

CS163: Deep Learning for Computer Vision

Lecture 14: Detection + Segmentation

Announcement

- Assignment 3 dues on Nov 24 (this Sunday)
- Assignment 4 will be released today
- Next Wednesday afternoon's lecture (right before thanksgiving break) will be independent study to watch records of research spotlight talks
 - We invite several student researchers (PhD and undergrad students) to talk about their research projects on computer vision
 - Each talk will be 5 – 10 mins, which you will watch and answer one question about the specific project, also you need to raise a question

Last Time: Computer Vision Tasks

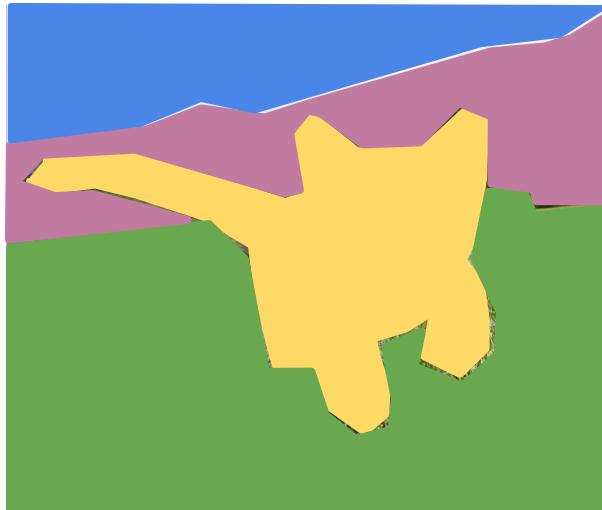
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Objects

Instance Segmentation



DOG, DOG, CAT

[This image](#) is CCO public domain

Last Time: Object Detection

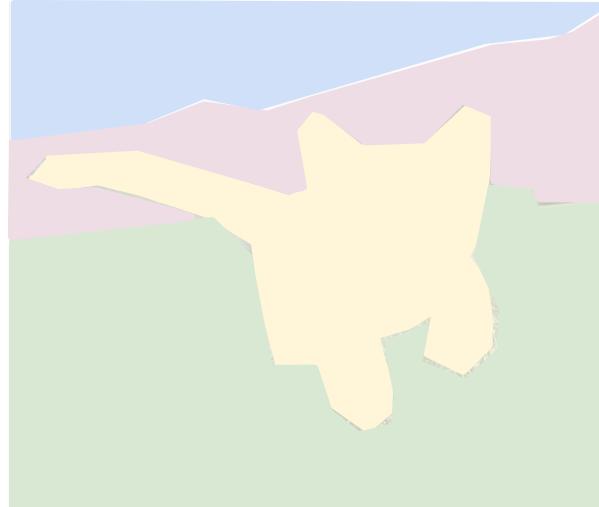
Classification



CAT

No spatial extent

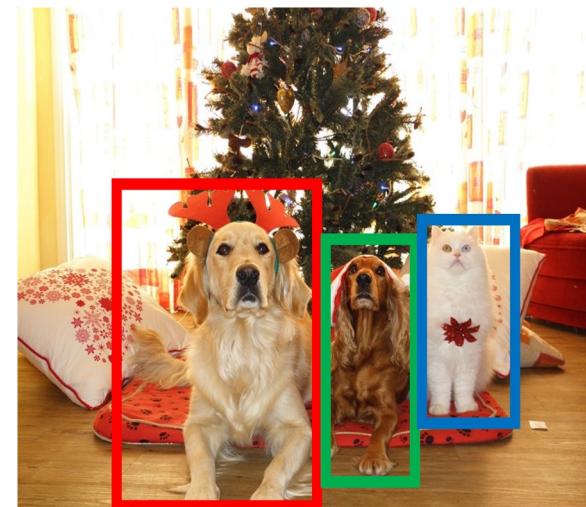
Semantic Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Objects

Instance Segmentation



DOG, DOG, CAT

Object Detection: Impact of Deep Learning

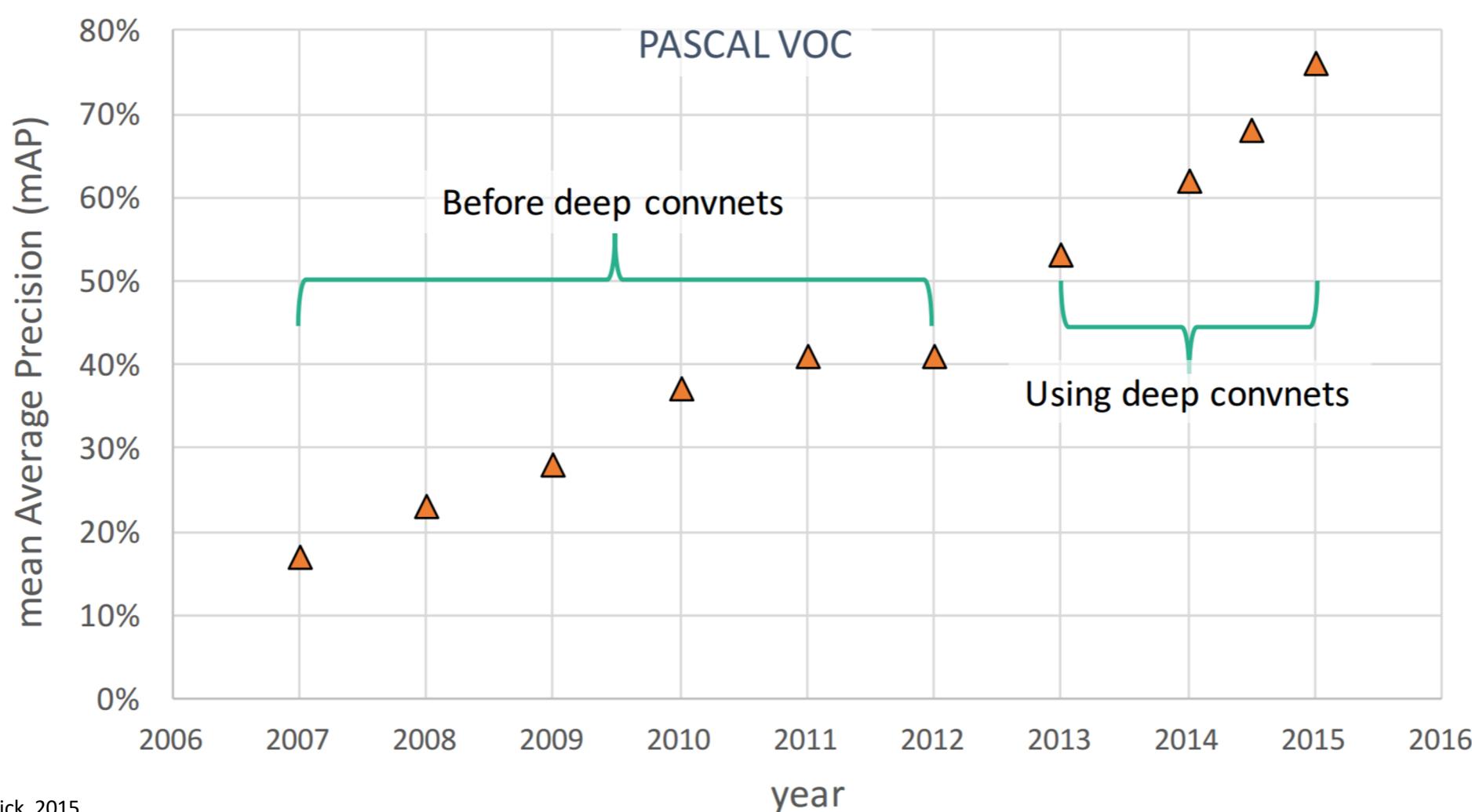
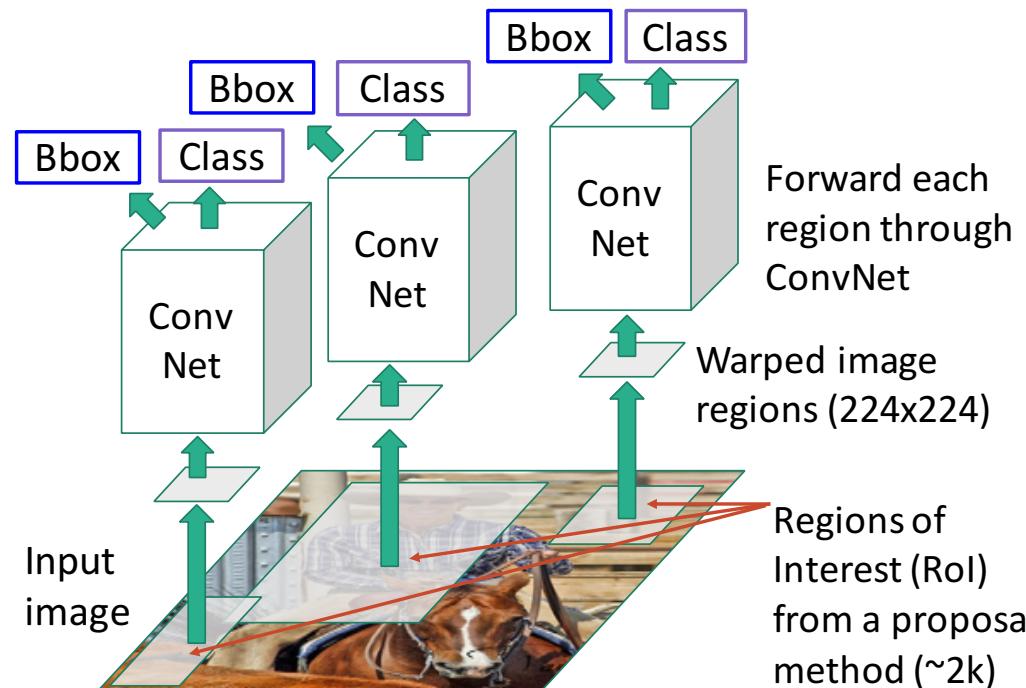


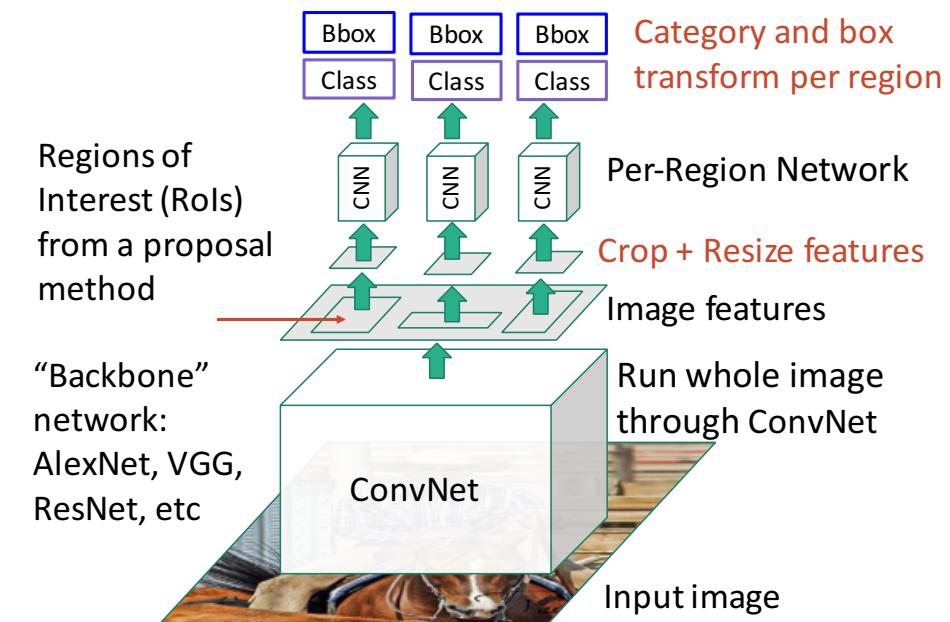
Figure copyright Ross Girshick, 2015.
Reproduced with permission.

Last Time: Object Detection Methods

“Slow” R-CNN: Run CNN independently for each region

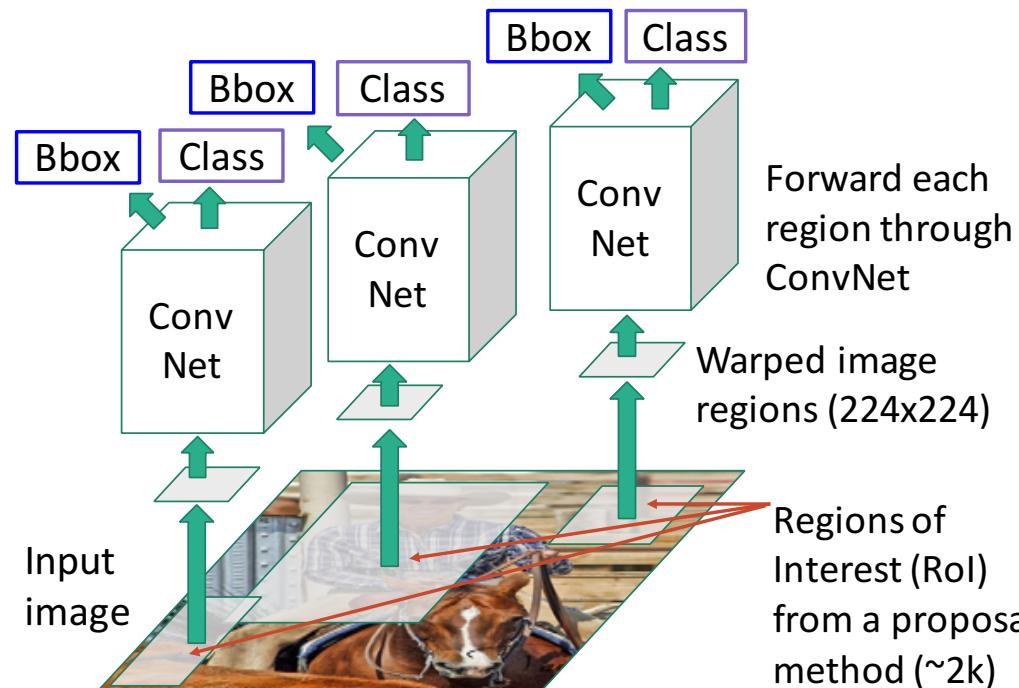


Fast R-CNN: Apply differentiable cropping to shared image features

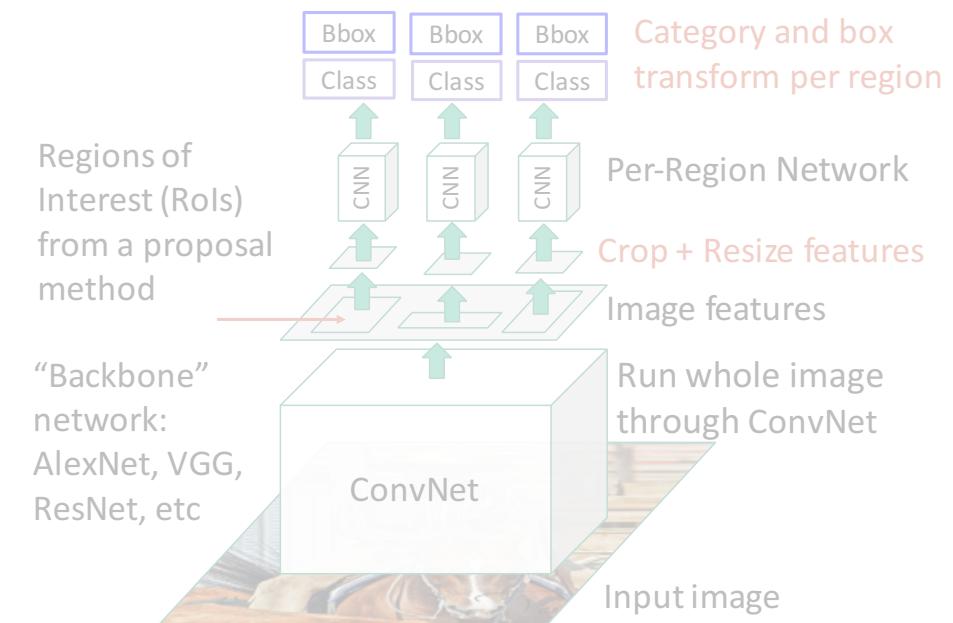


Recap: Slow R-CNN Training

“Slow” R-CNN: Run CNN independently for each region

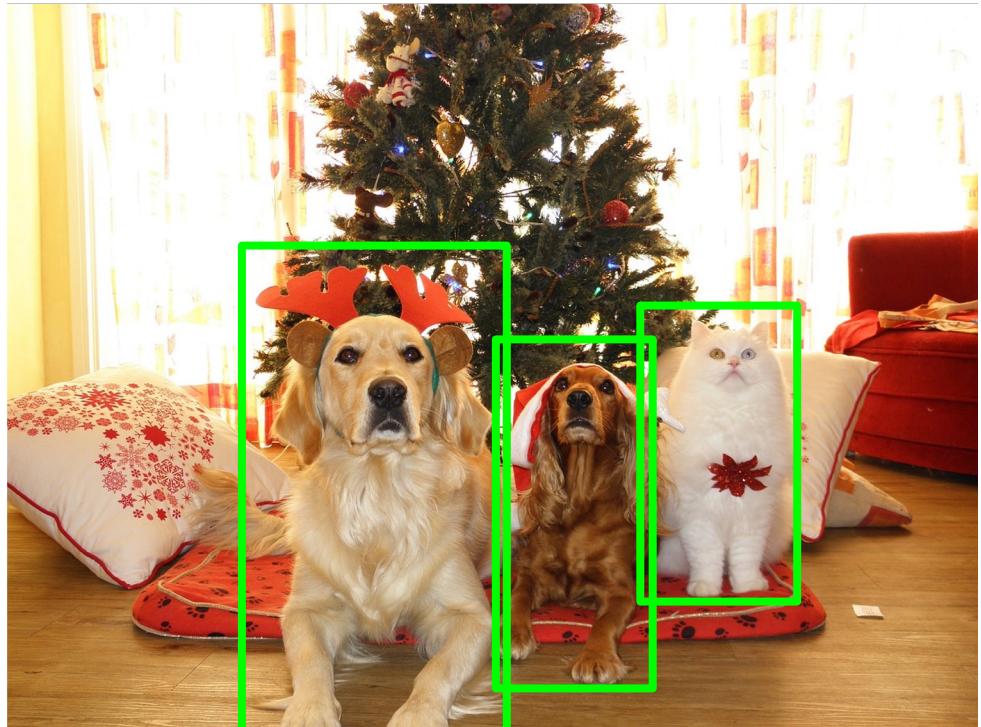


Fast R-CNN: Apply differentiable cropping to shared image features



“Slow” R-CNN Training

Input Image

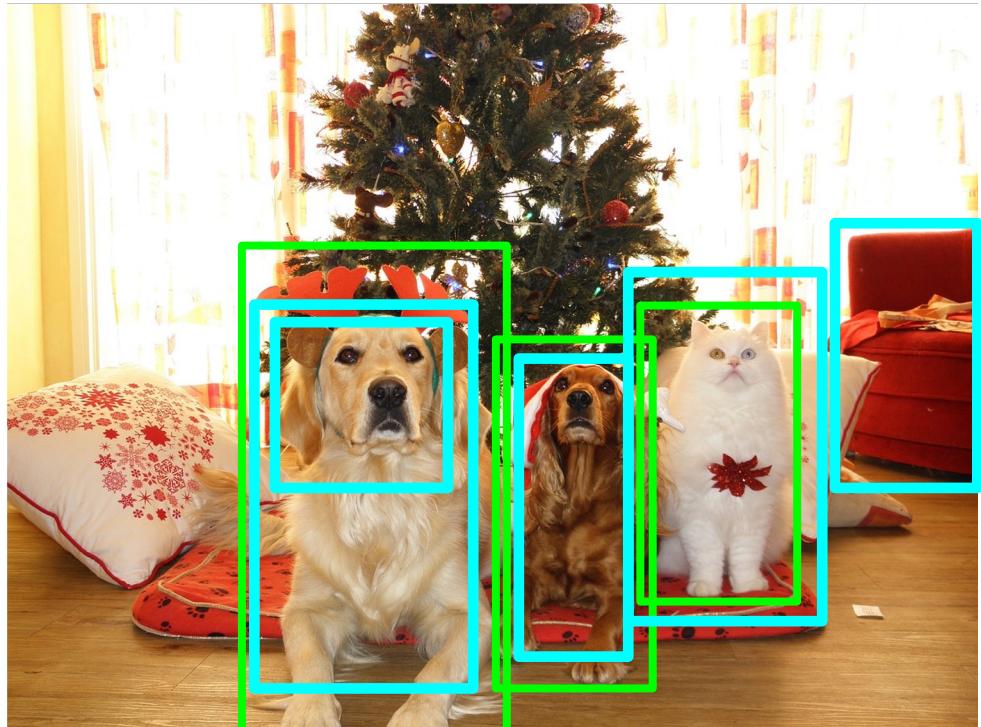


Ground-Truth boxes

[This image](#) is CC0 public domain

“Slow” R-CNN Training

Input Image



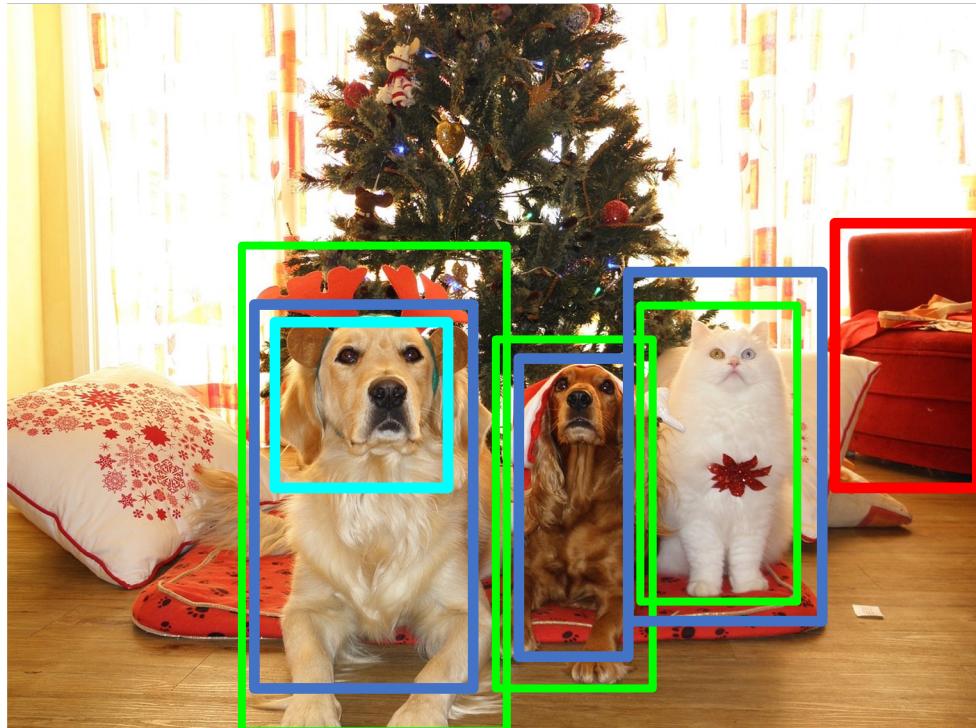
Ground-Truth boxes

Region Proposals

[This image](#) is CC0 public domain

“Slow” R-CNN Training

Input Image



Categorize each region proposal as positive, negative, or neutral based on overlap with ground-truth boxes

GT Boxes

Positive

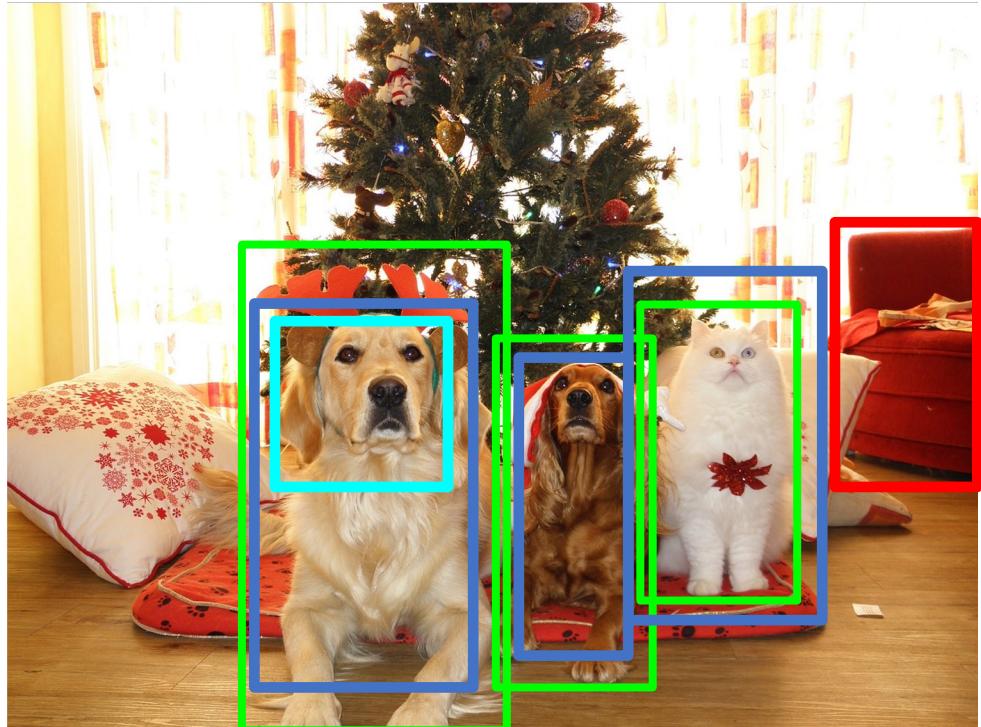
Neutral

Negative

[This image](#) is CC0 public domain

“Slow” R-CNN Training

Input Image



GT Boxes

Positive

Neutral

Negative

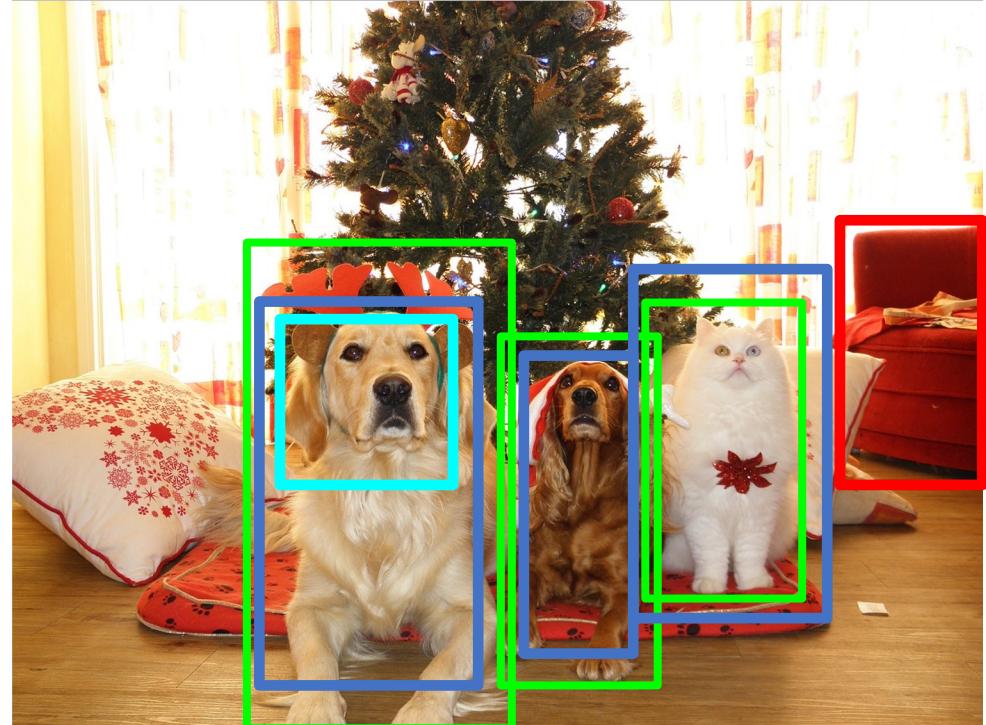


Crop pixels from each positive and negative proposal, resize to 224 x 224

[This image](#) is CC0 public domain

“Slow” R-CNN Training

Input Image



GT Boxes

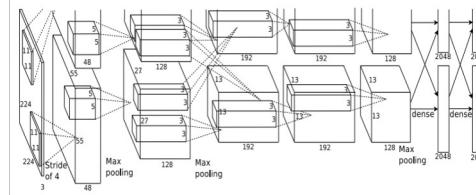
Positive

Neutral

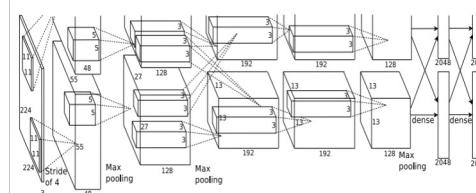
Negative



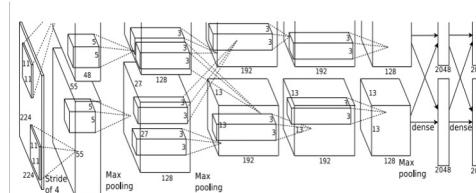
Run each region through CNN. For positive boxes predict class and box offset; for negative boxes just predict background class



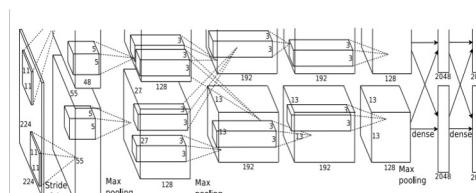
Class target: Dog
Box target: →



Class target: Cat
Box target: →



Class target: Dog
Box target: →



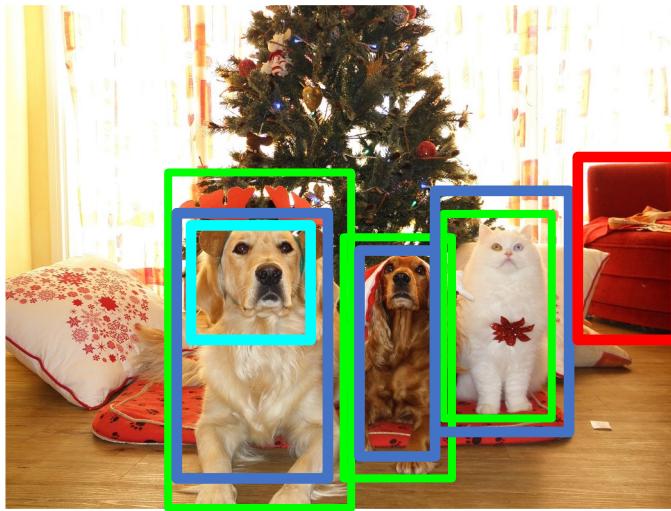
Class target: Background
Box target: None

This image is CCO public domain

Fast R-CNN Training

Crop features for each region, use them to predict class and box targets per region

Input Image



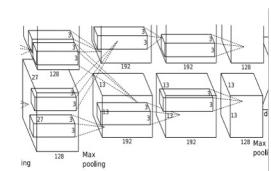
GT Boxes

Positive

Neutral

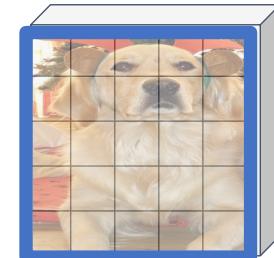
Negative

Image Features

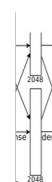


A photograph of a light-colored dog sitting on a wooden floor in front of a decorated Christmas tree. The image is overlaid with a 4x6 grid. Several objects are highlighted with colored bounding boxes: a red and white patterned object on the left (green box), the dog's head and upper body (blue box), the Christmas tree (green box), the dog's lower body and the floor (blue box), and a red object on the right (red box). A small white tag is visible near the bottom right corner of the image.

