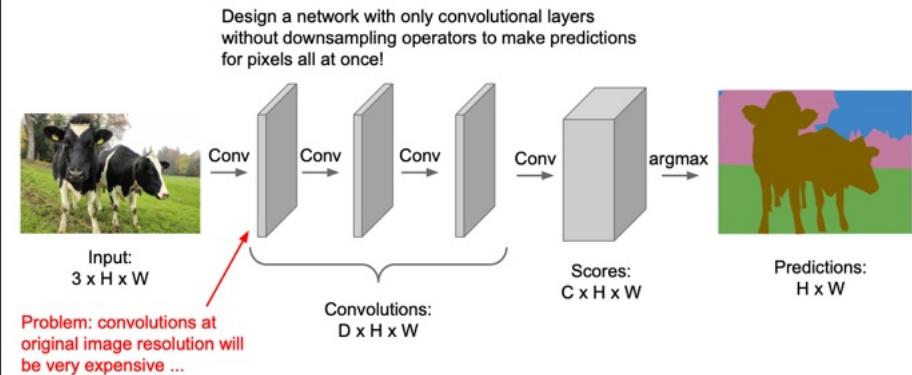


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 12

May 9, 2023

Semantic Segmentation Idea: Fully Convolutional



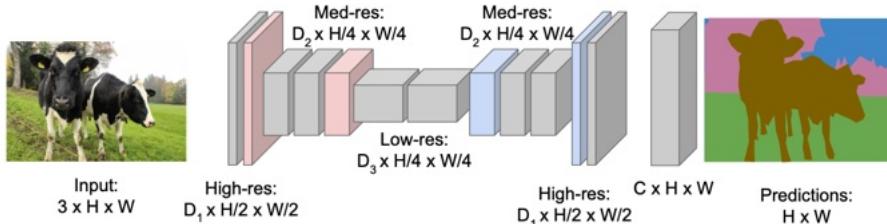
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 13

May 9, 2023

Semantic Segmentation Idea: Fully Convolutional

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



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al., "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 14

May 9, 2023

Semantic Segmentation Idea: Fully Convolutional

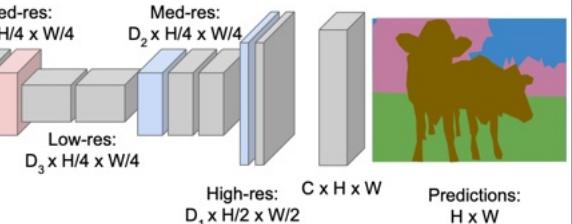
Downsampling:
Pooling, strided convolution



Input:
 $3 \times H \times W$

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

Upsampling:
???



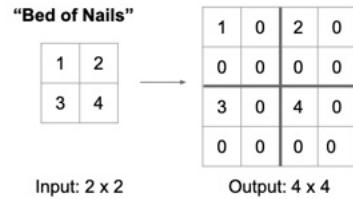
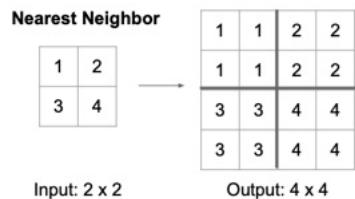
Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al., "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 15

May 9, 2023

In-Network upsampling: “Unpooling”



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 16

May 9, 2023

In-Network upsampling: “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

5	6
7	8

Output: 2 x 2

Max Unpooling

Use positions from pooling layer

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

Input: 2 x 2 Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers

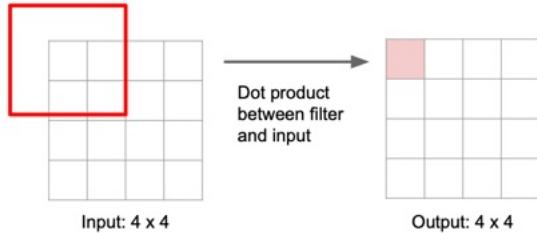
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 17

May 9, 2023

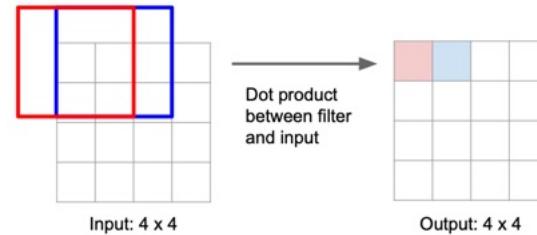
Learnable Upsampling

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



Learnable Upsampling

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



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 19

May 9, 2023

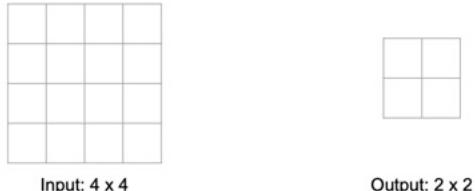
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 20

May 9, 2023

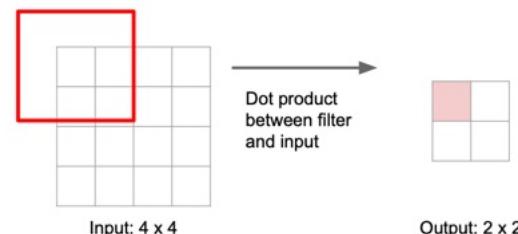
Learnable Upsampling

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Learnable Upsampling

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 21

May 9, 2023

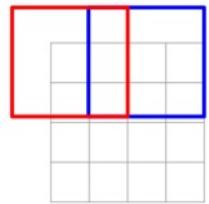
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 22

May 9, 2023

Learnable Upsampling

Recall: Normal 3×3 convolution, stride 2 pad 1



Input: 4×4

Dot product
between filter
and input



Output: 2×2

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

Stride gives ratio between
movement in input and
output

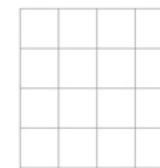
We can interpret strided
convolution as "learnable
downsampling".

Learnable Upsampling: Transposed Convolution

3×3 transposed convolution, stride 2 pad 1



Input: 2×2



Output: 4×4

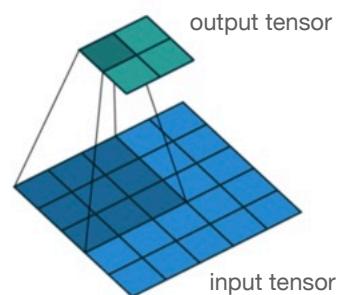


Figure 1: Normal Convolution

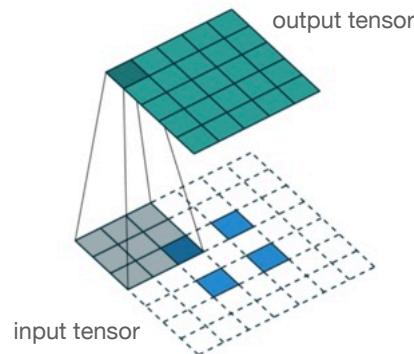
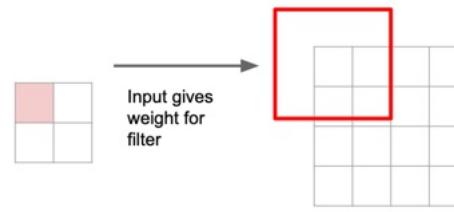


Figure 2: Transposed convolution.

Learnable Upsampling: Transposed Convolution

3×3 transposed convolution, stride 2 pad 1

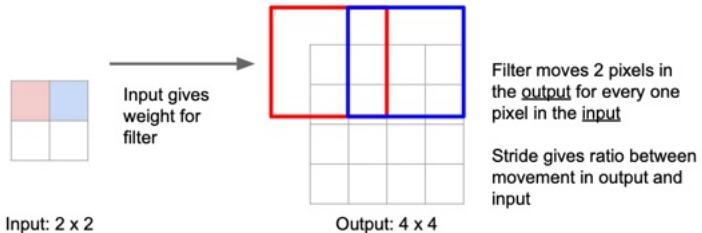


Input: 2×2

Output: 4×4

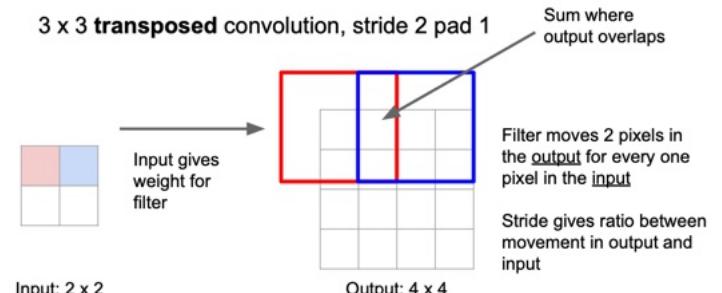
Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1

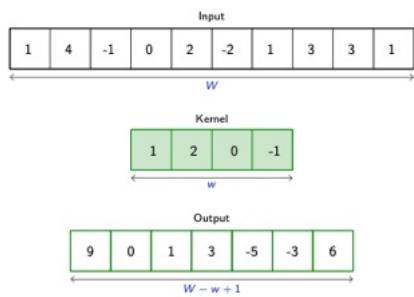


Learnable Upsampling: Transposed Convolution

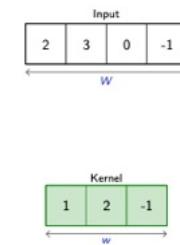
3 x 3 transposed convolution, stride 2 pad 1

