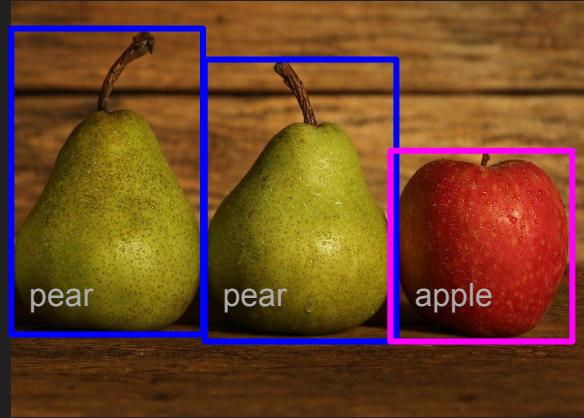


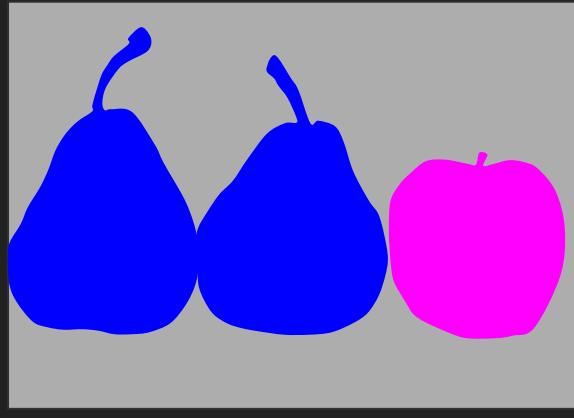
Convolutional neural networks (convnets)

Part 2 - Segmentation & Detection

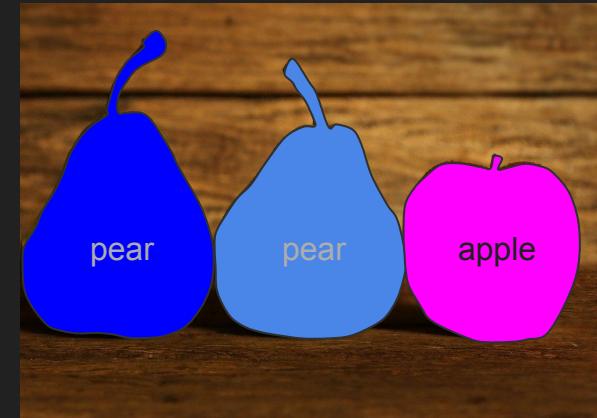
Computer vision tasks



Object detection



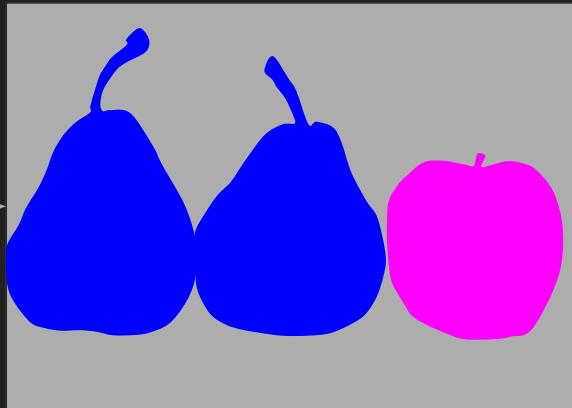
Semantic segmentation
(e.g. pear, apple, background)



Instance
segmentation

Semantic segmentation

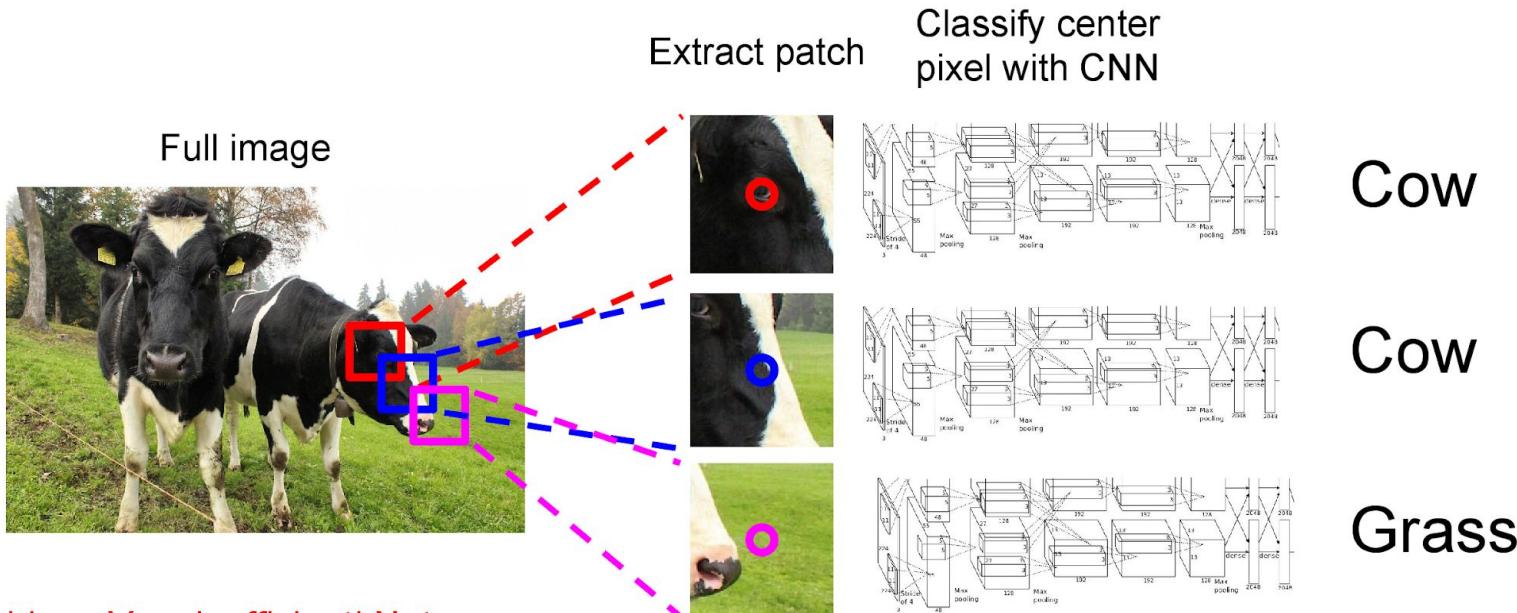
Semantic segmentation



Semantic segmentation
(e.g. **pear**, **apple**, **background**)

- Assign a class to each pixel.
- No separation between instances of objects of same class.
- May be viewed as the classification problem on the pixel level.

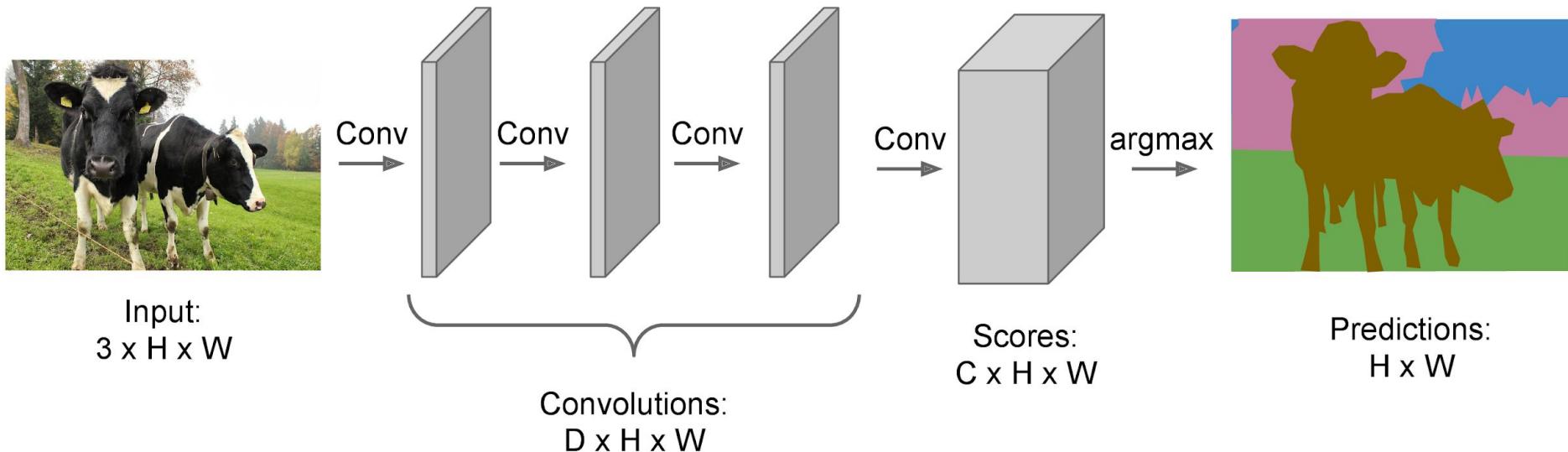
Semantic Segmentation Idea: Sliding Window



Problem: Very inefficient! Not reusing shared features between overlapping patches

Semantic Segmentation Idea: Fully Convolutional

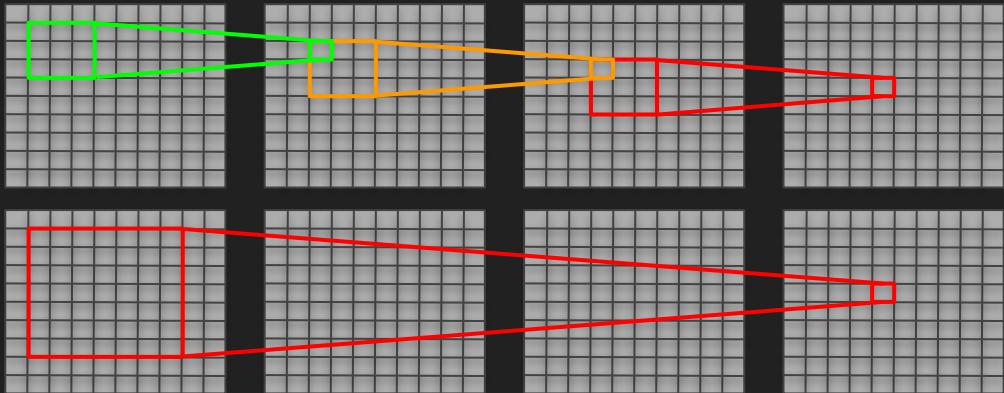
Design a network as a bunch of convolutional layers
to make predictions for pixels all at once!