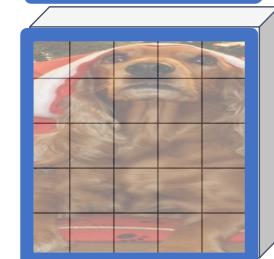
Class target: Dog
Box target: 



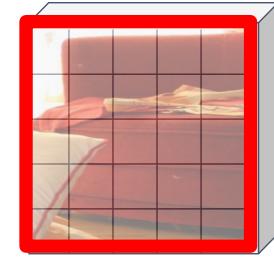
Class target: Cat
Box target: →



Class target: Dog
Box target: ➔



Class target: Background
Box target: None



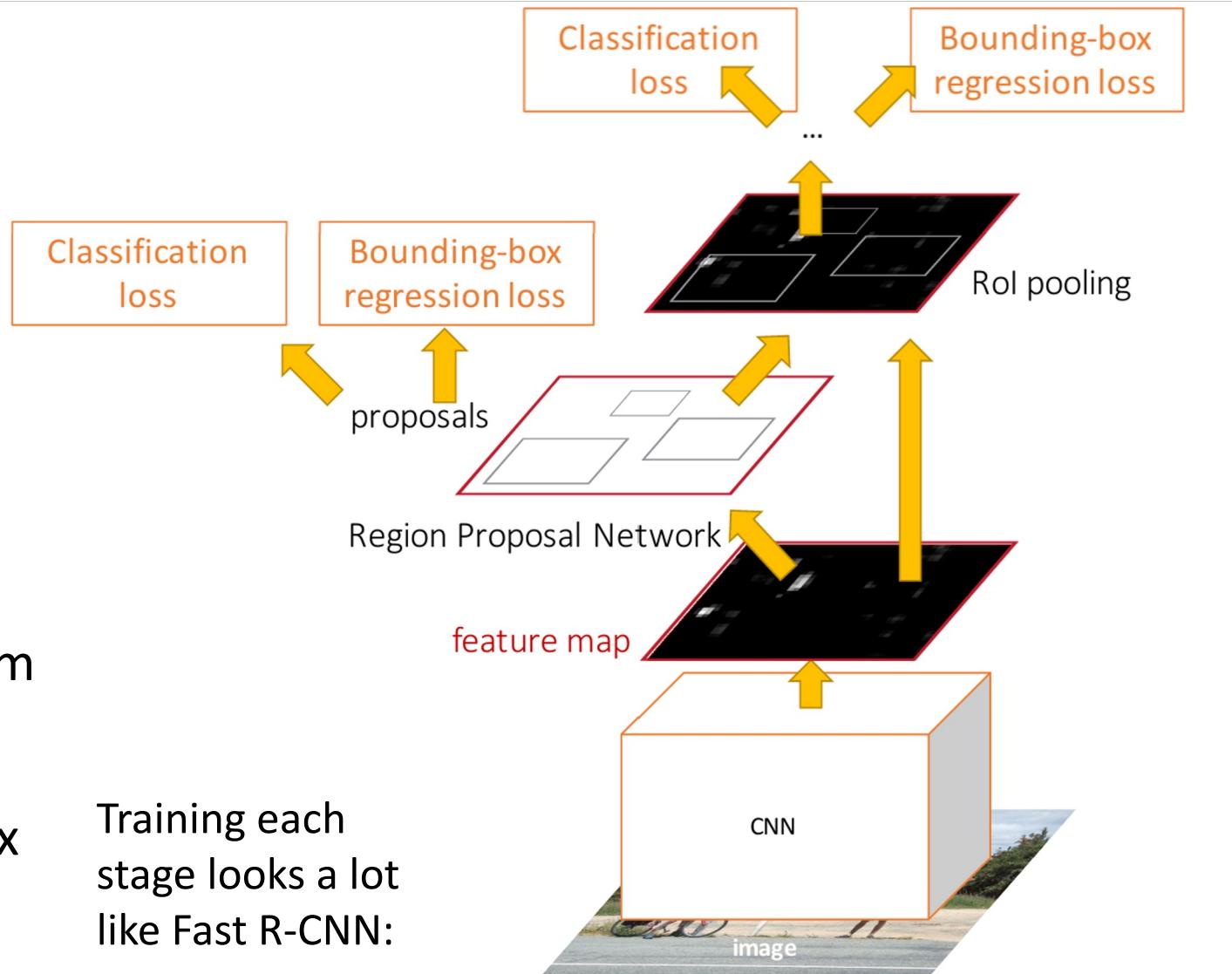
This image is CC0 public domain

Faster R-CNN: Learnable Region Proposals

Jointly train with 4 losses:

1. **RPN classification:** anchor box is object / not an object
 2. **RPN regression:** predict transform from anchor box to proposal box
 3. **Object classification:** classify proposals as background / object class
 4. **Object regression:** predict transform from proposal box to object box

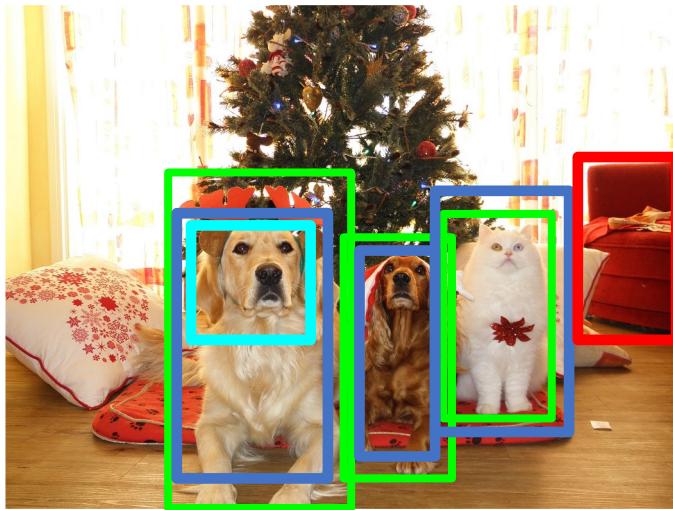
Training each stage looks a lot like Fast R-CNN:



Faster R-CNN Training: RPN Training

RPN predicts Object / Background for each anchor, as well as regresses from anchor to object box

Input Image



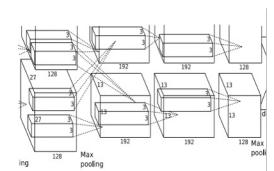
GT Boxes

Positive

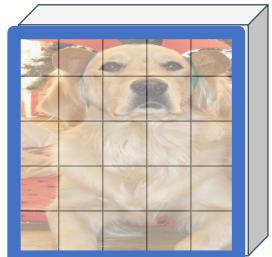
Neutral

Negative

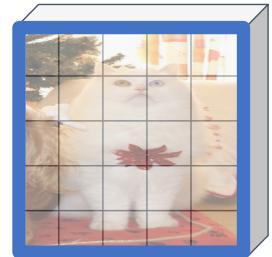
Image Features



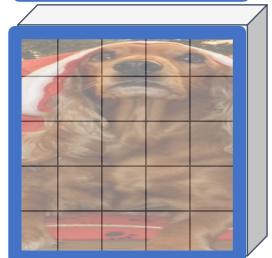
RPN gives lots of **anchors** which we classify as pos / neg / neutral by matching with ground-truth



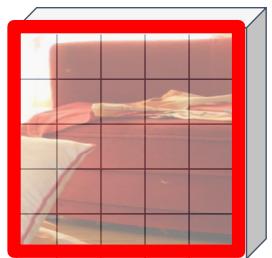
Class target: Obj
Box target: 



Class target: Obj
Box target: 



Class target: Obj
Box target: —————→

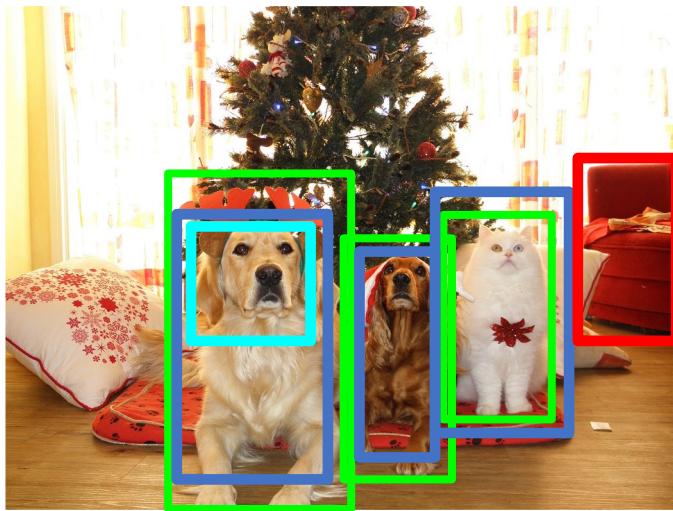


Class target: Background
Box target: None

This image is CC0 public domain

Faster R-CNN Training: Stage 2

Input Image



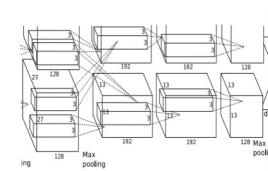
GT Boxes

Positive

Neutral

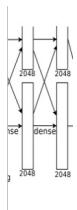
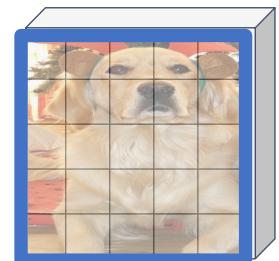
Negative

Image Features

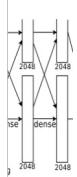
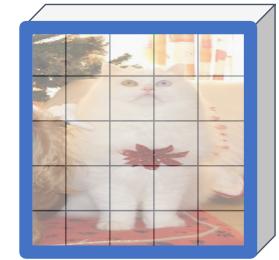


Now proposals come from RPN
rather than selective search,
but otherwise this works the
same as Fast R-CNN training

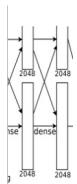
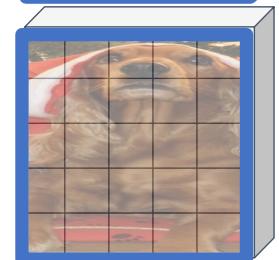
Crop features for each proposal, use them to predict class and box targets per region



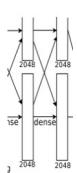
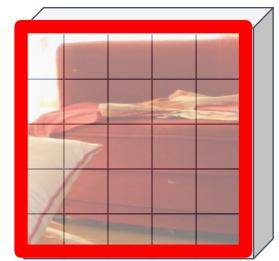
Class target: Dog
Box target: 



Class target: Cat
Box target: →



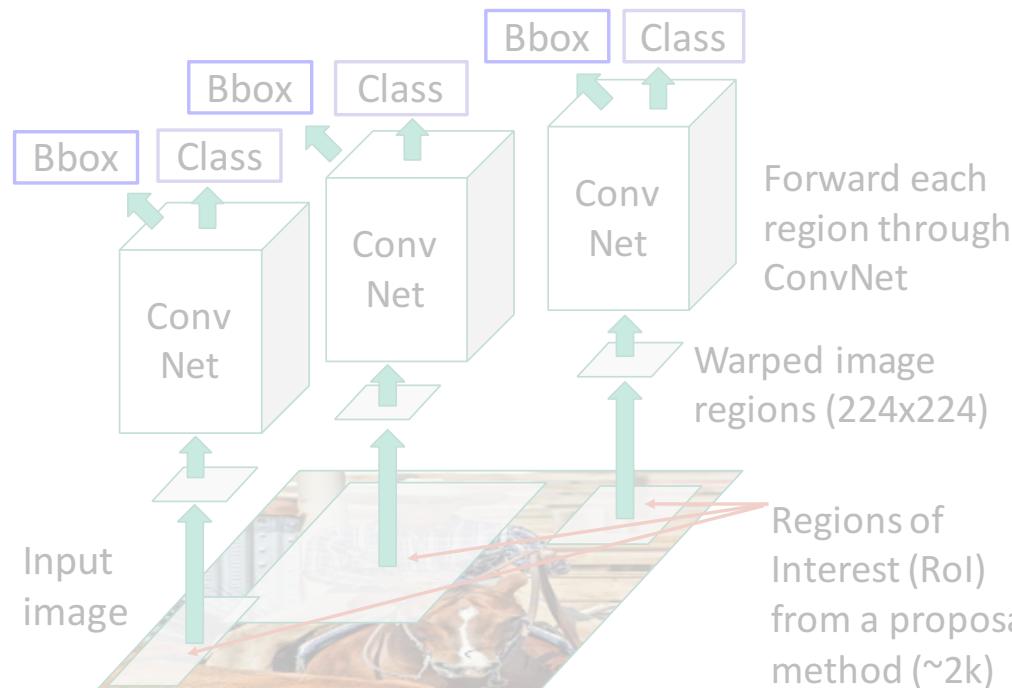
Class target: Dog
Box target: →



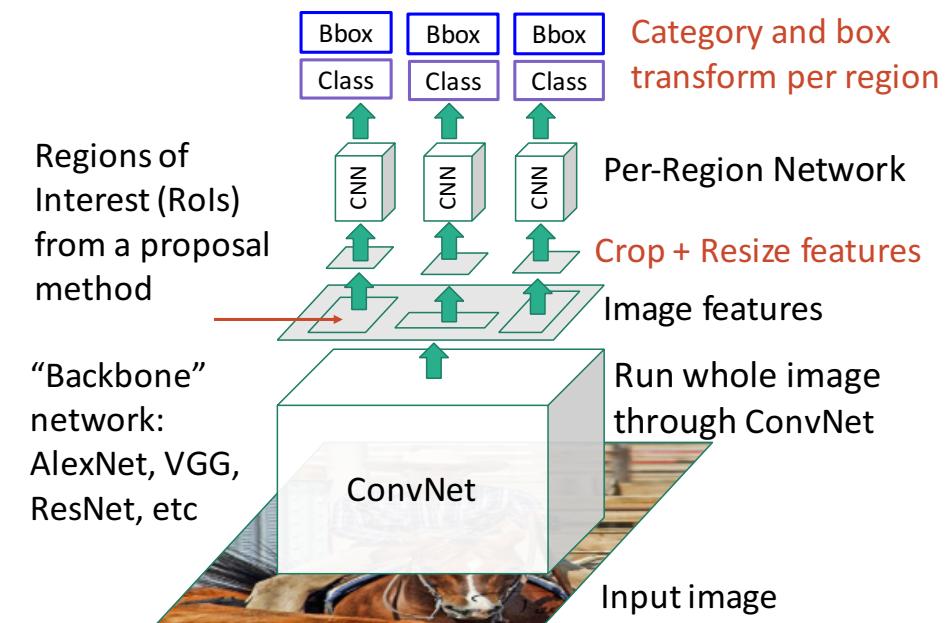
Class target: Background
Box target: None

Recap: Fast R-CNN Feature Cropping

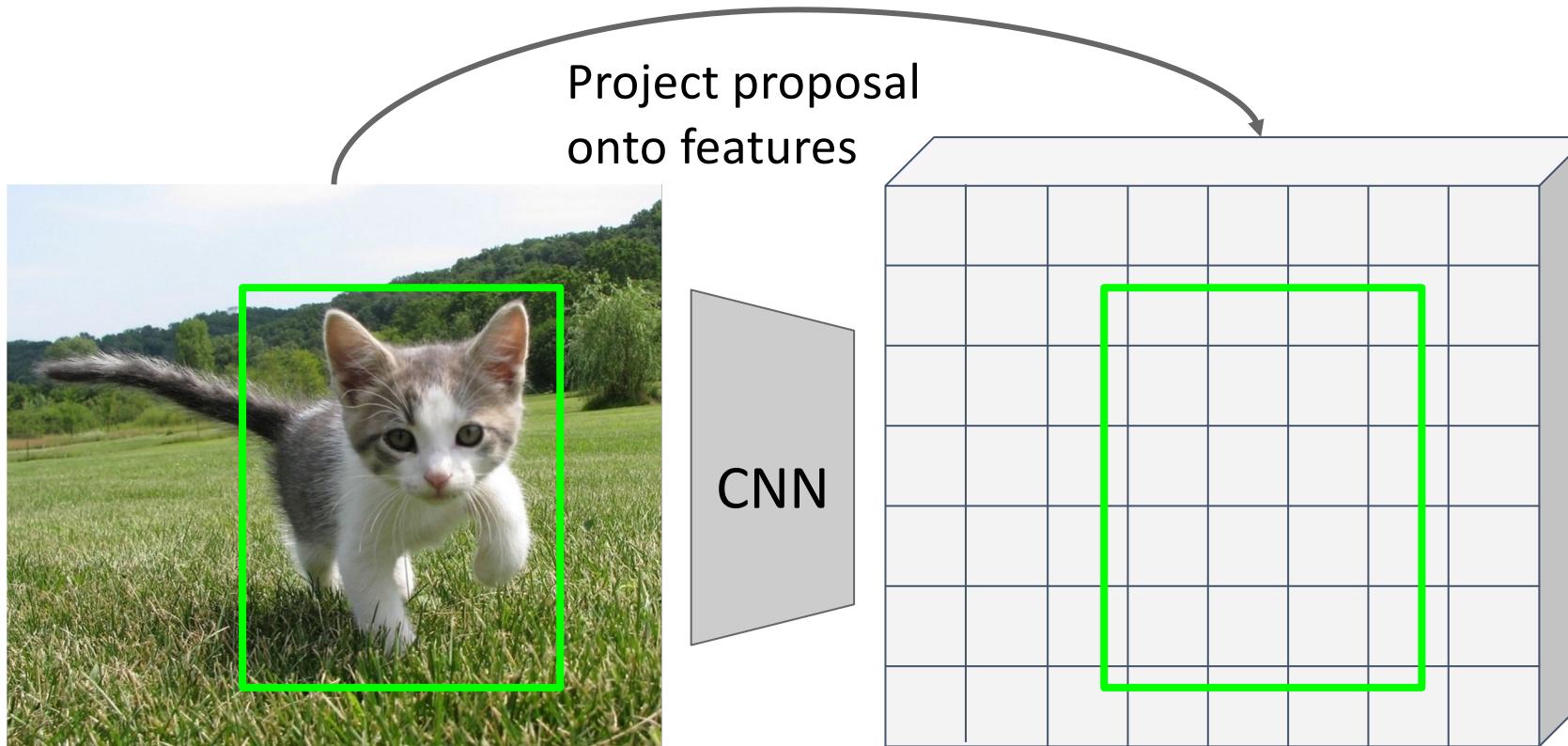
“Slow” R-CNN: Run CNN independently for each region



Fast R-CNN: Apply differentiable cropping to shared image features



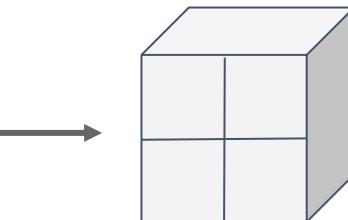
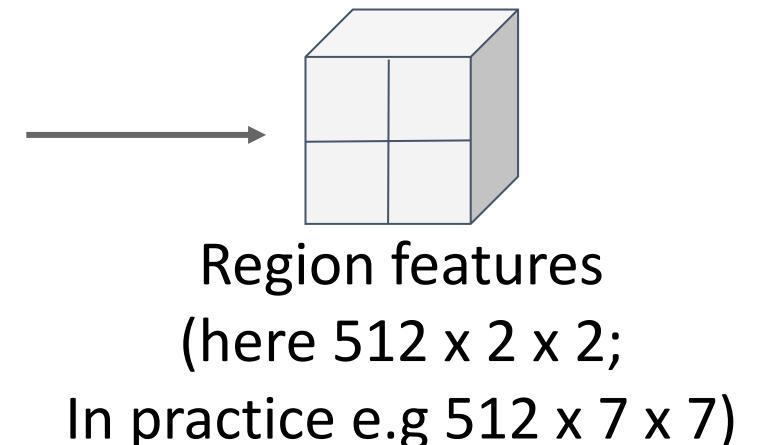
Cropping Features



Input Image
(e.g. $3 \times 640 \times 480$)

Image features
(e.g. $512 \times 20 \times 15$)

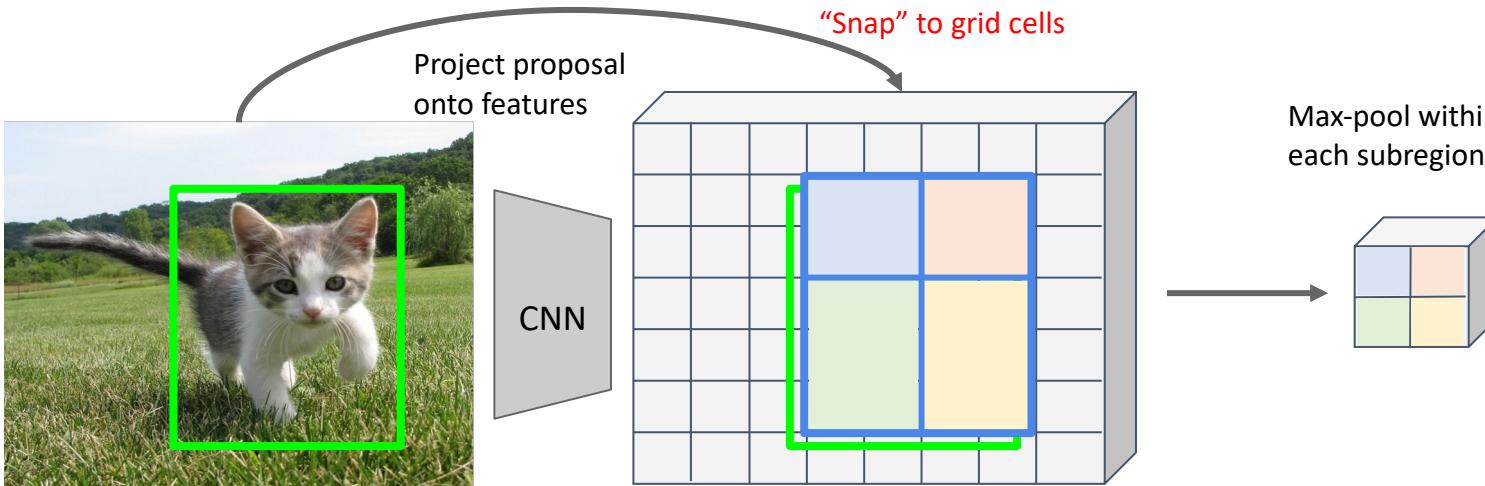
Goal: Crop features for region proposal, and resize to a fixed size for downstream processing, in a differentiable way



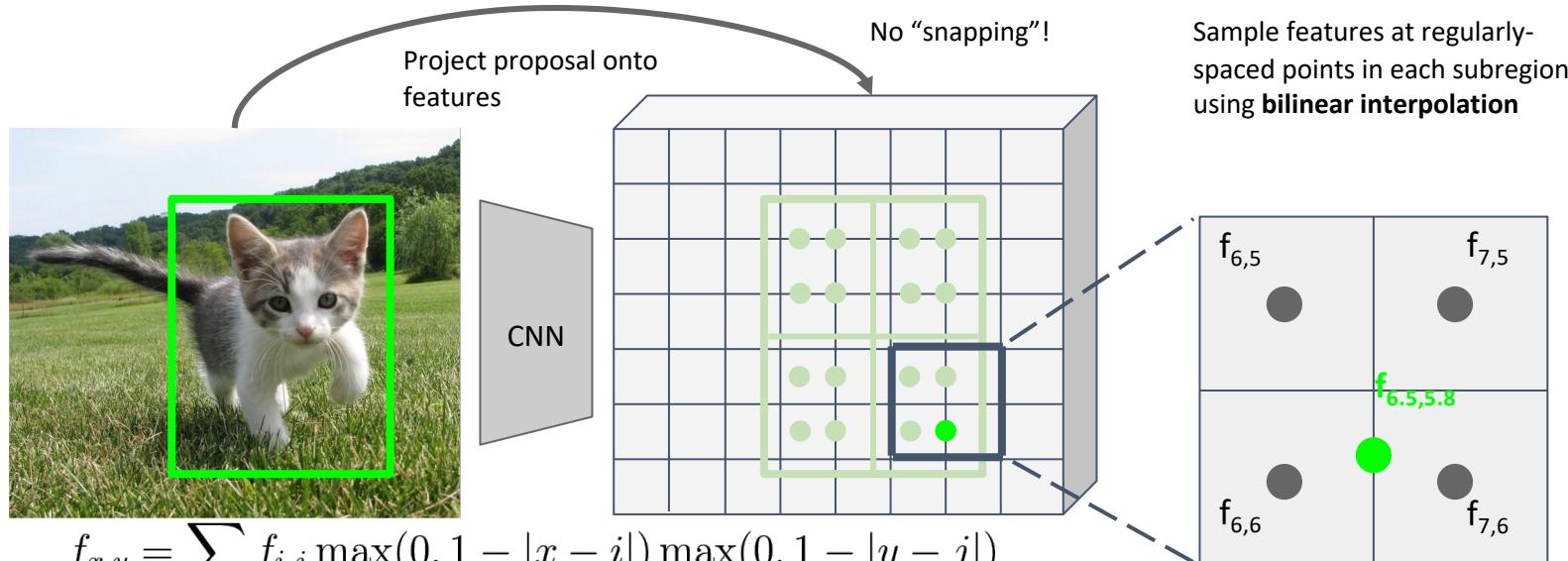
Region features
(here $512 \times 2 \times 2$;
In practice e.g $512 \times 7 \times 7$)

Cropping Features

RoI Pool



RoI Align

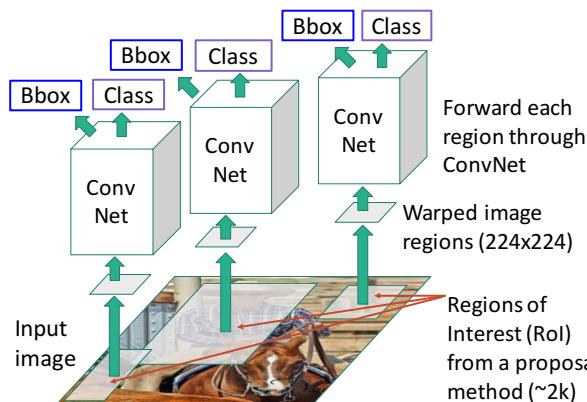


$$f_{x,y} = \sum f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$$

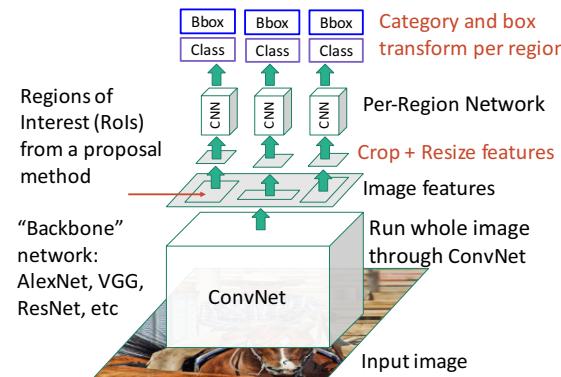
$$\begin{aligned} f_{6.5,5.8} = & (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ & + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

Object Detection Methods

“Slow” R-CNN: Run CNN independently for each region

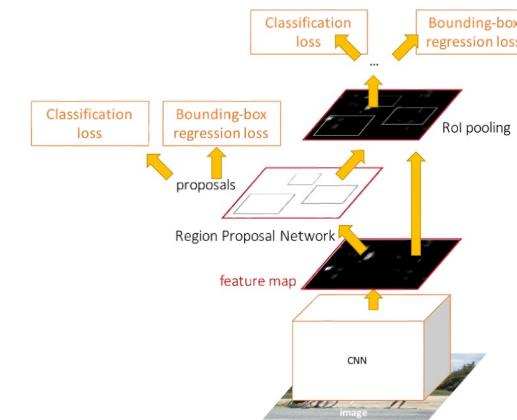


Fast R-CNN: Apply differentiable cropping to shared image features

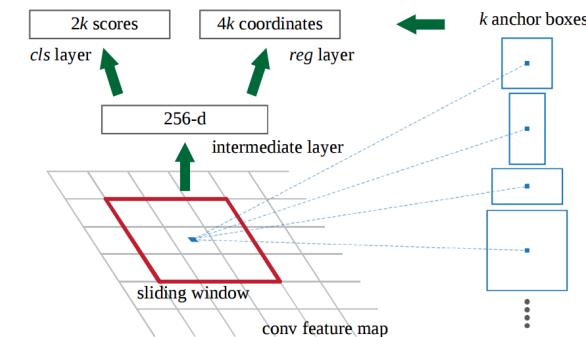


It relies on anchor boxes. Can we do detection without anchors?

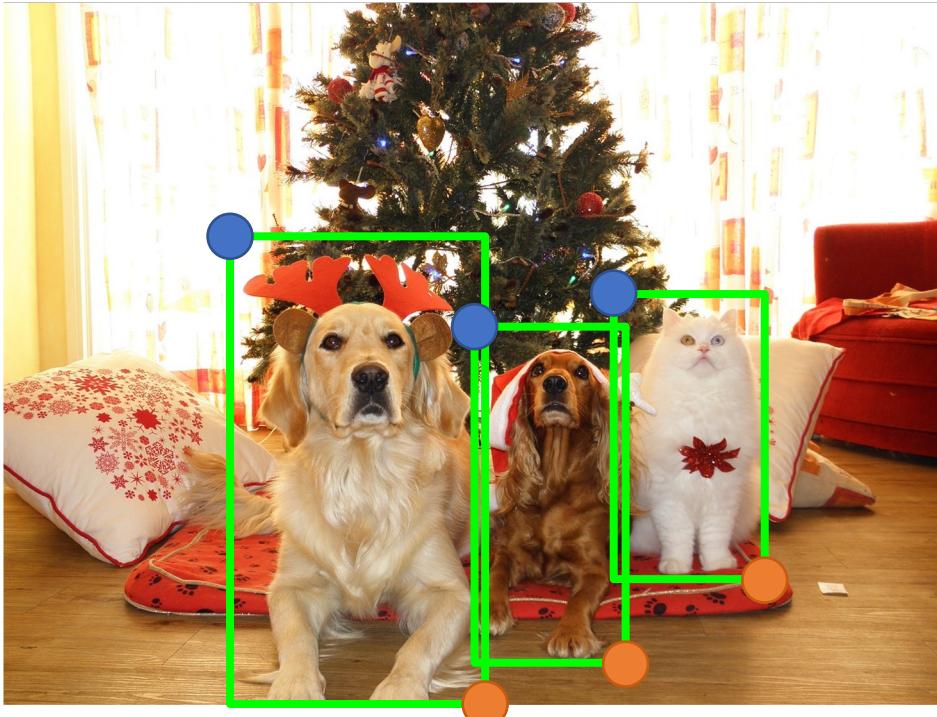
Faster R-CNN: Compute proposals with CNN



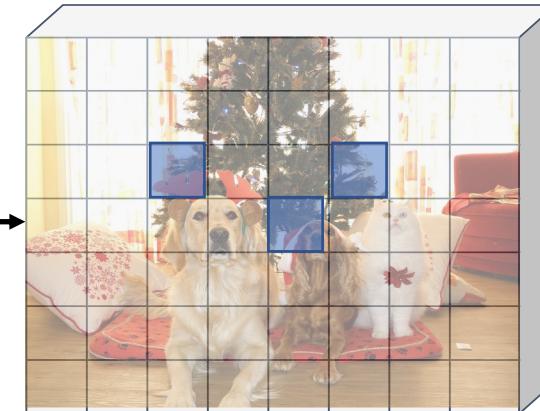
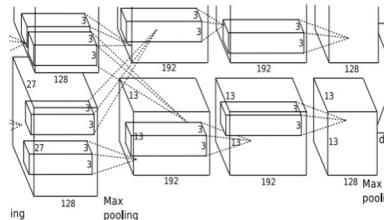
Single-Stage (SSD, YOLO): Fully convolutional detector



Detection without Anchors (anchor-free methods): CornerNet



Backbone CNN



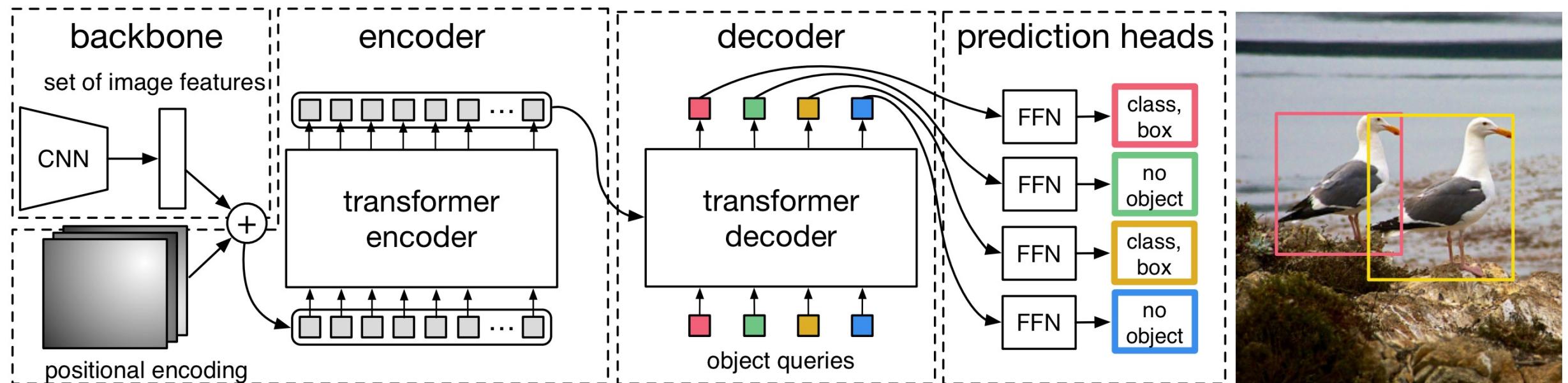
Upper left corners
Heatmap: $C \times H \times W$
Embeddings: $D \times H \times W$



Matching corners are trained to have similar embeddings

Lower right corners
Heatmap: $C \times H \times W$
Embeddings: $D \times H \times W$

Detection without Anchors (anchor-free methods): Transformer-based DETR



Carion et al, "DETR: End-to-End Object Detection with Transformers", ECCV 2020

Computer Vision Tasks: Object Detection

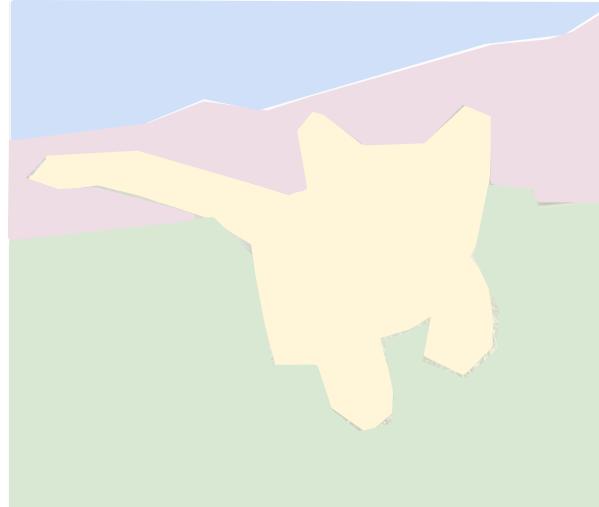
Classification



CAT

No spatial extent

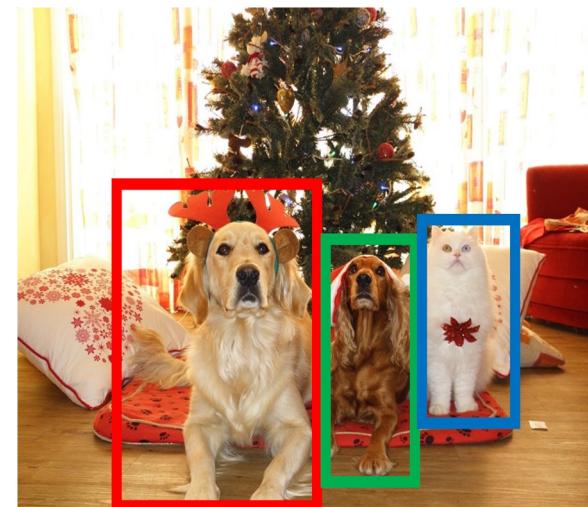
Semantic
Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Objects

Instance
Segmentation



DOG, DOG, CAT

Computer Vision Tasks: Semantic Segmentation

Classification



CAT

No spatial extent

Semantic
Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Objects

Instance
Segmentation



DOG, DOG, CAT

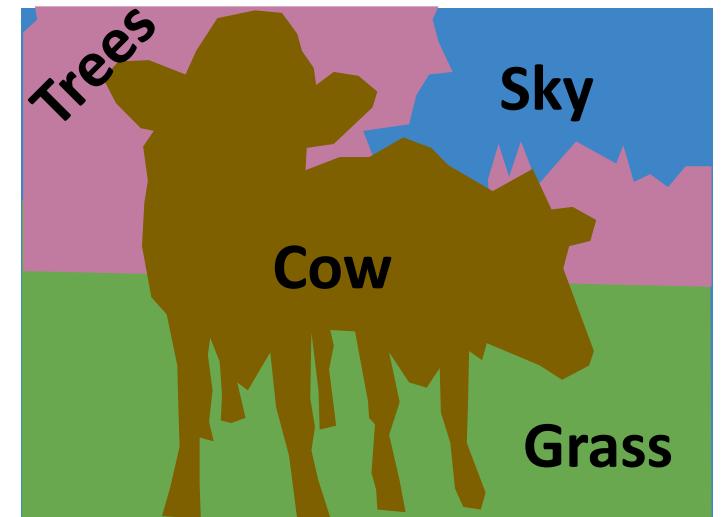
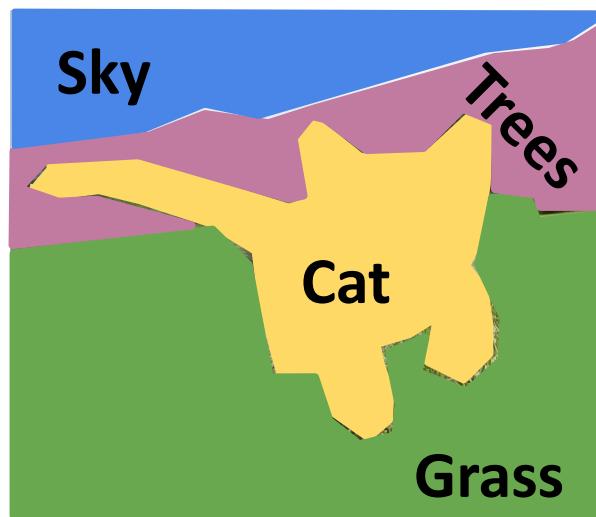
Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



[This image is CC0 public domain](#)



Things and Stuff

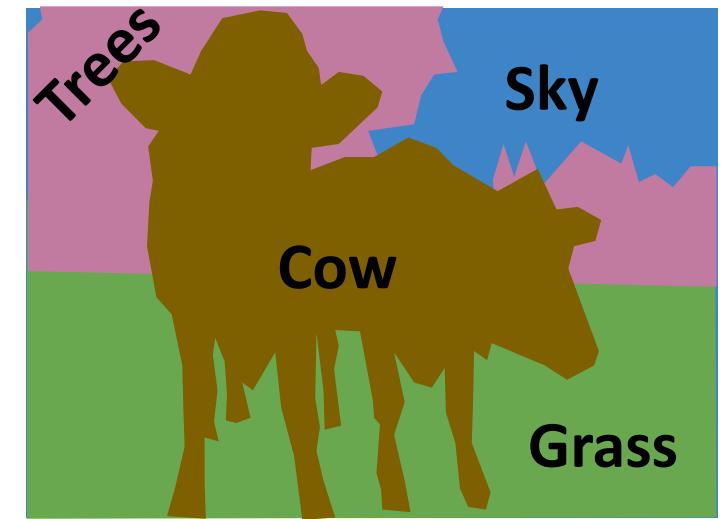
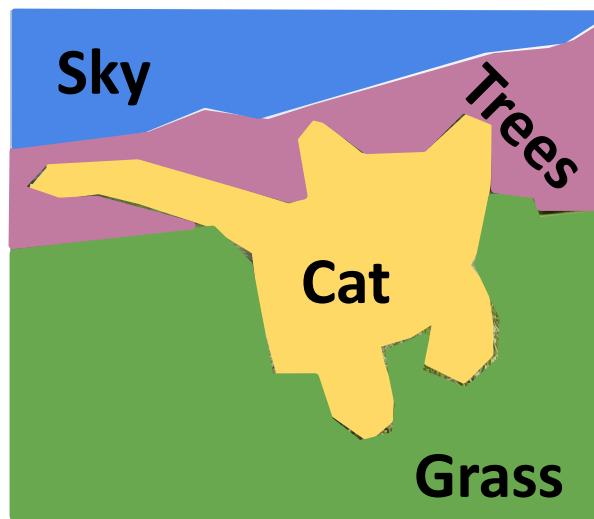
Things: Object categories
that can be separated into
object instances
(e.g. cats, cars, person)



[This image is CC0 public domain](#)



Stuff: Object categories
that cannot be separated
into instances
(e.g. sky, grass, water, trees)