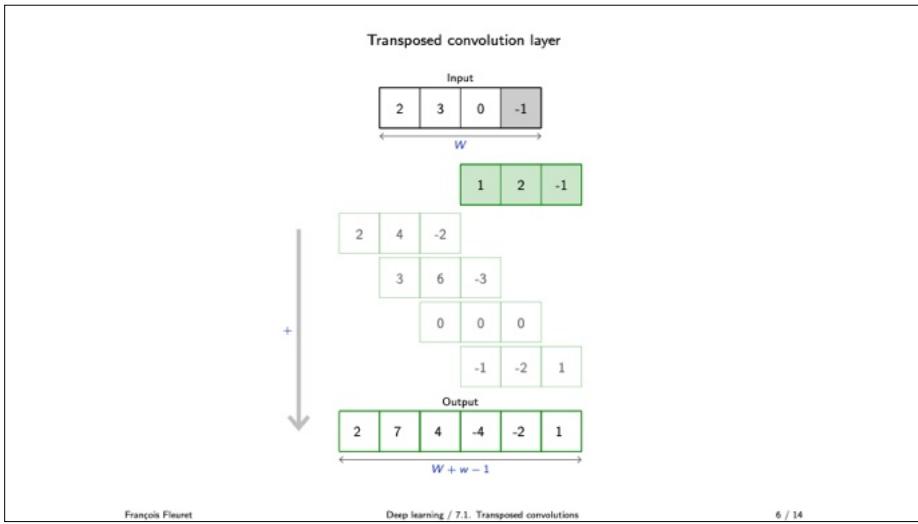
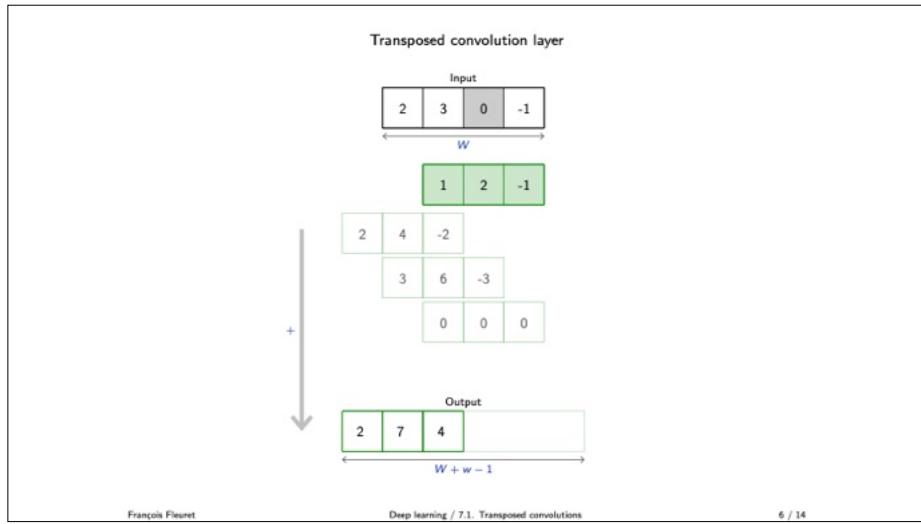
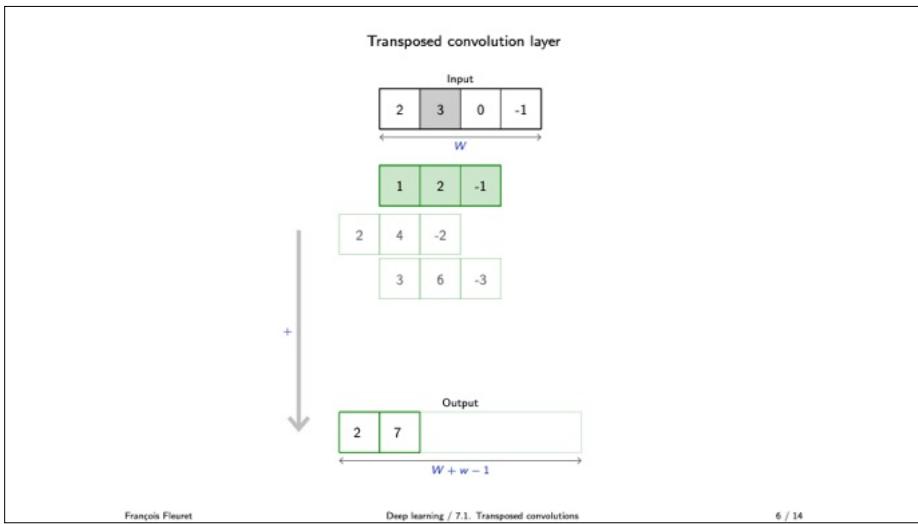
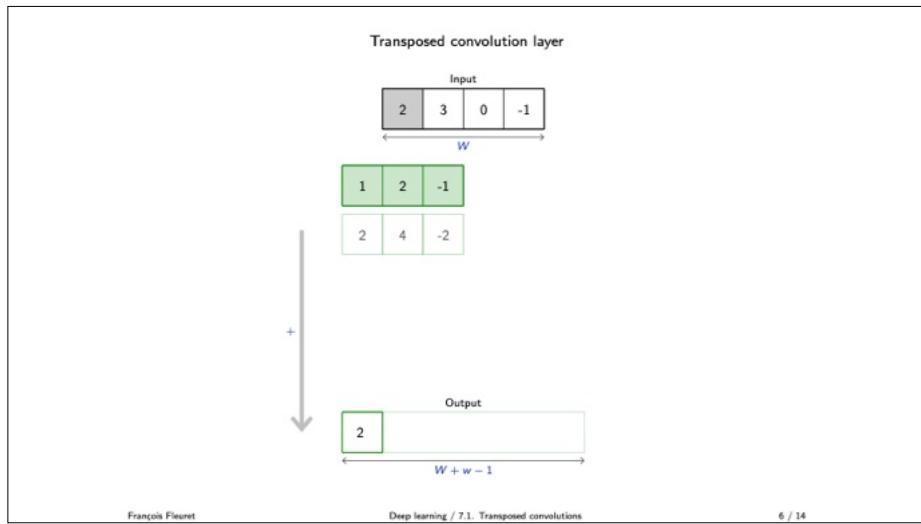


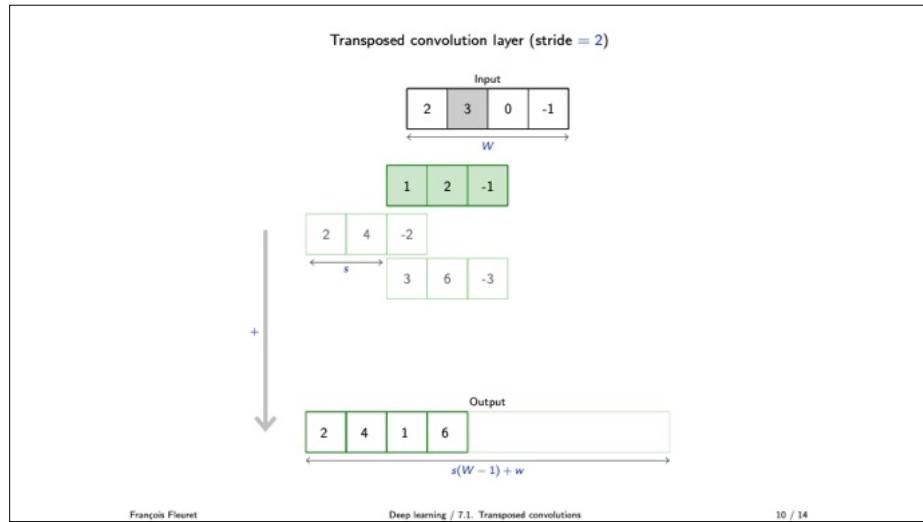
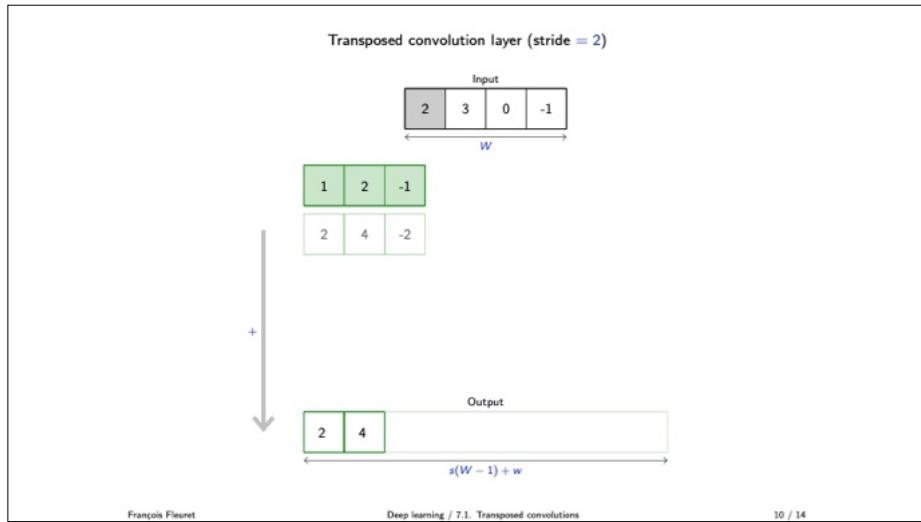
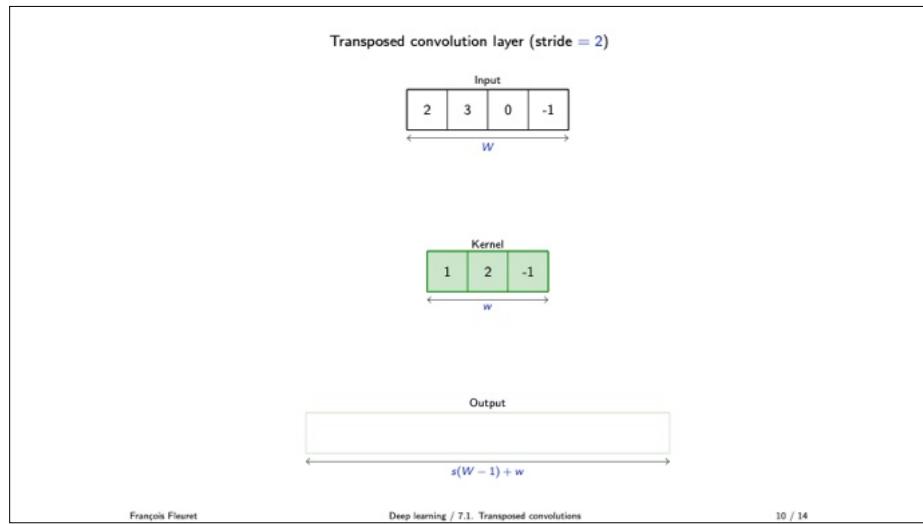
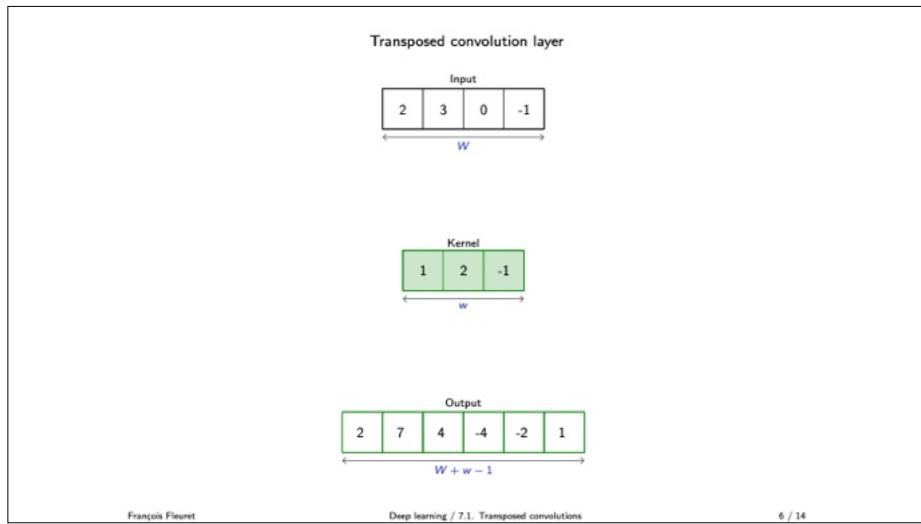
Convolution layer