Semantic segmentation - keeping same size

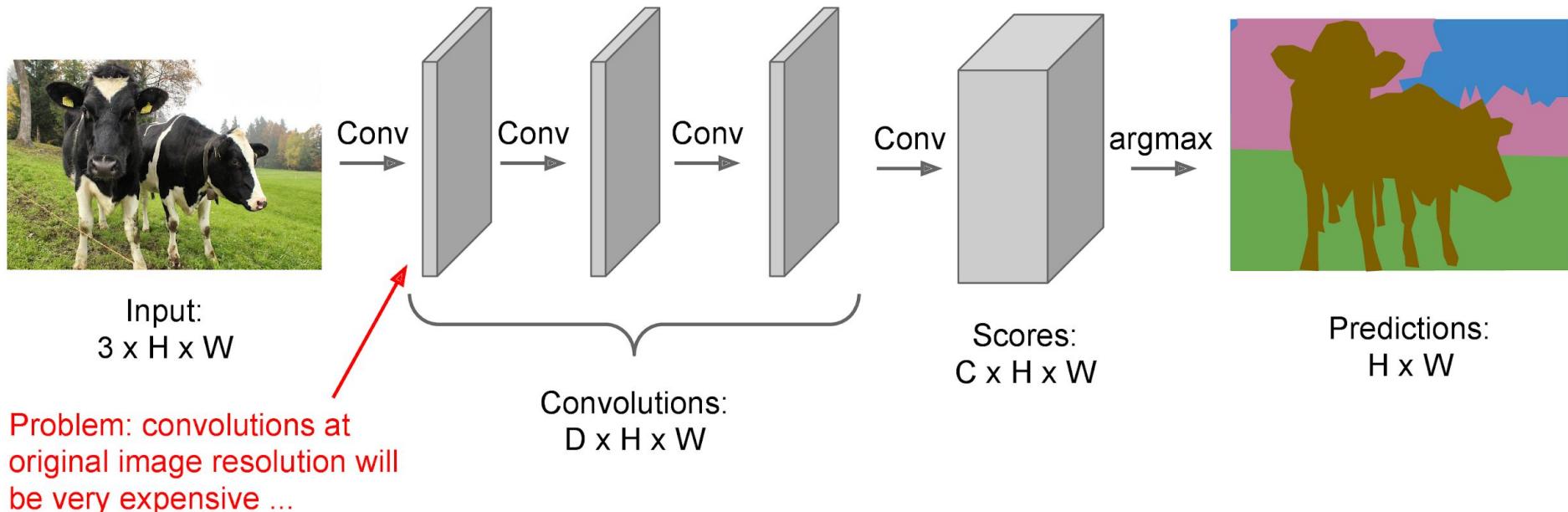
If $w \times h$ is unchanged:

- Receptive field is small.
 - Note: dilation.
- If the number of filters(channels) is large, the memory and computation cost is large.



Semantic Segmentation Idea: Fully Convolutional

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

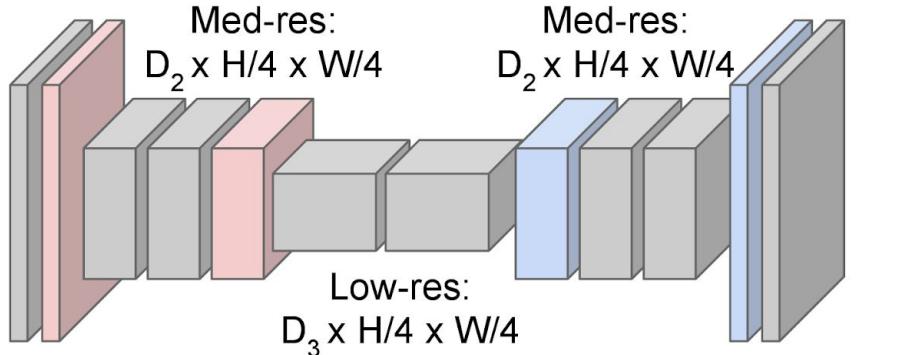


Semantic Segmentation Idea: Fully Convolutional

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

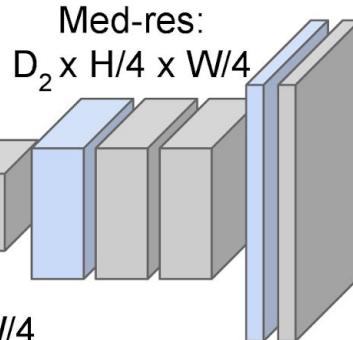


Input:
 $3 \times H \times W$



High-res:
 $D_1 \times H/2 \times W/2$

Low-res:
 $D_3 \times H/4 \times W/4$



High-res:
 $D_1 \times H/2 \times W/2$



Predictions:
 $H \times W$

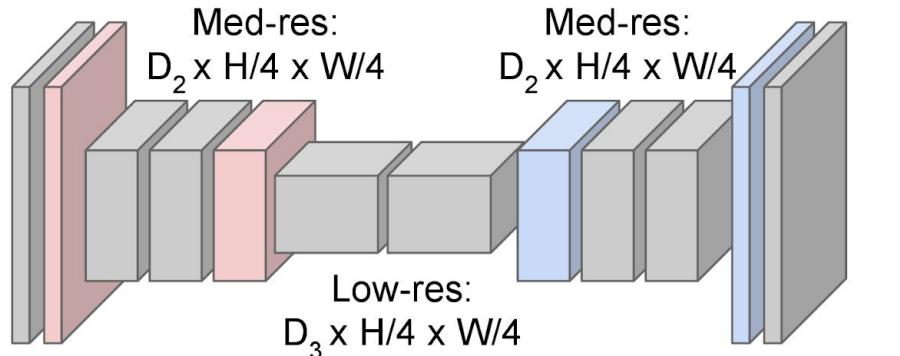
Semantic Segmentation Idea: Fully Convolutional

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!



High-res:
 $D_1 \times H/2 \times W/2$

Upsampling:
???