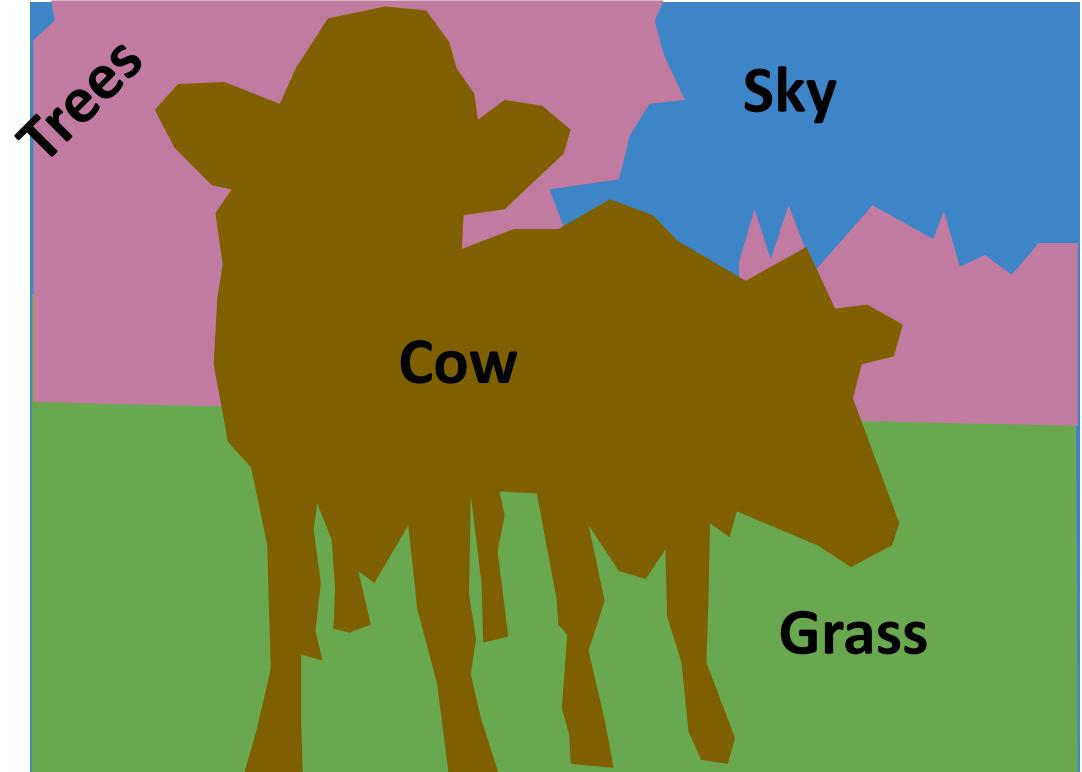


Object Detection vs Semantic Segmentation

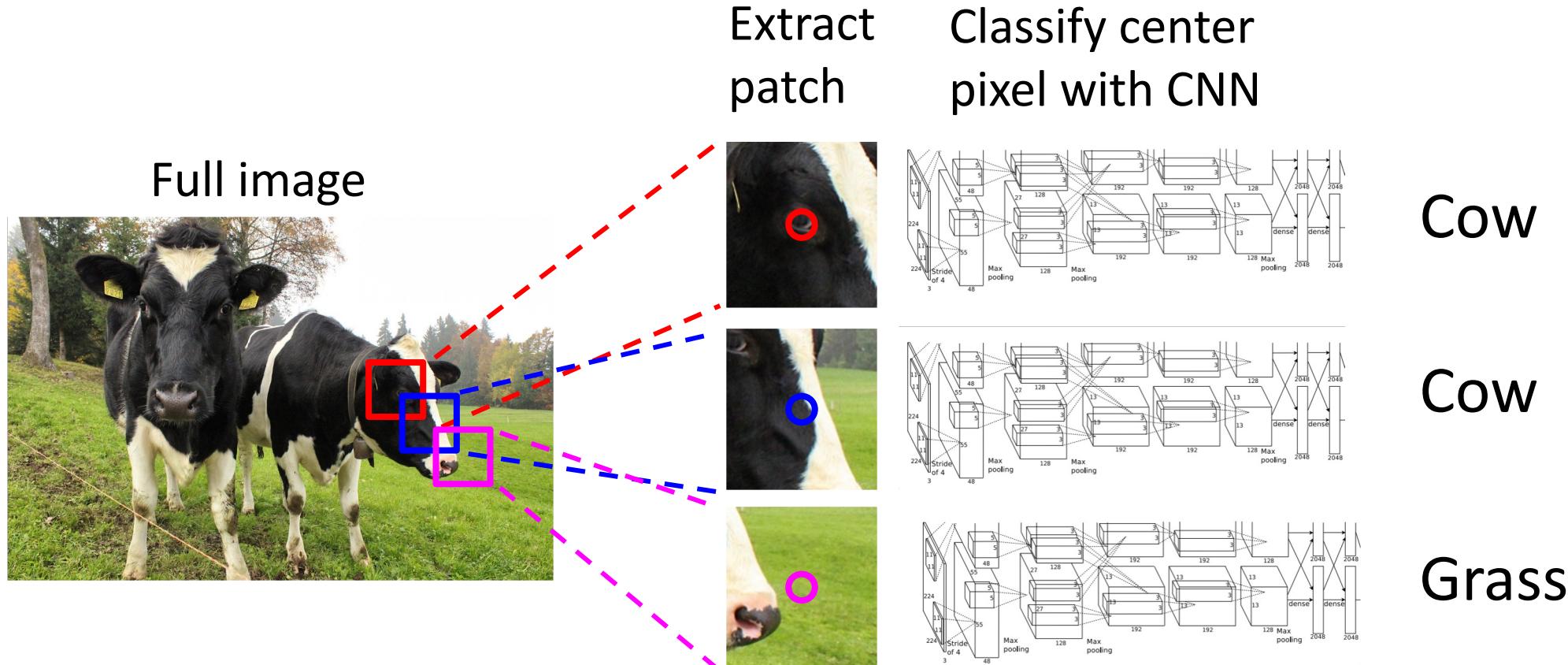
Object Detection: Detects individual object instances, but only gives box
(Only things!)



Semantic Segmentation: Gives per-pixel labels, but merges instances (Both things and stuff) where each pixel has a label

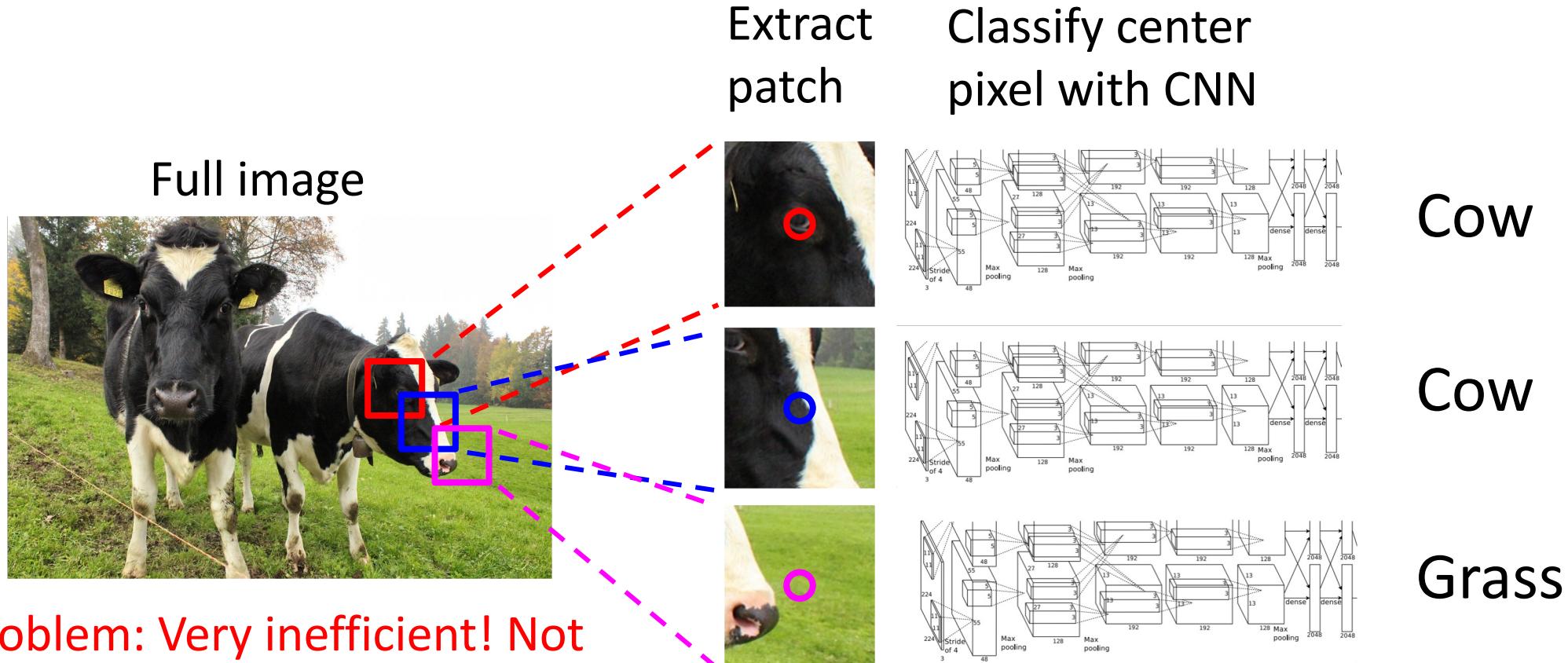


Semantic Segmentation Idea: Sliding Window



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

Semantic Segmentation Idea: Sliding Window

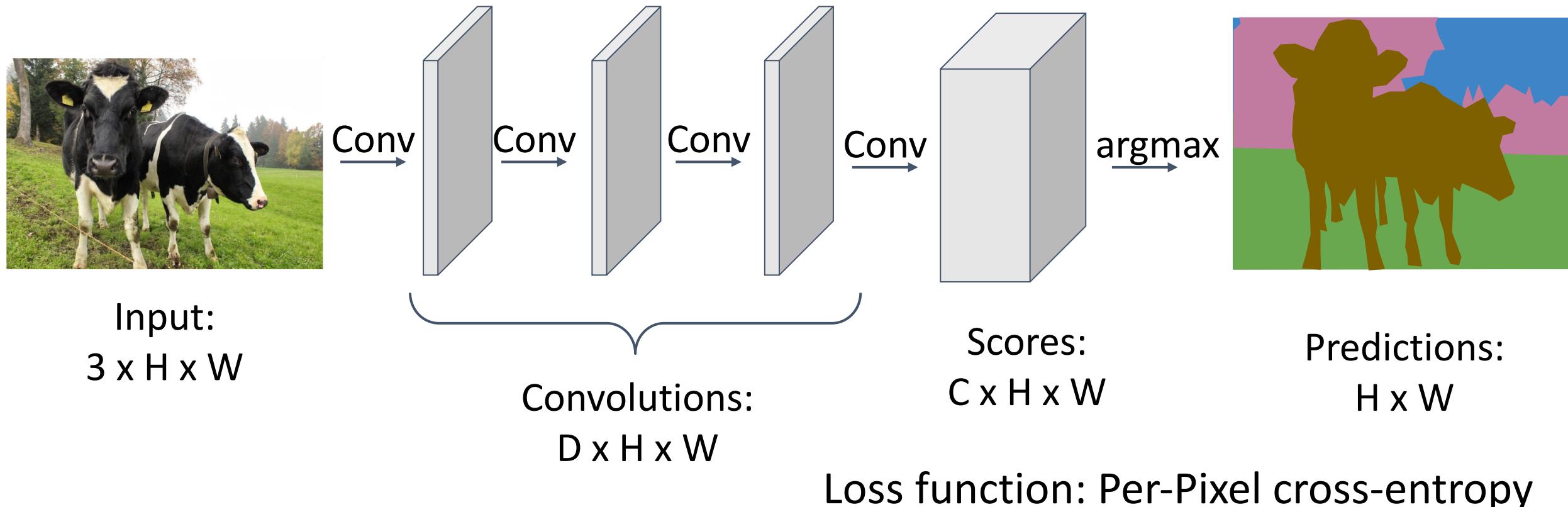


Problem: Very inefficient! Not reusing shared features between overlapping patches

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 Network

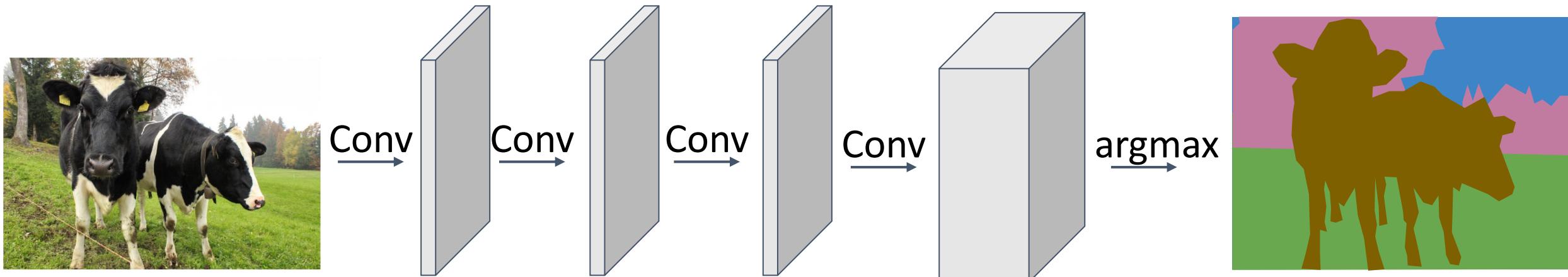
Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Semantic Segmentation: Fully Convolutional Network

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

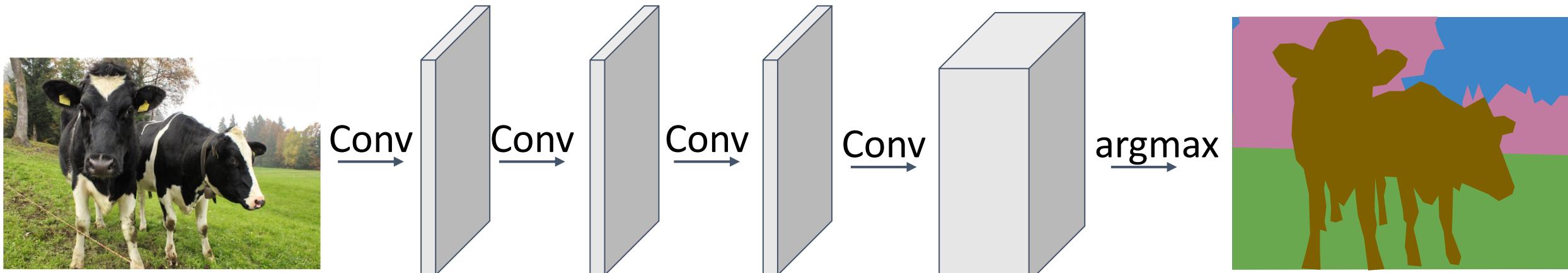


Input: **Problem #1:** Effective receptive
 $3 \times H \times W$ field size is linear in number of
conv layers: With L 3×3 conv
layers, receptive field is $1+2L$

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Semantic Segmentation: Fully Convolutional Network

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



Input: $3 \times H \times W$ **Problem #1:** Effective receptive field size is linear in number of conv layers: With L 3×3 conv layers, receptive field is $1+2L$

Problem #2: Convolution on high res images is expensive!
Recall ResNet stem aggressively downsamples

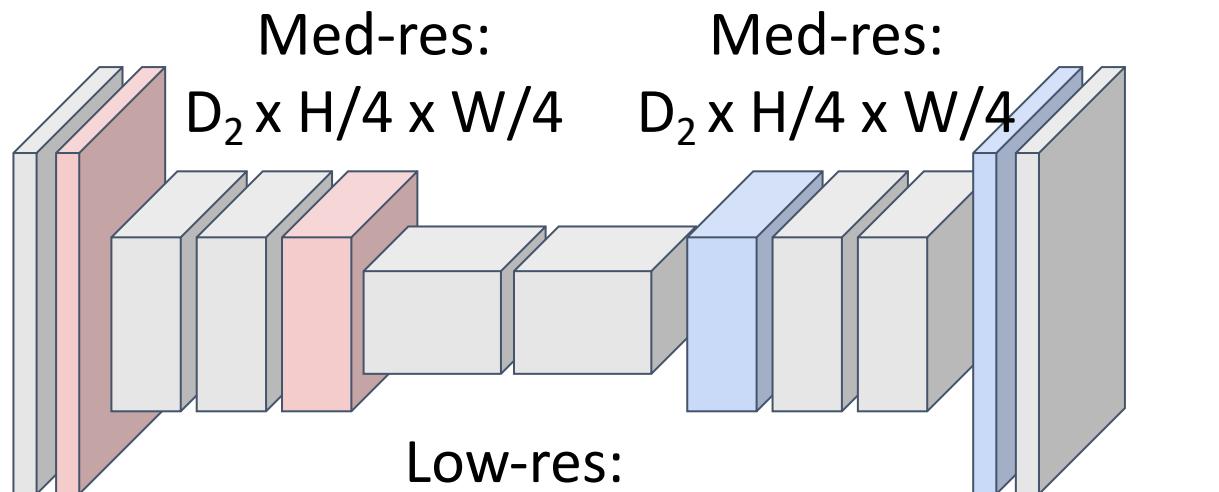
Semantic Segmentation: Fully Convolutional Network

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



Input:
 $3 \times H \times W$

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



Predictions:
 $H \times W$

Or it is also called Encoder-Decoder Architecture

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

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

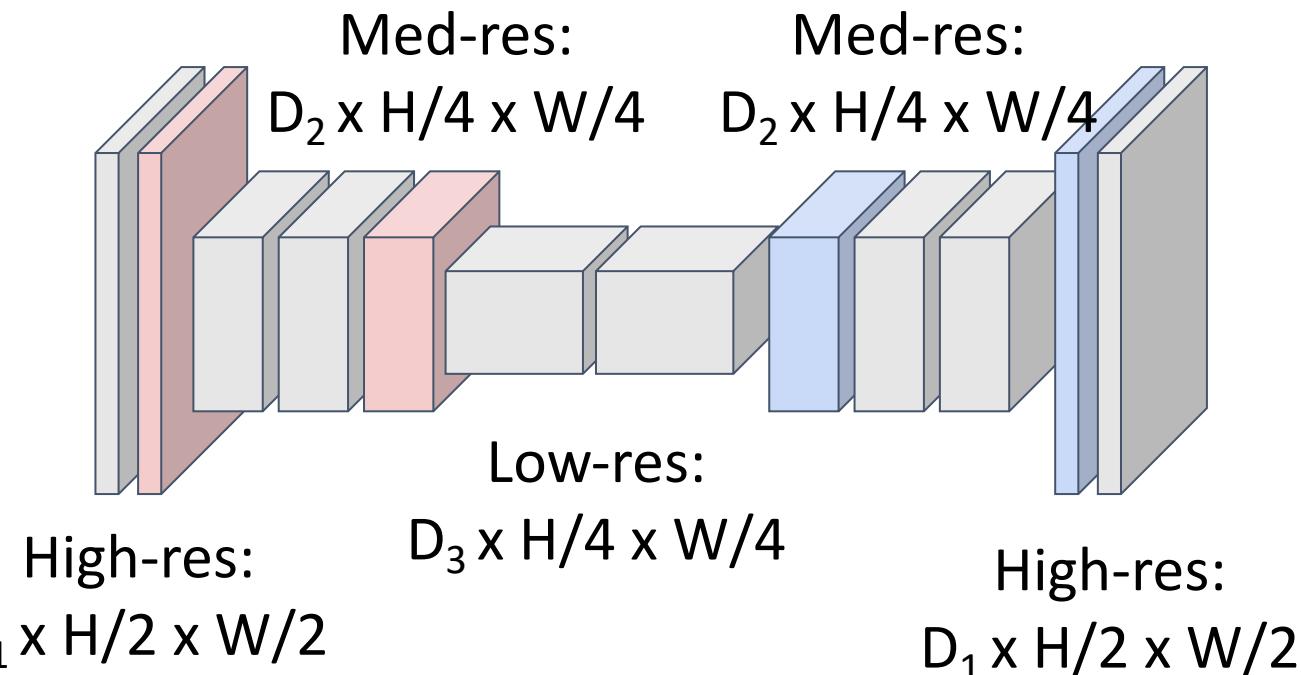
Semantic Segmentation: Fully Convolutional Network

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

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



Or it is also called Encoder-Decoder Architecture

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

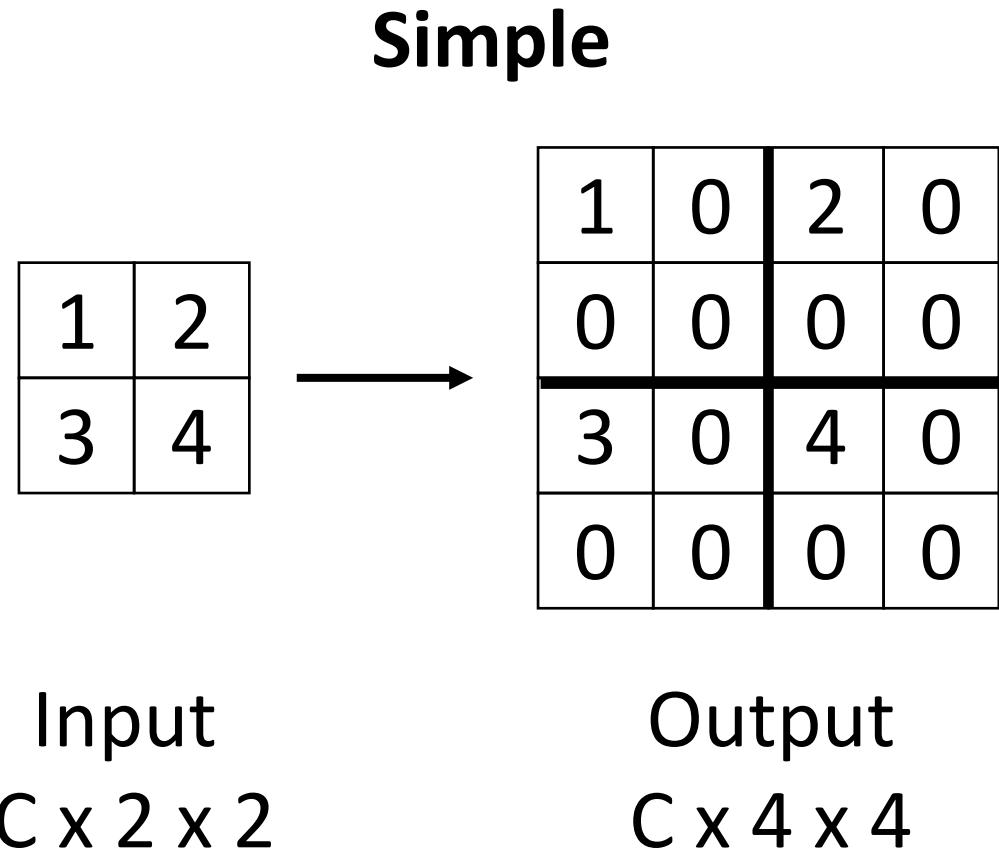
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Upsampling:
???



Predictions:
 $H \times W$

In-Network Upsampling: “Unpooling”



In-Network Upsampling: “Unpooling”

Simple

1	2
3	4



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

Nearest Neighbor

1	2
3	4



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

Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

In-Network Upsampling: Bilinear Interpolation

1		2	
3		4	



1.00	1.25	1.75	2.00
1.50	1.75	2.25	2.50
2.50	2.75	3.25	3.50
3.00	3.25	3.75	4.00

Input: $C \times 2 \times 2$

Output: $C \times 4 \times 4$

$$f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$$

Use two closest neighbors in x and y
to construct linear approximations

$$j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$$

In-Network Upsampling: Bicubic Interpolation

1	2
3	4



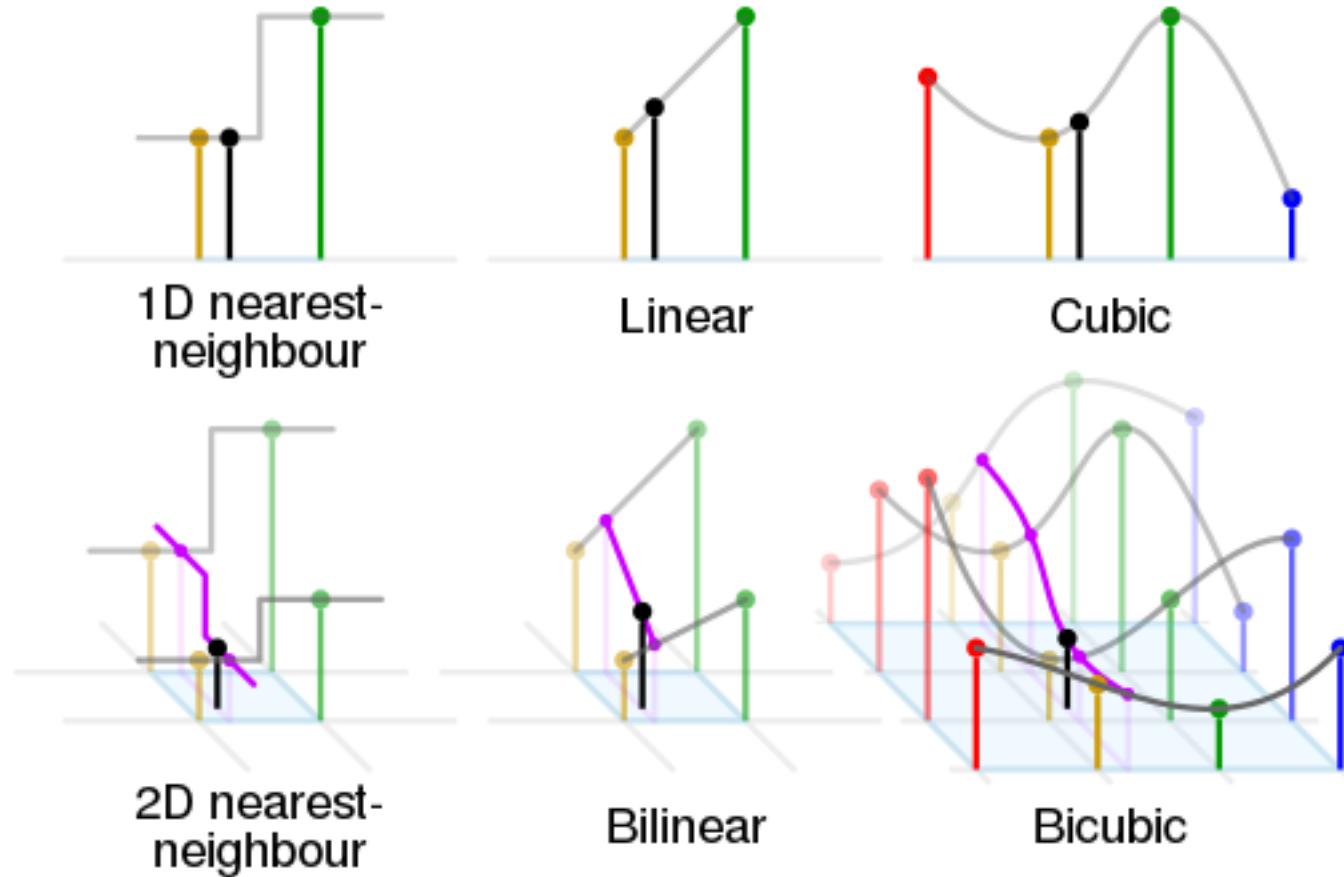
0.68	1.02	1.56	1.89
1.35	1.68	2.23	2.56
2.44	2.77	3.32	3.65
3.11	3.44	3.98	4.32

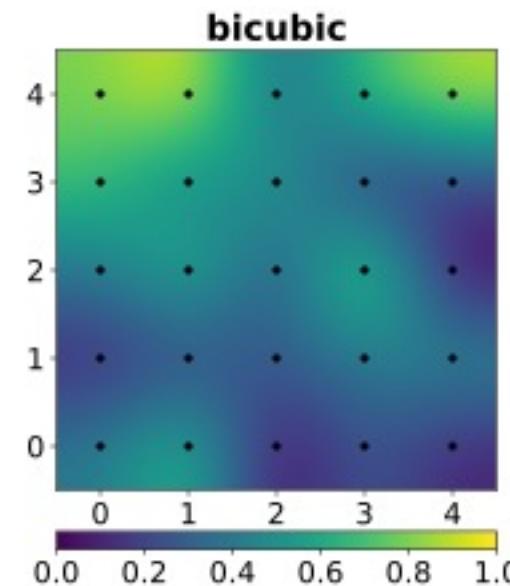
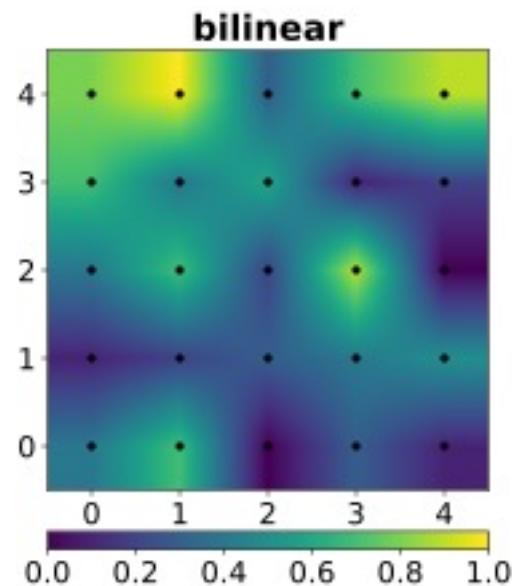
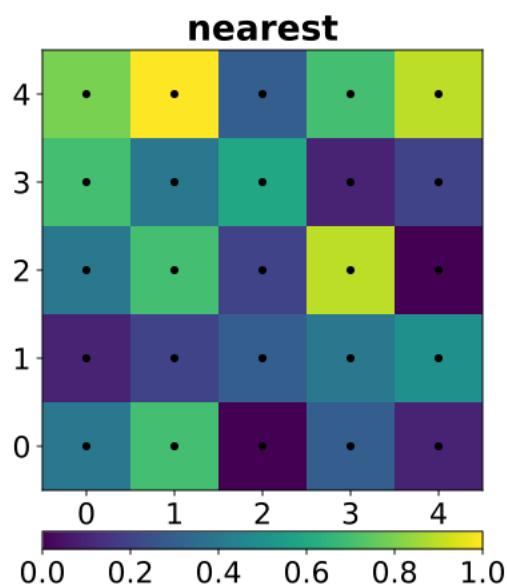
Input: $C \times 2 \times 2$

Output: $C \times 4 \times 4$

Use **three** closest neighbors in x and y to
construct **cubic** approximations
(This is how we normally resize images!)

Different interpolation methods





In-Network Upsampling: “Max Unpooling”

Max Pooling: Remember which position had the max

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

→

5	6
7	8

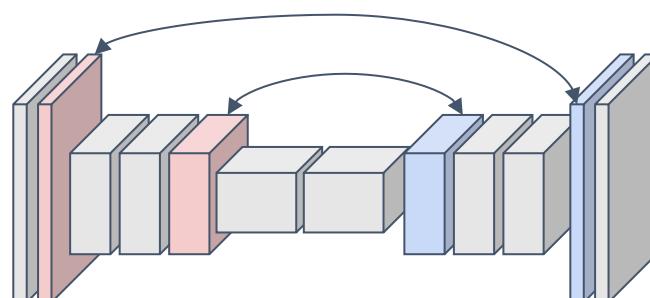
Rest
of
net

→

1	2
3	4

→

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

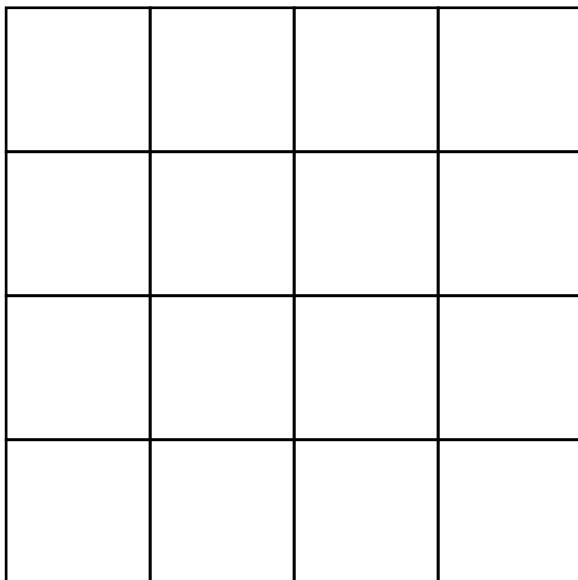


Pair each downsampling layer with an upsampling layer

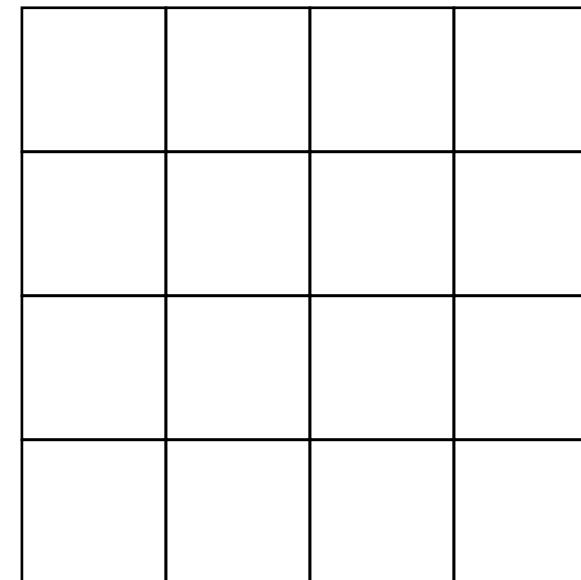
Noh et al, “Learning Deconvolution Network for Semantic Segmentation”, ICCV 2015

Learnable Upsampling: Transposed Convolution

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



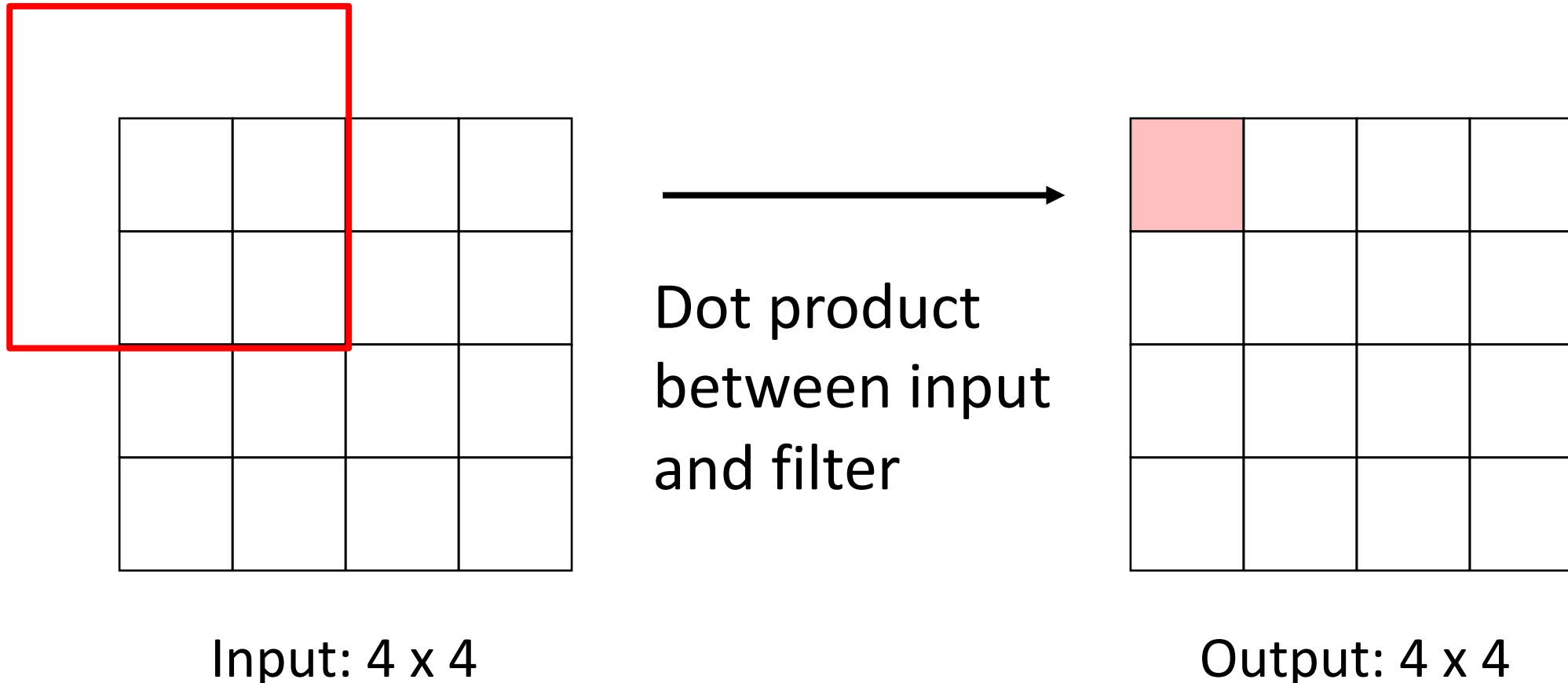
Input: 4×4



Output: 4×4

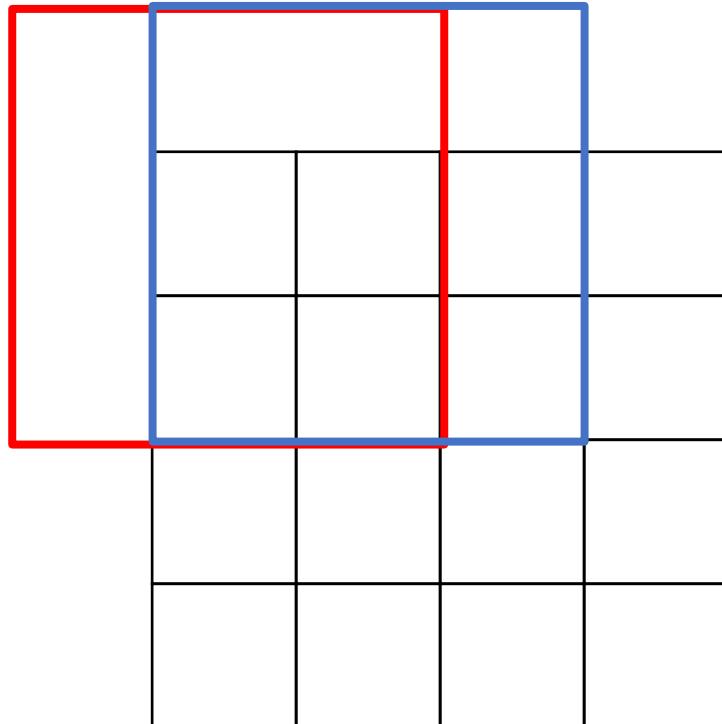
Learnable Upsampling: Transposed Convolution