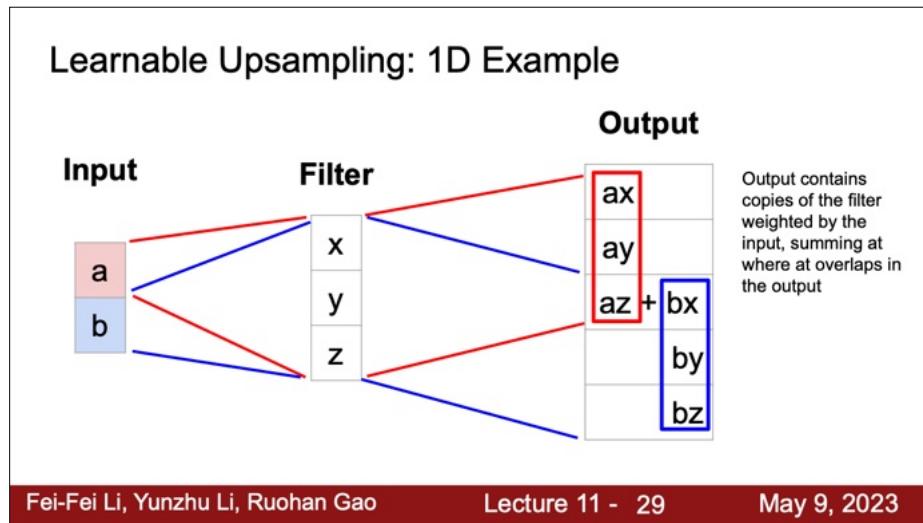
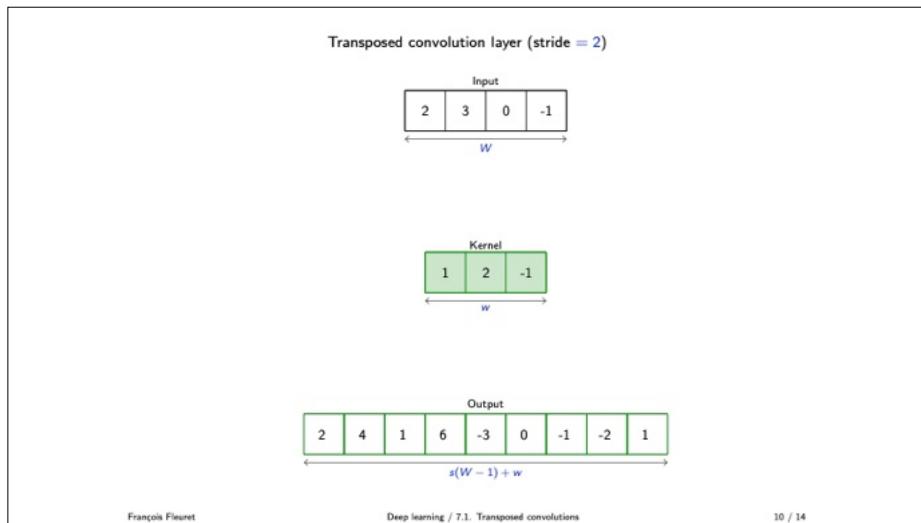
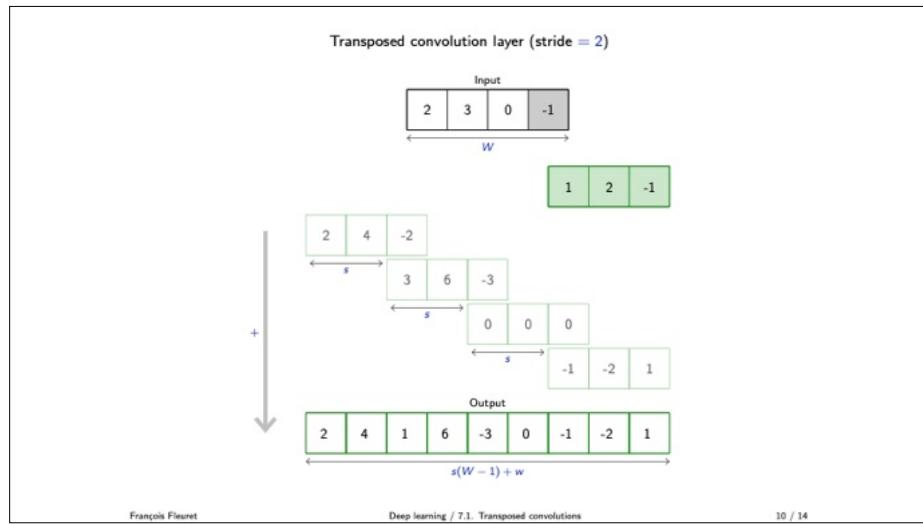
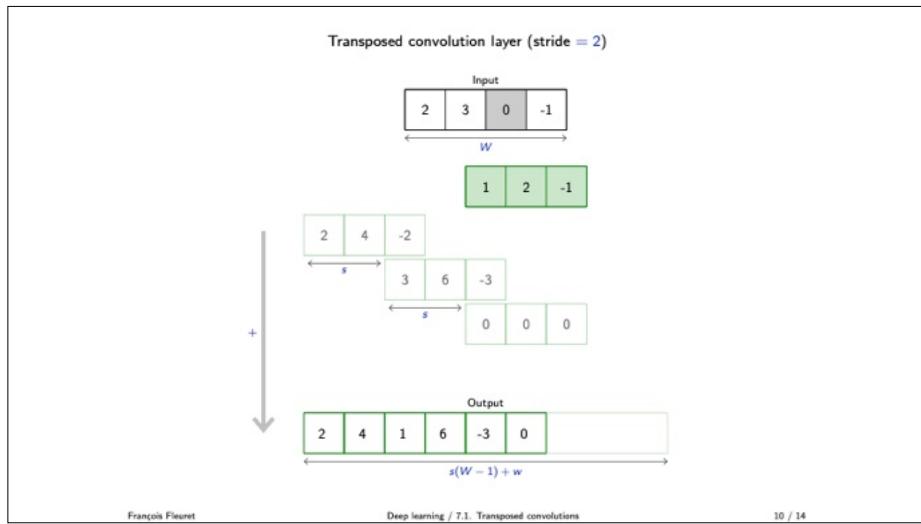


Transposed convolution layer









Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 30

May 9, 2023

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Transposed convolution multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

Example: 1D transposed conv, kernel size=3, stride=2, padding=0

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 31

May 9, 2023

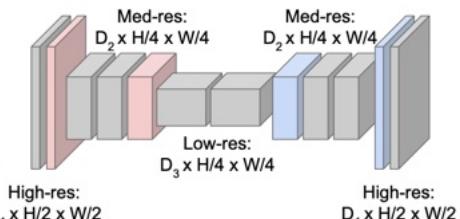
Semantic Segmentation Idea: Fully Convolutional

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

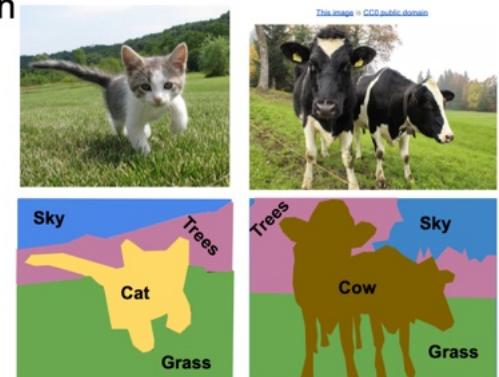


Predictions:
 $H \times W$

Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 32

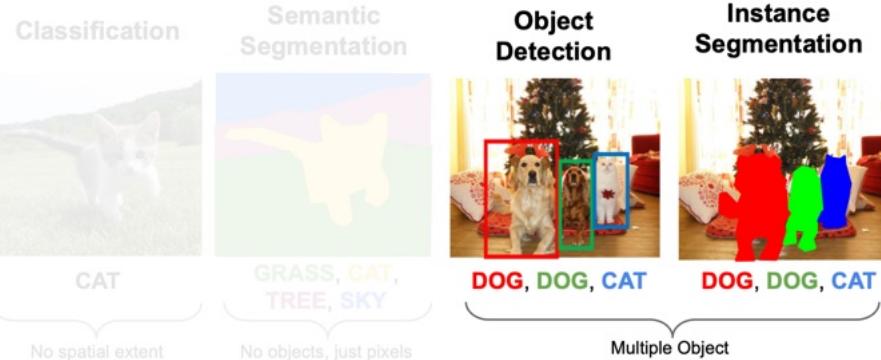
May 9, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 34

May 9, 2023

Object Detection

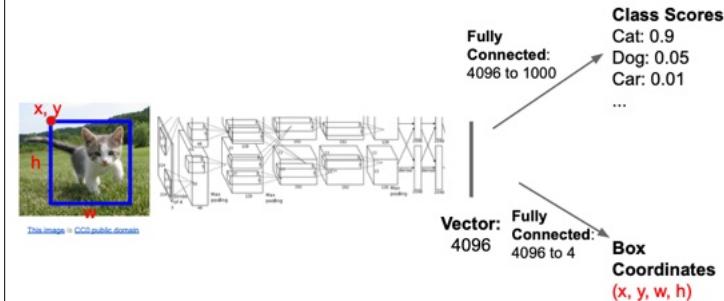


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 36

May 9, 2023

Object Detection: Single Object (Classification + Localization)

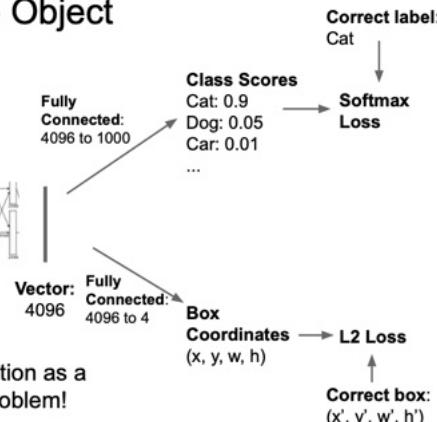
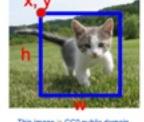


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 37

May 9, 2023

Object Detection: Single Object (Classification + Localization)

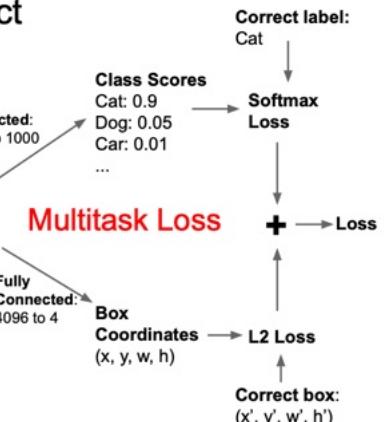
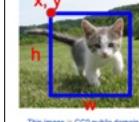


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 38

May 9, 2023

Object Detection: Single Object (Classification + Localization)



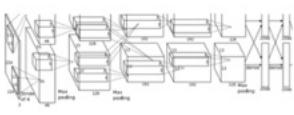
Treat localization as a regression problem!