Predictions:
 $H \times W$

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

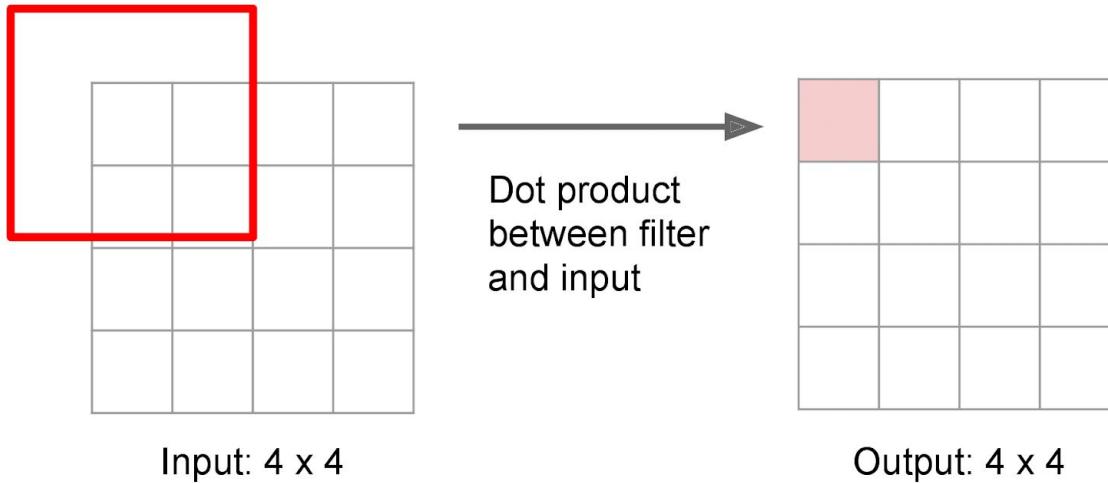
Learnable upsampling - explanation

Typical 3×3 convolution, stride 1 pad 1

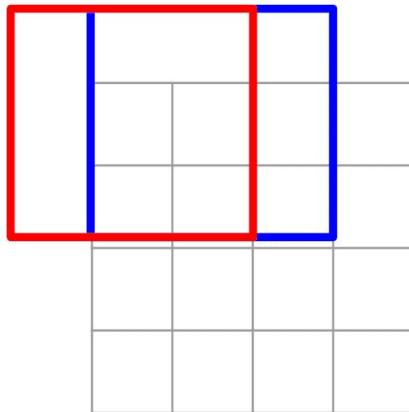
Input: 4×4

Output: 4×4

Typical 3×3 convolution, stride 1 pad 1

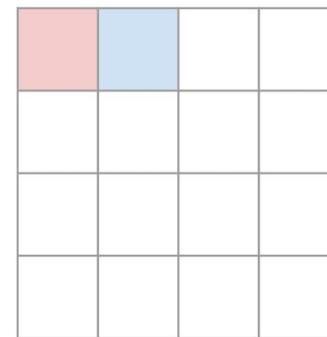


Typical 3×3 convolution, stride 1 pad 1



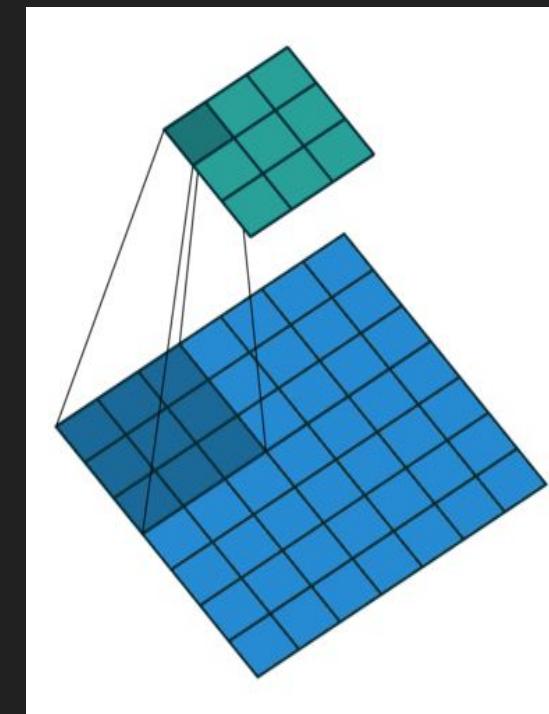
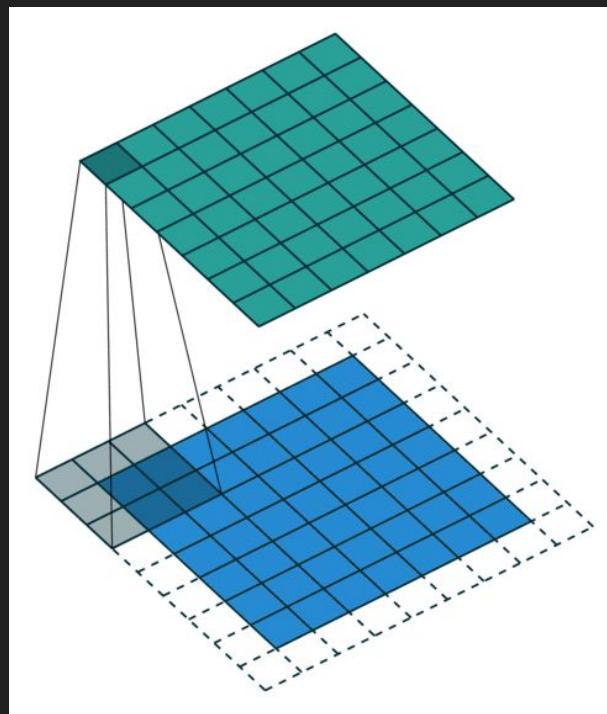
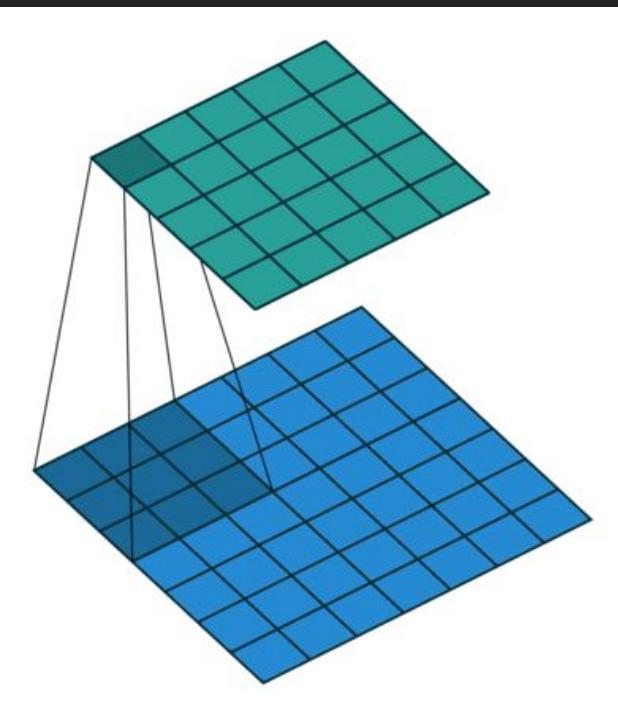
Input: 4×4

Dot product
between filter
and input

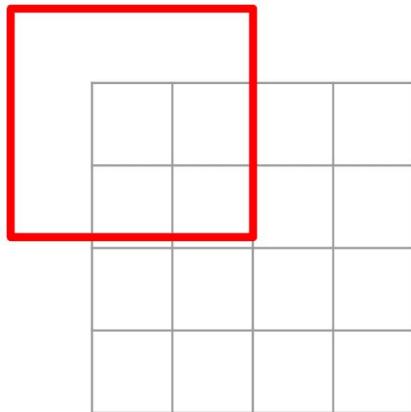


Output: 4×4

Regular convolution

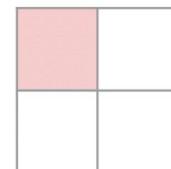


Typical 3×3 convolution, stride 2 pad 1



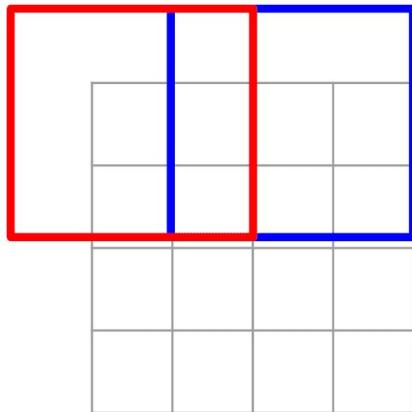
Input: 4×4

Dot product
between filter
and input



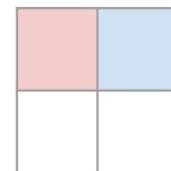
Output: 2×2

Typical 3×3 convolution, stride 2 pad 1



Input: 4×4

Dot product
between filter
and input



Output: 2×2

Learnable Upsampling: “Deconvolution”

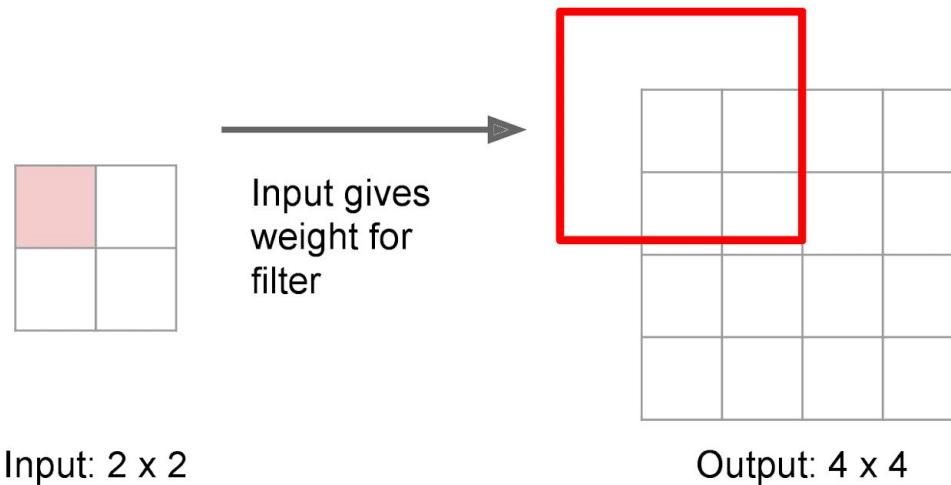
3 x 3 “deconvolution”, stride 2 pad 1

Input: 2 x 2

Output: 4 x 4

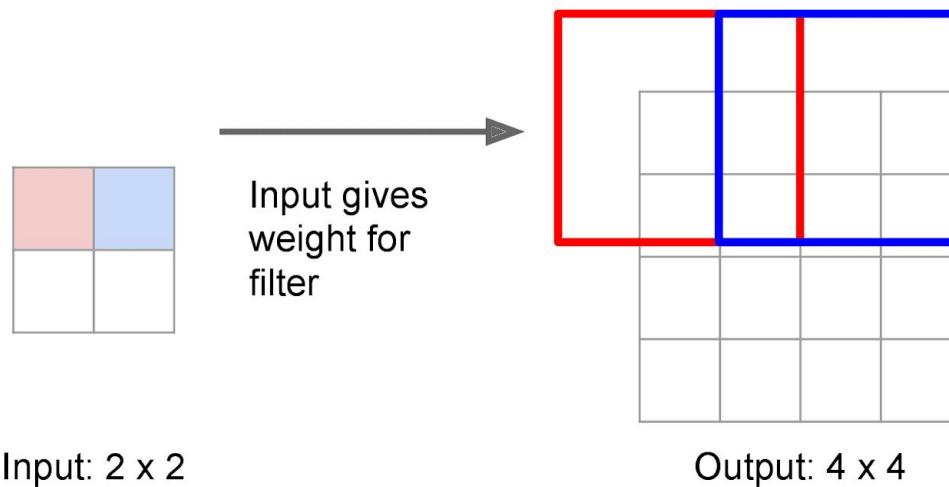
Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1

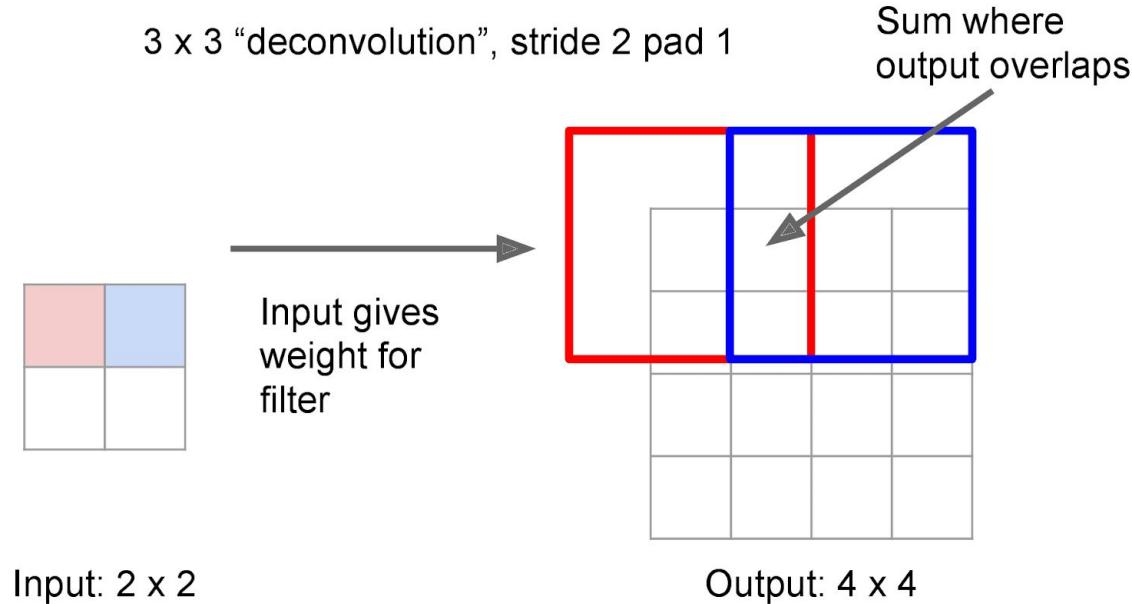


Learnable Upsampling: “Deconvolution”