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



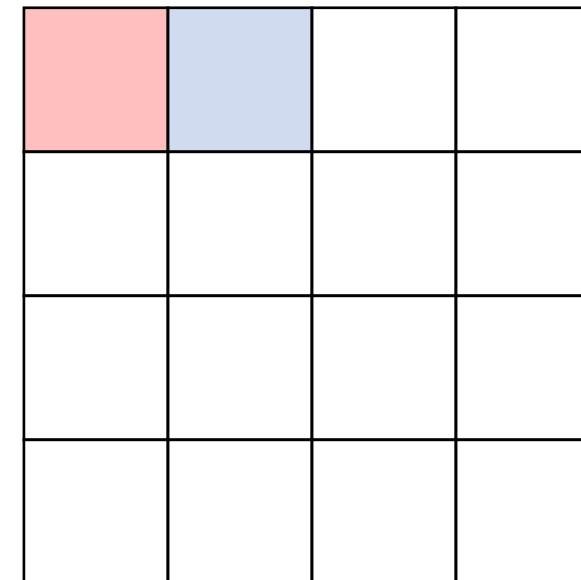
Learnable Upsampling: Transposed Convolution

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



Input: 4×4

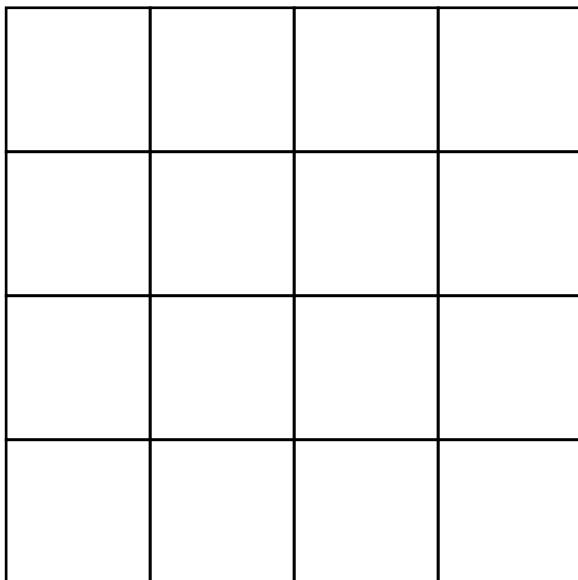
Dot product
between input
and filter



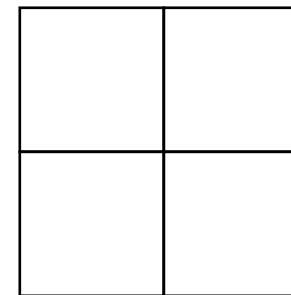
Output: 4×4

Learnable Upsampling: Transposed Convolution

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



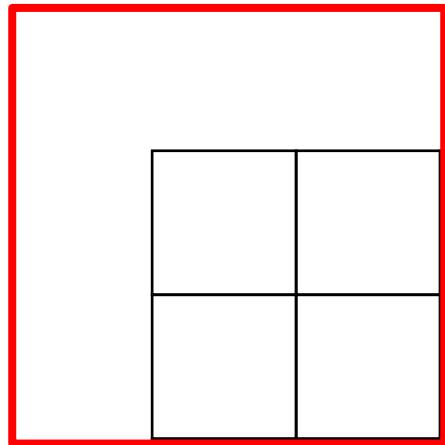
Input: 4×4



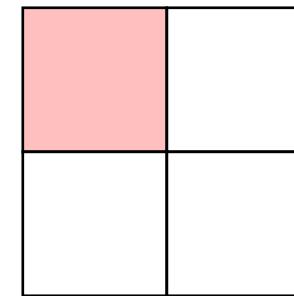
Output: 2×2

Learnable Upsampling: Transposed Convolution

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



Dot product
between input
and filter

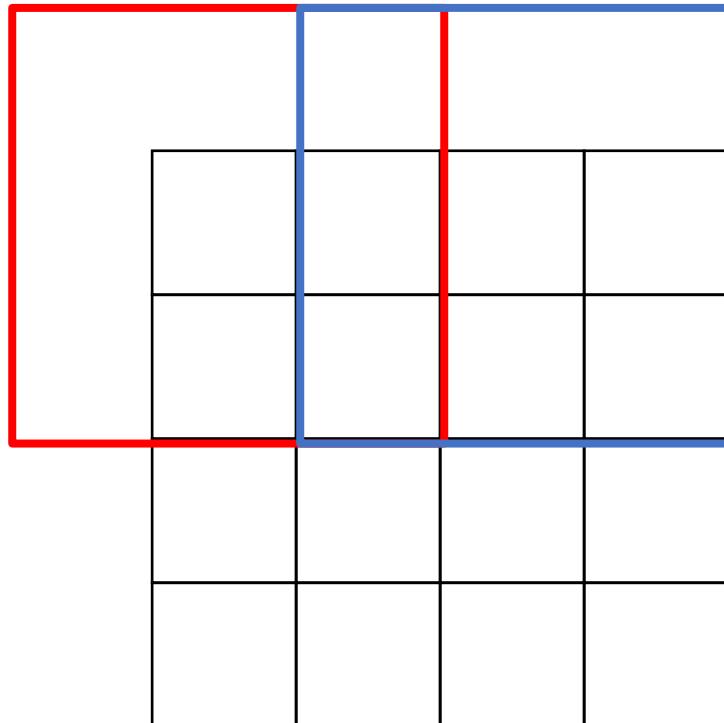


Input: 4×4

Output: 2×2

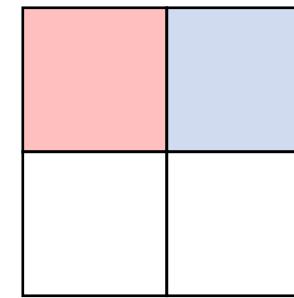
Learnable Upsampling: Transposed Convolution

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



Input: 4×4

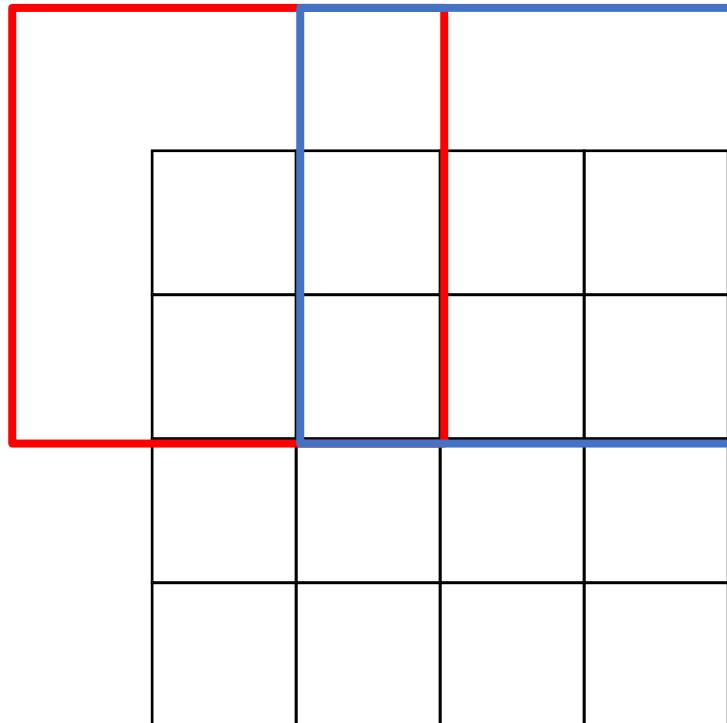
Dot product
between input
and filter



Output: 2×2

Learnable Upsampling: Transposed Convolution

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

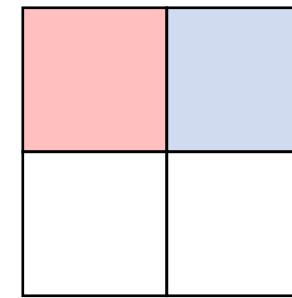


Input: 4×4

Convolution with stride > 1 is “Learnable Downsampling”
How about “Learnable Upsampling”?



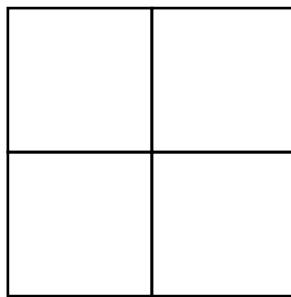
Dot product
between input
and filter



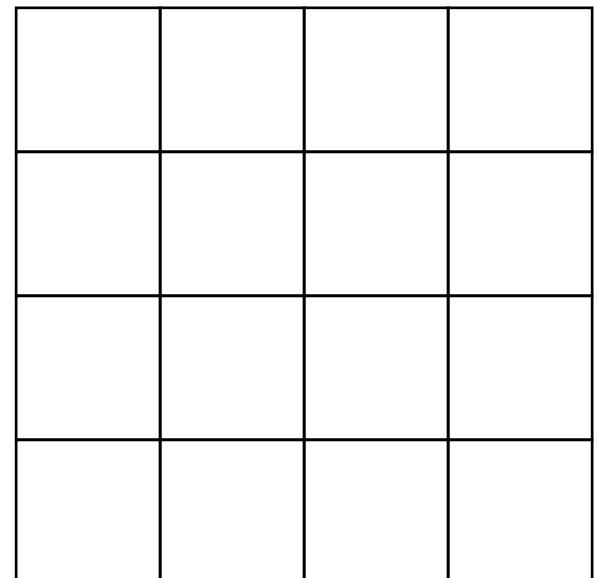
Output: 2×2

Learnable Upsampling: Transposed Convolution

3 x 3 convolution transpose, stride 2



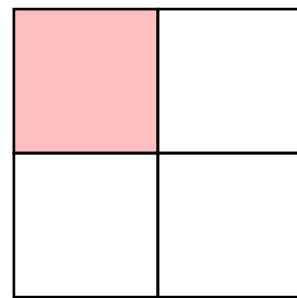
Input: 2 x 2



Output: 4 x 4

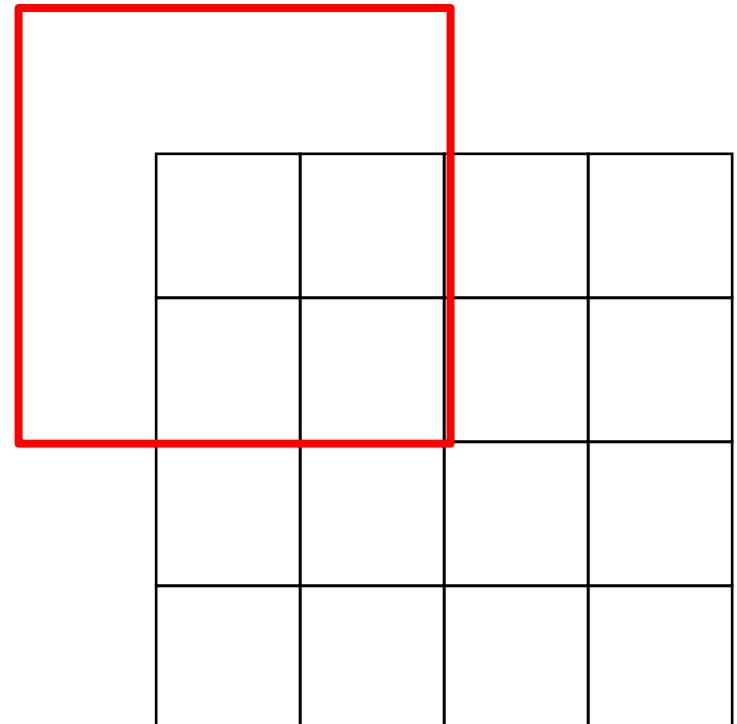
Learnable Upsampling: Transposed Convolution

3 x 3 convolution transpose, stride 2



Input: 2 x 2

→
Weight filter by
input value and
copy to output

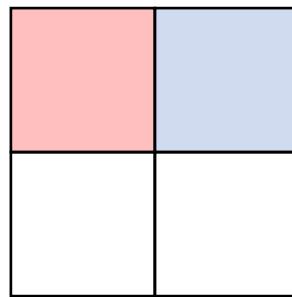


Output: 4 x 4

Learnable Upsampling: Transposed Convolution

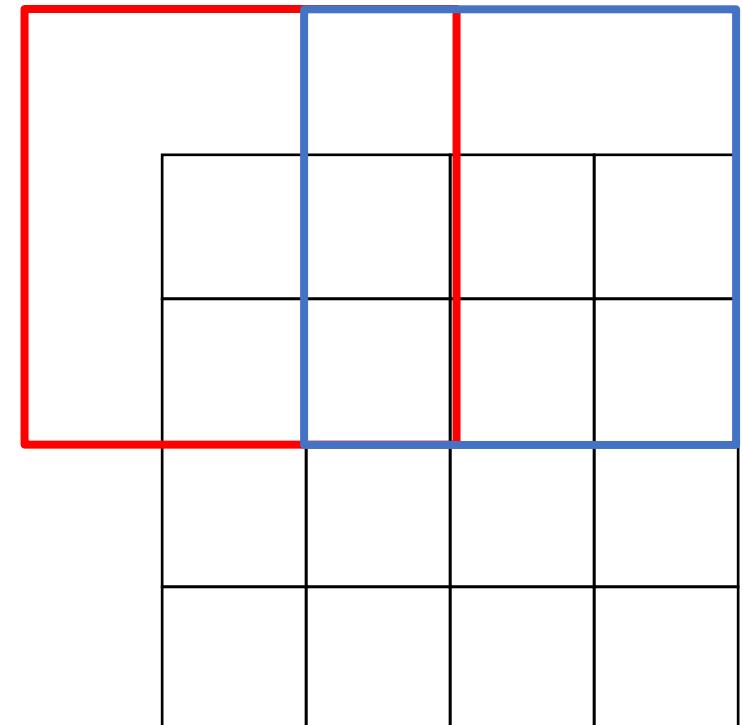
3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in output
for every 1 pixel in input



Input: 2 x 2

Weight filter by
input value and
copy to output

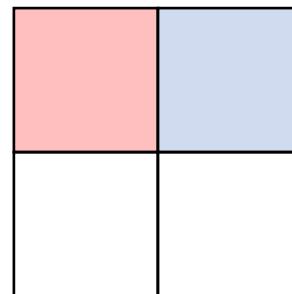


Output: 4 x 4

Learnable Upsampling: Transposed Convolution

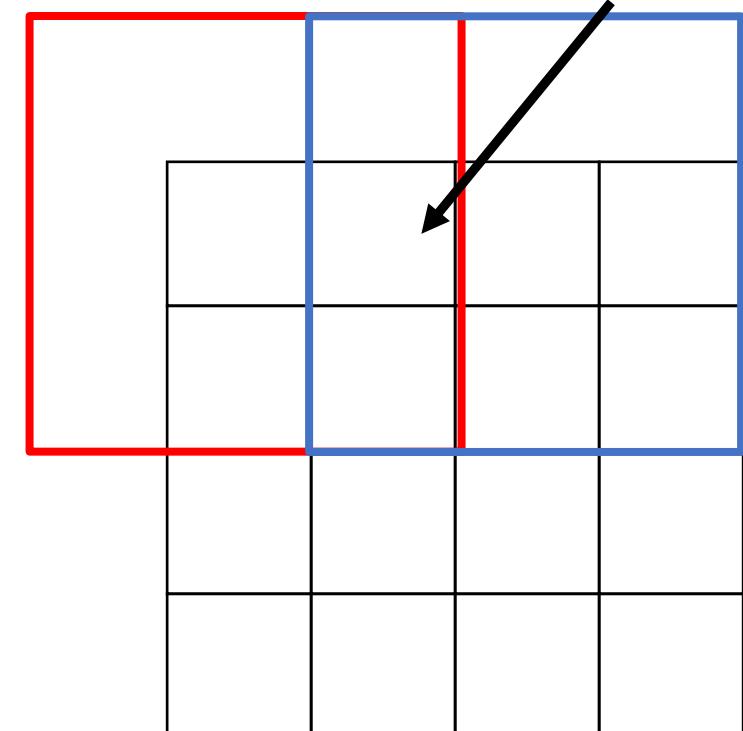
3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in output
for every 1 pixel in input



Input: 2 x 2

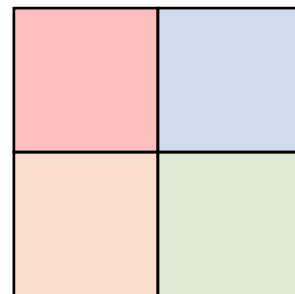
Weight filter by
input value and
copy to output



Learnable Upsampling: Transposed Convolution

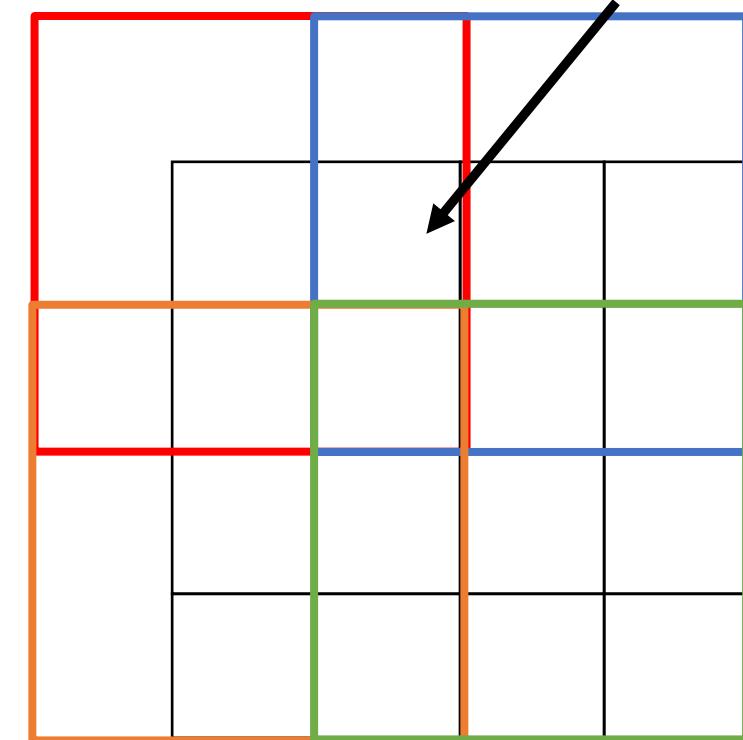
3 x 3 convolution transpose, stride 2

This gives 5x5 output – need to trim one pixel from top and left to give 4x4 output



Input: 2 x 2

Weight filter by
input value and
copy to output



Transposed Convolution: 1D example

Input

a
b

Filter

x
y
z

Output

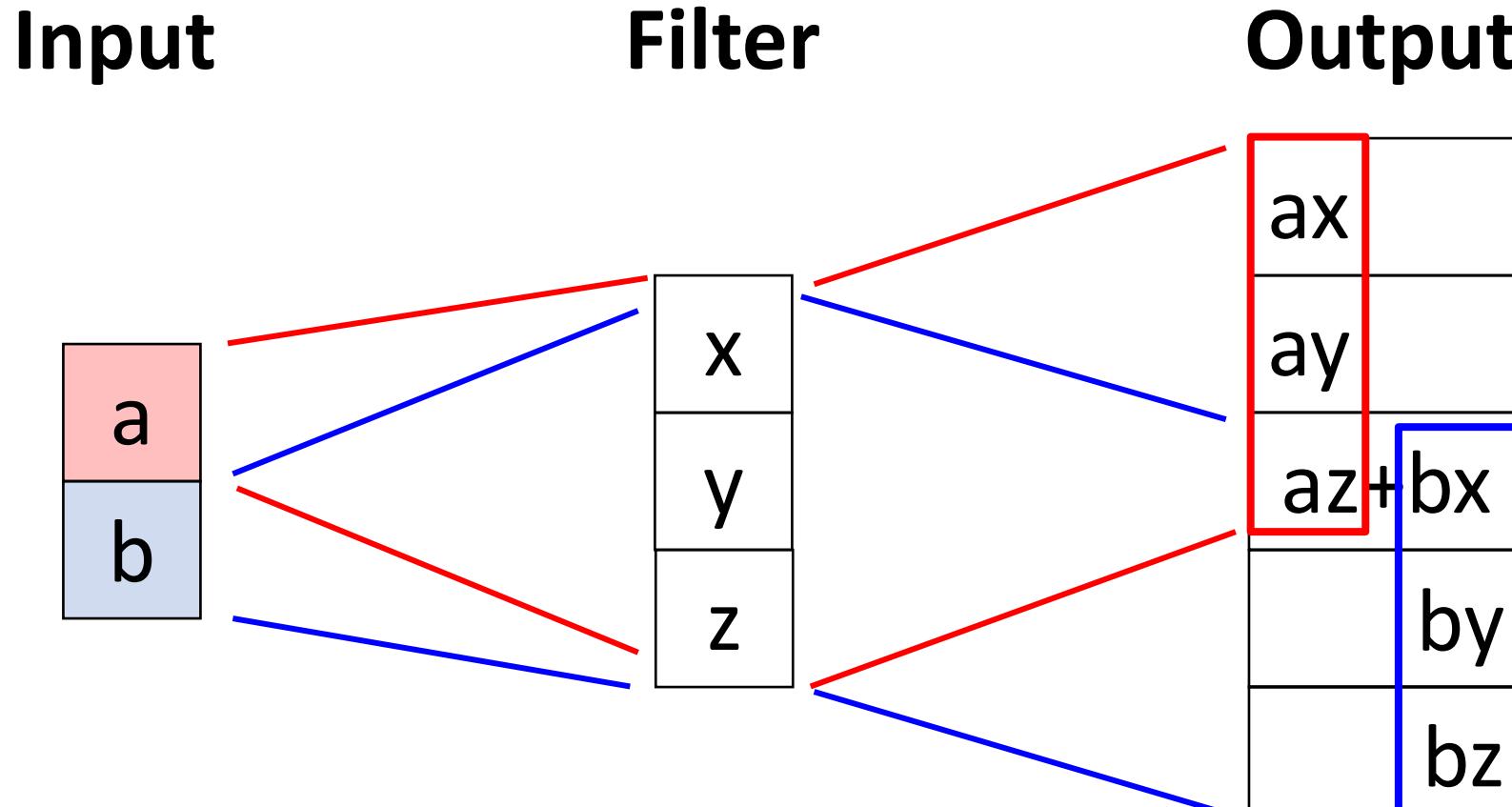
ax
ay
az+bx
by
bz

Output has copies of filter weighted by input

Stride 2: Move 2 pixels output for each pixel in input

Sum at overlaps

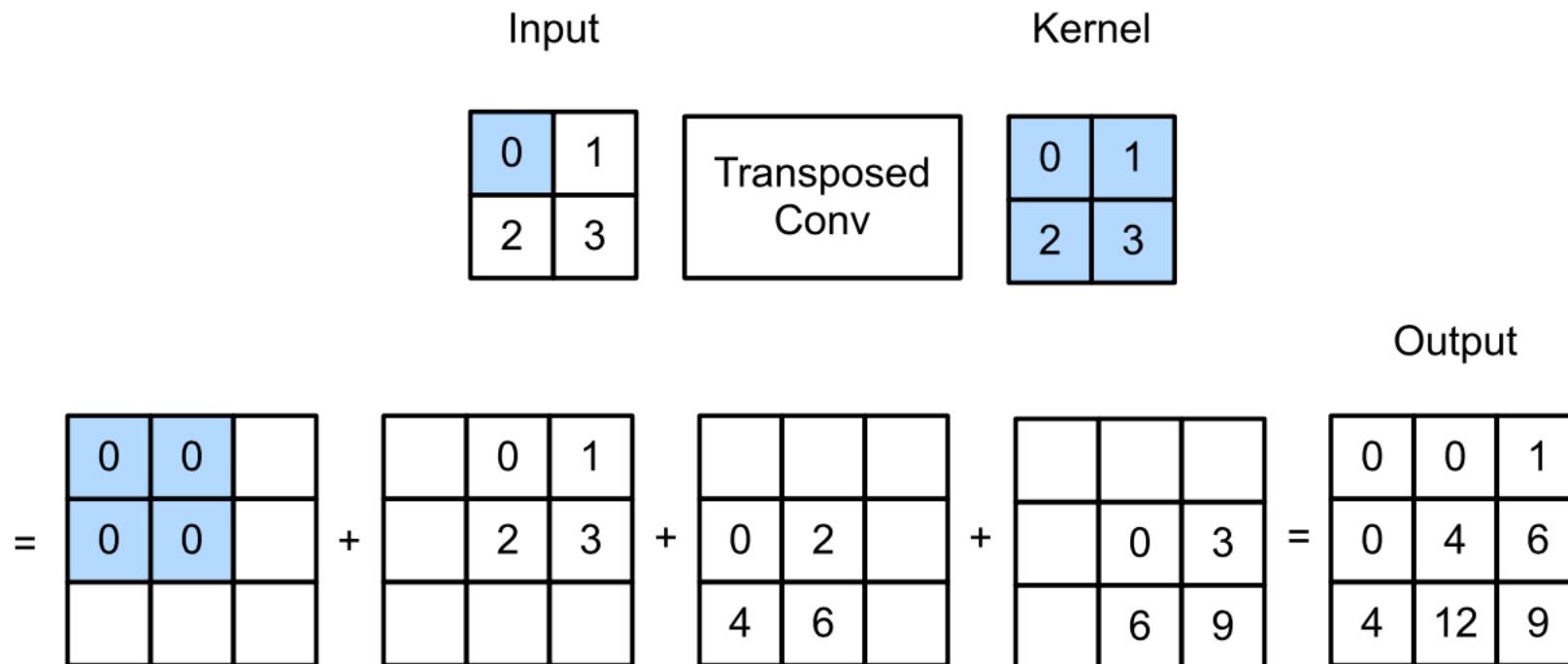
Transposed Convolution: 1D example



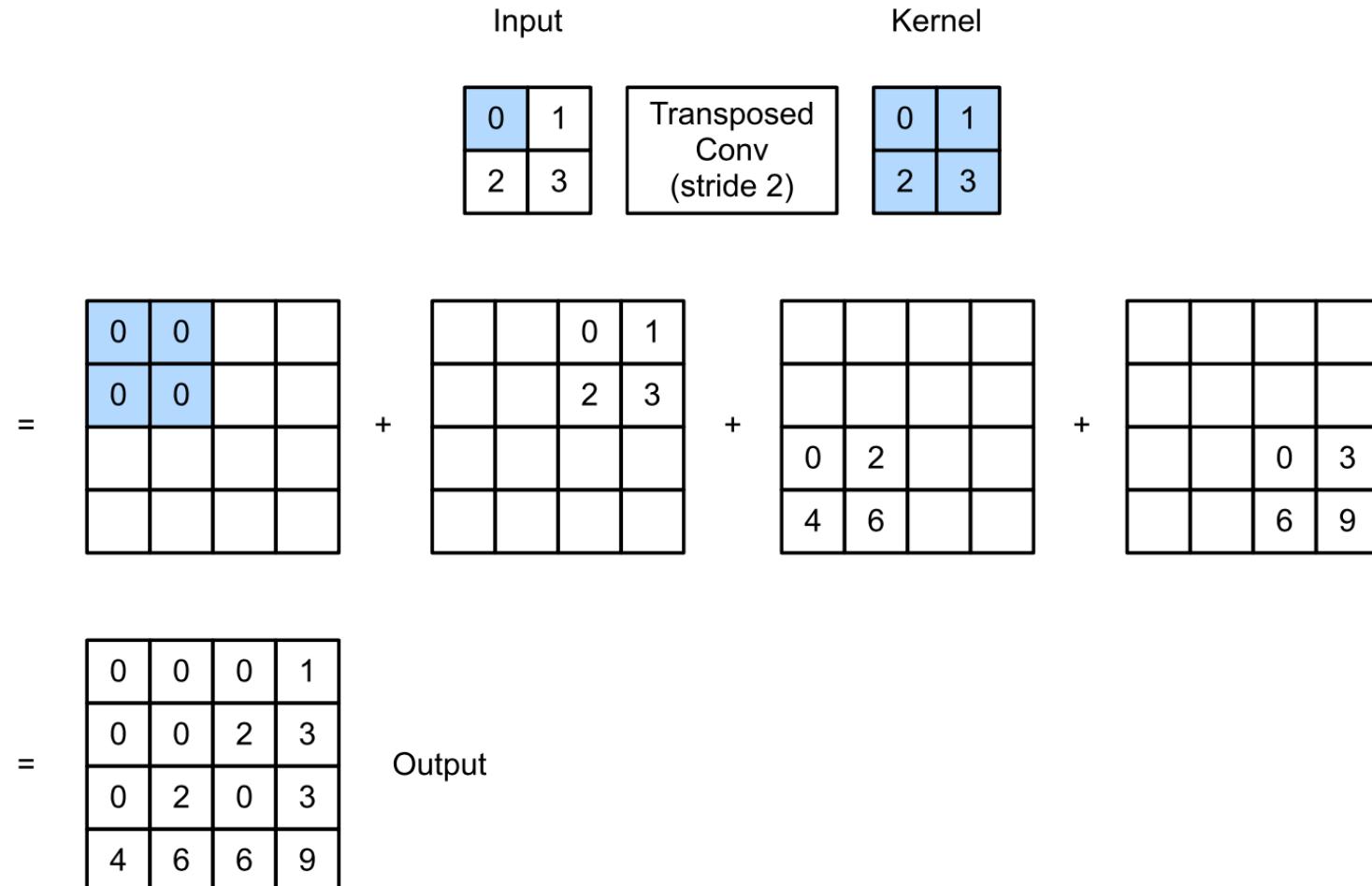
This has many names:

- Deconvolution (bad)!
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution
- Transposed Convolution (best name)

Transposed Convolution: 2D example



Transposed Convolution: 2D example



Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

Transposed convolution multiplies by the transpose of the same matrix:

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

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

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

When stride=1, transposed conv is just a regular conv (with different padding rules)

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

Transposed convolution multiplies by the transpose of the same matrix:

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

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

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

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:

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

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

When stride>1, transposed convolution cannot be expressed as normal conv

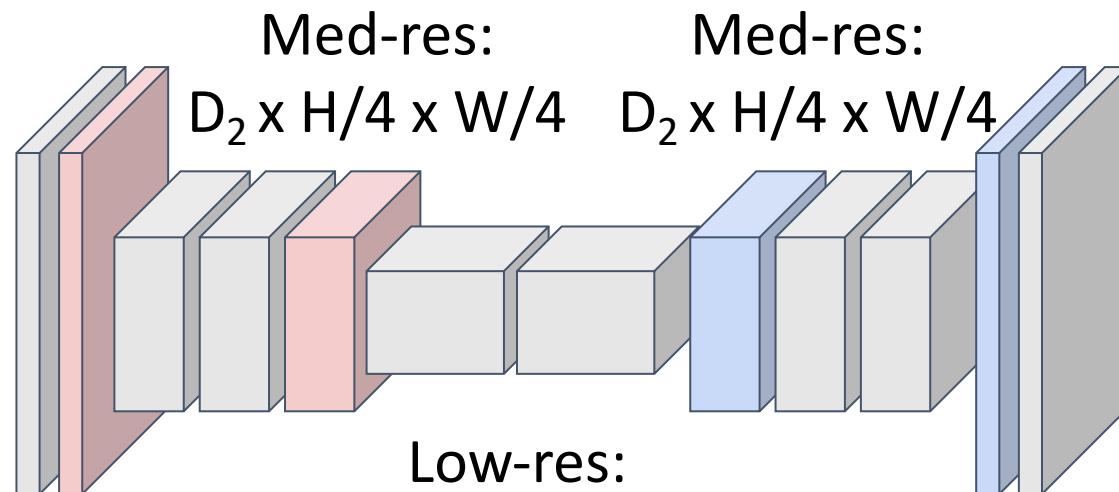
Semantic Segmentation: Fully Convolutional Network

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

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



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

Upsampling:
Interpolation,
transposed conv

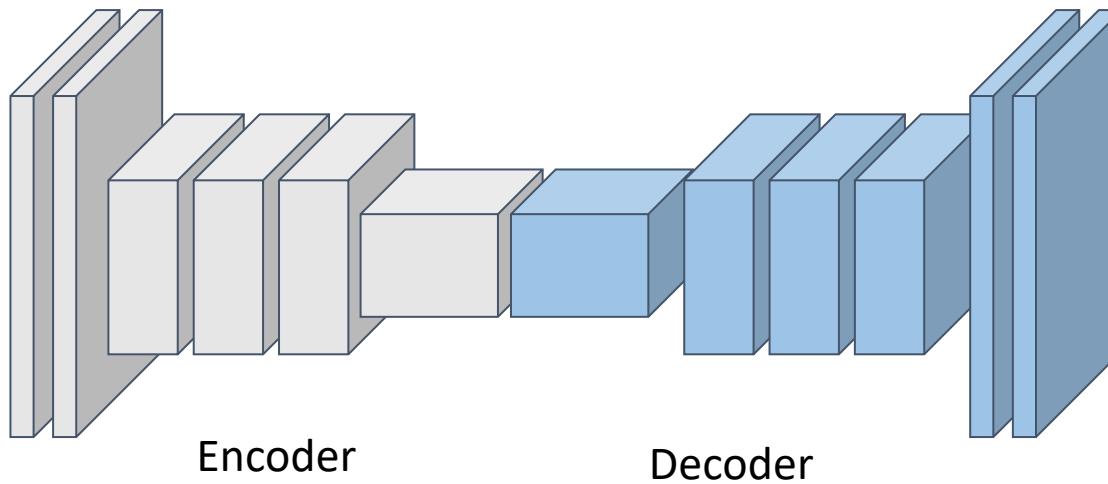


Loss function: Per-Pixel cross-entropy

Semantic segmentation networks (encoder+decoder)

Encoder side: How you can better encode spatial context information?

Decoder side: How you can maintain the spatial structure of the mask?