Fei-Fei Li, Yunzhu Li, Ruohan Gao

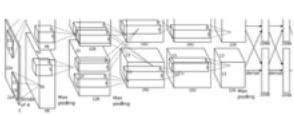
Lecture 11 - 39

May 9, 2023

Object Detection: Multiple Objects



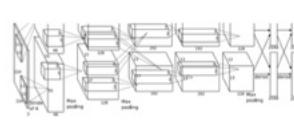
CAT: (x, y, w, h)



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)

....

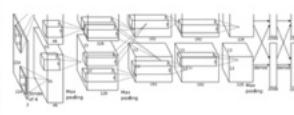
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 40

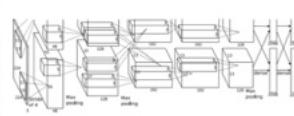
May 9, 2023

Each image needs a
different number of outputs!

Object Detection: Multiple Objects



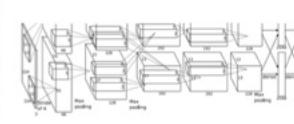
CAT: (x, y, w, h) 4 numbers



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



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

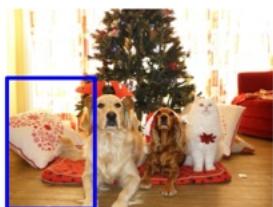
DUCK: (x, y, w, h) numbers!

Fei-Fei Li, Yunzhu Li, Ruohan Gao

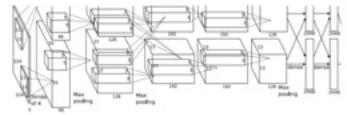
Lecture 11 - 41

May 9, 2023

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

Fei-Fei Li, Yunzhu Li, Ruohan Gao

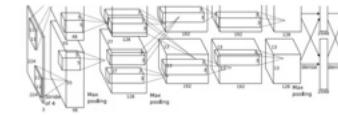
Lecture 11 - 42

May 9, 2023

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

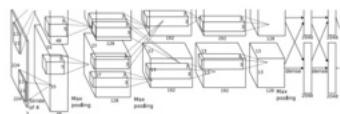
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 43

May 9, 2023

Object Detection: Multiple Objects

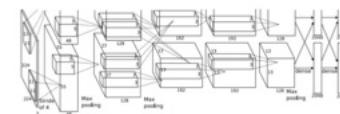
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection: Multiple Objects

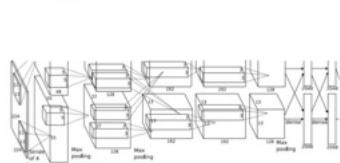
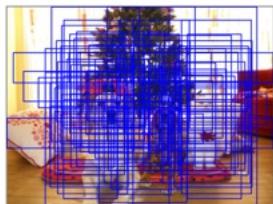
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

Object Detection: Multiple Objects

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

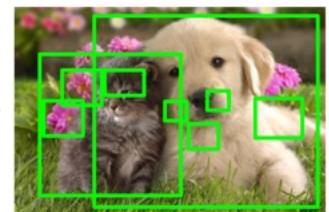


Dog? NO
Cat? YES
Background? NO

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

Region Proposals: Selective Search

- 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



Alexe et al., "Measuring the objectness of image windows", TPAMI 2012
Uijlings et al., "Selective Search for Object Recognition", IJCV 2013
Cheng et al., "BING: Binarized normal gradients for objectness estimation at 300fps", CVPR 2014
Zitnick and Dollár, "Edge boxes: Locating object proposals from edges", ECCV 2014

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.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 48

May 9, 2023

R-CNN



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

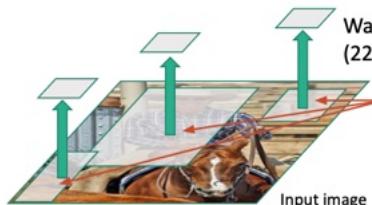
Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 49

May 9, 2023

R-CNN



Warped image regions
(224x224 pixels)
Regions of Interest
(RoI) from a proposal
method (~2k)

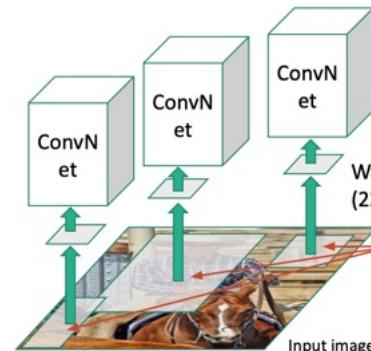
Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 50

May 9, 2023

R-CNN



Forward each region
through ConvNet
(ImageNet-pretrained)

Warped image regions
(224x224 pixels)
Regions of Interest
(RoI) from a proposal
method (~2k)

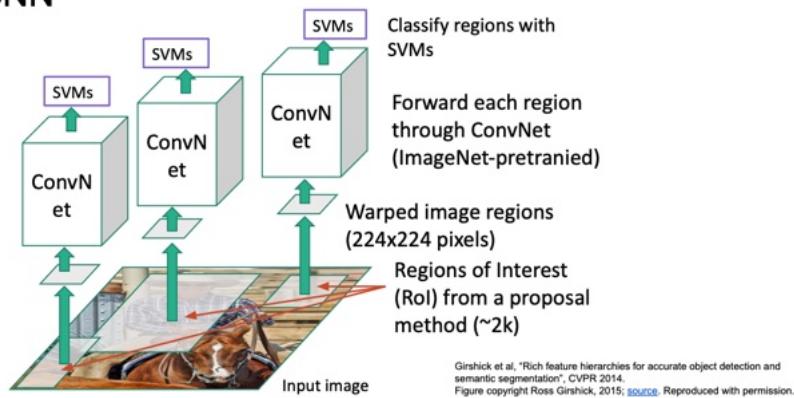
Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 51