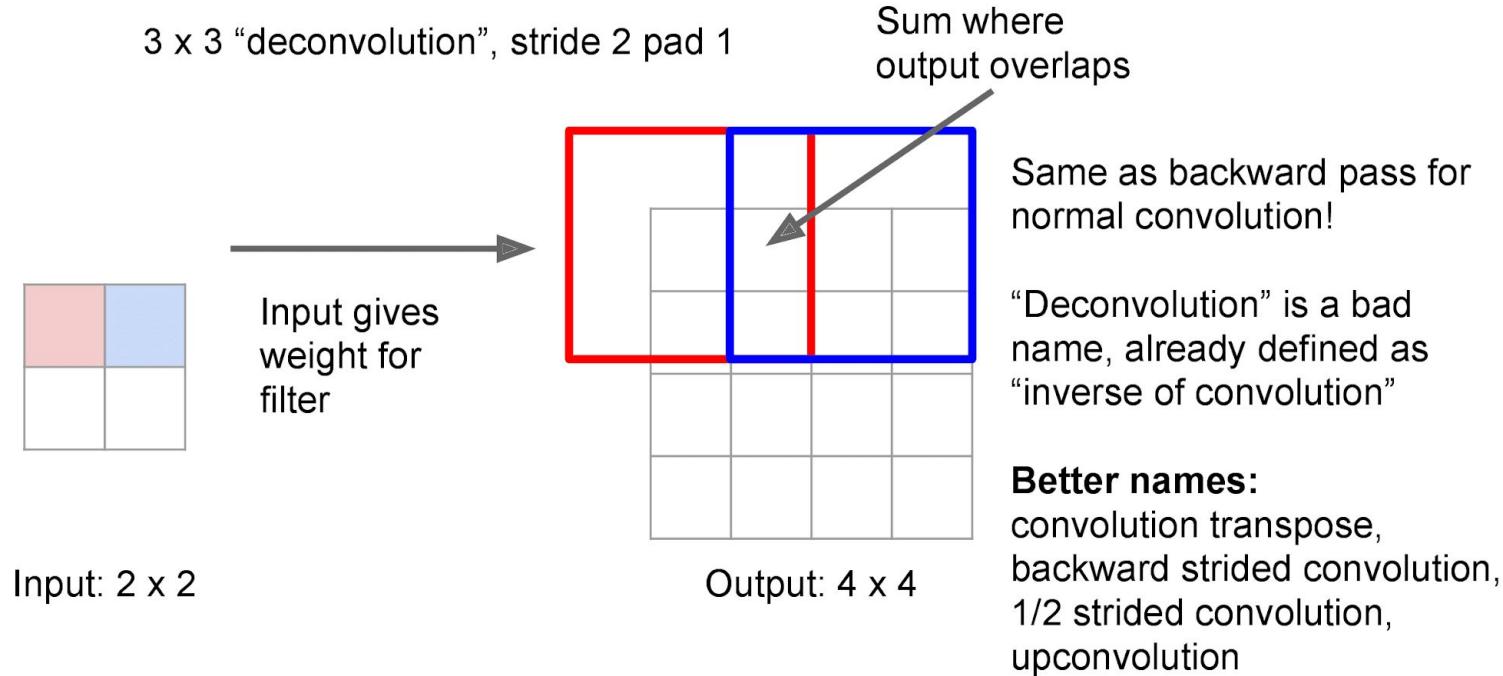
3 x 3 “deconvolution”, stride 2 pad 1



Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”

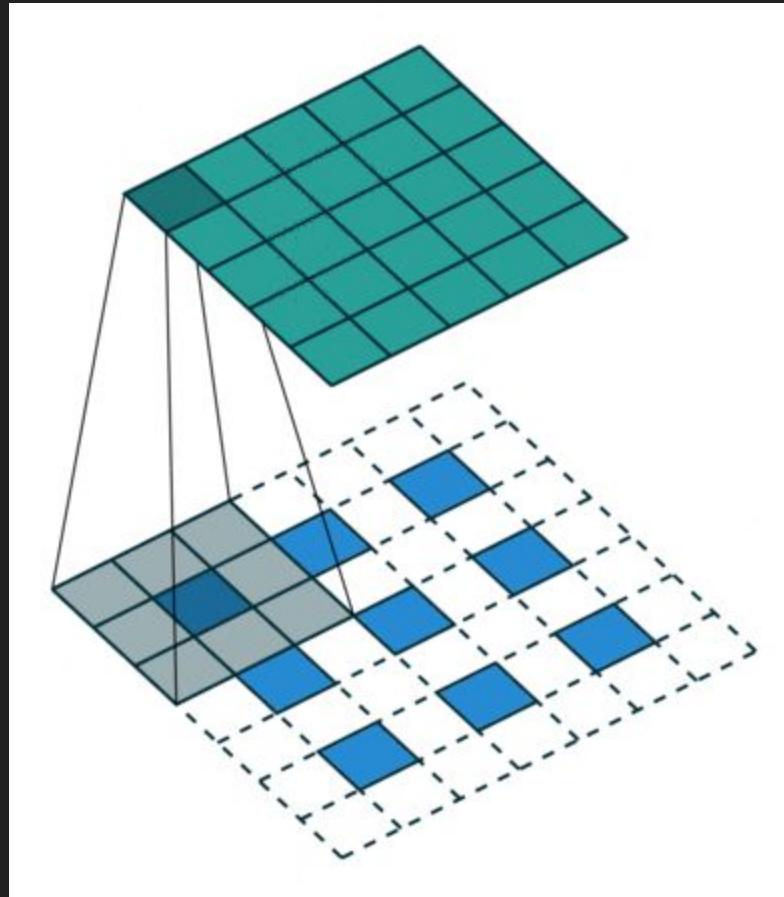




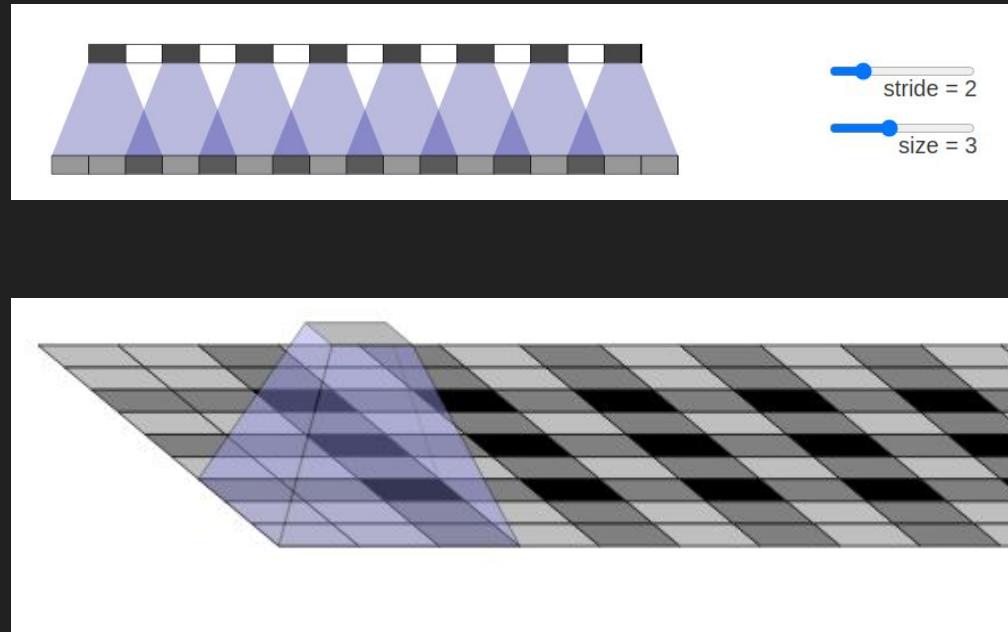
How many input pixel an output pixel depends on?

ⓘ Start presenting to display the poll results on this slide.

Learnable upsampling, alternative view

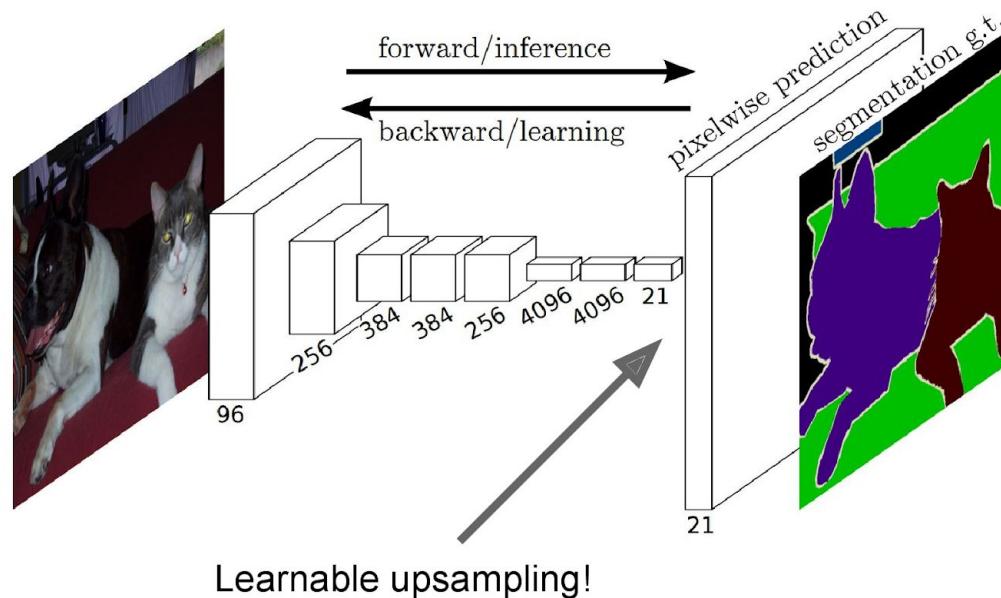


Checkerboard pattern in images generated by convolution transpose



Learnable upsampling - explanation end

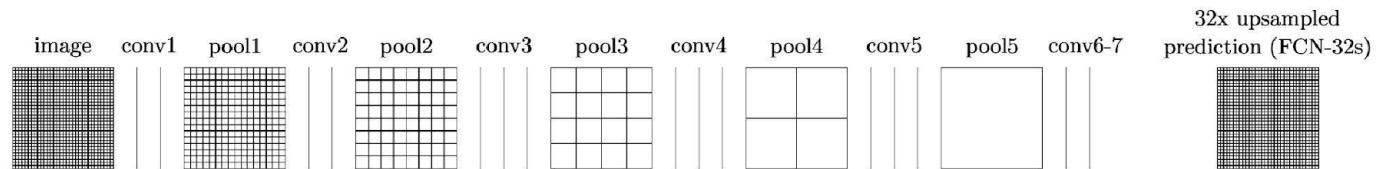
Semantic Segmentation: Upsampling



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Slide taken from CS231n@stanford