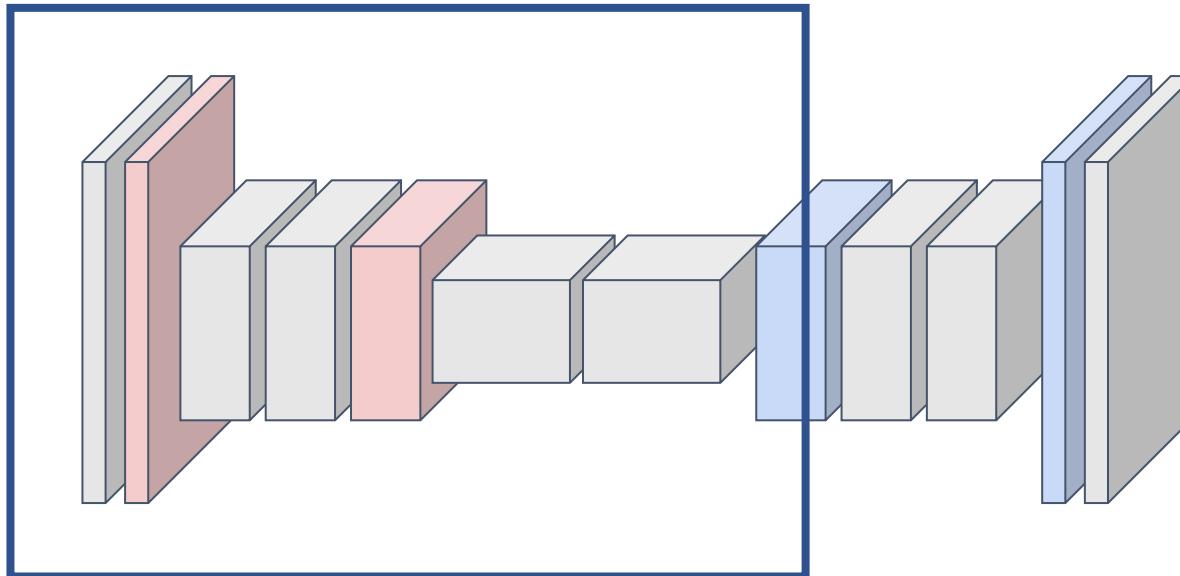


See the list of recent models at <https://github.com/open-mmlab/mmsegmentation>

Innovation on encoder: How to better encode the scene context



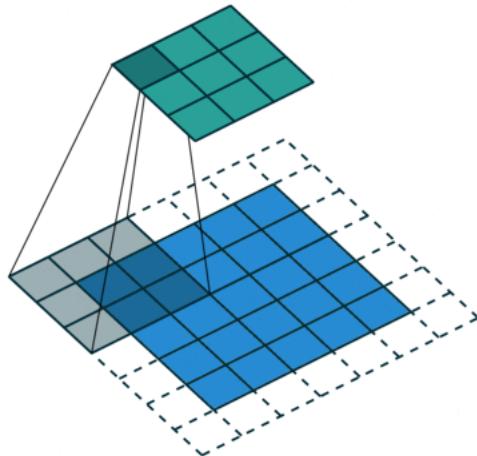
Input:
 $3 \times H \times W$



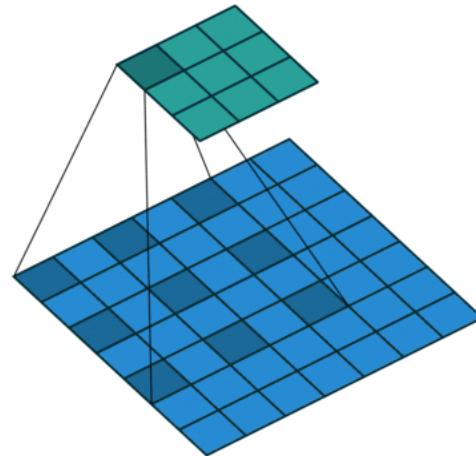
Predictions:
 $H \times W$

Dilated Convolution / Astrous Convolution

Standard Convolution

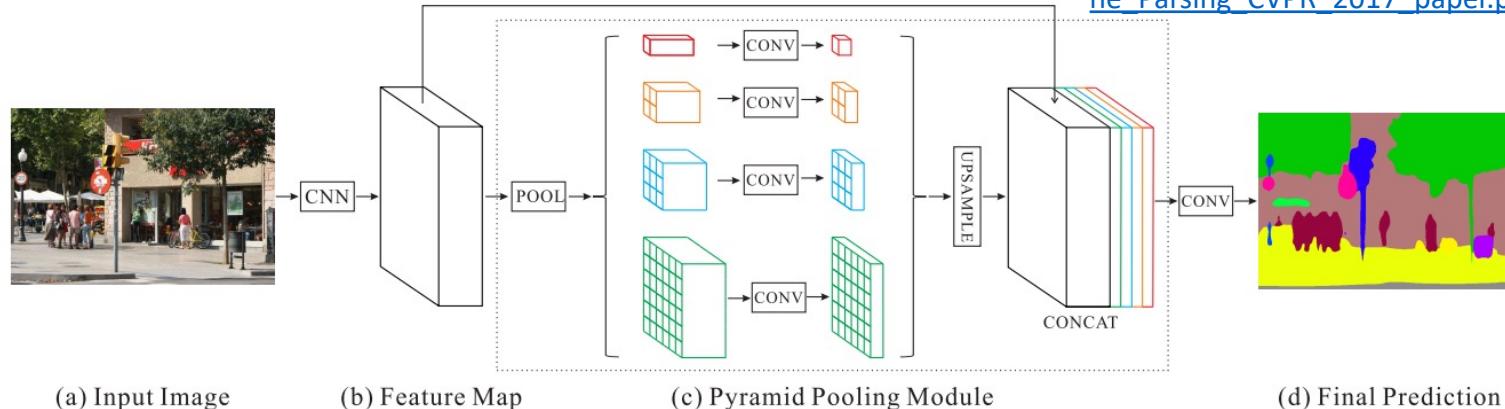


Dilated Convolution: a larger receptive field



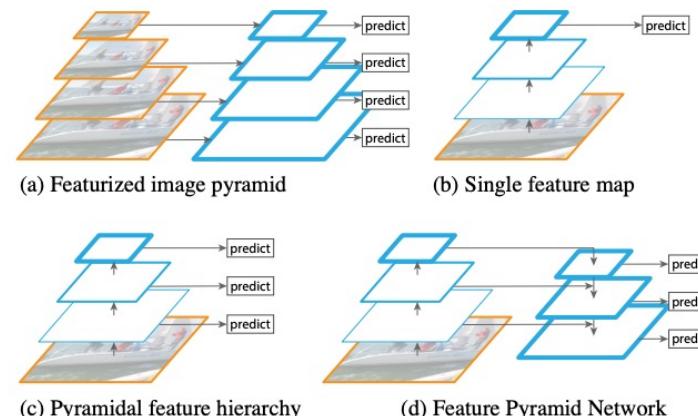
Feature Pyramid Structures

Pyramid Scene Parsing Network (PSPnet)



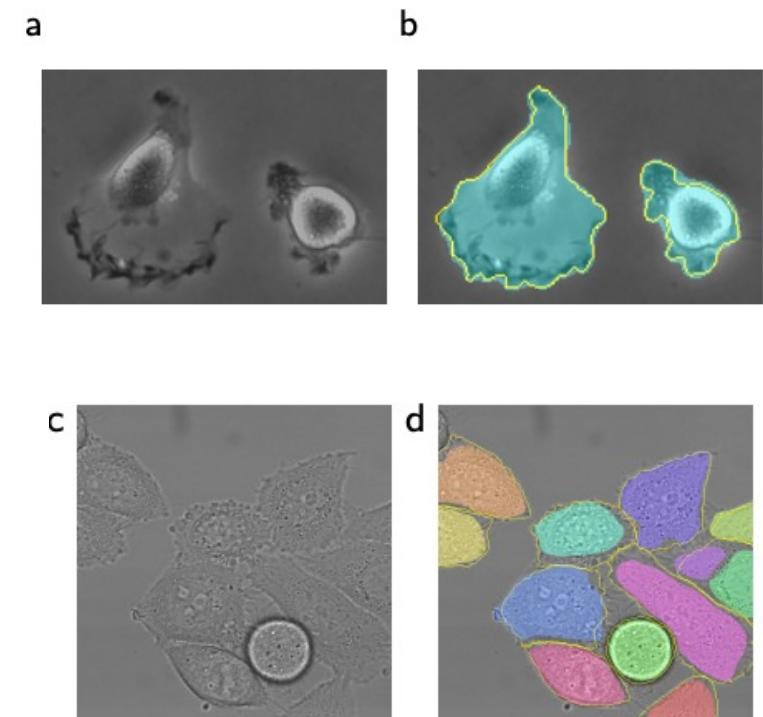
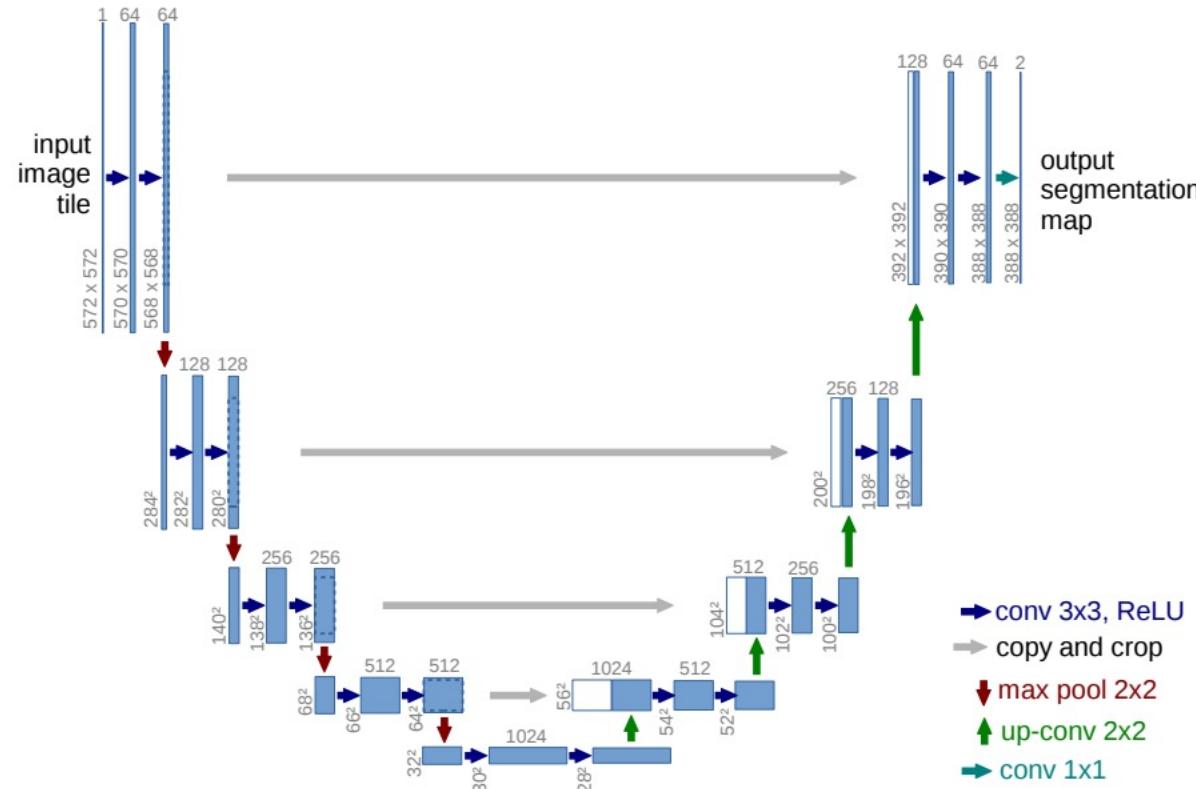
https://openaccess.thecvf.com/content_cvpr_2017/papers/Zhao_Pyramid_Scene_Parsing_CVPR_2017_paper.pdf

Feature Pyramid Networks (FPN)



<https://arxiv.org/pdf/1612.03144.pdf>

U-Net: a popular network for biomedical image segmentation (>38,000 citations)



U-Net: Convolutional Networks for Biomedical Image Segmentation <https://arxiv.org/pdf/1505.04597.pdf>
Olaf Ronneberger, Philipp Fischer, Thomas Brox

Semantic segmentation datasets

Cityscapes Dataset
(30 classes and 5 K images in driving scenes)



<https://www.cityscapes-dataset.com/>

MIT ADE20K
(150 classes, 25K training images)

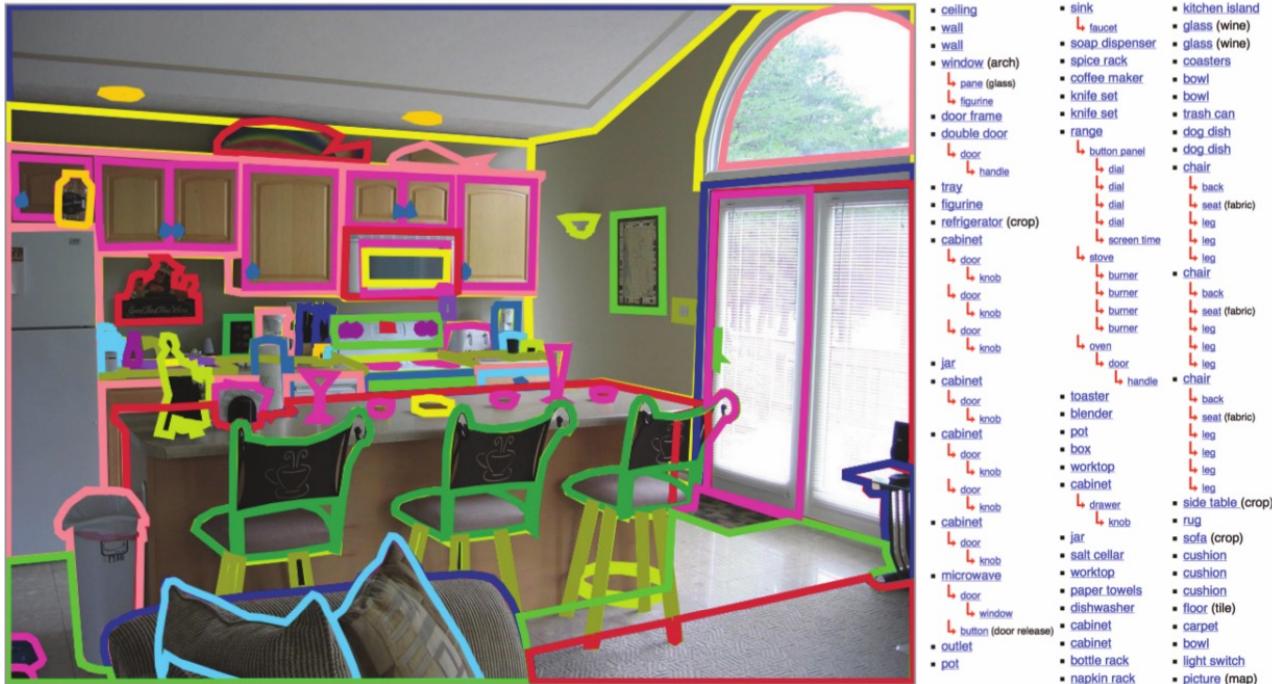


Scene Parsing through ADE20K Dataset. B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso and A. Torralba. Computer Vision and Pattern Recognition (CVPR), 2017.

Constructing ADE20K Dataset

Ms. Adela Barriuso

- Annotating each object instances in a scene
- Single expert annotator for a few years of work



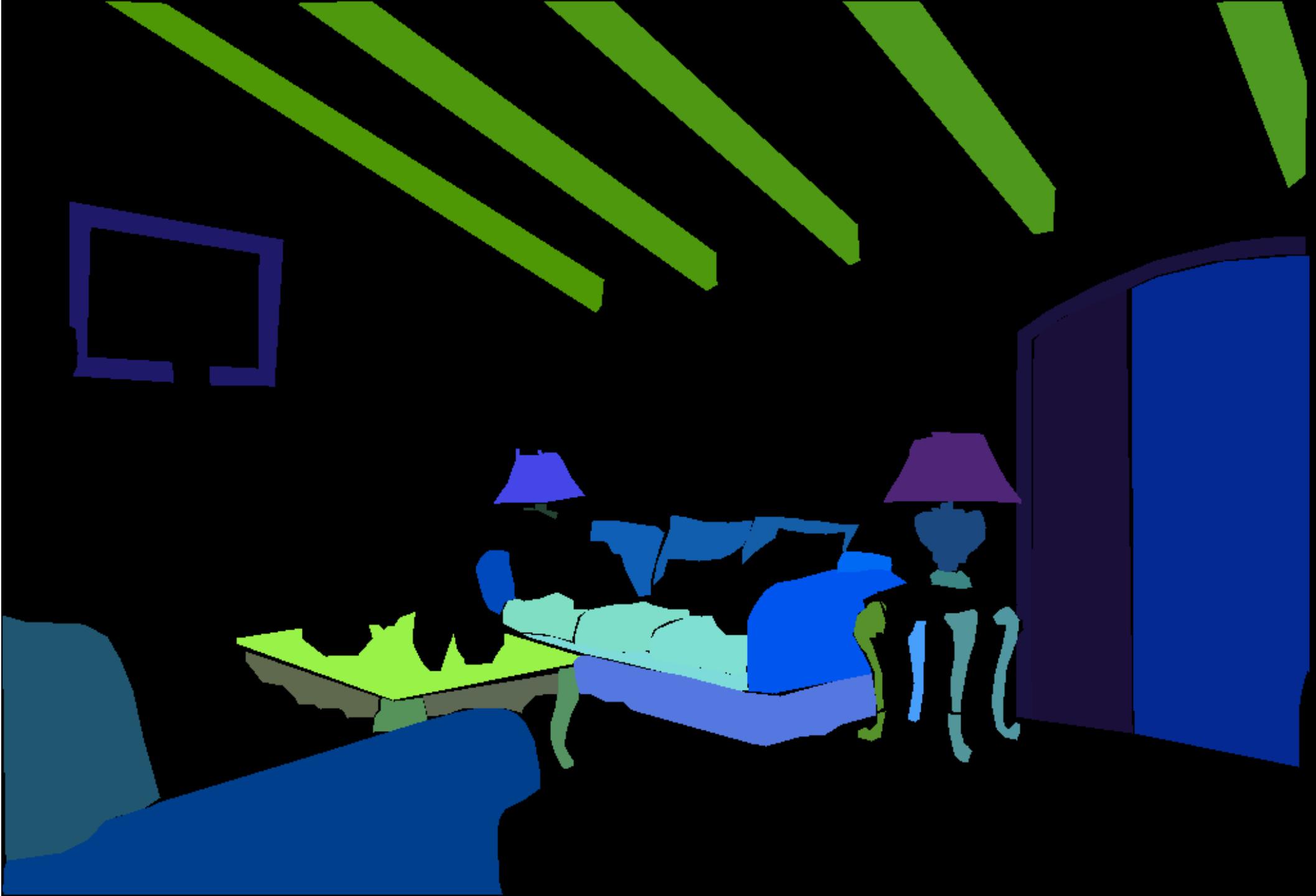
Labelme Annotation Tool

<http://groups.csail.mit.edu/vision/datasets/ADE20K/>





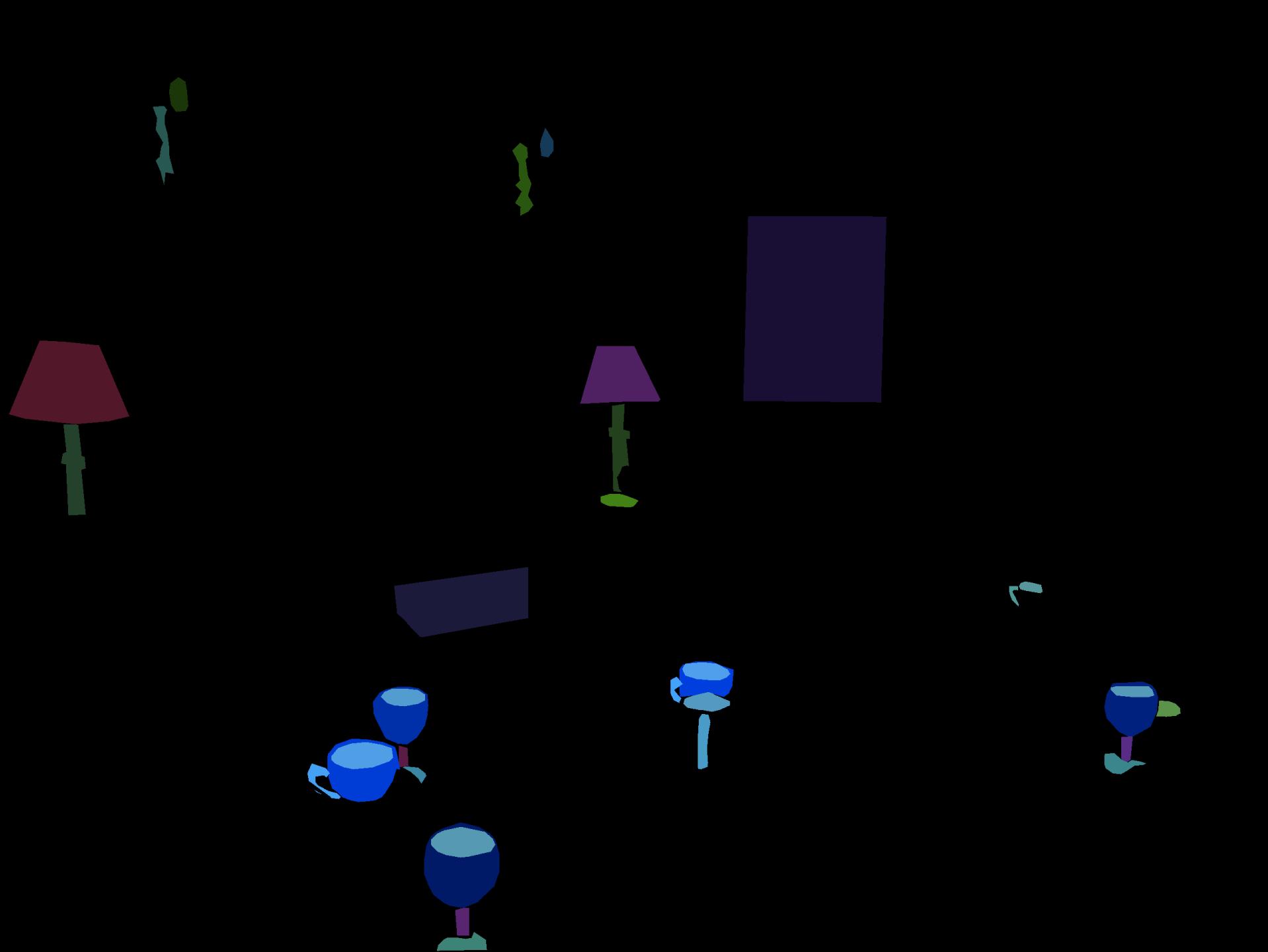




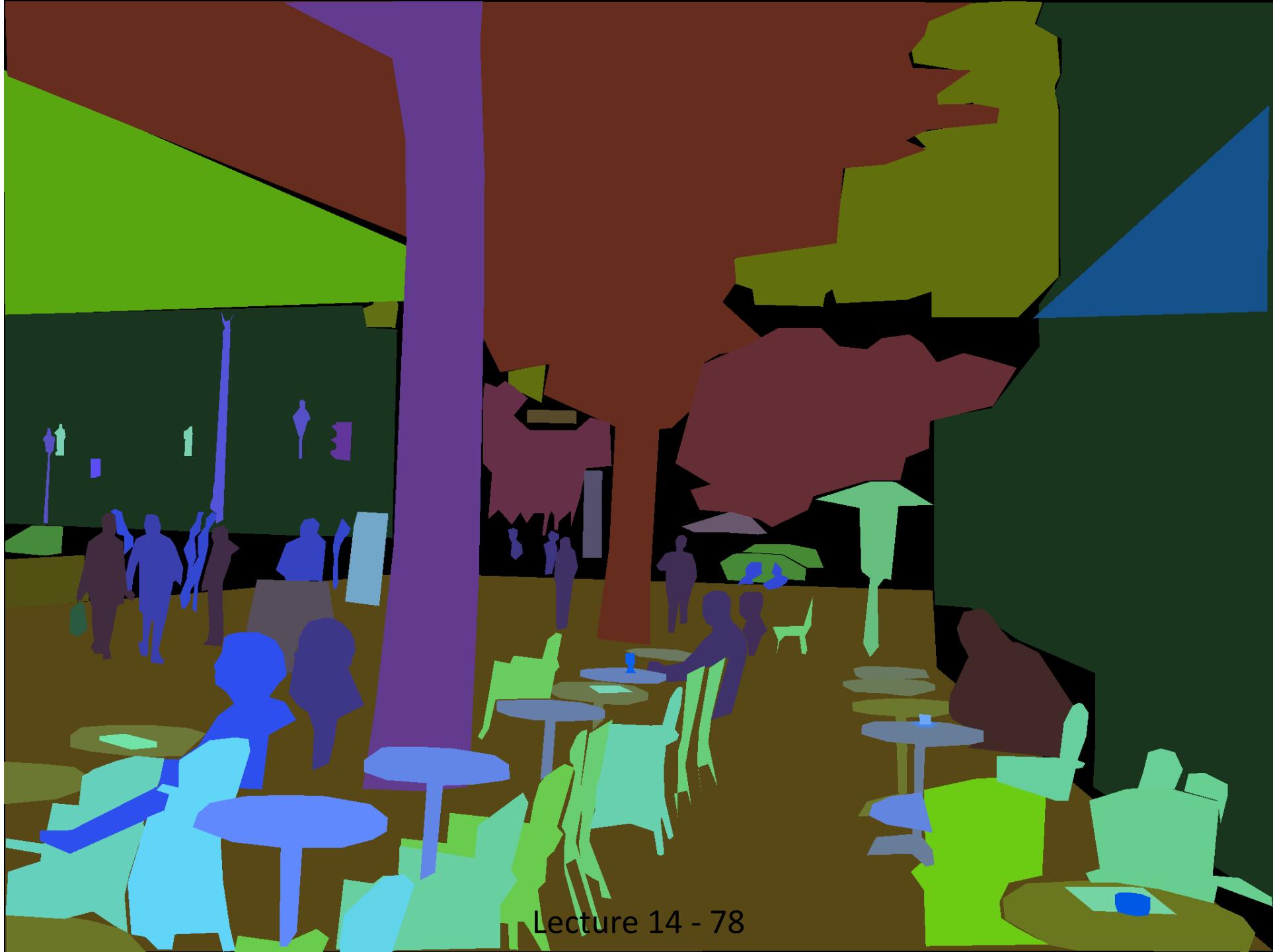


Lecture 14 - 74

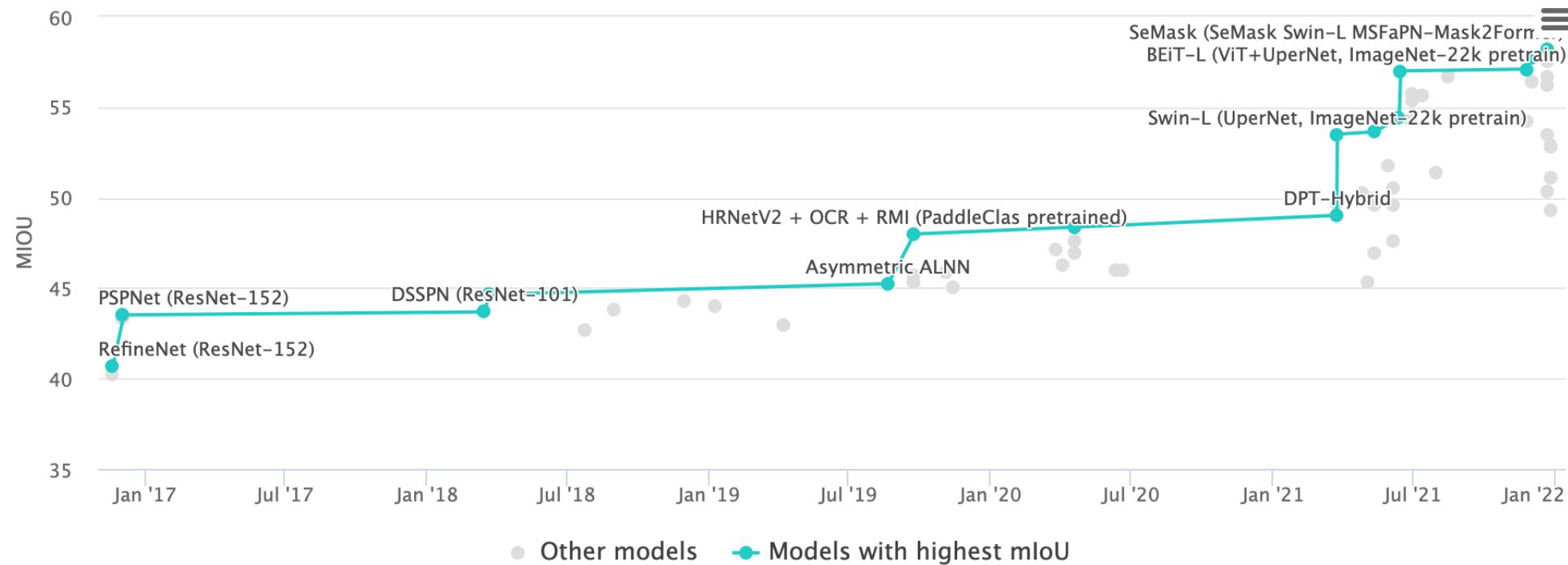








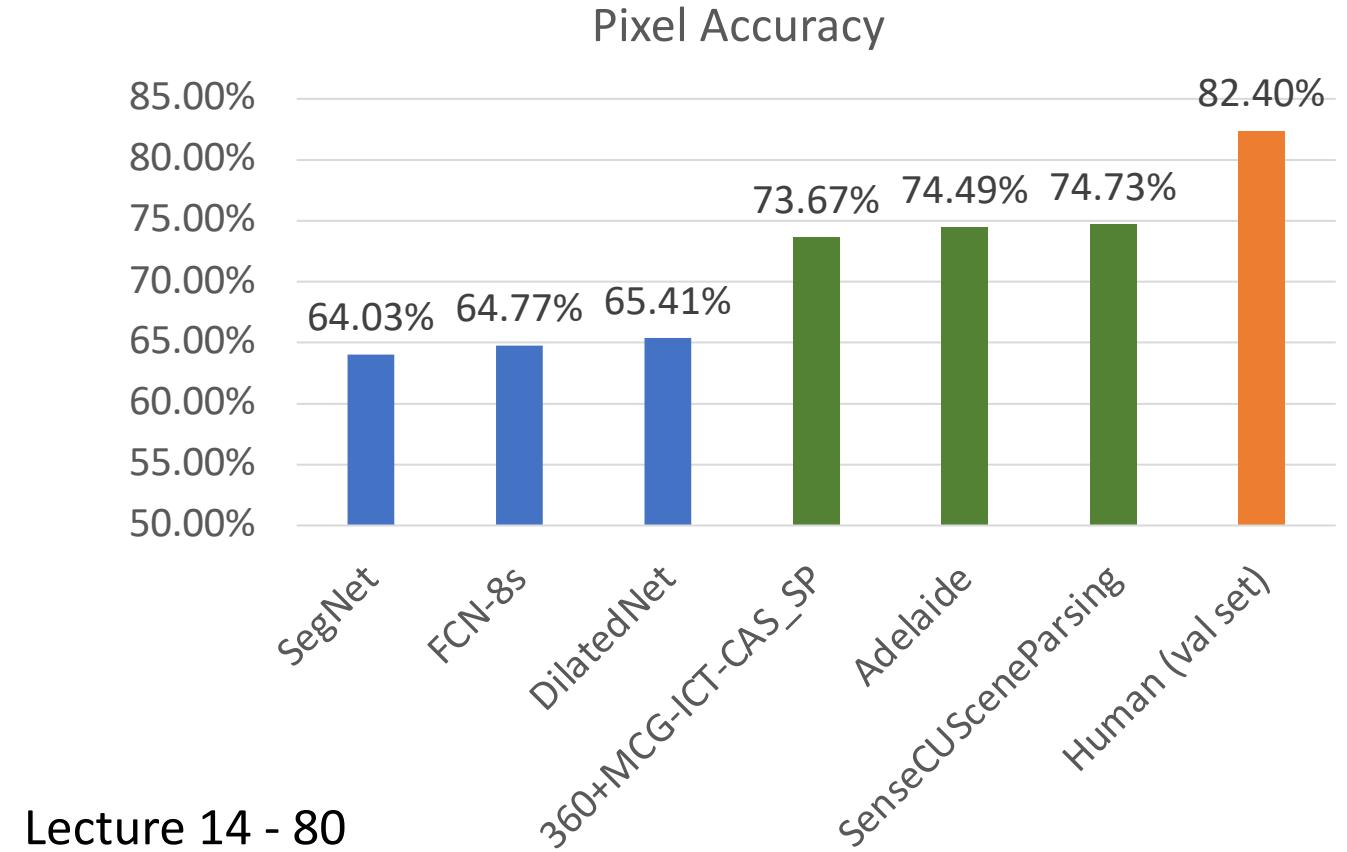
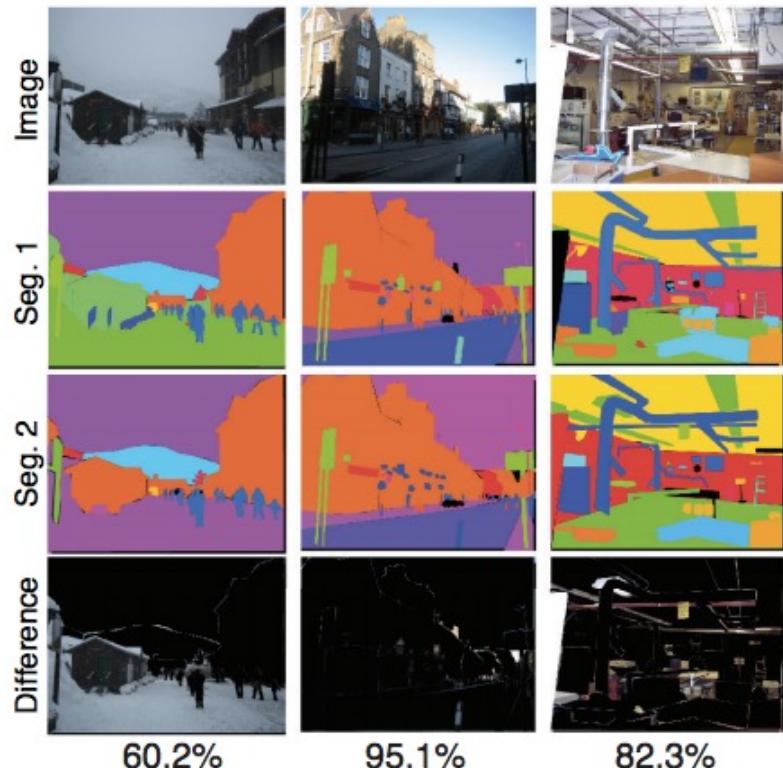
Progress of semantic segmentation on ADE20K



<https://paperswithcode.com/sota/semantic-segmentation-on-ade20k-val>

Data Consistency and Human Performance

- 61 images from val set are re-annotated after 6 months.
82.4% pixels got the same label.



Image

Ground-truth

SenseCU...

Adelaide

360+MCG...