May 9, 2023

R-CNN

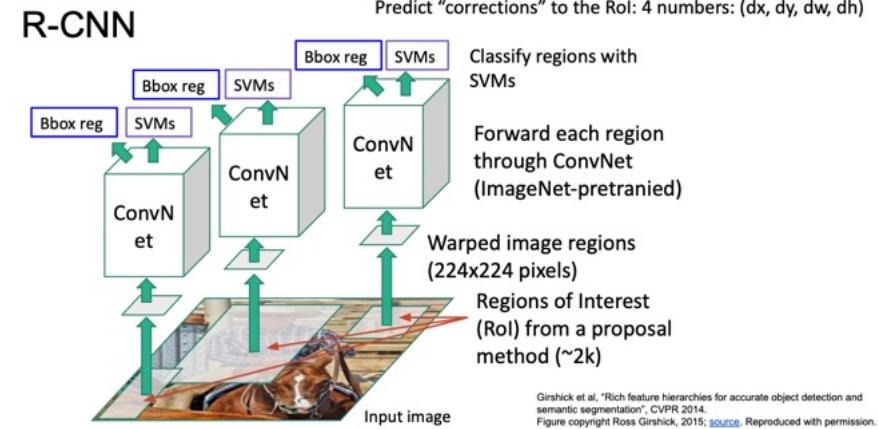


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 52

May 9, 2023

R-CNN

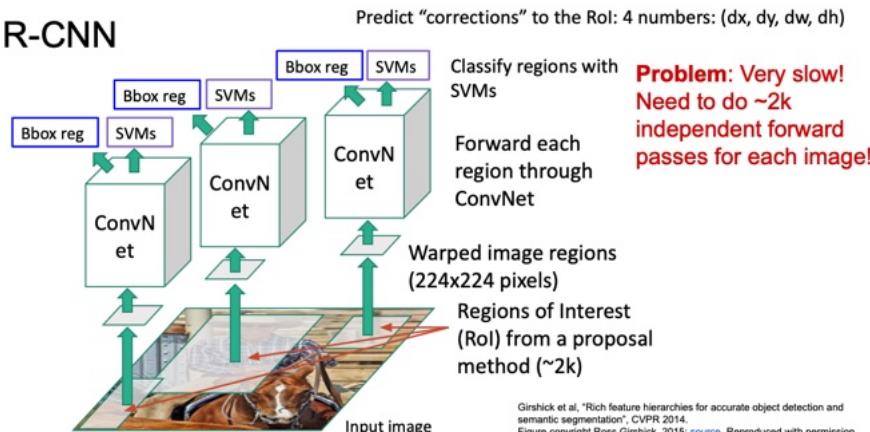


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 53

May 9, 2023

R-CNN

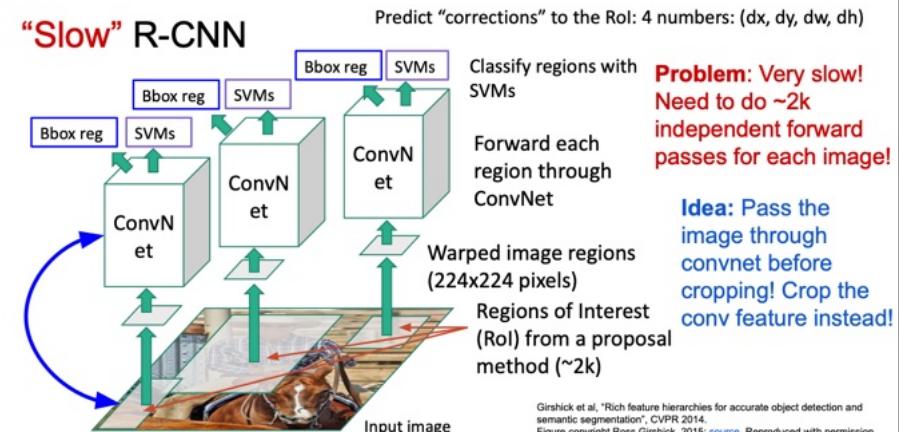


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 54

May 9, 2023

"Slow" R-CNN



Fei-Fei Li, Yunzhu Li, Ruohan Gao

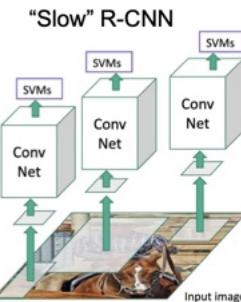
Lecture 11 - 55

May 9, 2023

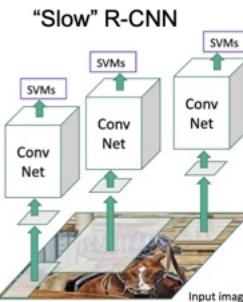
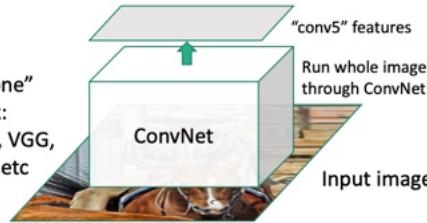
Fast R-CNN



Input image



Fast R-CNN



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 56

May 9, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

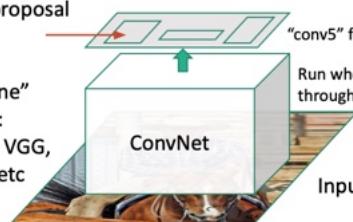
Lecture 11 - 57

May 9, 2023

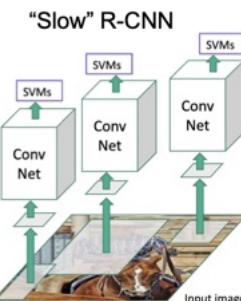
Fast R-CNN

Regions of Interest (Rois)
from a proposal
method

"Backbone"
network:
AlexNet, VGG,
ResNet, etc



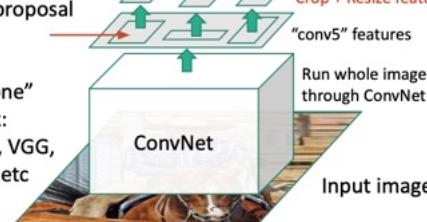
"conv5" features
Run whole image through ConvNet
Input image



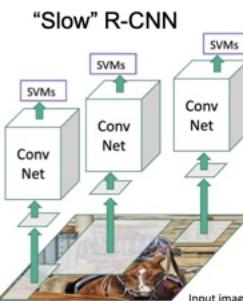
Fast R-CNN

Regions of
Interest (Rois)
from a proposal
method

"Backbone"
network:
AlexNet, VGG,
ResNet, etc



Crop + Resize features
"conv5" features
Run whole image through ConvNet
Input image



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 58

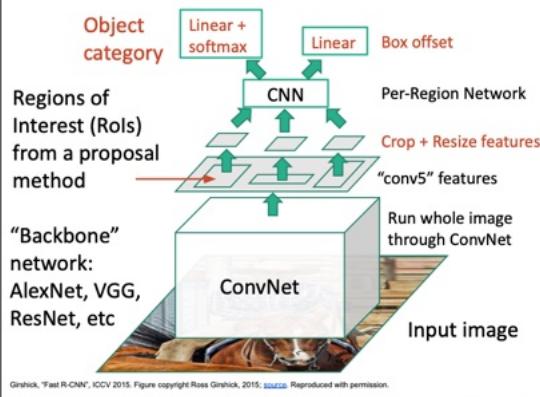
May 9, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 59

May 9, 2023

Fast R-CNN

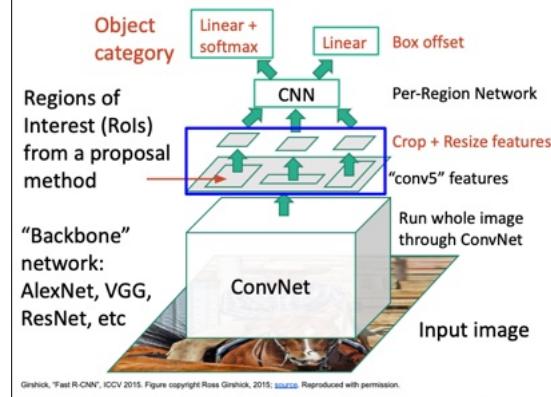


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 60

May 9, 2023

Fast R-CNN

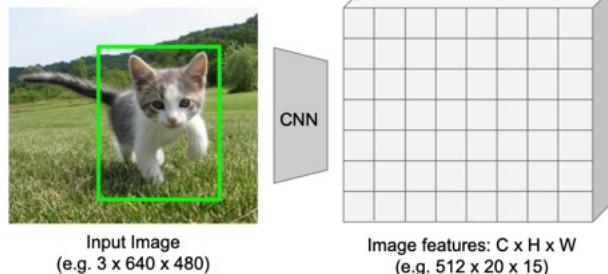


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 61

May 9, 2023

Cropping Features: RoI Pool



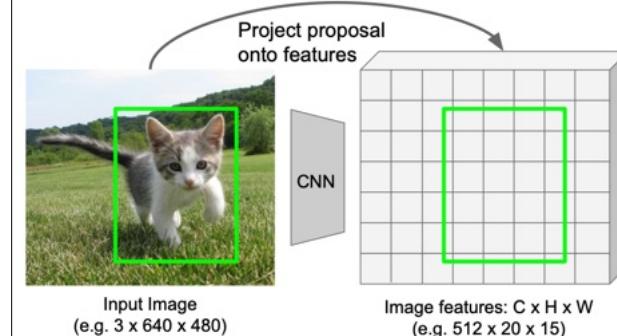
Girshick, “Fast R-CNN”, ICCV 2015.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 62

May 9, 2023

Cropping Features: RoI Pool



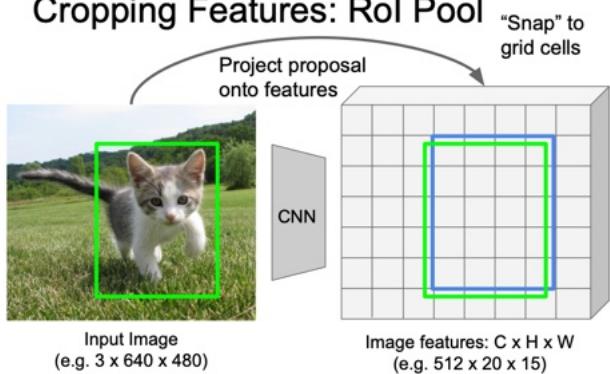
Girshick, “Fast R-CNN”, ICCV 2015.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

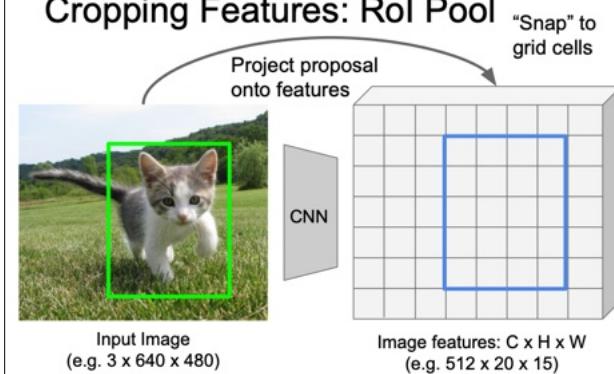
Lecture 11 - 63

May 9, 2023

Cropping Features: RoI Pool



Cropping Features: RoI Pool



Fei-Fei Li, Yunzhu Li, Ruohan Gao

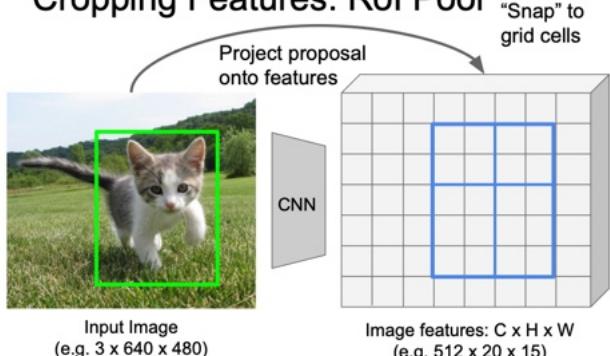
Lecture 11 - 64

May 9, 2023

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

May 9, 2023

Cropping Features: RoI Pool



Divide into 2x2 grid of (roughly) equal subregions

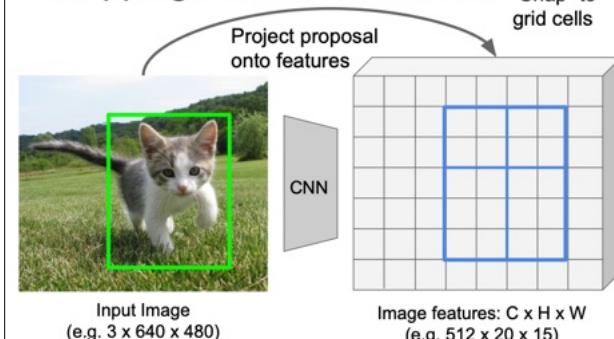
Max-pool within each subregion

Region features (here 512 x 2 x 2; In practice e.g. 512 x 7 x 7)

Divide into 2x2 grid of (roughly) equal subregions

Region features always the same size even if input regions have different sizes!