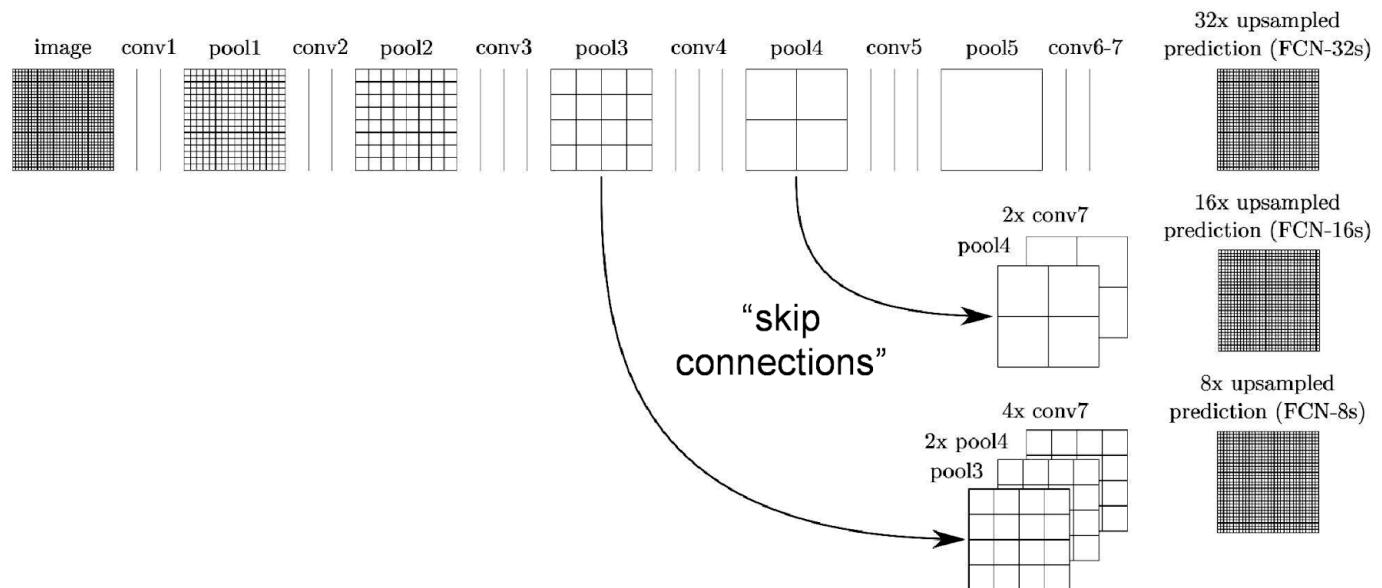
Semantic Segmentation: Upsampling



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

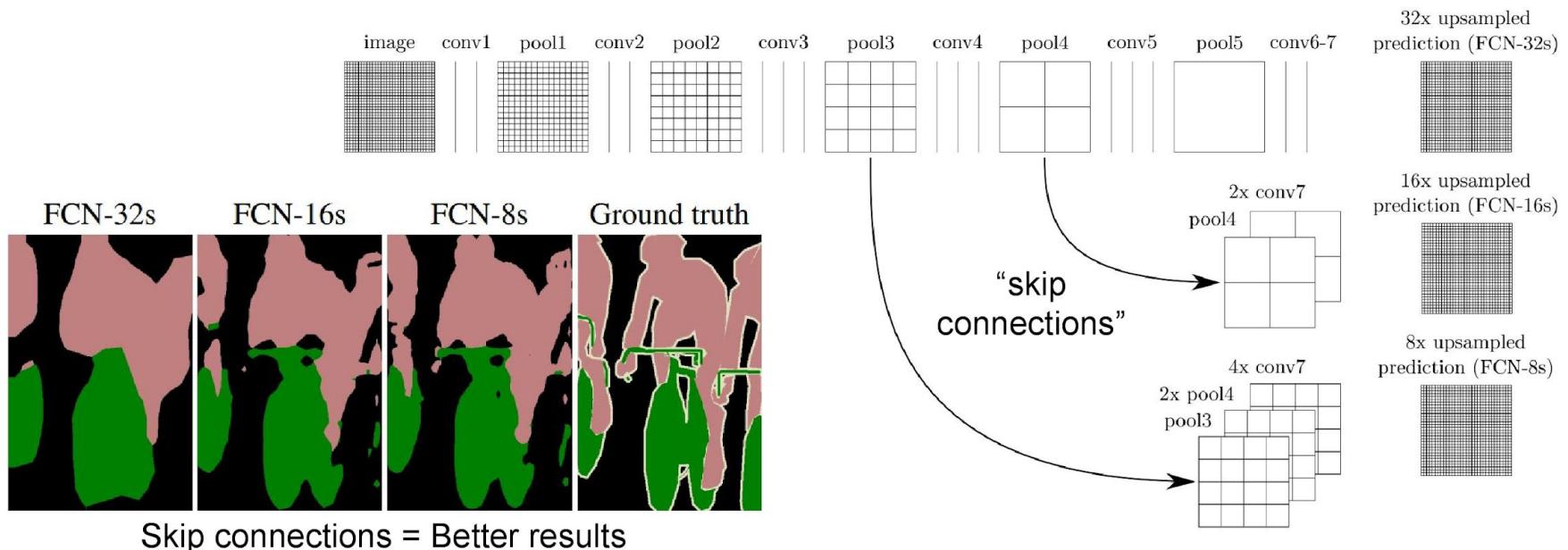
Slide taken from CS231n@stanford

Semantic Segmentation: Upsampling



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

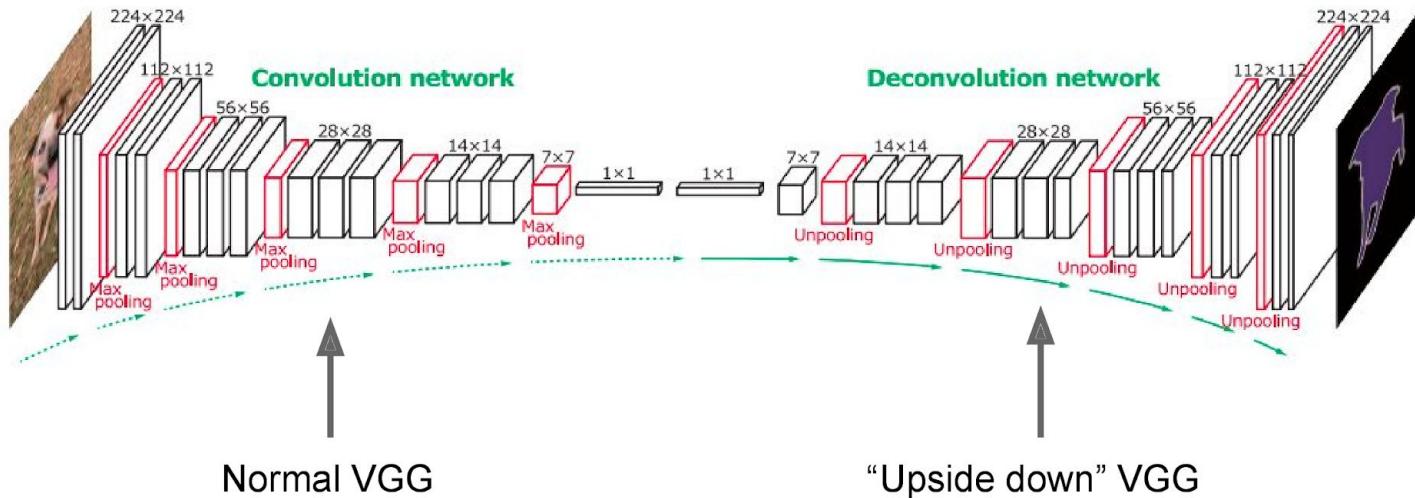
Semantic Segmentation: Upsampling



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Slide taken from CS231n@stanford

Semantic Segmentation: Upsampling



Noh et al, “Learning Deconvolution Network for Semantic Segmentation”, ICCV 2015

6 days of training on Titan X...