SegModel

CASIA_IVA



96.55%



96.36%



96.27%



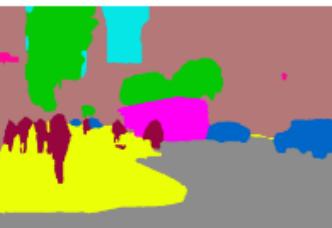
91.96%



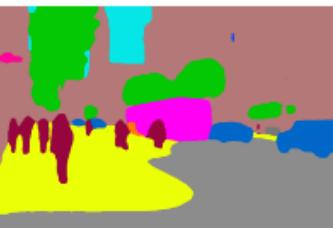
91.04%



92.73%



91.88%



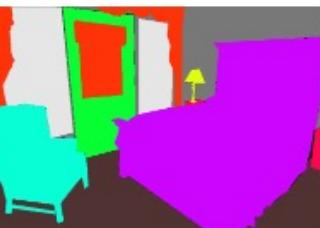
93.02%



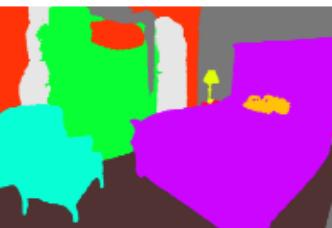
92.62%



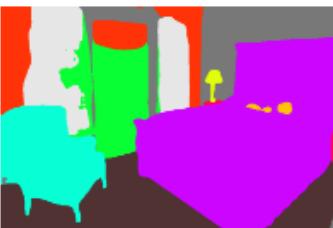
91.47%



92.25%



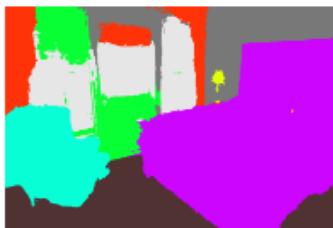
85.66%



90.47%



79.81%



84.85%



94.14%



95.63%



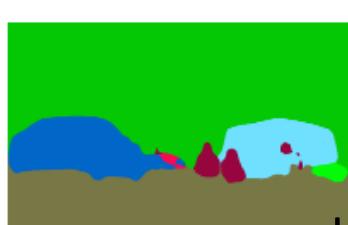
96.55%



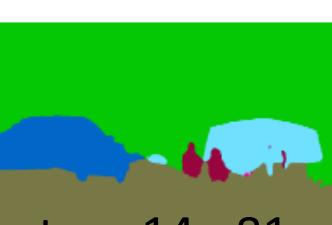
94.19%



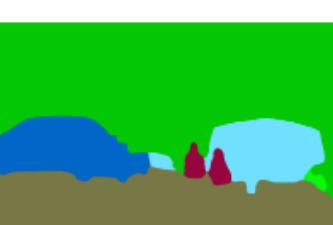
89.59%



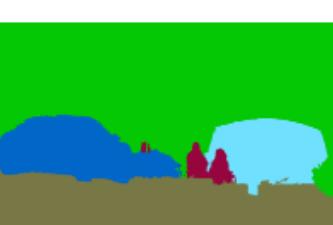
93.71%



93.51%



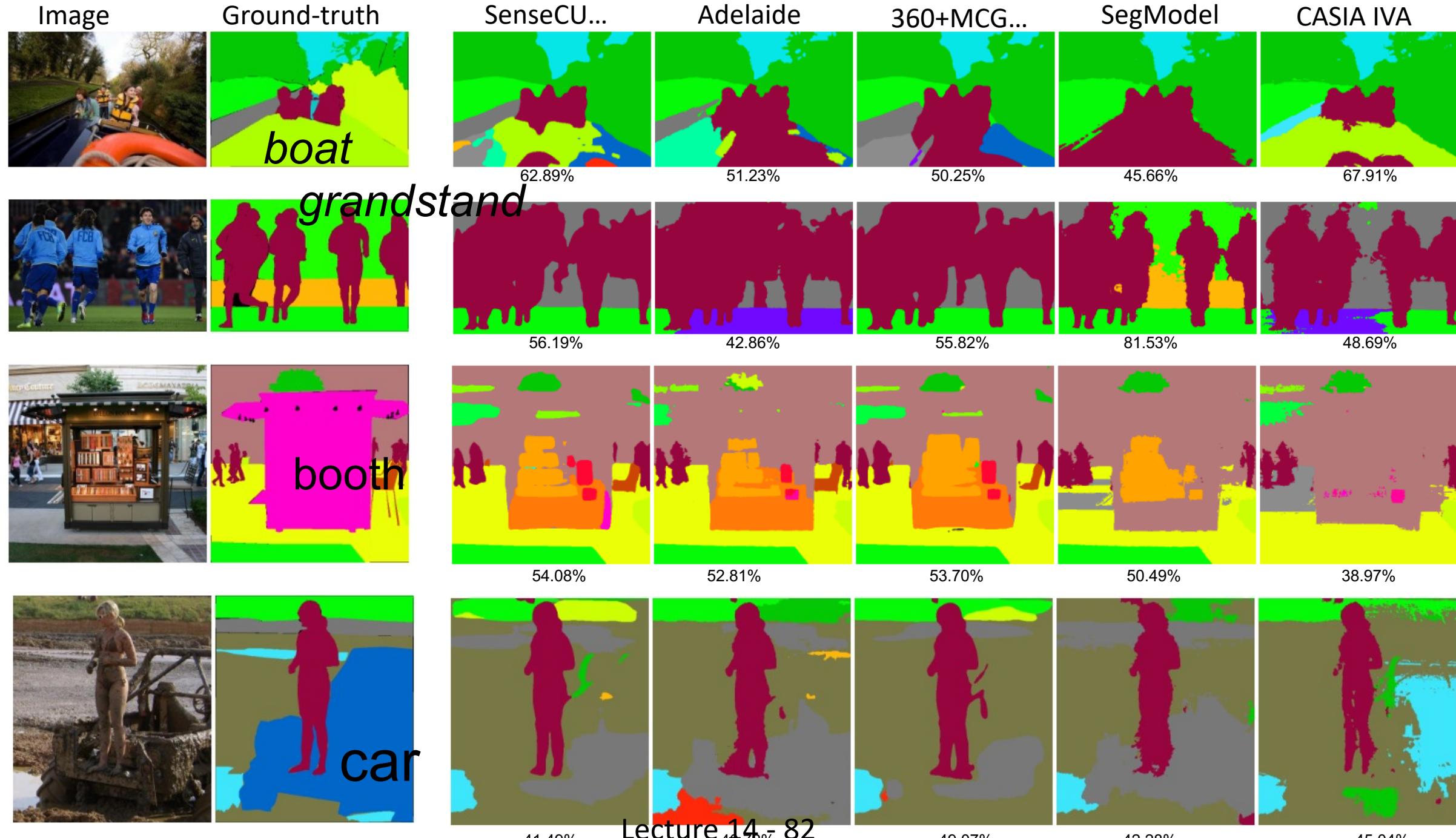
94.84%



94.89%



94.54%



Semantic segmentation libraries

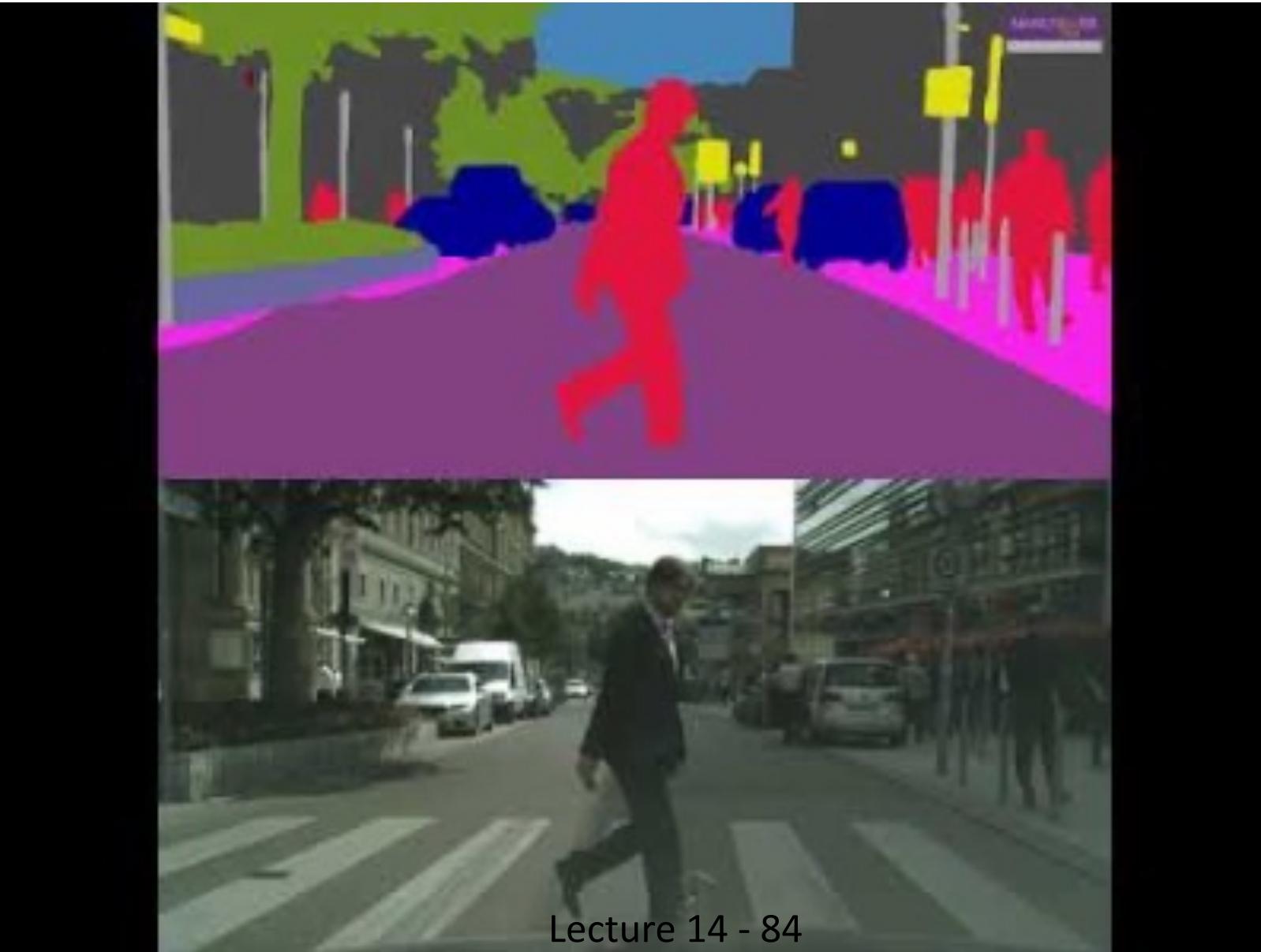
- Semantic segmentation on MIT ADE20K:
 - Link: <https://github.com/CSAILVision/semantic-segmentation-pytorch>
 - Colab: <https://colab.research.google.com/github/CSAILVision/semantic-segmentation-pytorch/blob/master/notebooks/DemoSegmenter.ipynb>
- mmSegmentattion: a zoo of any segmentation NNs
 - Link: <https://github.com/open-mmlab/mmsegmentation>

Supported backbones:
<input checked="" type="checkbox"/> ResNet (CVPR'2016)
<input checked="" type="checkbox"/> ResNeXt (CVPR'2017)
<input checked="" type="checkbox"/> HRNet (CVPR'2019)
<input checked="" type="checkbox"/> ResNeSt (ArXiv'2020)
<input checked="" type="checkbox"/> MobileNetV2 (CVPR'2018)
<input checked="" type="checkbox"/> MobileNetV3 (ICCV'2019)
<input checked="" type="checkbox"/> Vision Transformer (ICLR'2021)
<input checked="" type="checkbox"/> Swin Transformer (ICCV'2021)
<input checked="" type="checkbox"/> Twins (NeurIPS'2021)

Supported methods:
<input checked="" type="checkbox"/> FCN (CVPR'2015/TPAMI'2017)
<input checked="" type="checkbox"/> ERFNet (T-ITS'2017)
<input checked="" type="checkbox"/> UNet (MICCAI'2016/Nat. Methods'2019)
<input checked="" type="checkbox"/> PSPNet (CVPR'2017)
<input checked="" type="checkbox"/> DeepLabV3 (ArXiv'2017)
<input checked="" type="checkbox"/> BiSeNetV1 (ECCV'2018)
<input checked="" type="checkbox"/> PSANet (ECCV'2018)
<input checked="" type="checkbox"/> DeepLabV3+ (CVPR'2018)
<input checked="" type="checkbox"/> UPerNet (ECCV'2018)
<input checked="" type="checkbox"/> ICNet (ECCV'2018)
<input checked="" type="checkbox"/> NonLocal Net (CVPR'2018)
<input checked="" type="checkbox"/> EncNet (CVPR'2018)
<input checked="" type="checkbox"/> Semantic FPN (CVPR'2019)
<input checked="" type="checkbox"/> DANet (CVPR'2019)
<input checked="" type="checkbox"/> APCNet (CVPR'2019)
<input checked="" type="checkbox"/> EMANet (ICCV'2019)
<input checked="" type="checkbox"/> GCNet (ICCVW'2019/TPAMI'2020)
<input checked="" type="checkbox"/> FastFCN (ArXiv'2019)
<input checked="" type="checkbox"/> Fast-SCNN (ArXiv'2019)
<input checked="" type="checkbox"/> ISANet (ArXiv'2019/IJCV'2021)
<input checked="" type="checkbox"/> OCRNet (ECCV'2020)
<input checked="" type="checkbox"/> DNLNet (ECCV'2020)
<input checked="" type="checkbox"/> PointRend (CVPR'2020)
<input checked="" type="checkbox"/> CGNet (TIP'2020)
<input checked="" type="checkbox"/> BiSeNetV2 (IJCV'2021)
<input checked="" type="checkbox"/> STDC (CVPR'2021)
<input checked="" type="checkbox"/> SETR (CVPR'2021)
<input checked="" type="checkbox"/> DPT (ArXiv'2021)
<input checked="" type="checkbox"/> Segmenter (ICCV'2021)
<input checked="" type="checkbox"/> SegFormer (NeurIPS'2021)

Supported datasets:
<input checked="" type="checkbox"/> Cityscapes
<input checked="" type="checkbox"/> PASCAL VOC
<input checked="" type="checkbox"/> ADE20K
<input checked="" type="checkbox"/> Pascal Context
<input checked="" type="checkbox"/> COCO-Stuff 10k
<input checked="" type="checkbox"/> COCO-Stuff 164k
<input checked="" type="checkbox"/> CHASE_DB1
<input checked="" type="checkbox"/> DRIVE
<input checked="" type="checkbox"/> HRF
<input checked="" type="checkbox"/> STARE
<input checked="" type="checkbox"/> Dark Zurich
<input checked="" type="checkbox"/> Nighttime Driving
<input checked="" type="checkbox"/> LoveDA
<input checked="" type="checkbox"/> Potsdam
<input checked="" type="checkbox"/> Vaihingen

Demo video of semantic segmentation on CityScapes



Computer Vision Tasks: Instance Segmentation

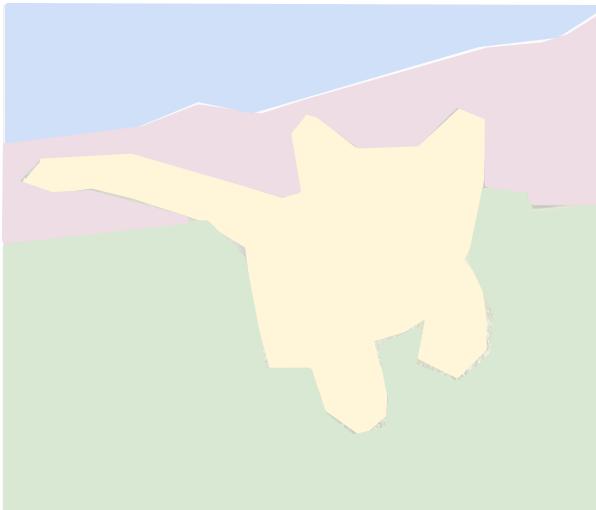
Classification



CAT

No spatial extent

Semantic
Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Objects

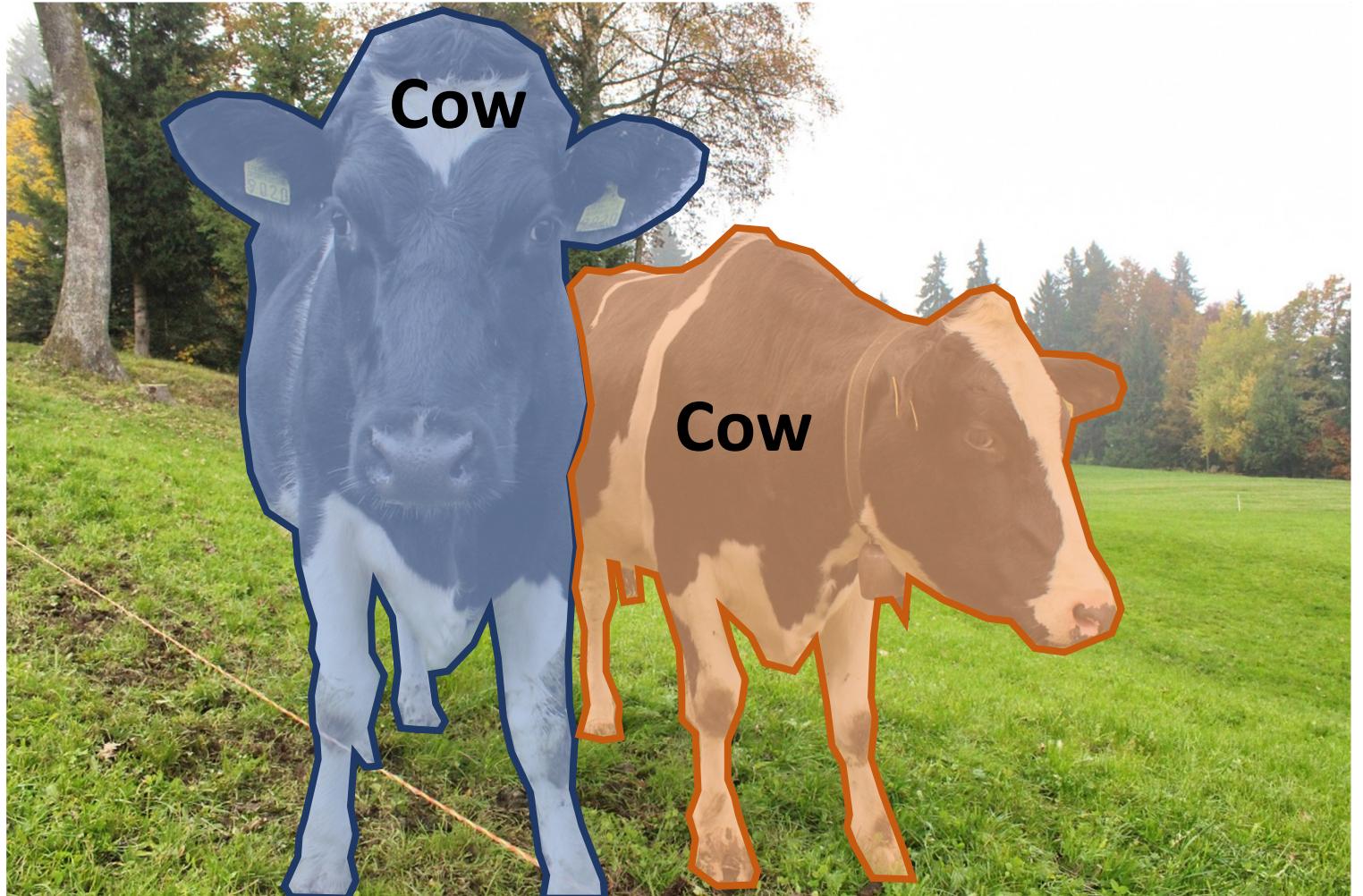
Instance
Segmentation



DOG, DOG, CAT

Computer Vision Tasks: Instance Segmentation

Instance Segmentation:
Detect all objects in the image, and identify the pixels that belong to each object (Only things!)



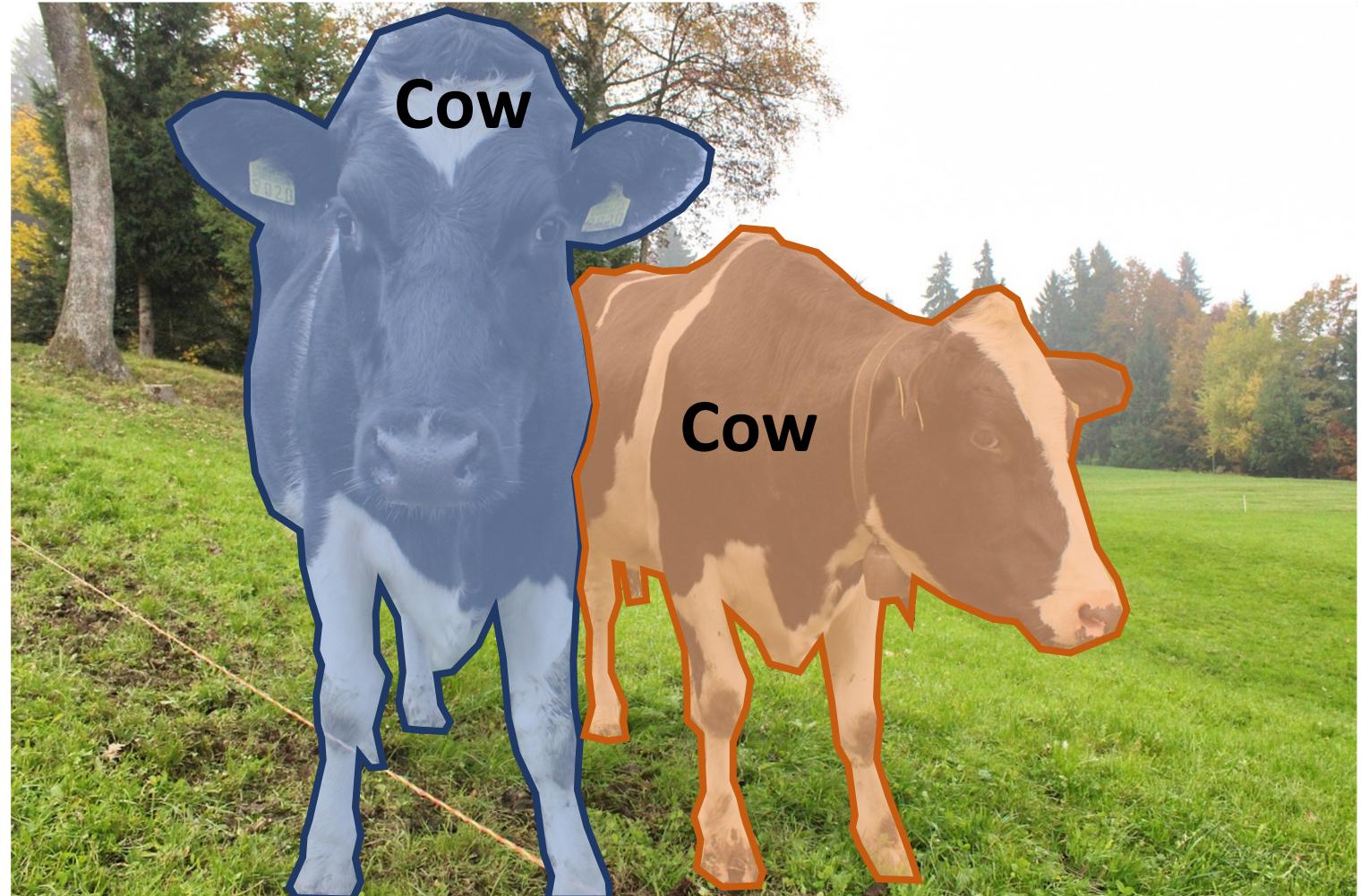
This image is CC0 public domain

Computer Vision Tasks: Instance Segmentation

Instance Segmentation:

Detect all objects in the image, and identify the pixels that belong to each object (Only things!)

Approach: Perform object detection, then predict a segmentation mask for each object!



This image is CC0 public domain

Object Detection: Faster R-CNN

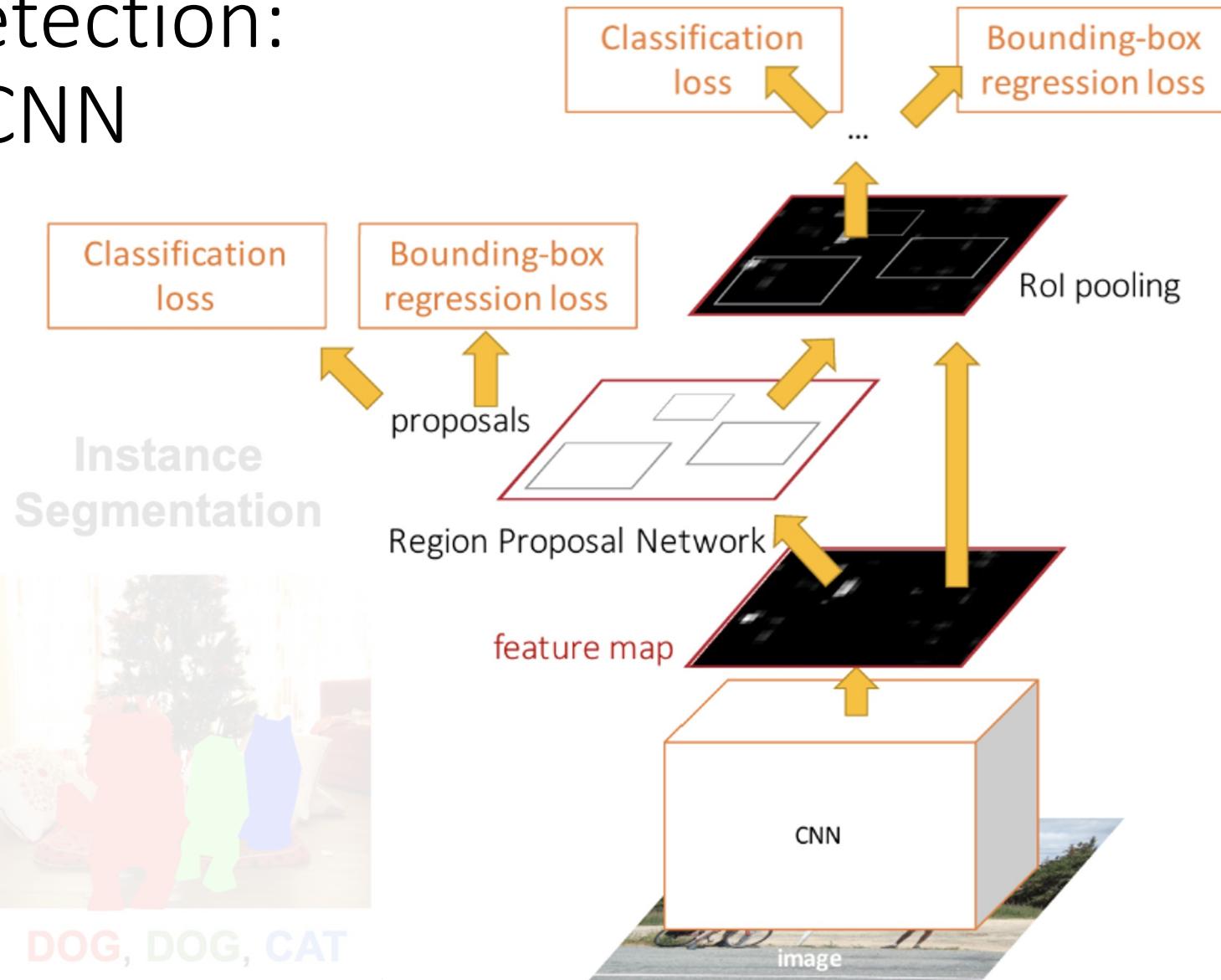
Object Detection



DOG, DOG, CAT



DOG, DOG, CAT



Instance Segmentation: Mask R-CNN

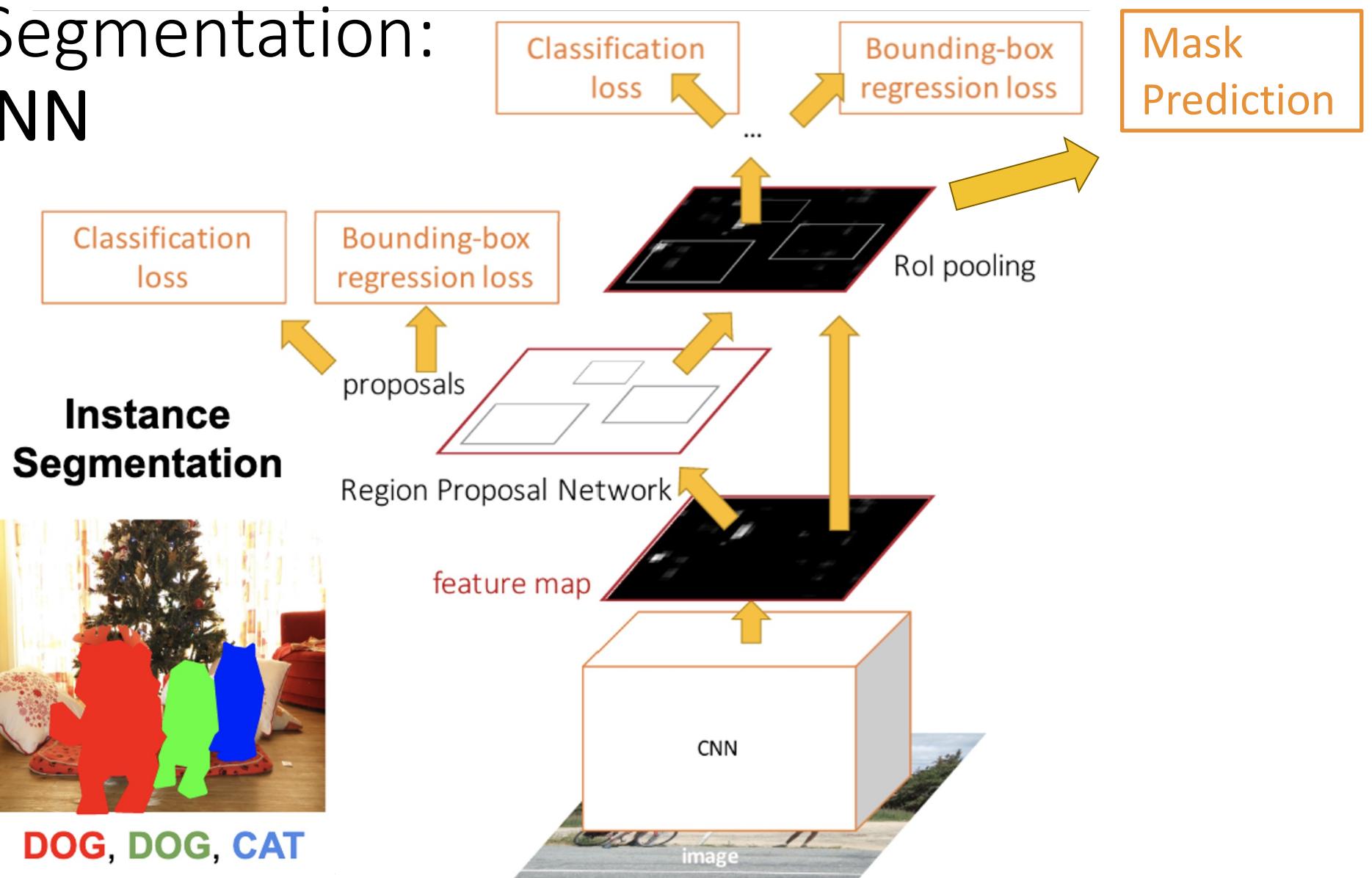
Object
Detection



DOG, DOG, CAT



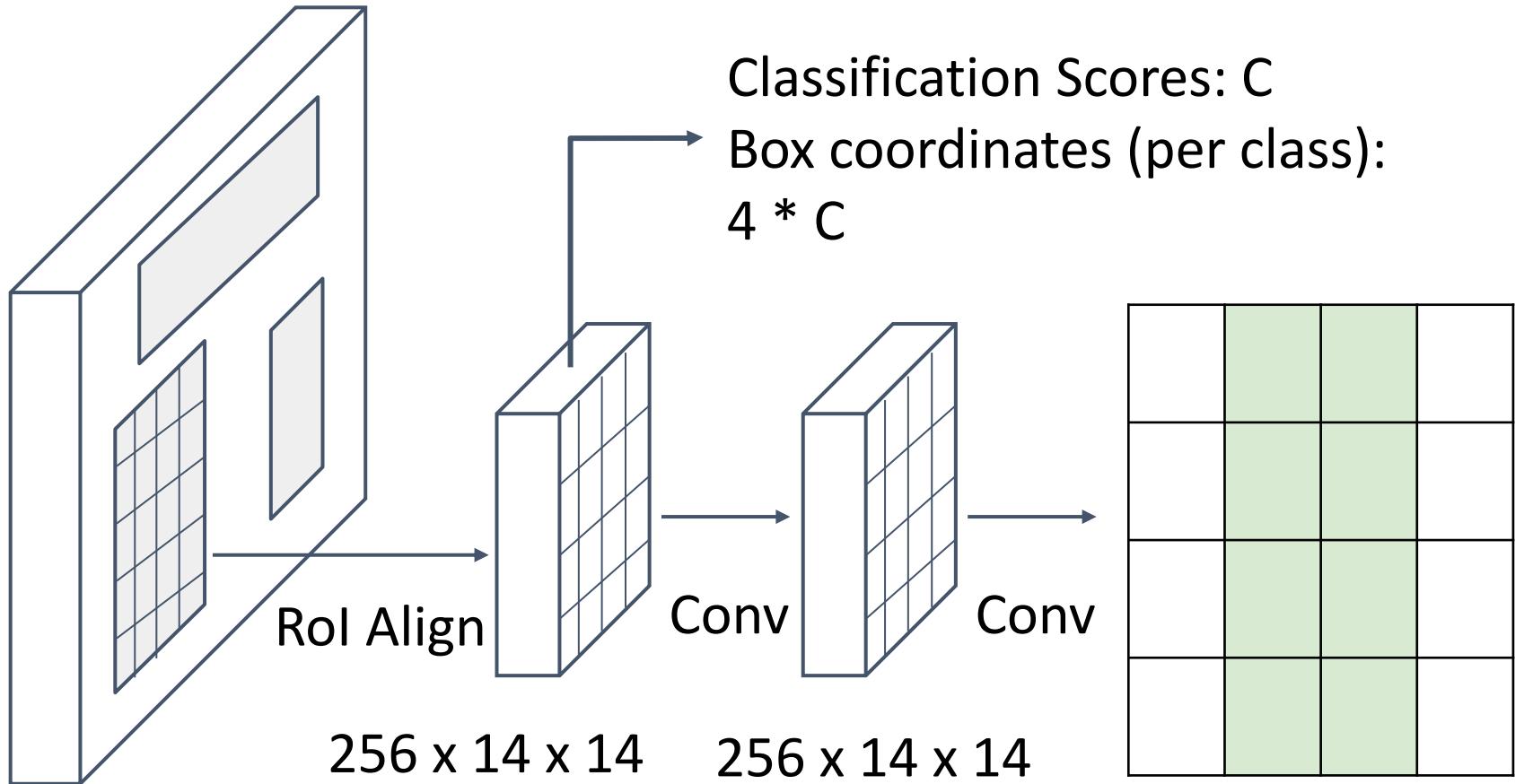
DOG, DOG, CAT



Mask R-CNN



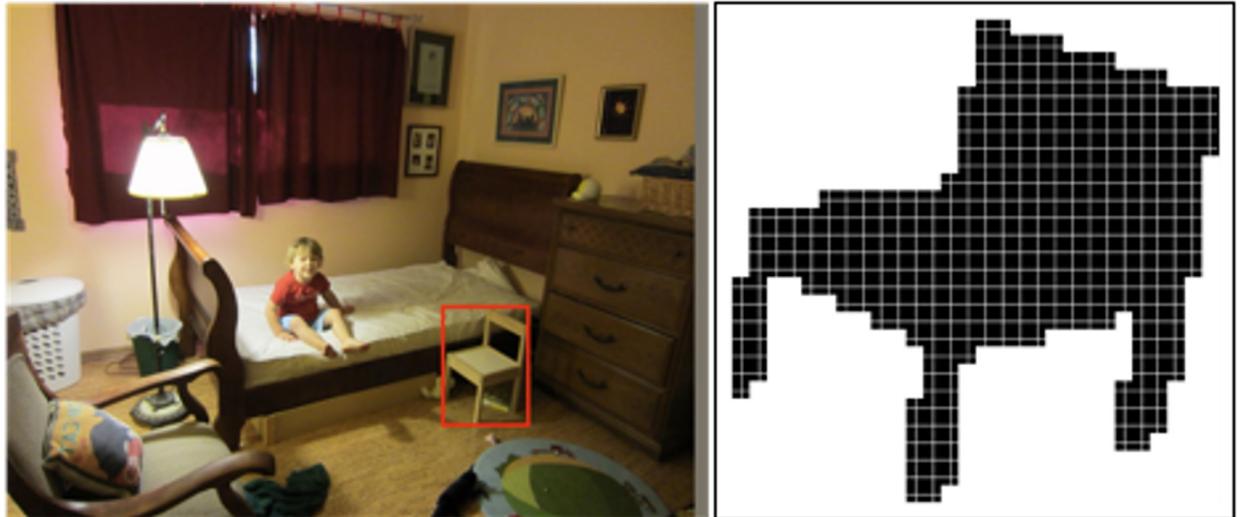
CNN
+RPN



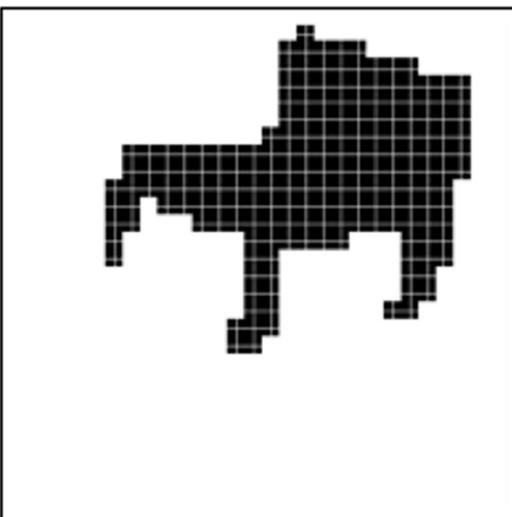
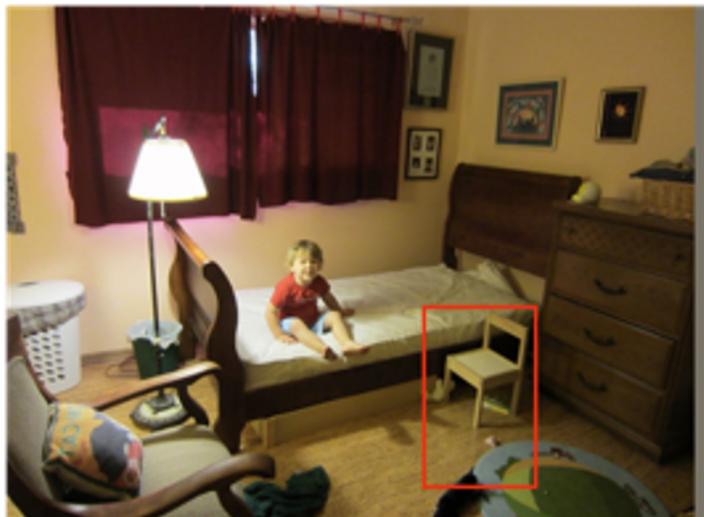
Classification Scores: C
Box coordinates (per class):
 $4 * C$