Cropping Features: RoI Pool



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 66

May 9, 2023

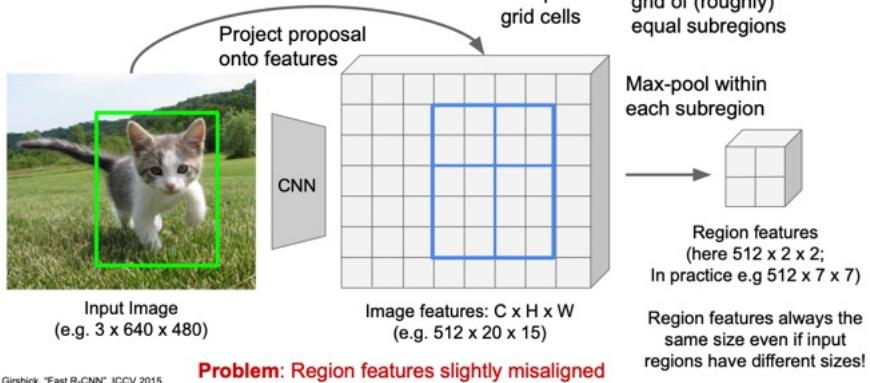
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 67

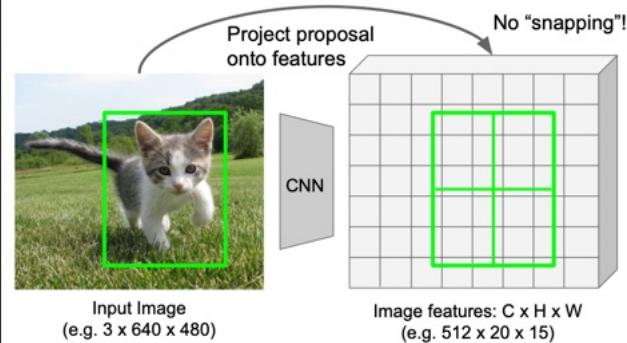
May 9, 2023

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

Cropping Features: RoI Pool



Cropping Features: RoI Align



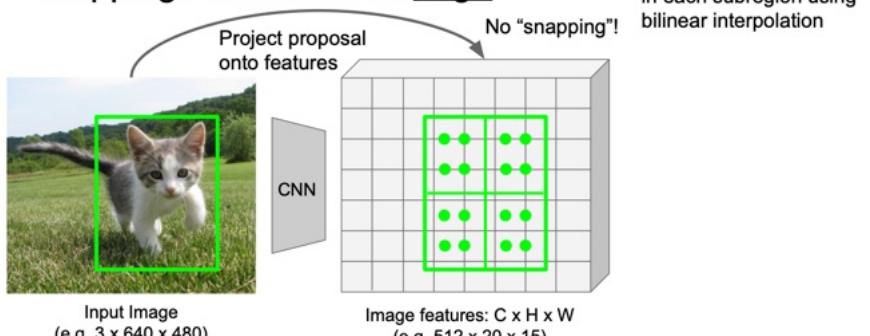
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 68

May 9, 2023

May 9, 2023

Cropping Features: RoI Align

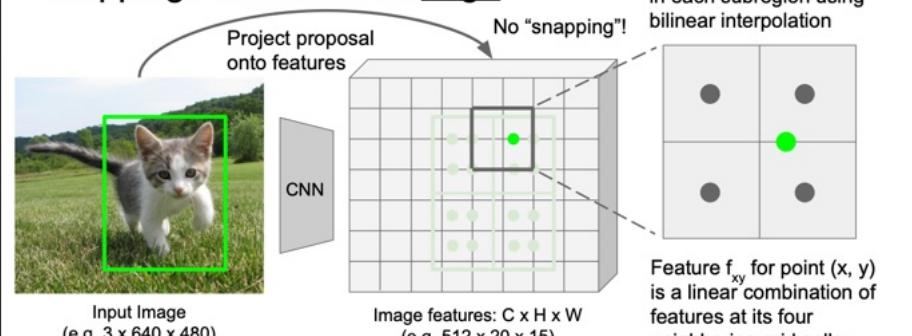


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 70

May 9, 2023

Cropping Features: RoI Align



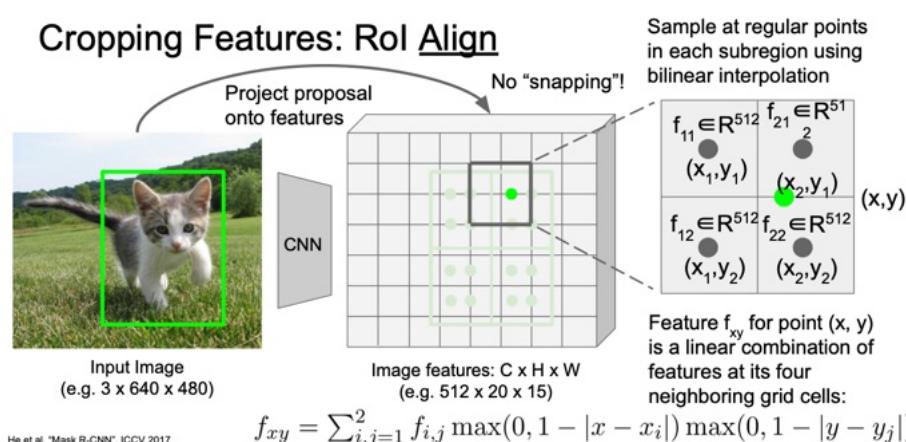
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 71

May 9, 2023

May 9, 2023

Cropping Features: RoI Align

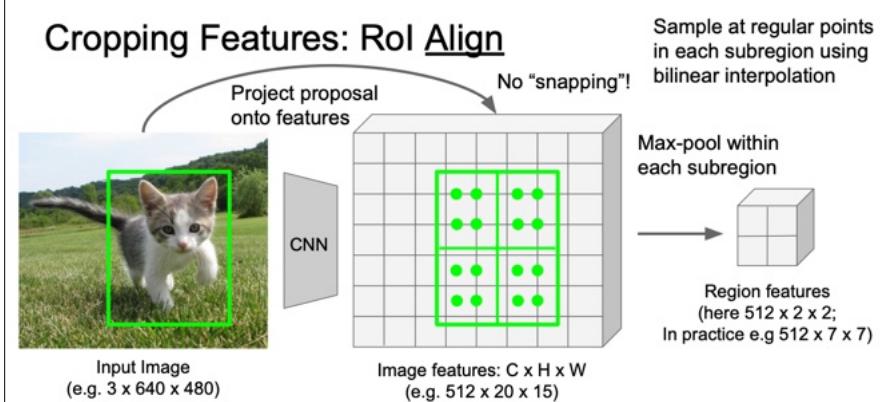


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 72

May 9, 2023

Cropping Features: RoI Align



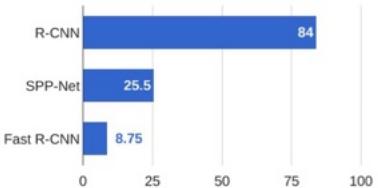
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 73

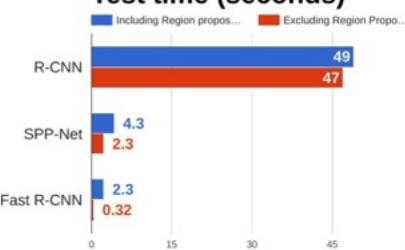
May 9, 2023

R-CNN vs Fast R-CNN

Training time (Hours)



Test time (seconds)

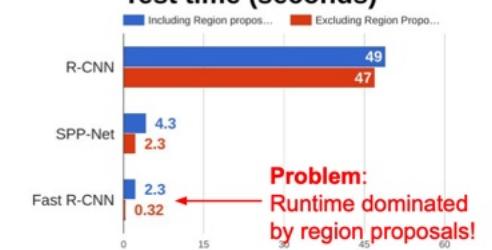


R-CNN vs Fast R-CNN

Training time (Hours)



Test time (seconds)



Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al., "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 74

May 9, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 75

May 9, 2023

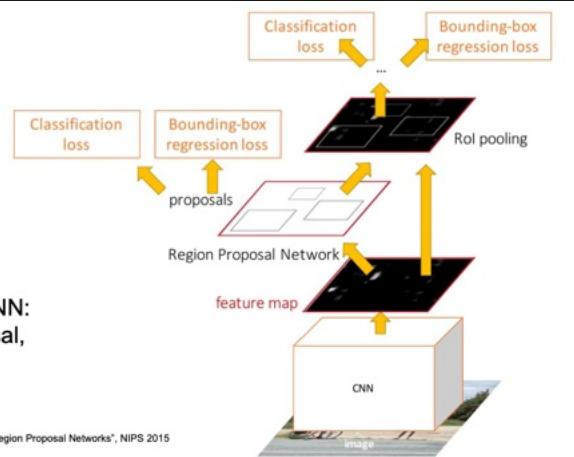
Faster R-CNN:

Make CNN do proposals!

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

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one

Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

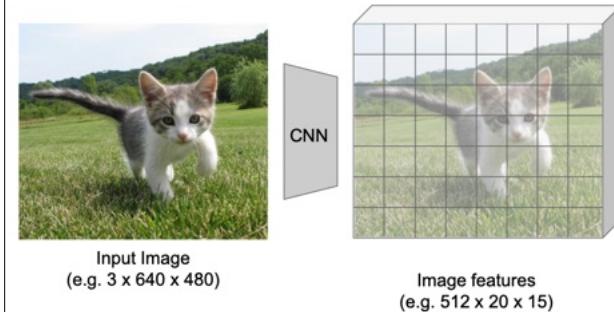


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 76

May 9, 2023

Region Proposal Network



Fei-Fei Li, Yunzhu Li, Ruohan Gao

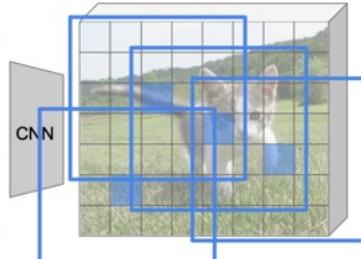
Lecture 11 - 77

May 9, 2023

Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)



Imagine an **anchor box**
of fixed size at each
point in the feature map

Fei-Fei Li, Yunzhu Li, Ruohan Gao

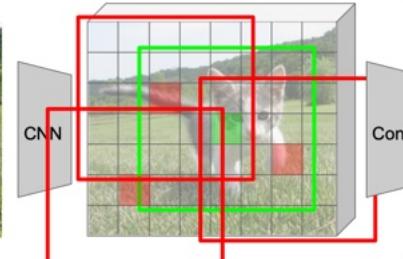
Lecture 11 - 78

May 9, 2023

Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)



Imagine an **anchor box**
of fixed size at each
point in the feature map

Anchor is an object?
1 x 20 x 15

At each point, predict
whether the corresponding
anchor contains an object
(binary classification)