In-Network 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

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

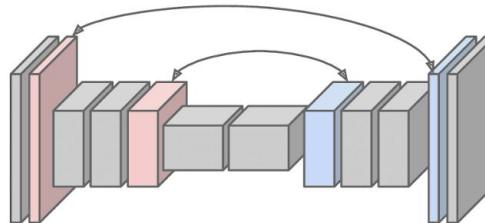
1	2
3	4

Rest of the network

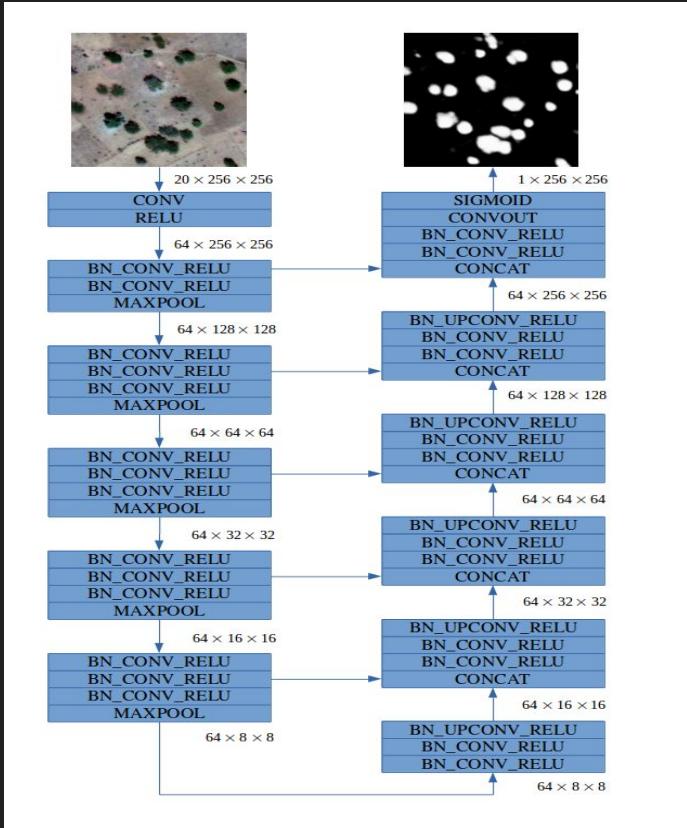
0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers



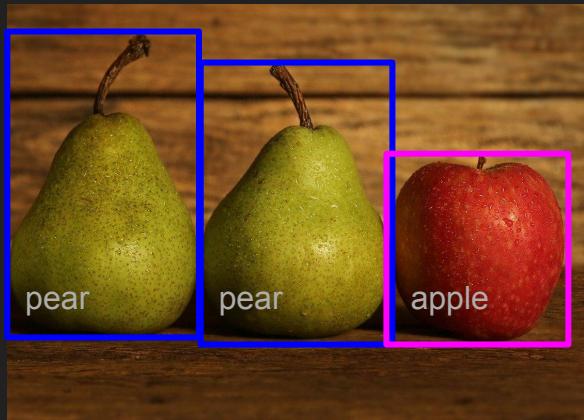
Semantic segmentation, take-home-message



- Use U-net as a default starting point for semantic segmentation (sample architecture on the left)
- Fully convolutional solution (no region proposals)

Object detection

Object detection

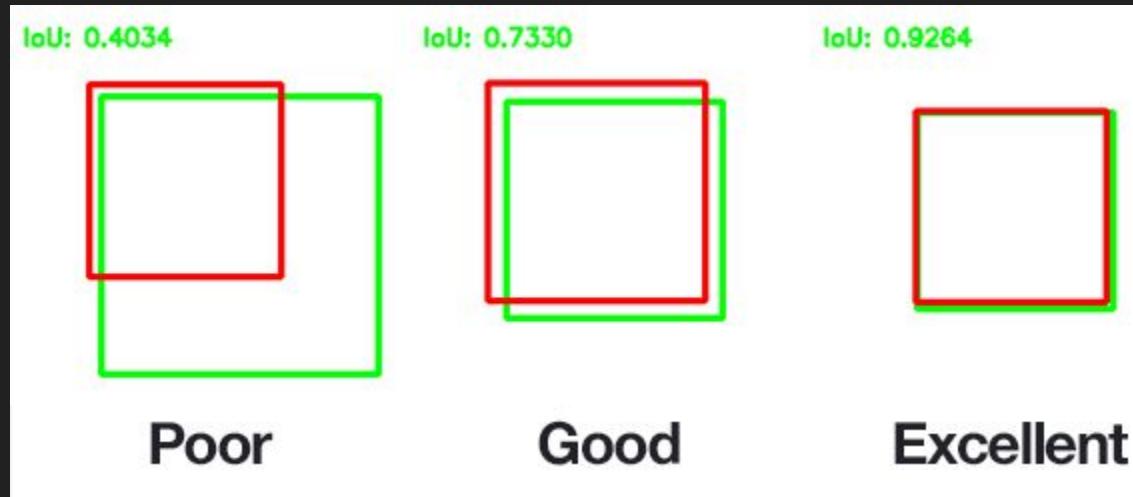
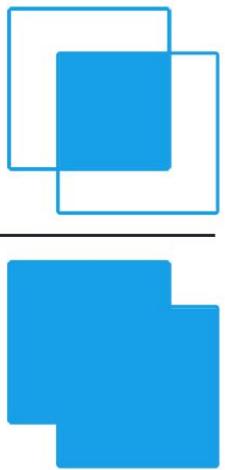


Object detection

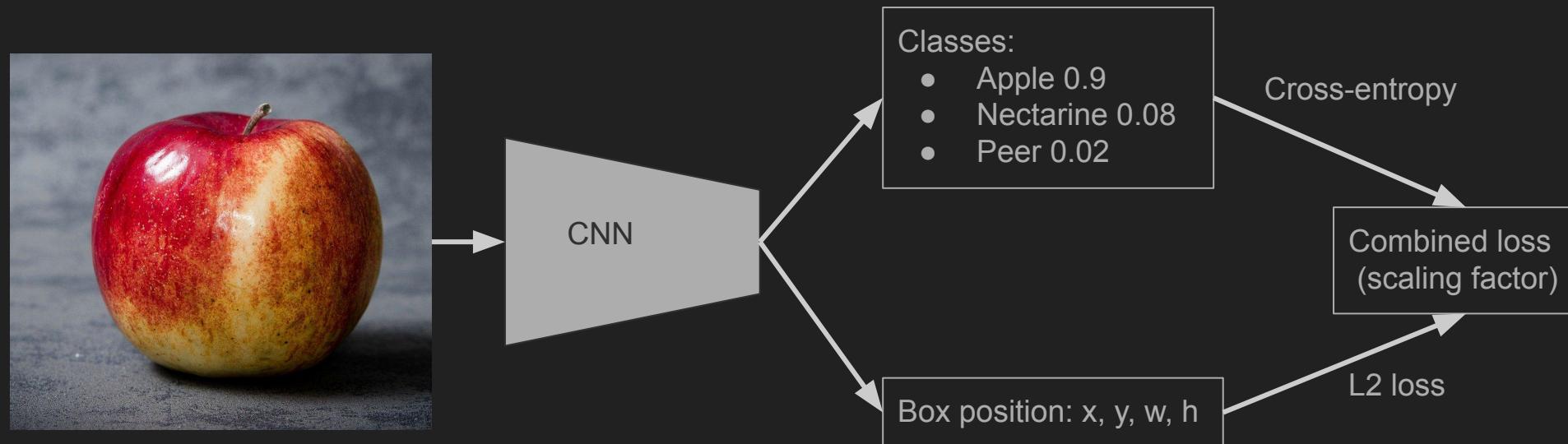
- Localize all objects with bounding boxes
- Axis aligned
- Class for each box.

Intersection over union (IoU)

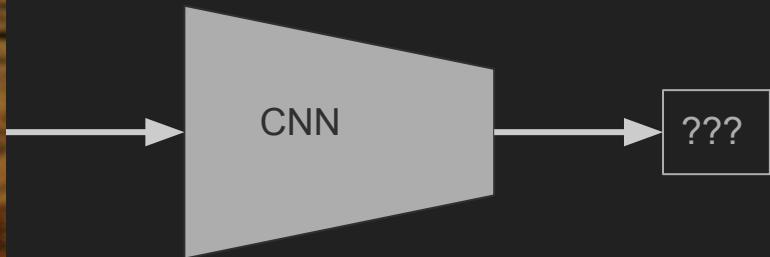
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