Predict a mask for
each of C classes:
 $C \times 28 \times 28$ (mask size)

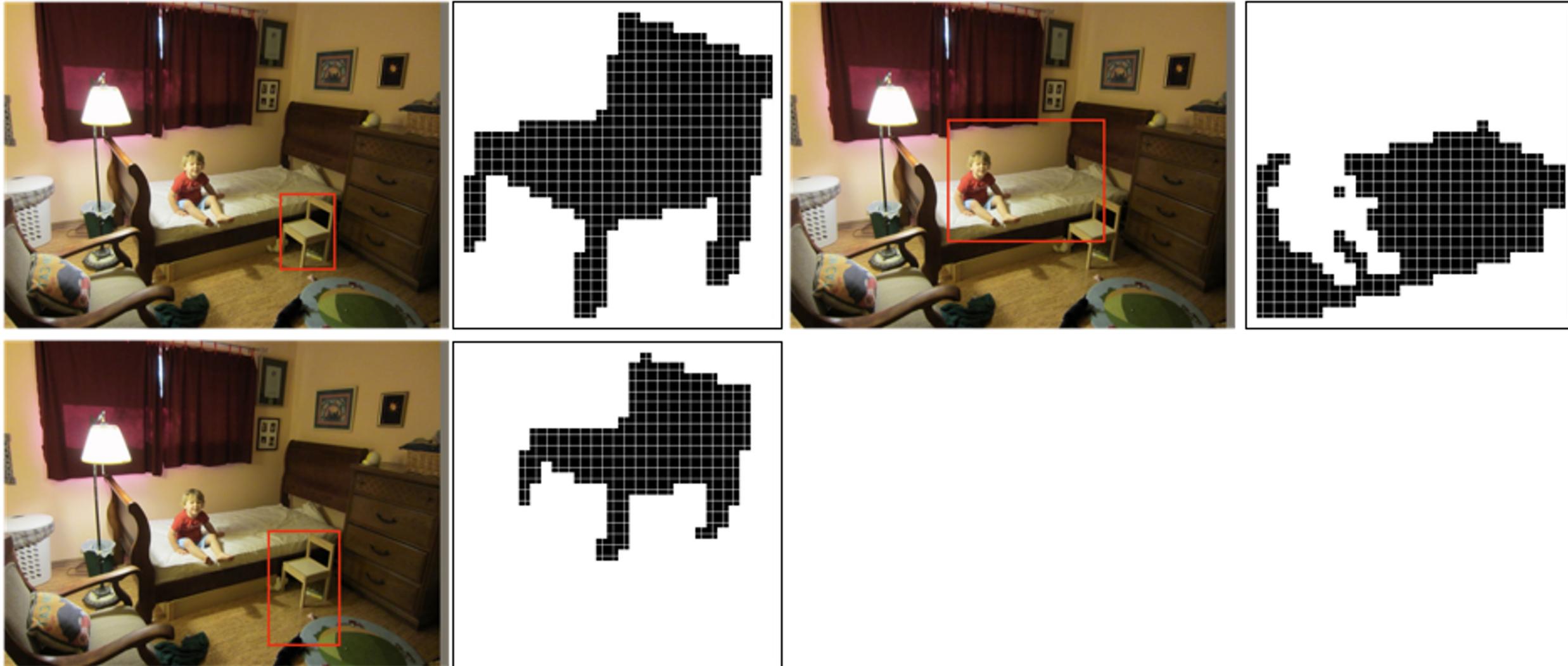
Mask R-CNN: Example Training Targets



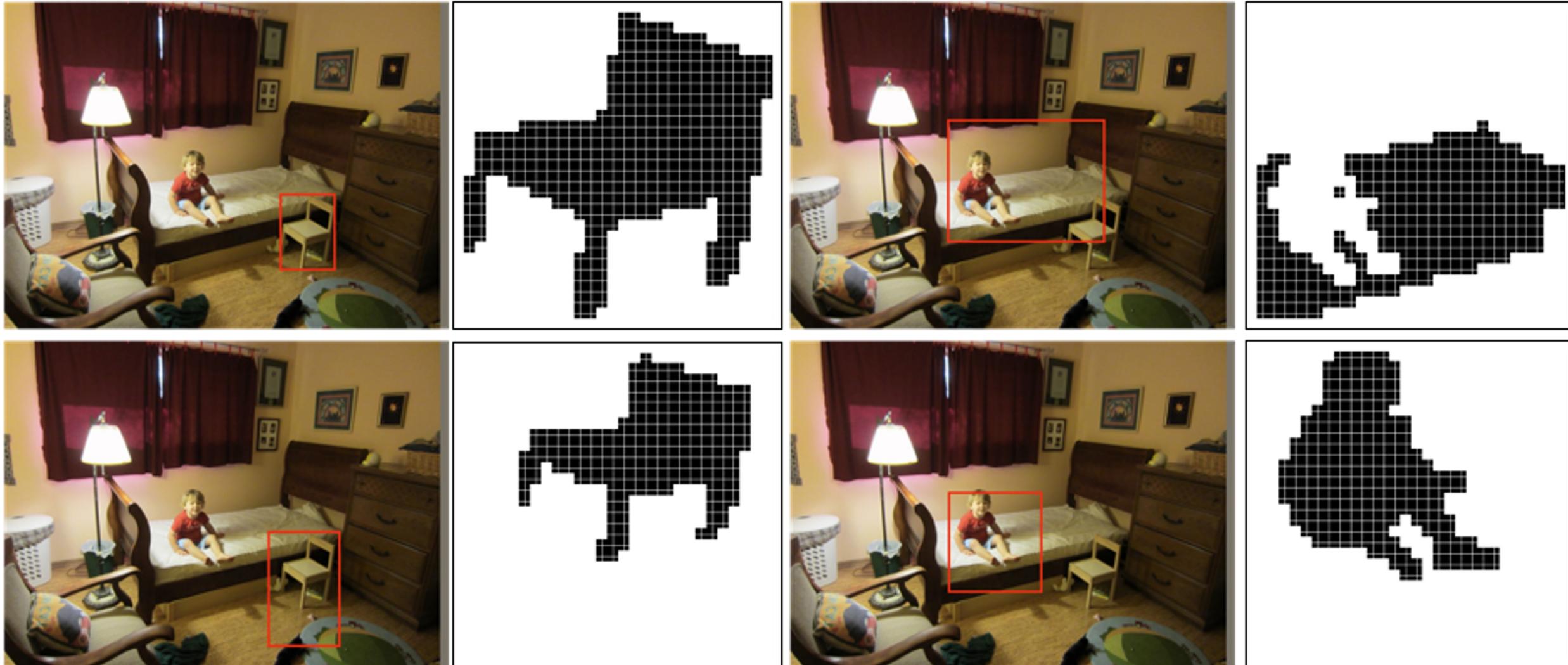
Mask R-CNN: Example Training Targets



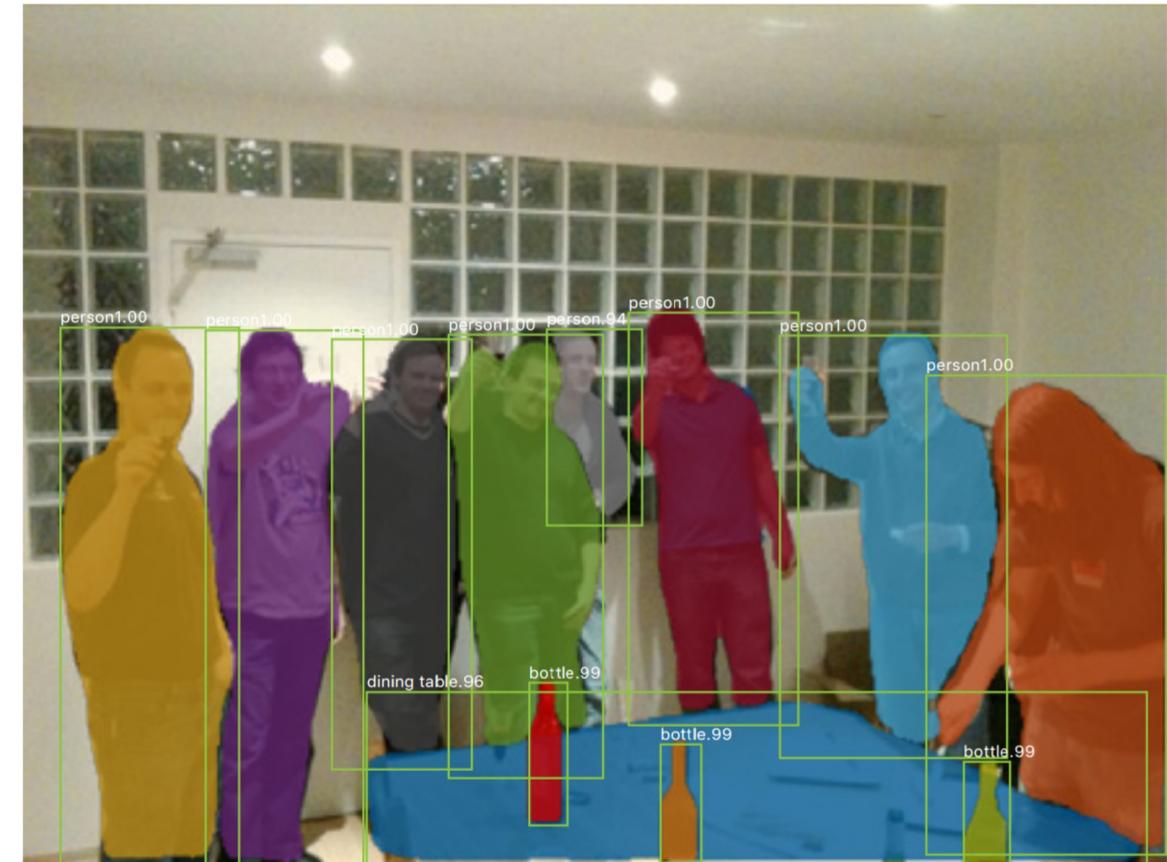
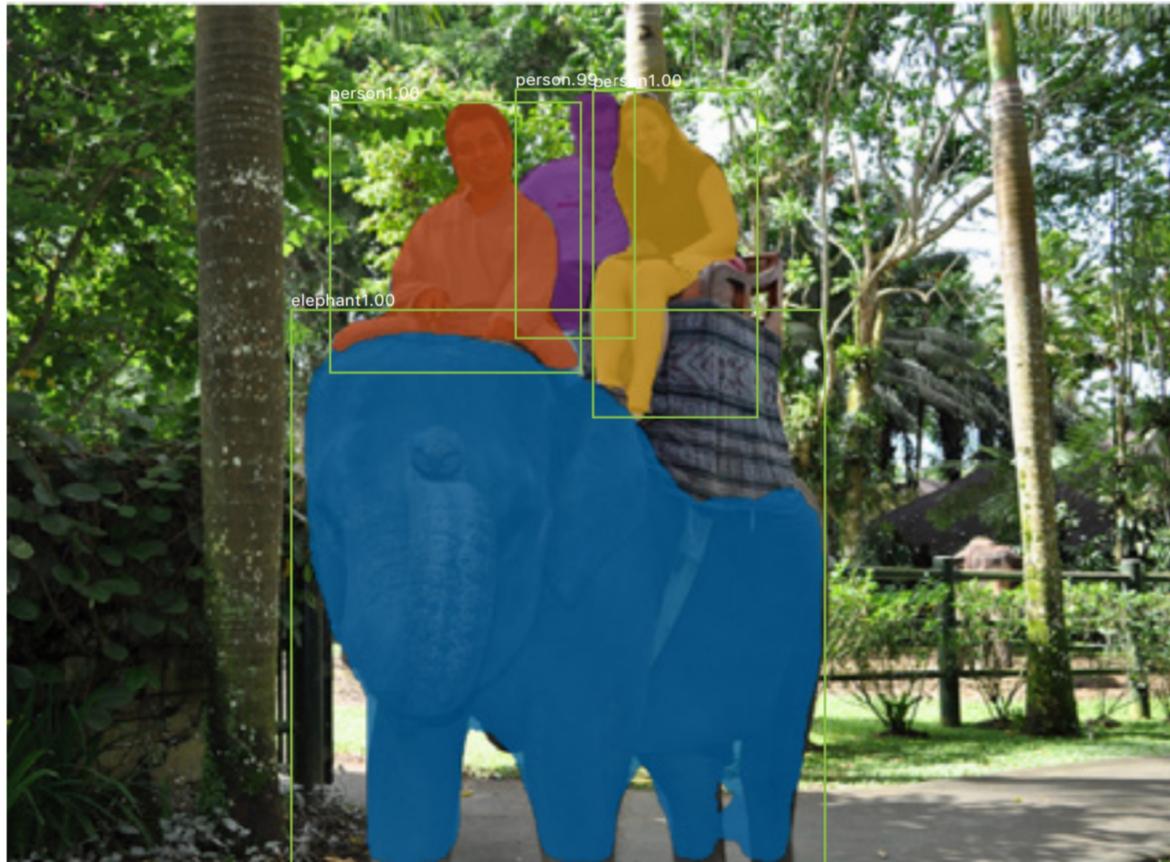
Mask R-CNN: Example Training Targets



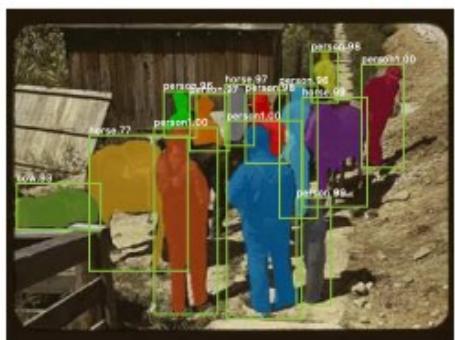
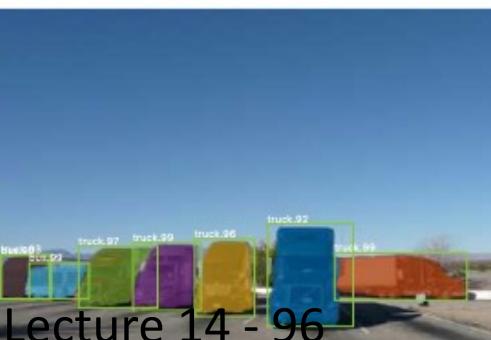
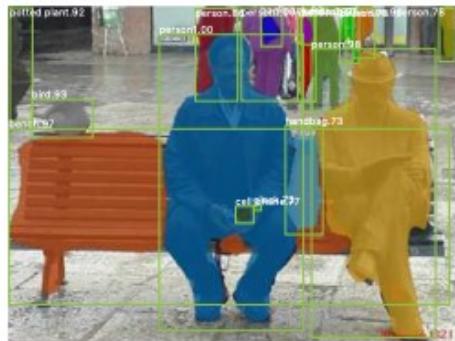
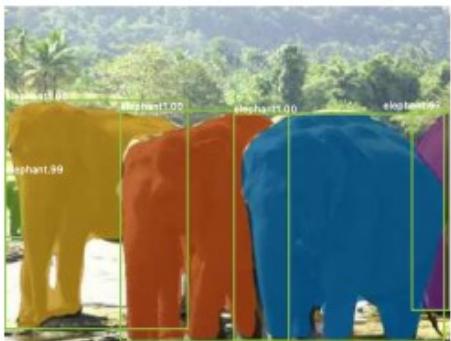
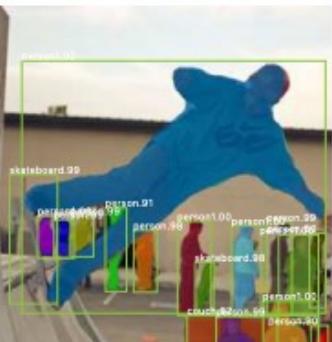
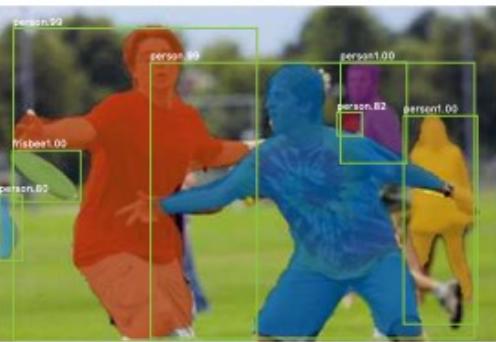
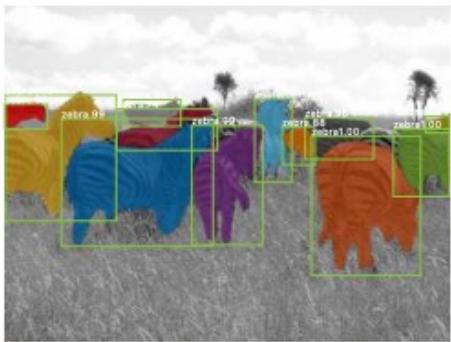
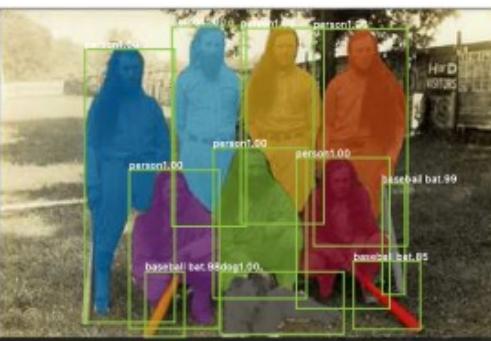
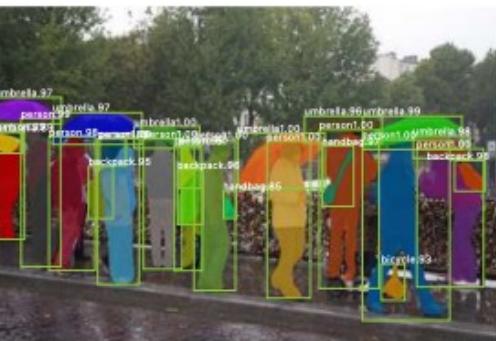
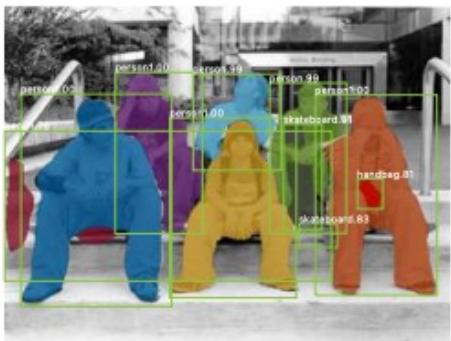
Mask R-CNN: Example Training Targets



Mask R-CNN: Very Good Results! Computer vision seems solved!



Mask R-CNN: by Kaiming He, Georgia Gkioxari, Piotr Dollar, Ross Girshick. ICCV'17, Best Paper Award:
https://openaccess.thecvf.com/content_ICCV_2017/papers/He_Mask_R-CNN_ICCV_2017_paper.pdf



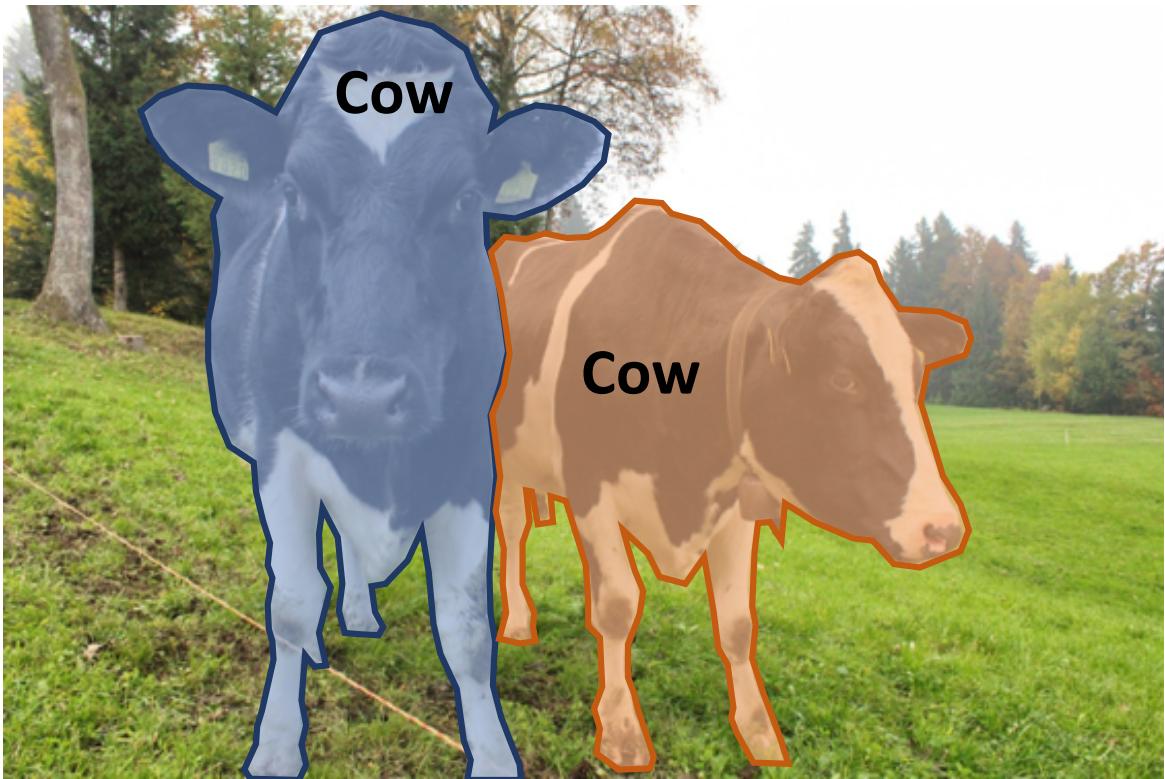
Mask R-CNN

- Object Detection
- Segmentation

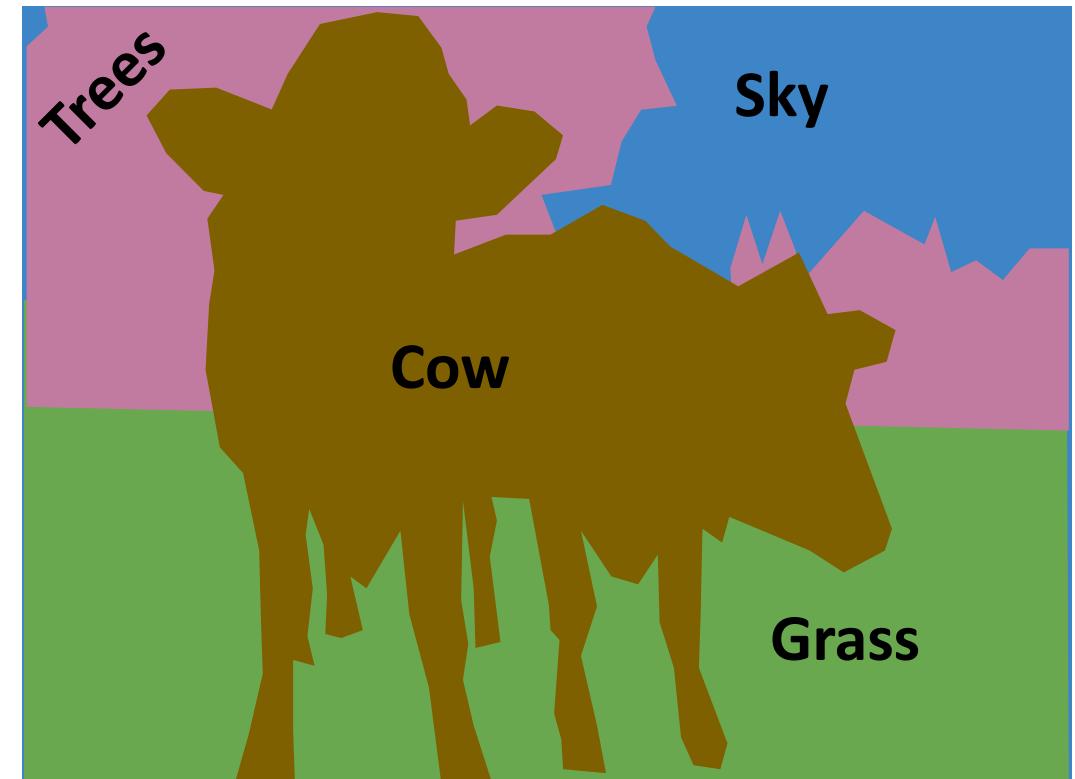


Beyond Instance Segmentation

Instance Segmentation: Separate object instances, but only things



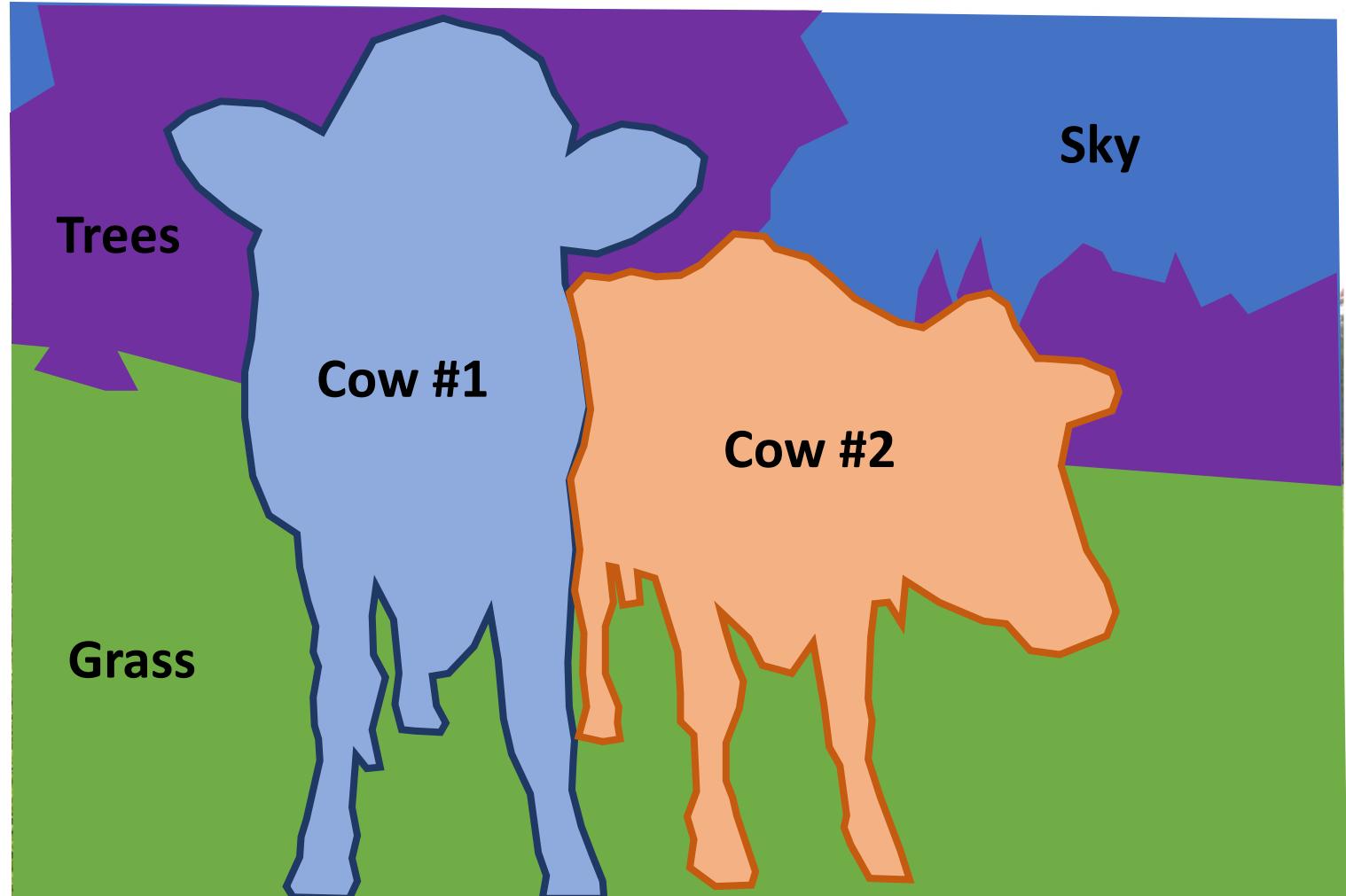
Semantic Segmentation: Identify both things and stuff, but doesn't separate instances



Beyond Instance Segmentation: Panoptic Segmentation

Label all pixels in the image (both things and stuff)

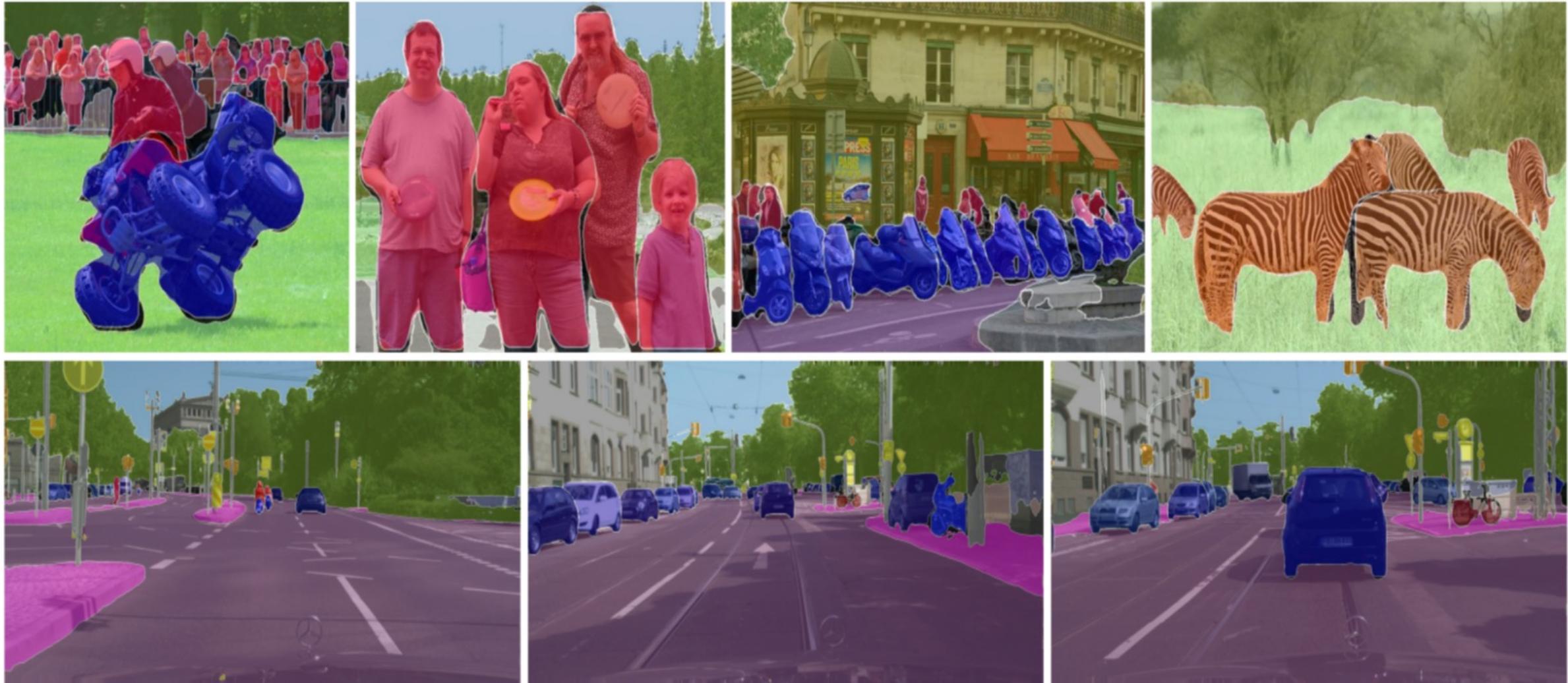
For “thing” categories also separate into instances



Kirillov et al, “Panoptic Segmentation”, CVPR 2019

Kirillov et al, “Panoptic Feature Pyramid Networks”, CVPR 2019

Beyond Instance Segmentation: Panoptic Segmentation



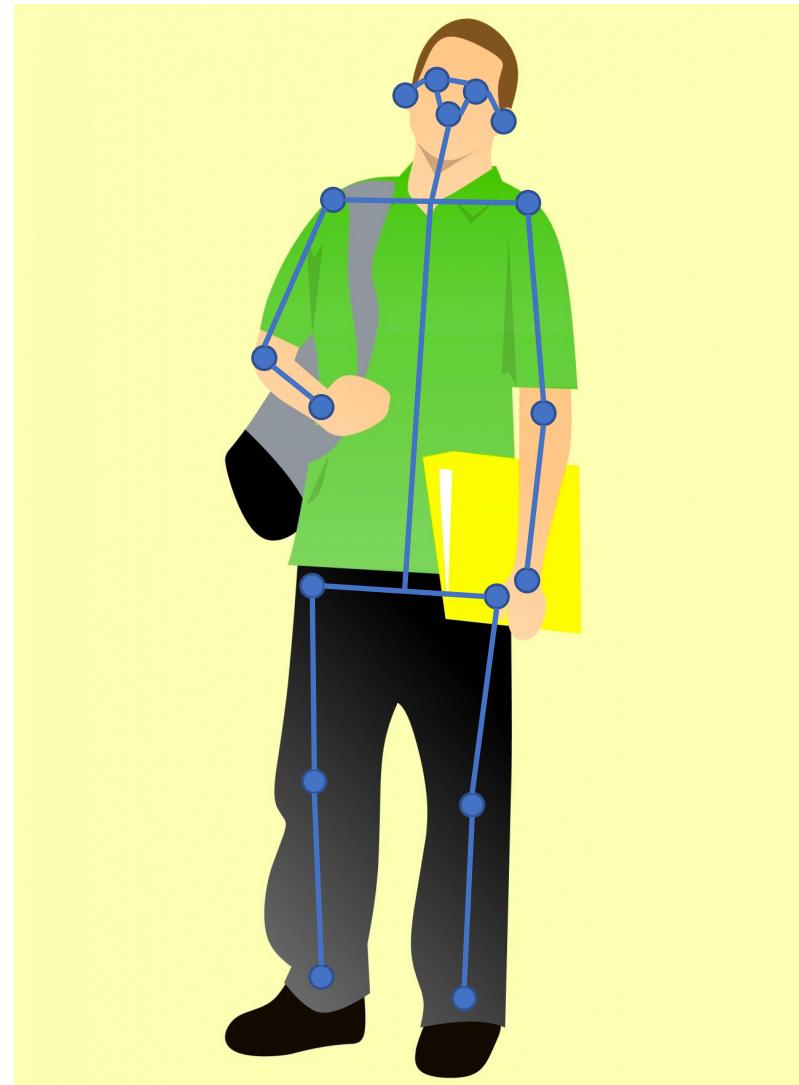
Kirillov et al, "Panoptic Feature Pyramid Networks", CVPR 2019