Fei-Fei Li, Yunzhu Li, Ruohan Gao

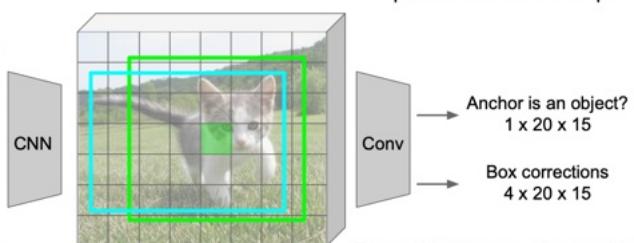
Lecture 11 - 79

May 9, 2023

Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)



Imagine an **anchor box**
of fixed size at each
point in the feature map

For positive boxes, also predict
a corrections from the anchor to
the ground-truth box (regress 4
numbers per pixel)

Fei-Fei Li, Yunzhu Li, Ruohan Gao

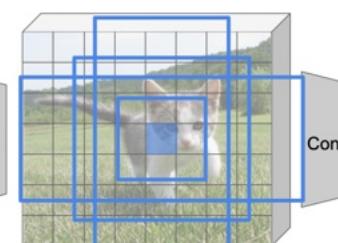
Lecture 11 - 80

May 9, 2023

Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)



In practice use K different
anchor boxes of different
size / scale at each point

Fei-Fei Li, Yunzhu Li, Ruohan Gao

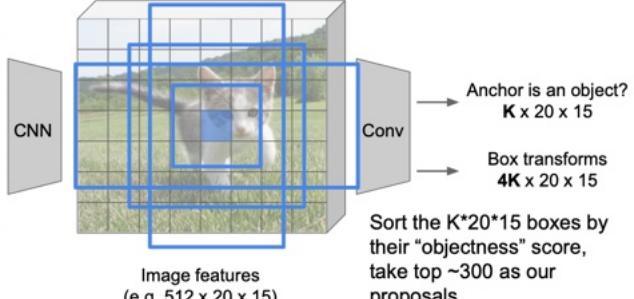
Lecture 11 - 81

May 9, 2023

Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)

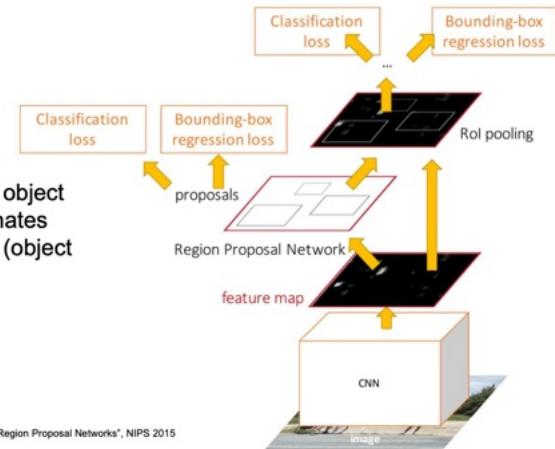


In practice use K different
anchor boxes of different
size / scale at each point

Faster R-CNN:

Make CNN do proposals!

- 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



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 82

May 9, 2023

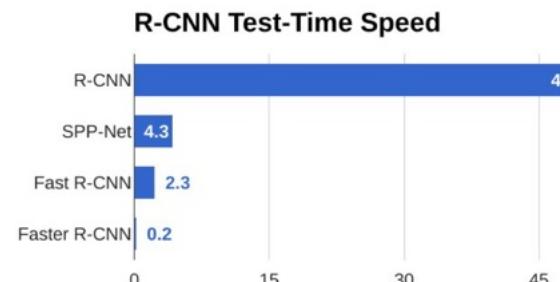
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 83

May 9, 2023

Faster R-CNN:

Make CNN do proposals!



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 84

May 9, 2023

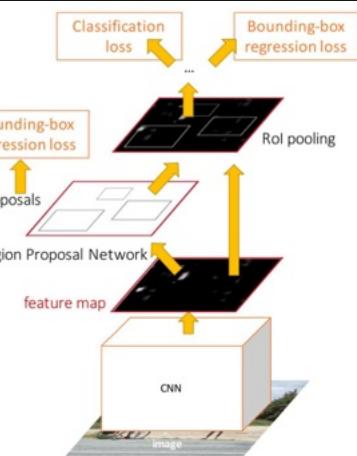
Faster R-CNN:

Make CNN do proposals!

Glossing over many details:

- Ignore overlapping proposals with **non-max suppression**
- How are anchors determined?
- How do we sample positive / negative samples for training the RPN?
- How to parameterize bounding box regression?

Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 85

May 9, 2023

Faster R-CNN:

Make CNN do proposals!

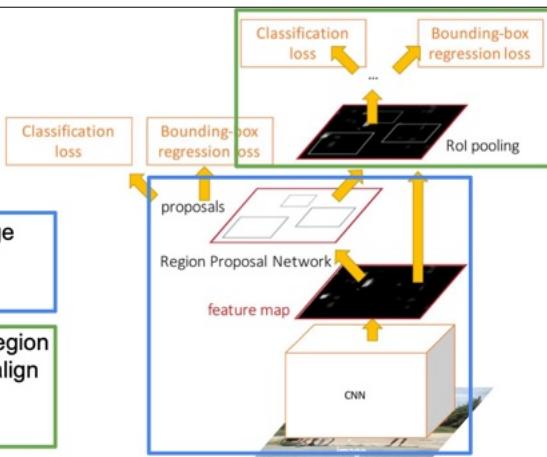
Faster R-CNN is a
Two-stage object detector

First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 86

May 9, 2023

Faster R-CNN:

Make CNN do proposals!

Do we really need
the second stage?

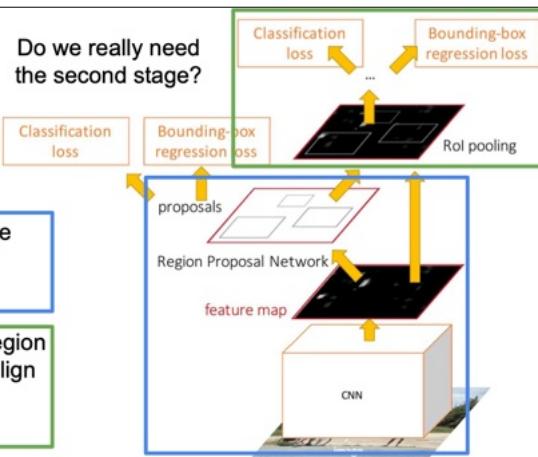
Faster R-CNN is a
Two-stage object detector

First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

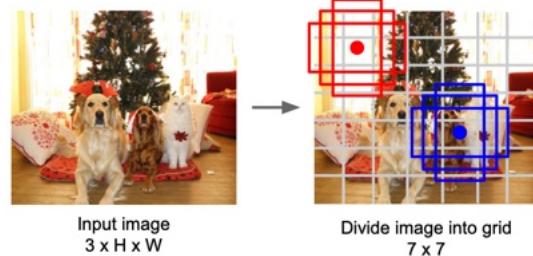


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 87

May 9, 2023

Single-Stage Object Detectors: YOLO / SSD / RetinaNet



- 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)
- Looks a lot like RPN, but category-specific!

Redmon et al., "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016
 Liu et al., "SSD: Single-Shot MultiBox Detector", ECCV 2016
 Lin et al., "Focal Loss for Dense Object Detection", ICCV 2017

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 88

May 9, 2023

Object Detection: Lots of variables ...

Backbone Network
VGG16
ResNet-101
Inception V2
Inception V3
Inception ResNet
MobileNet

"Meta-Architecture"
Two-stage: Faster R-CNN
Single-stage: YOLO / SSD
Hybrid: R-FCN

Image Size
Region Proposals
...

Takeaways
Faster R-CNN is slower but more accurate

SSD is much faster but not as accurate

Bigger / Deeper backbones work better

Huang et al., "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017
 Zou et al., "Object Detection in 20 Years: A Survey", arXiv 2019

R-FCN: Dai et al., "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016
 Inception-V2: Szegedy et al., "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015
 Inception-V3: Szegedy et al., "Rethinking the Inception Architecture for Computer Vision", arXiv 2016
 Inception ResNet: Szegedy et al., "Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016
 MobileNet: Howard et al., "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 90