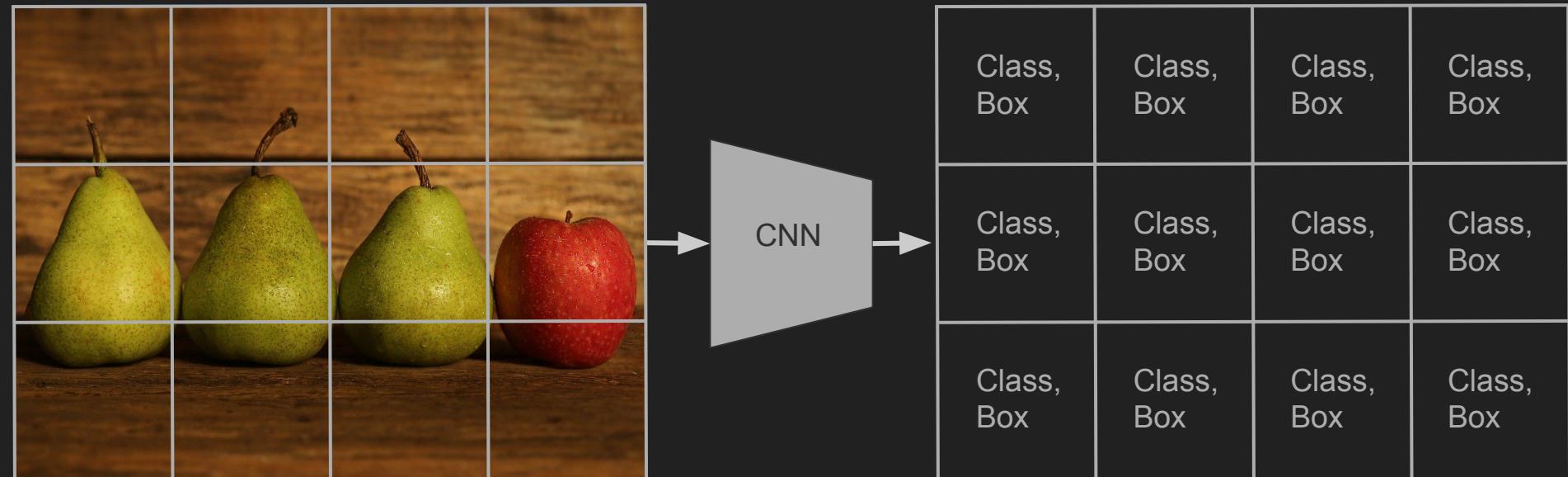
Single object case



Multi object case

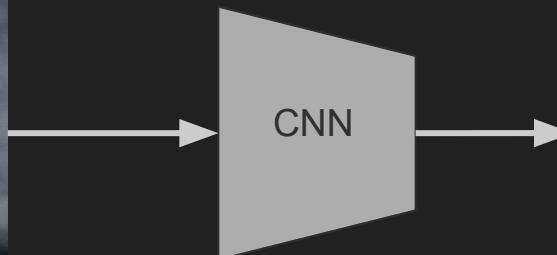
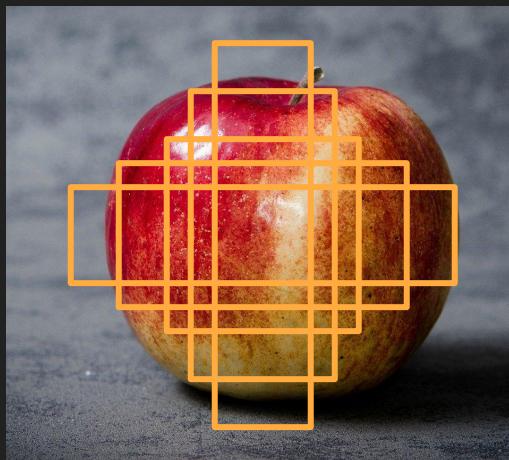


Idea 1: grid of predictions



Idea 2: Anchor boxes

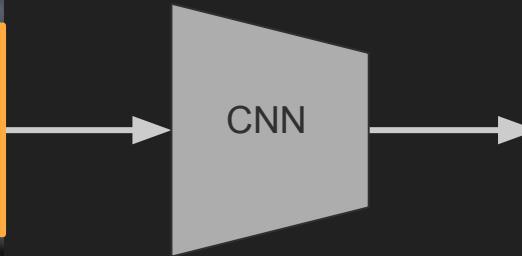
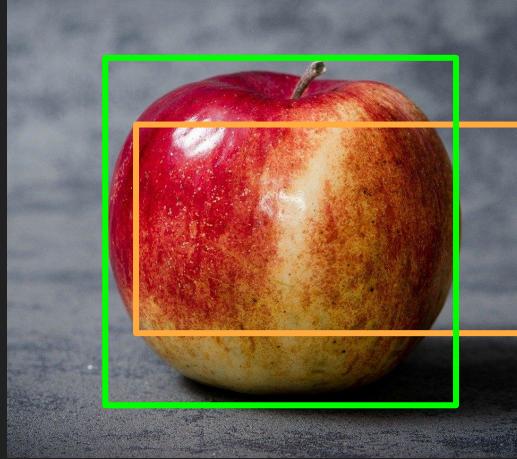
- Consider a set of predefined rectangles: **anchor boxes** (thousands or more).
- Always the same set of boxes, regardless of the image.



For each anchor box:

- Classify into $C+1$ classes
(C classes + background)
- Refinement: scale+offset

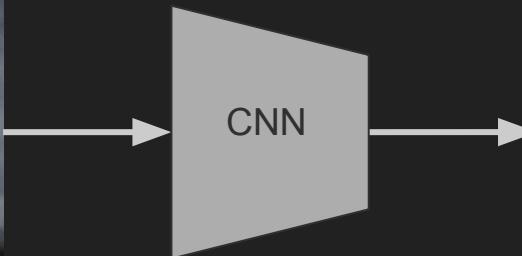
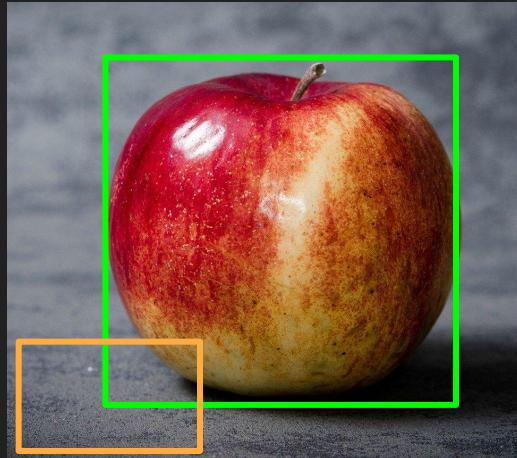
Idea 2: anchor boxes



Class: apple

Refinement:

- Scale width = 0.915
- Scale height = 1.66
- Offset x = -0.12
- Offset y = 0.02

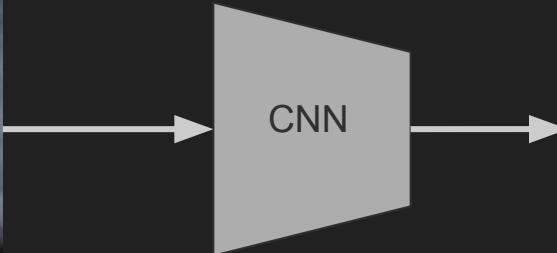
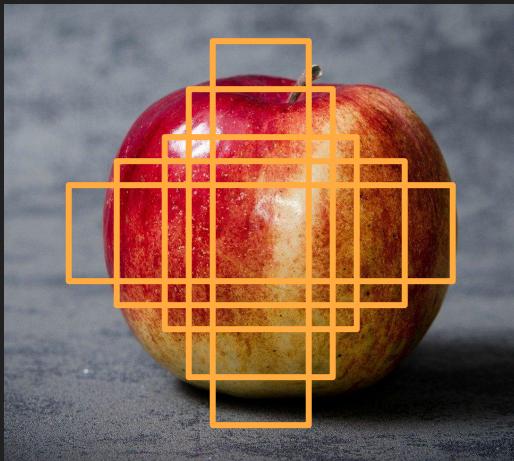


Class: background
(because IoU < threshold)

Refinement:

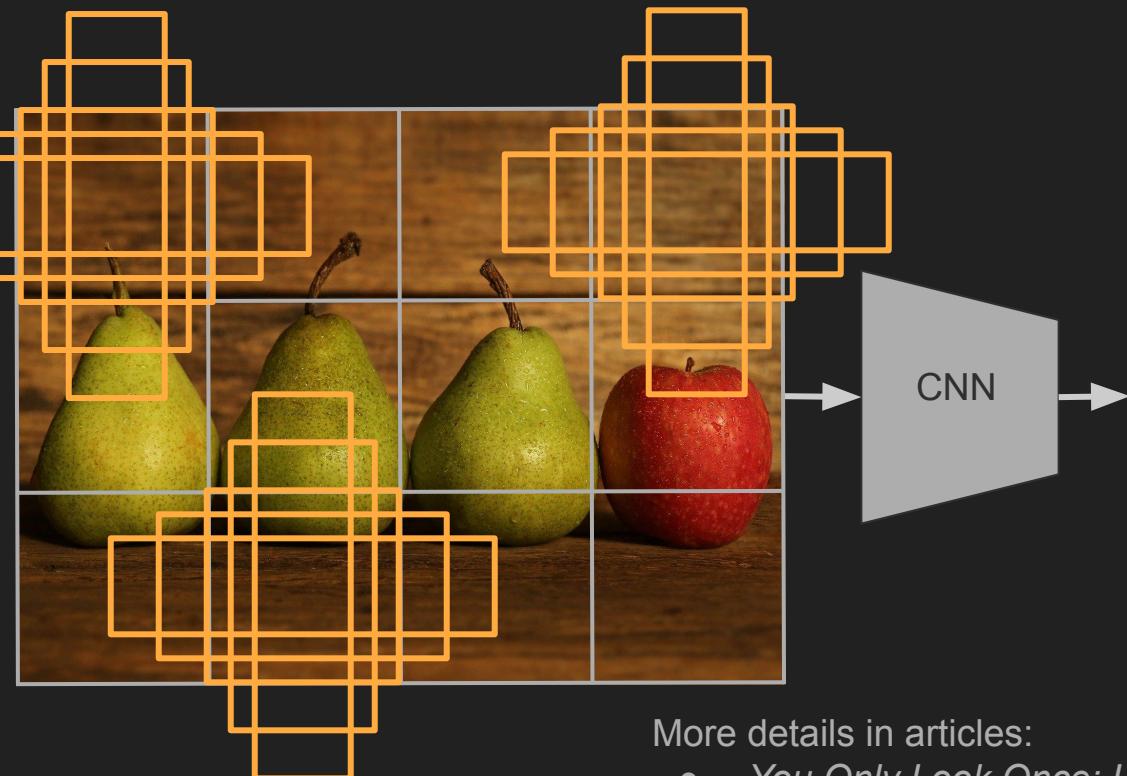
- Scale width = None
- Scale height = None
- Offset x = None
- Offset y = None

Idea 2: Anchor boxes



- Classify into $C+1$ classes
(C classes + background)
- Refinement: scale+offset

Idea 1 + Idea 2 = Single Stage Detector



More details in articles:

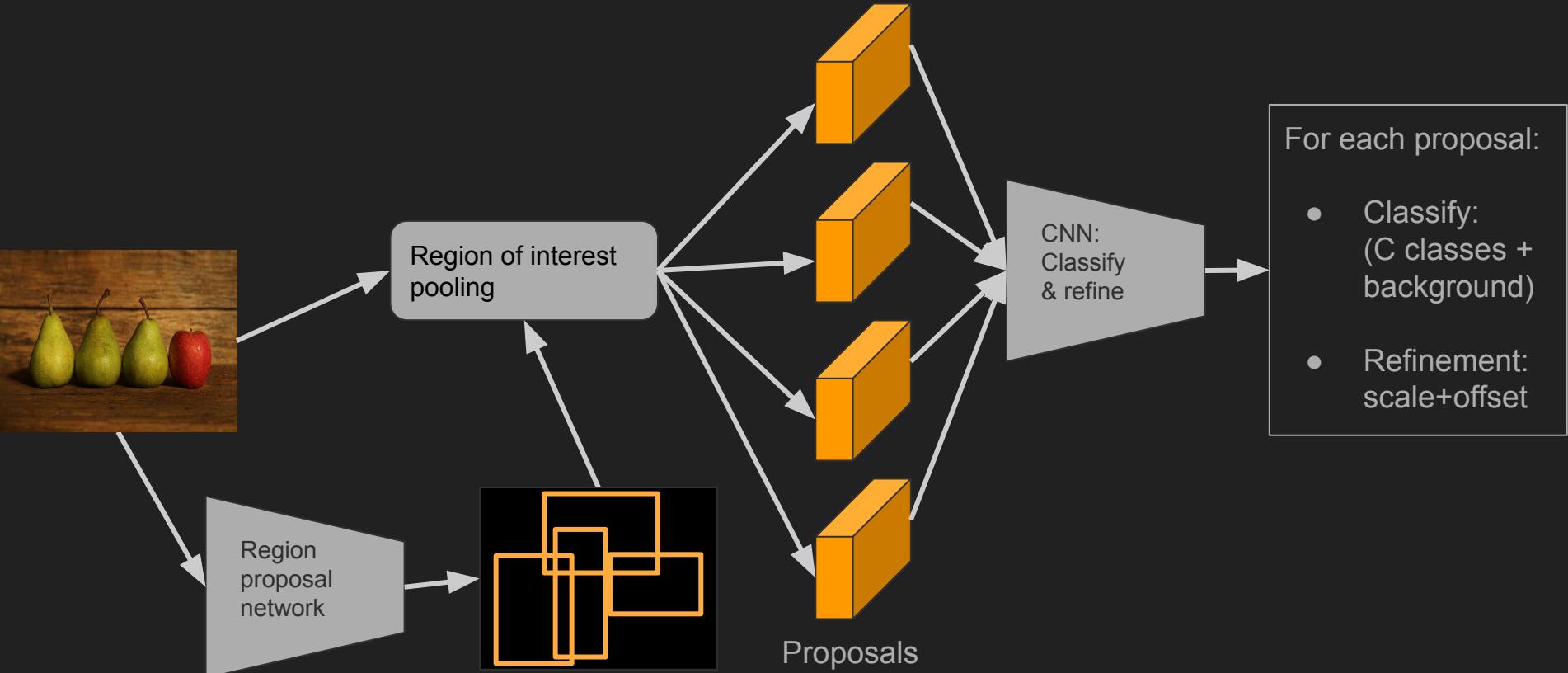
- *You Only Look Once: Unified, Real-Time Object Detection*
- *SSD: Single Shot MultiBox Detector*
- *Focal Loss for Dense Object Detection*



Assume: 5x5 grid, 10 anchor boxes, C=3 classes (apples, pears, oranges) What is the size of the output tensor for a single input image? Remember about the class representing background (not detection).

- ⓘ Start presenting to display the poll results on this slide.

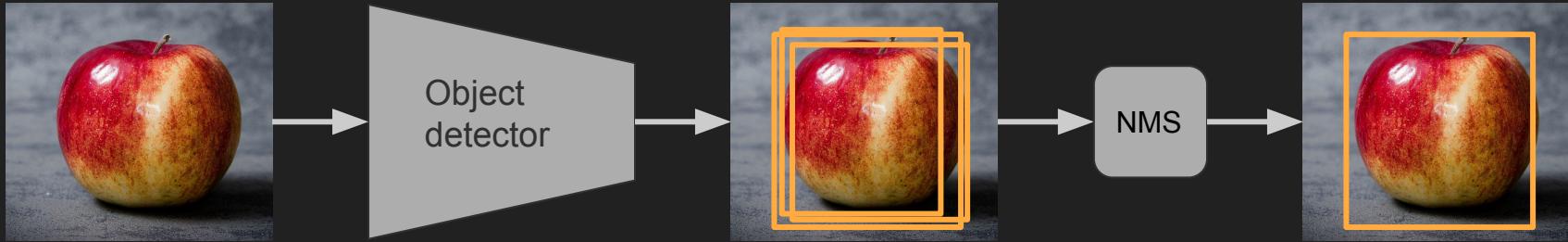
2-stage detection (sketch)



2-stage detection (sketch)

- Stage one: generate region proposals (object detection without classification)
- For each proposal run the network again on the patch:
 - Optimization: instead of running on the patch of the original image, run on a patch of features generated by the backbone.

Non-maximum suppression



NMS:

- Sort detection by decreasing prediction scores.
- Process detected boxes:
 - If a given object has $\text{IoU} > \text{threshold}$ with some unprocessed candidate, remove the candidate.
- In what scenario NMS is likely to remove valid predictions?

Non-maximum suppression



In what scenario NMS is likely to remove valid predictions?

- Bounding boxes of two disjoint objects have big IoU.

Anchor free detectors: pixel level box predictions

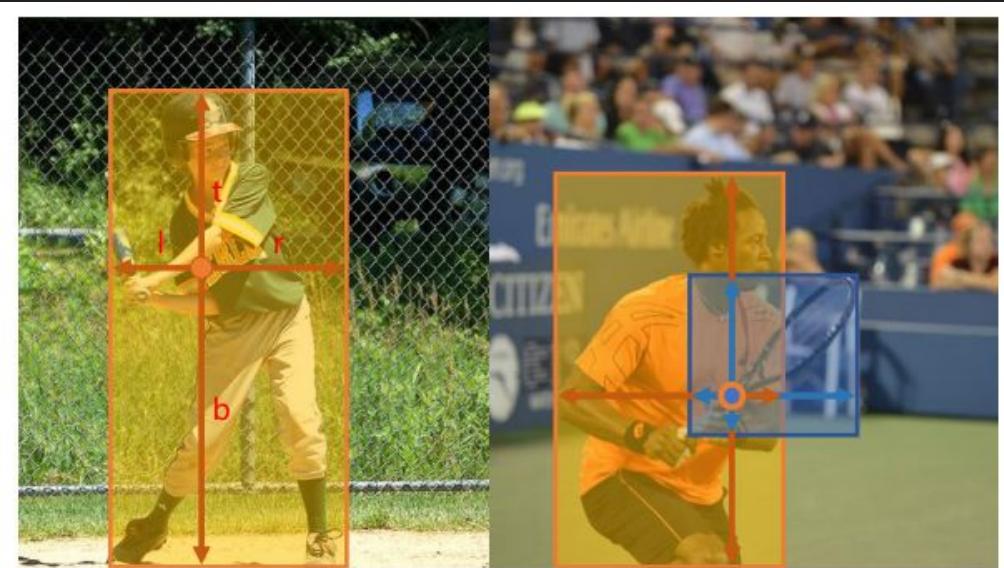
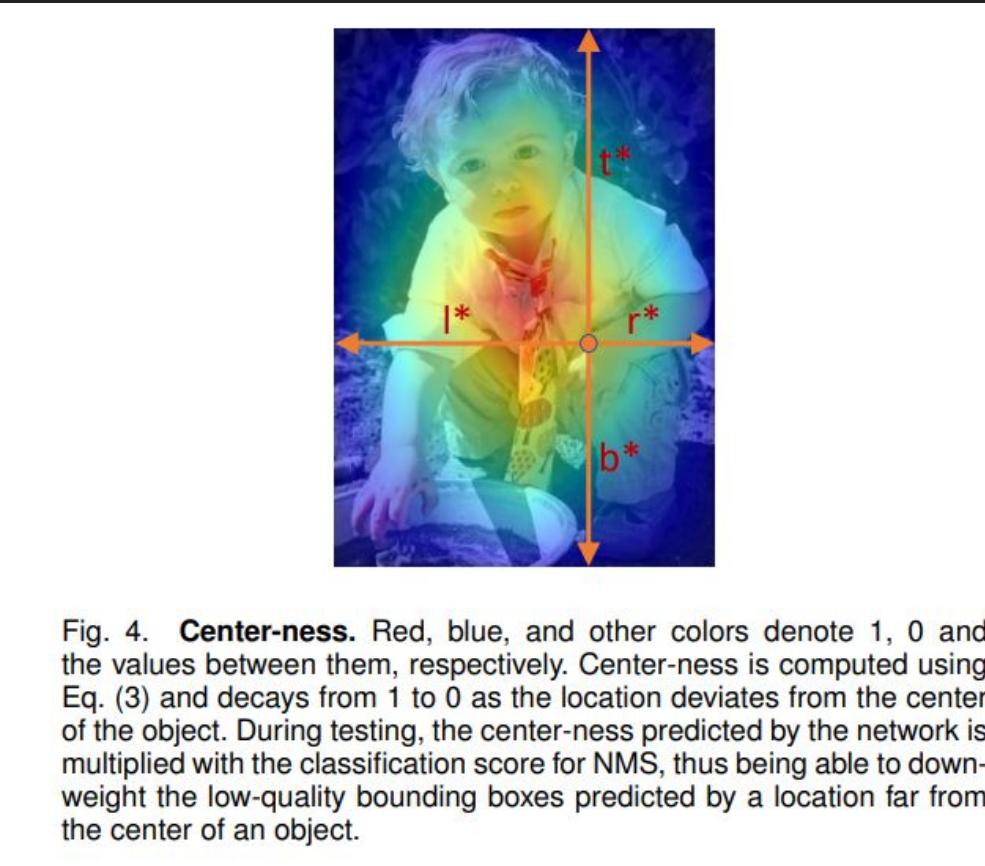


Fig. 1. Overall concept of FCOS. As shown in the left image, FCOS works by predicting a 4D vector (l, t, r, b) encoding the location of a bounding box at each foreground pixel (supervised by ground-truth bounding box information during training). The right plot shows that when a location residing in multiple bounding boxes, it can be ambiguous in terms of which bounding box this location should regress.

[FCOS paper](#)

Anchor free detectors: pixel level box predictions



[FCOS paper](#)

Object detection architectures - summary

- One stage detectors:
 - fast, less precise.
- Two stage detectors:
 - often give best results,
 - typically more computationally intensive.
- Anchor-free detectors (FCOS, CornerNet, CenterNet).

Feedback is a gift

<https://tinyurl.com/dnn-2025-11-12>