May 9, 2023

Instance Segmentation

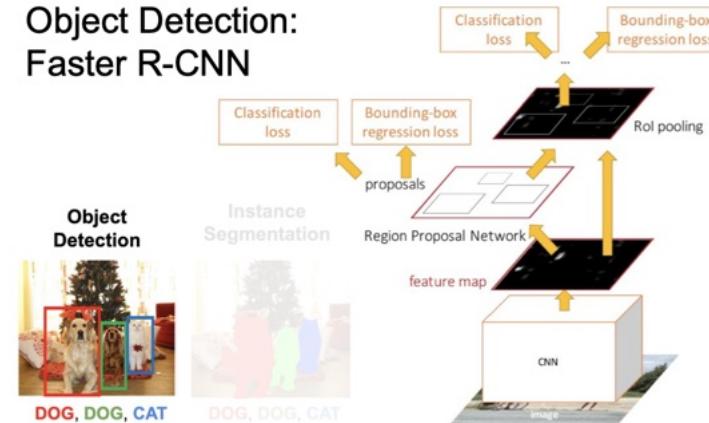


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 91

May 9, 2023

Object Detection: Faster R-CNN

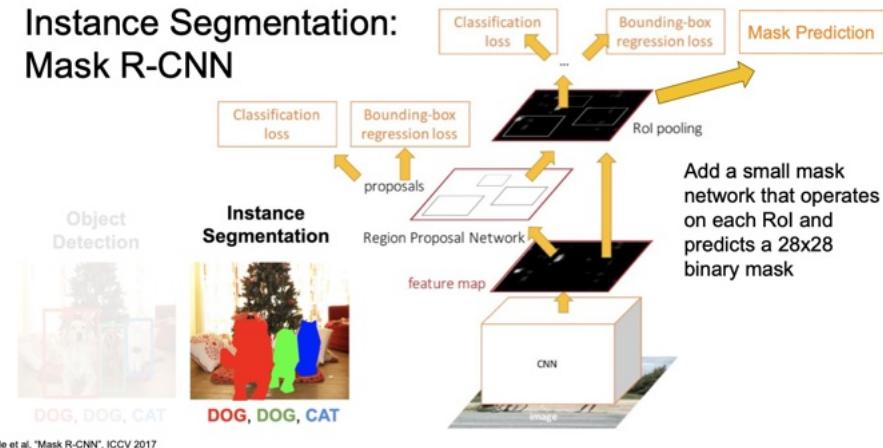


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 92

May 9, 2023

Instance Segmentation: Mask R-CNN

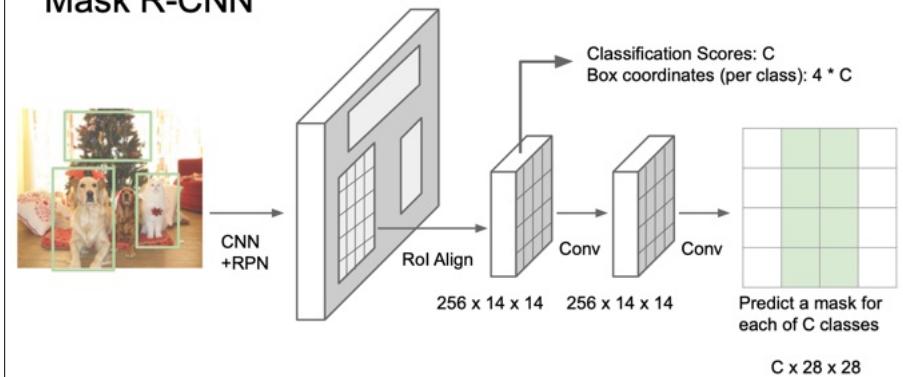


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 93

May 9, 2023

Mask R-CNN

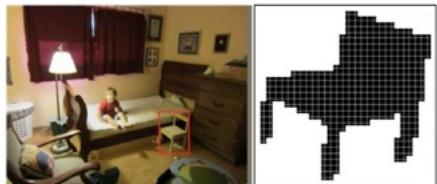


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 94

May 9, 2023

Mask R-CNN: Example Mask Training Targets

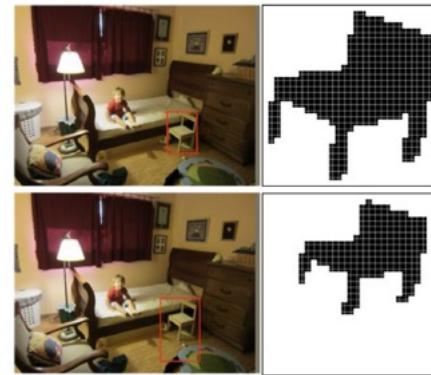


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 95

May 9, 2023

Mask R-CNN: Example Mask Training Targets

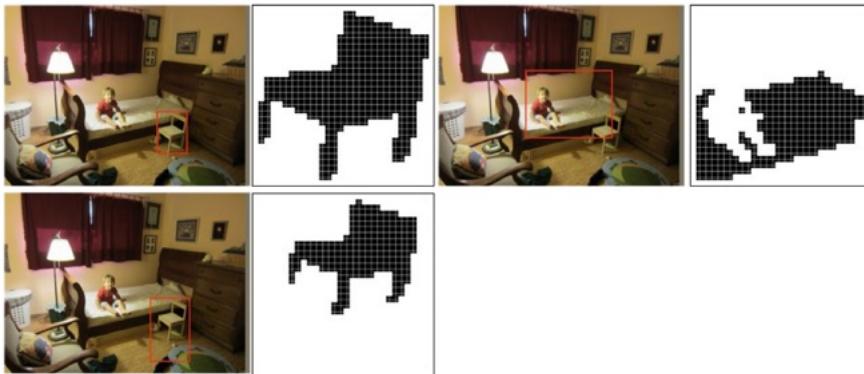


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 96

May 9, 2023

Mask R-CNN: Example Mask Training Targets

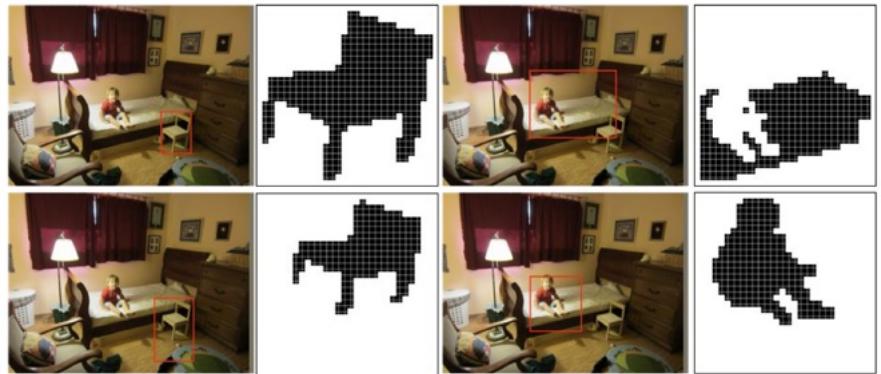


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 97

May 9, 2023

Mask R-CNN: Example Mask Training Targets

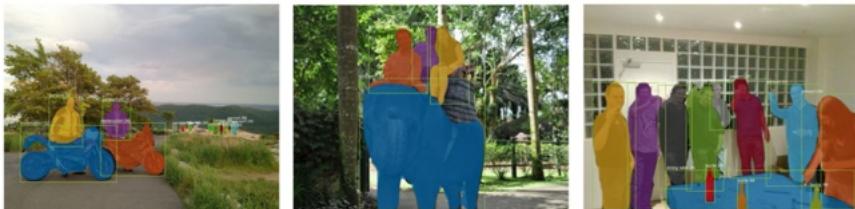


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 98

May 9, 2023

Mask R-CNN: Very Good Results!



He et al., "Mask R-CNN", ICCV 2017

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 99

May 9, 2023

Mask R-CNN Also does pose



He et al., "Mask R-CNN", ICCV 2017

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 100

May 9, 2023

Open Source Frameworks

Lots of good implementations on GitHub!

TensorFlow Detection API:

https://github.com/tensorflow/models/tree/master/research/object_detection

Faster RCNN, SSD, RFCN, Mask R-CNN, ...

Detectron2 (PyTorch)

<https://github.com/facebookresearch/detectron2>

Mask R-CNN, RetinaNet, Faster R-CNN, RPN, Fast R-CNN, R-FCN, ...

Finetune on your own dataset with pre-trained models

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 101

May 9, 2023

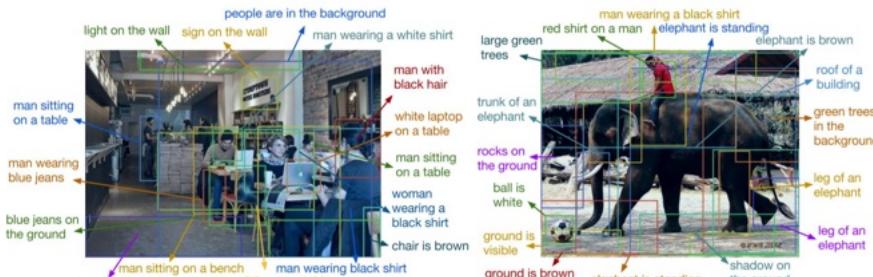
Beyond 2D Object Detection...

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 102

May 9, 2023

Object Detection + Captioning = Dense Captioning



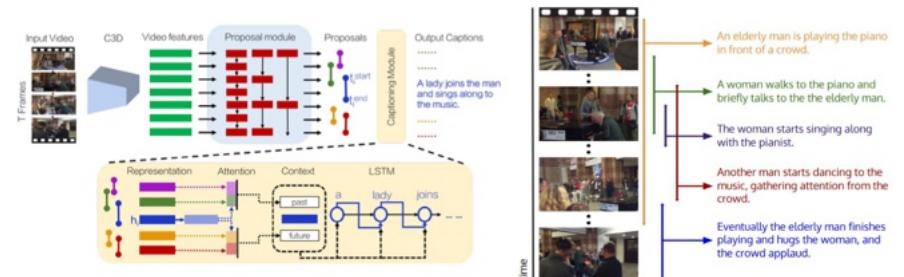
Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016.
Figure copyright IEEE, 2016. Reproduced for educational purposes.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 103

May 9, 2023

Dense Video Captioning



Ranjay Krishna et al., "Dense-Captioning Events in Videos", ICCV 2017.
Figure copyright IEEE, 2017. Reproduced with permission.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 105

May 9, 2023

Objects + Relationships = Scene Graphs

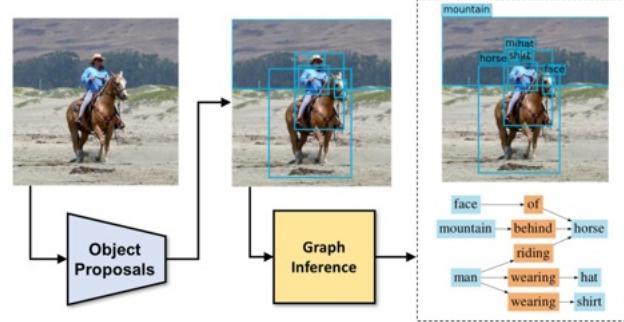


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 106

May 9, 2023

Scene Graph Prediction



Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017
Figure copyright IEEE, 2018. Reproduced for educational purposes.

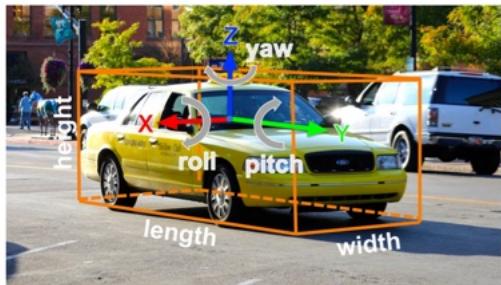
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 106

Lecture 11 - 107

May 9, 2023

3D Object Detection



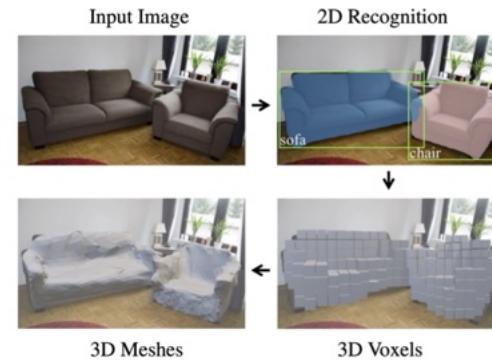
2D Object Detection:
2D bounding box
(x, y, w, h)

3D Object Detection:
3D oriented bounding box
(x, y, z, w, h, l, r, p, y)

Simplified bbox: no roll & pitch

Much harder problem than 2D object detection!

3D Shape Prediction: Mesh R-CNN



Gkioxari et al., Mesh RCNN, ICCV 2019

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 11 - 108

May 9, 2023