Beyond Instance Segmentation: Human Keypoints

Represent the pose of a human
by locating a set of **keypoints**

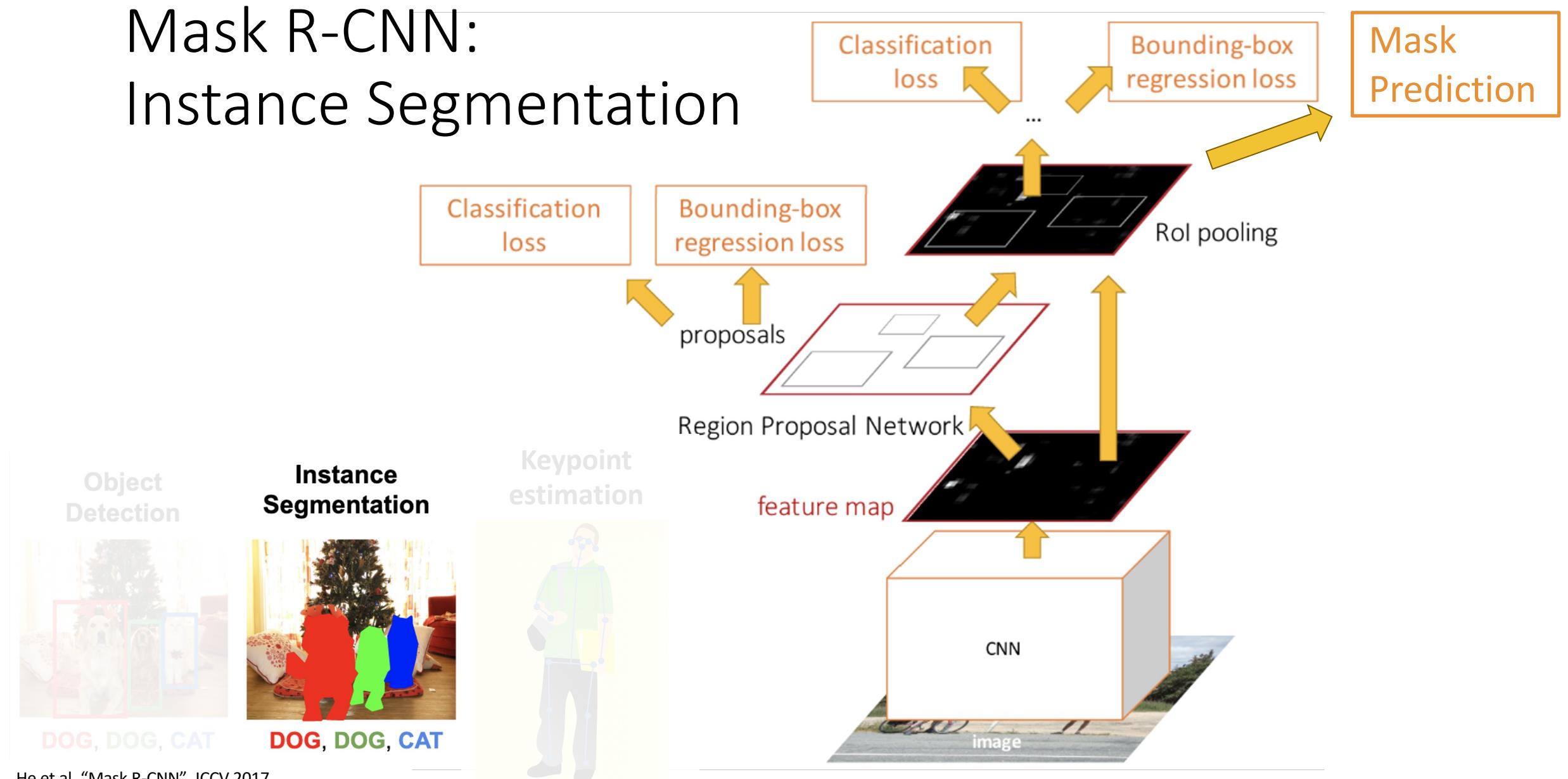
e.g. 17 keypoints:

- Nose
- Left / Right eye
- Left / Right ear
- Left / Right shoulder
- Left / Right elbow
- Left / Right wrist
- Left / Right hip
- Left / Right knee
- Left / Right ankle

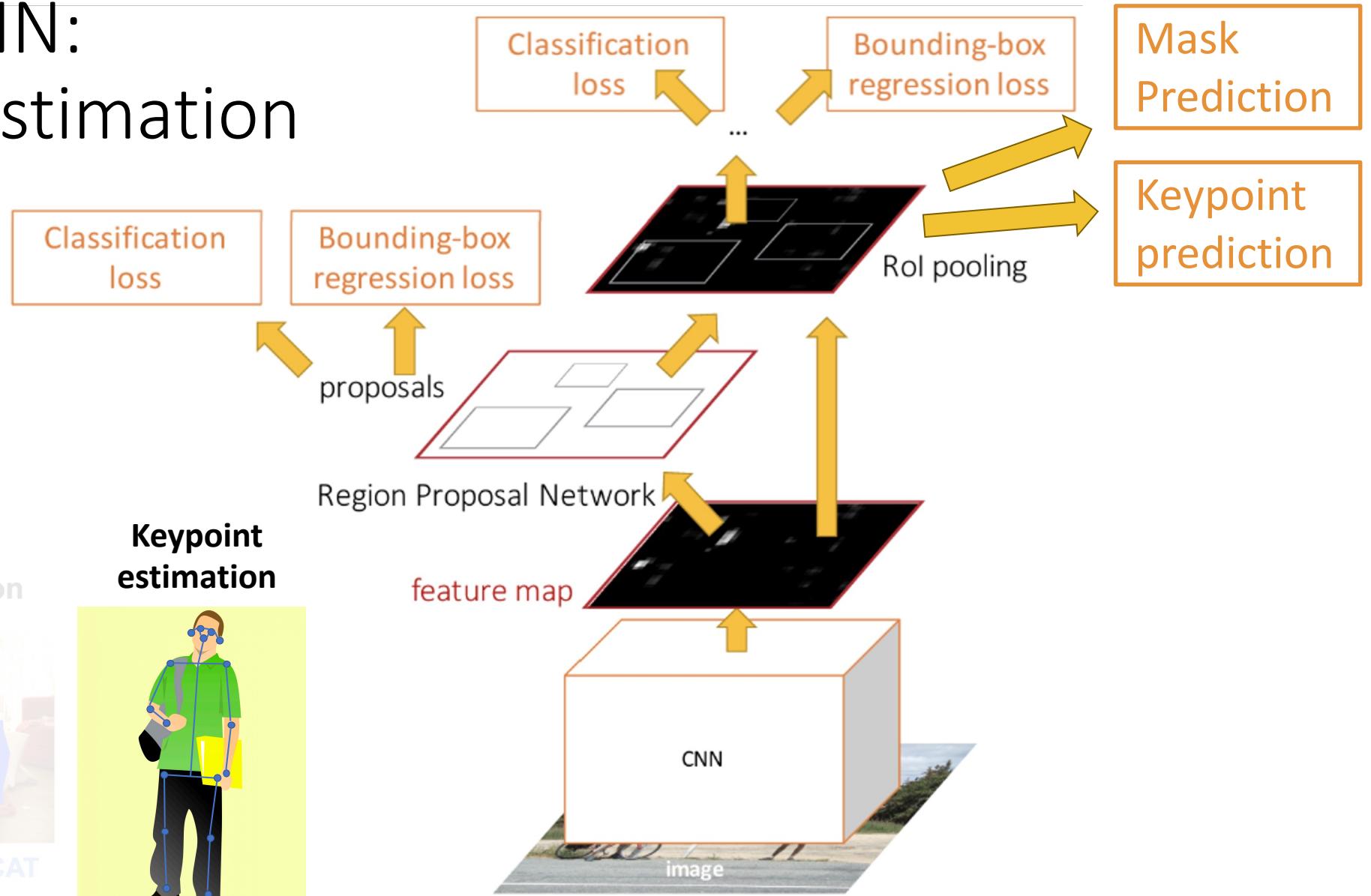
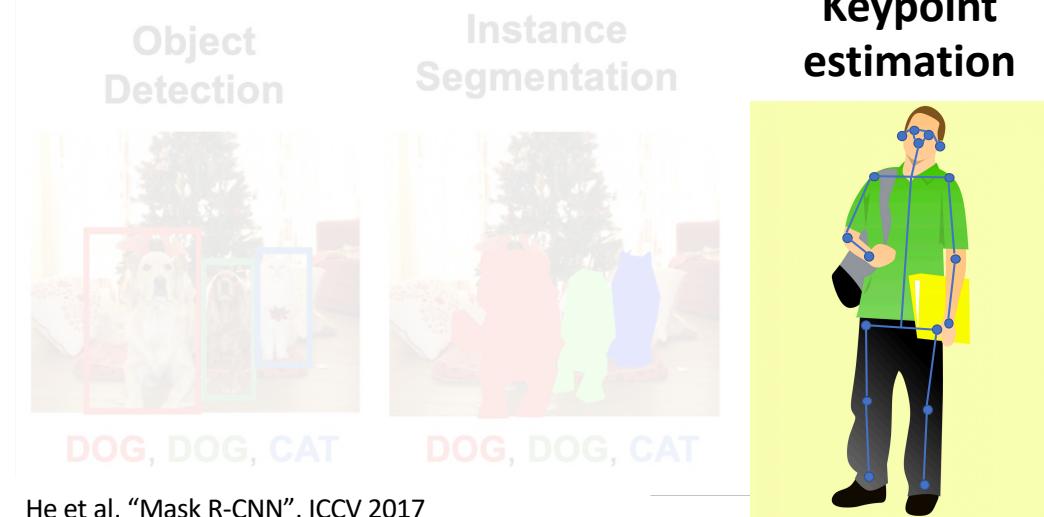


[Person image](#) is CC0 public domain

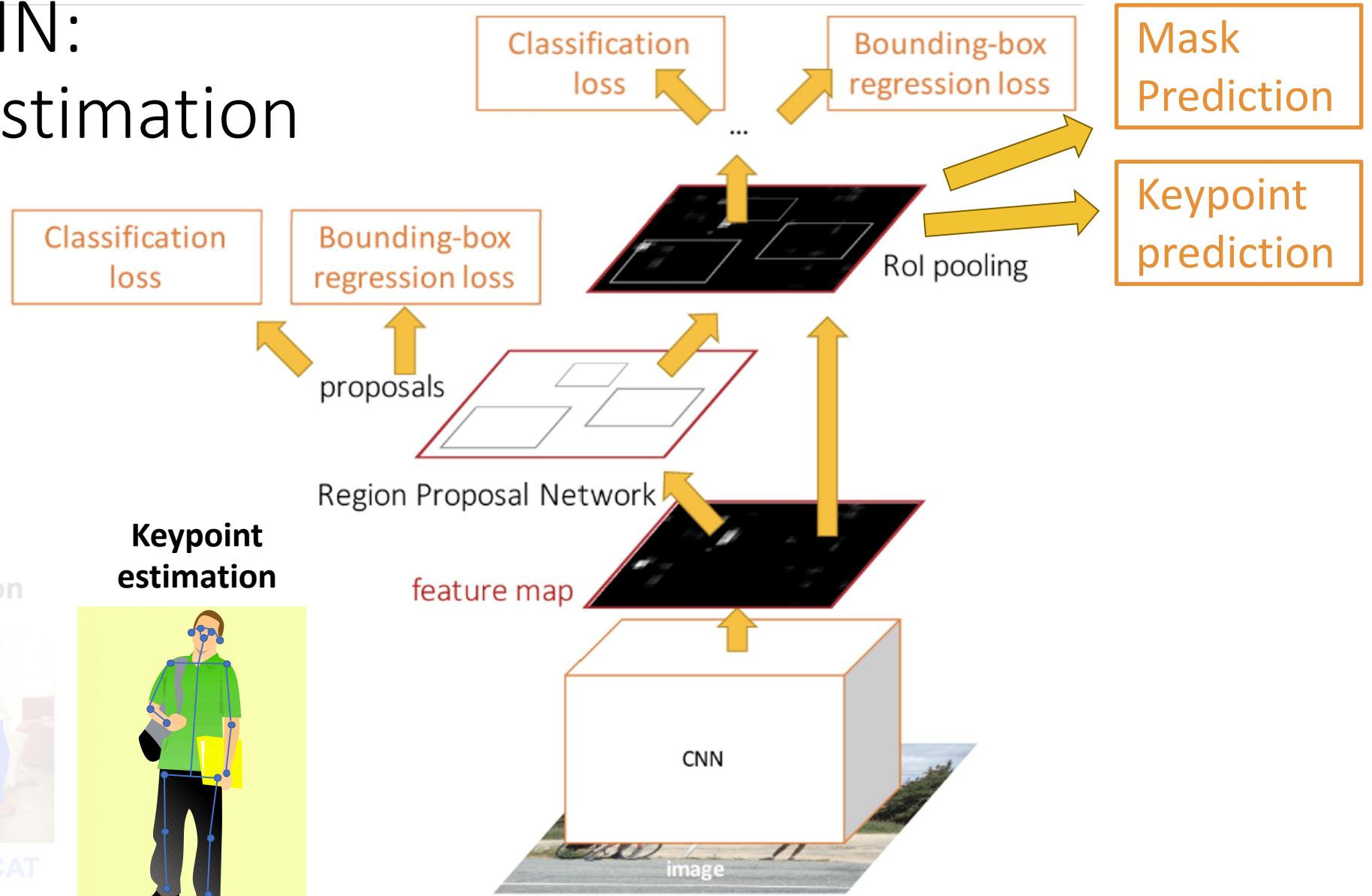
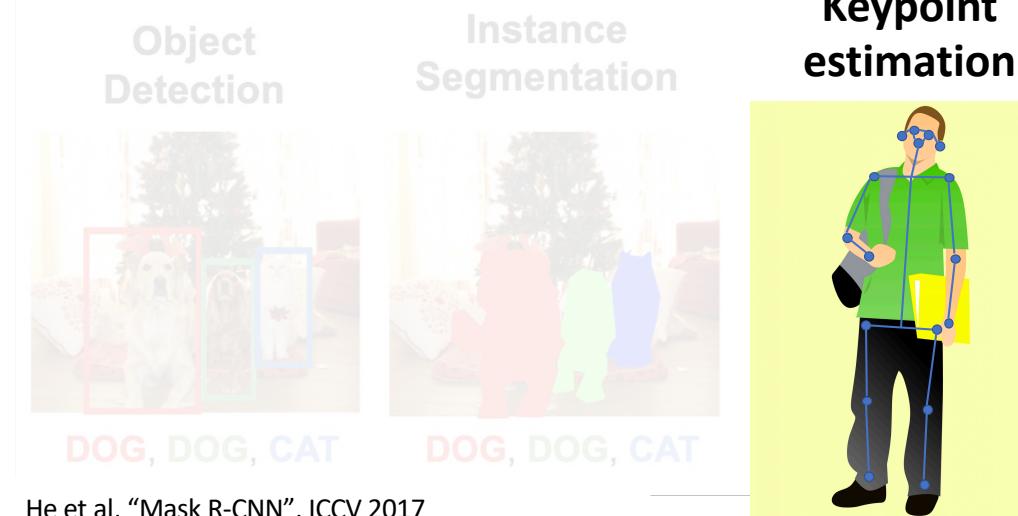
Mask R-CNN: Instance Segmentation



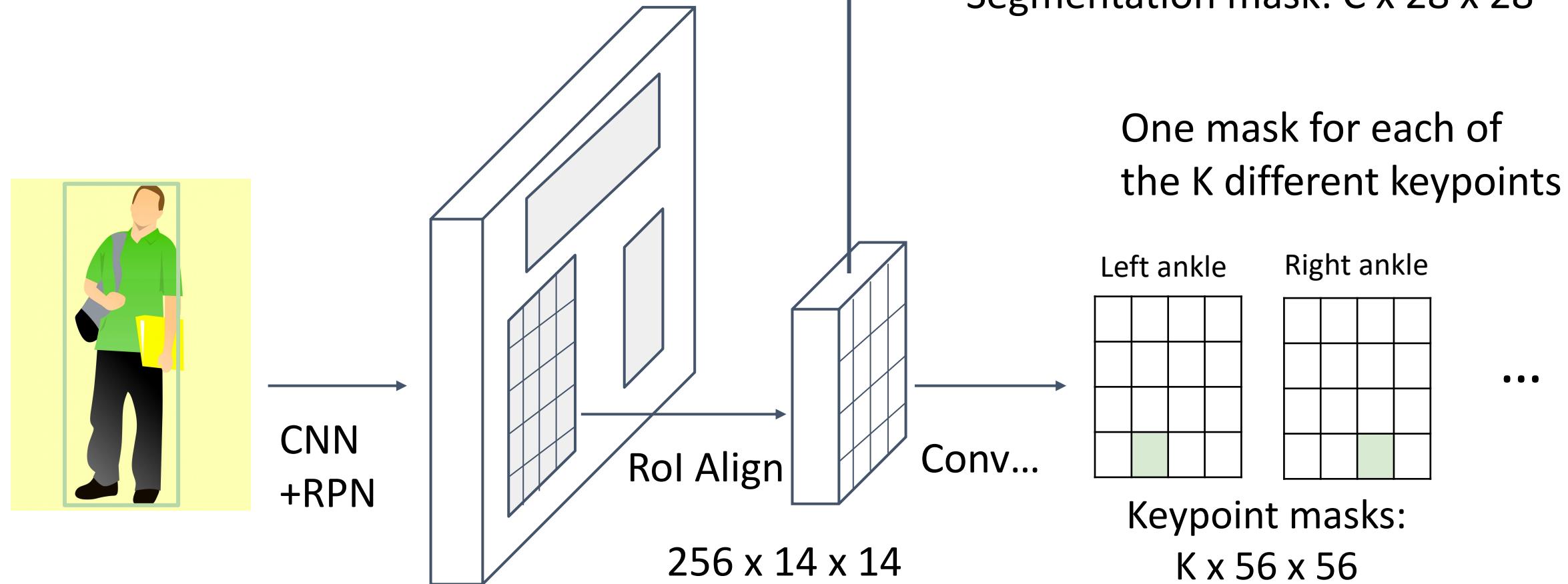
Mask R-CNN: Keypoint Estimation



Mask R-CNN: Keypoint Estimation

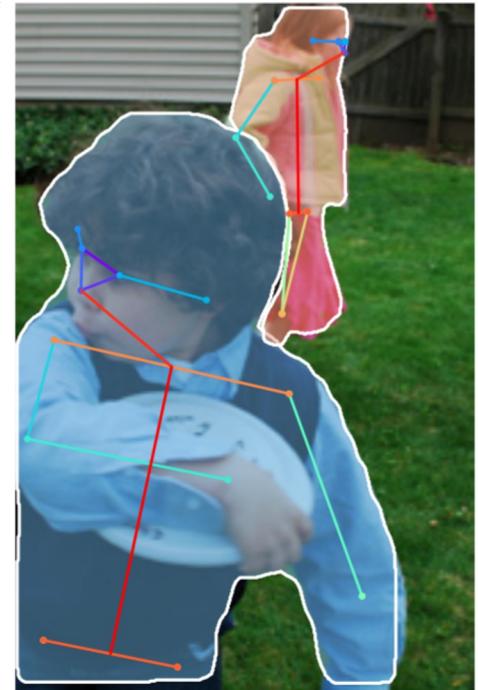


Mask R-CNN: Keypoints



Ground-truth has one “pixel” turned on per keypoint. Train with softmax loss

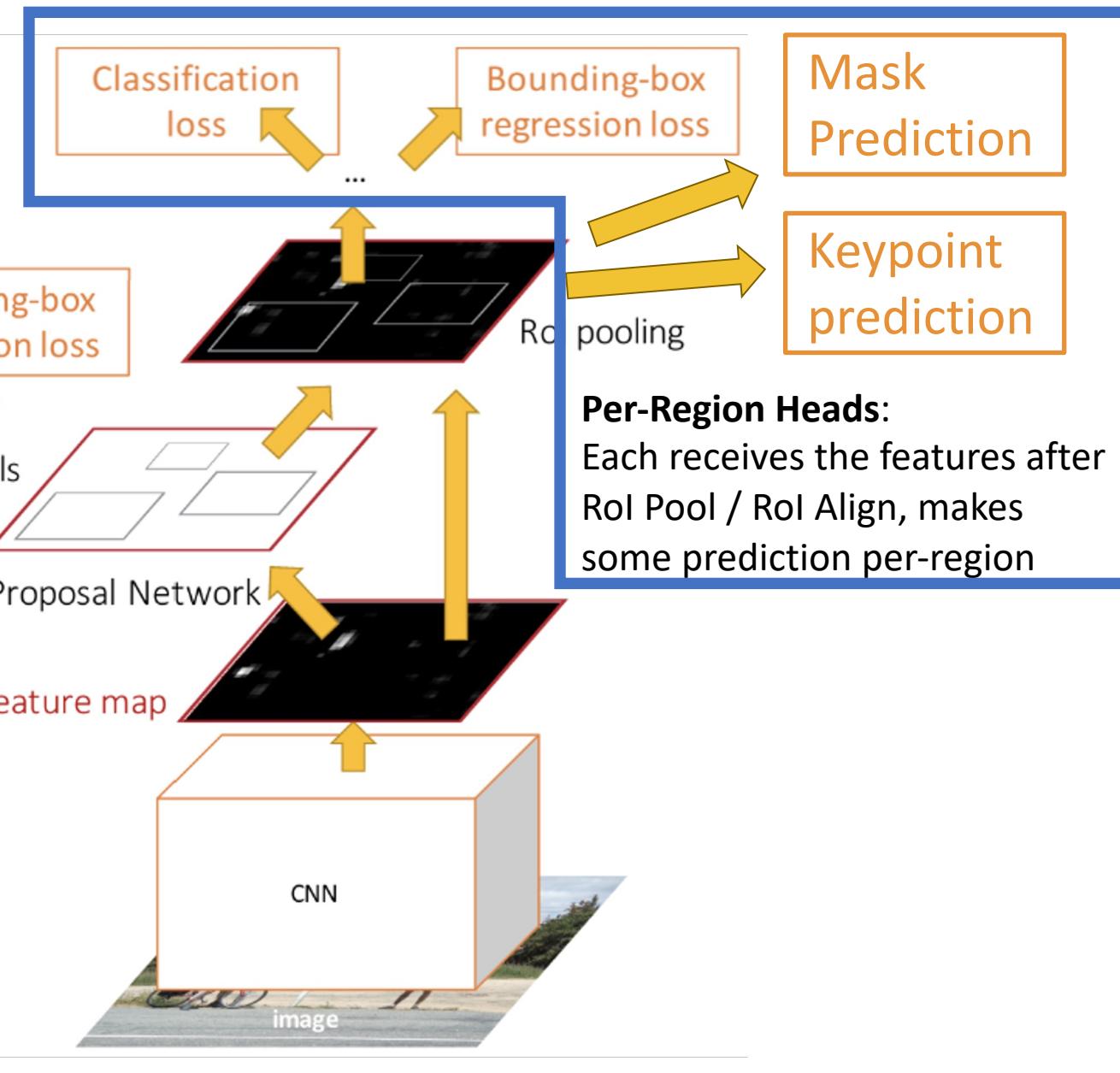
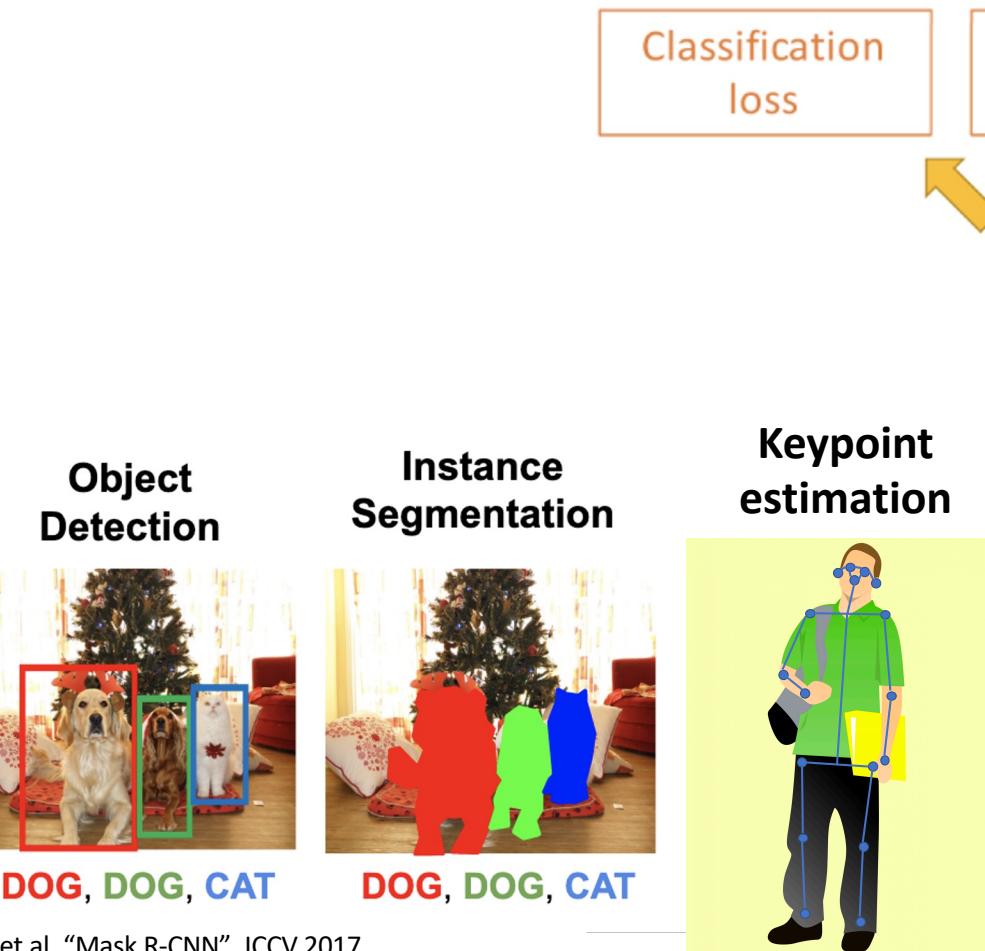
Joint Instance Segmentation and Pose Estimation



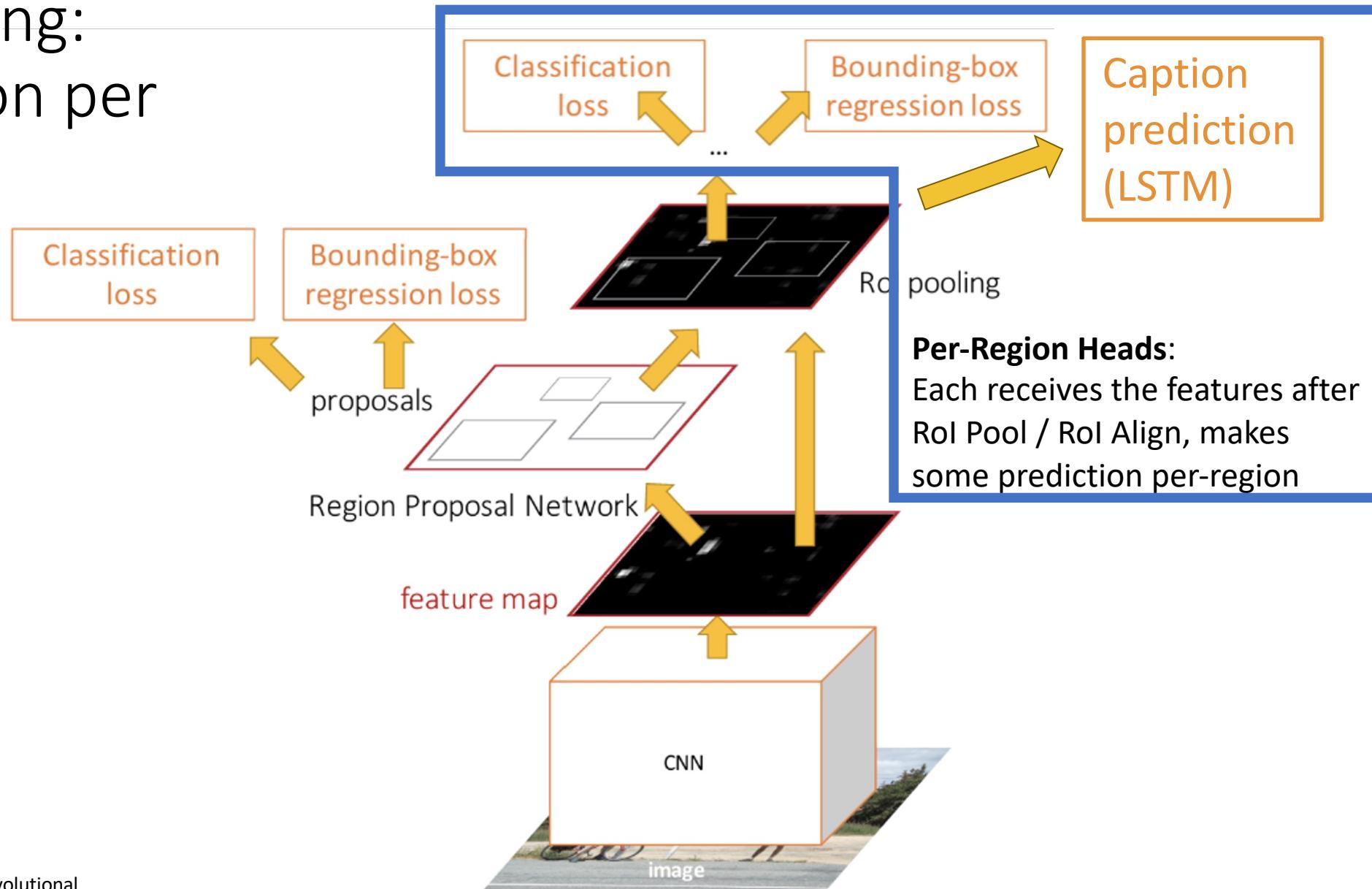
Mask R-CNN for pose estimation and instance segmentation demo video



General Idea: Add Per-Region “Heads” to Mask R-CNN!

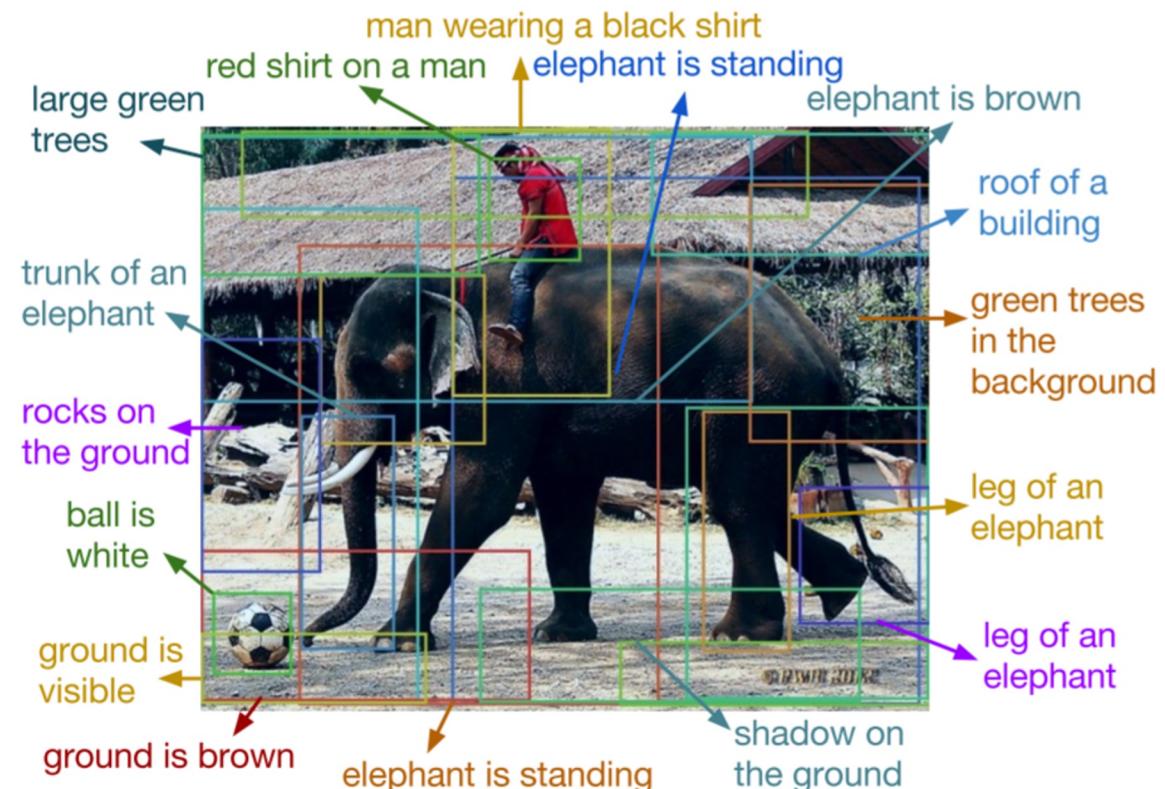
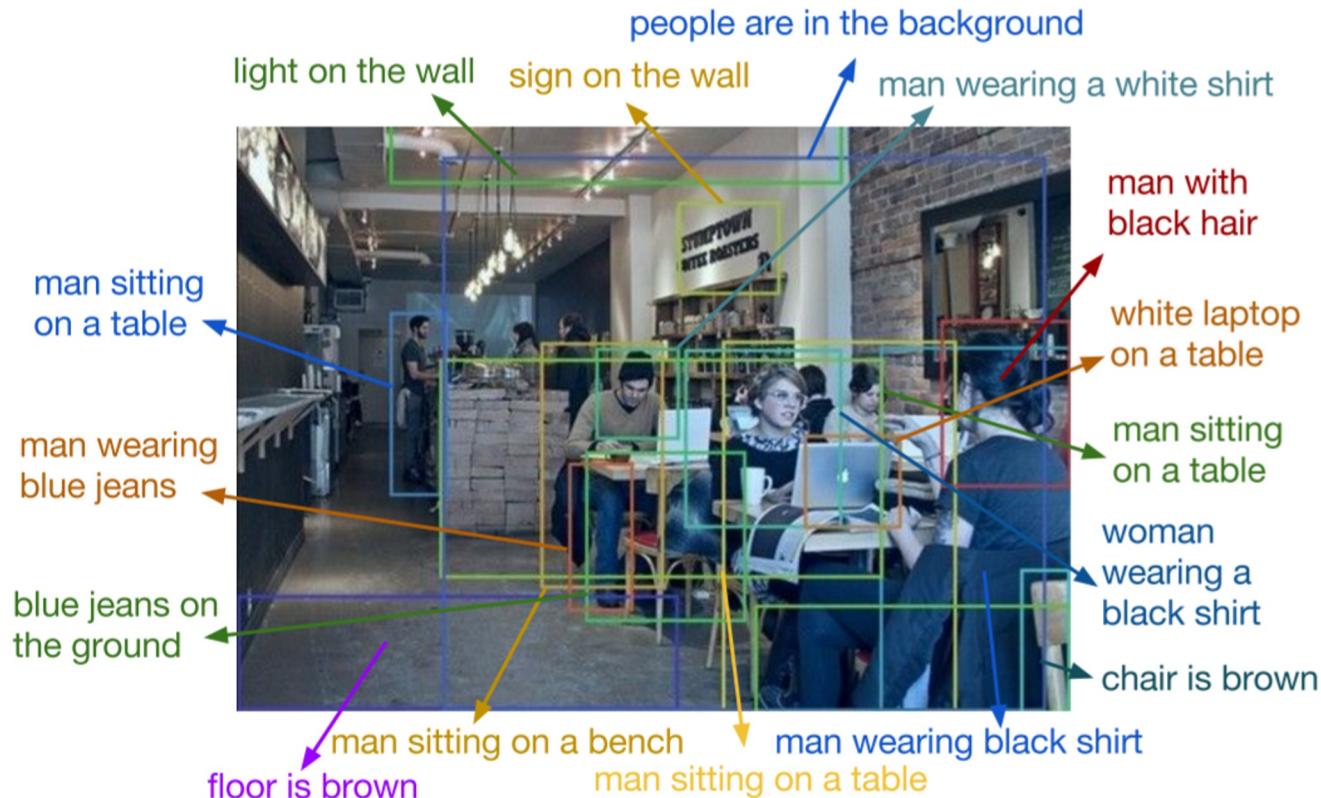


Dense Captioning: Predict a caption per region

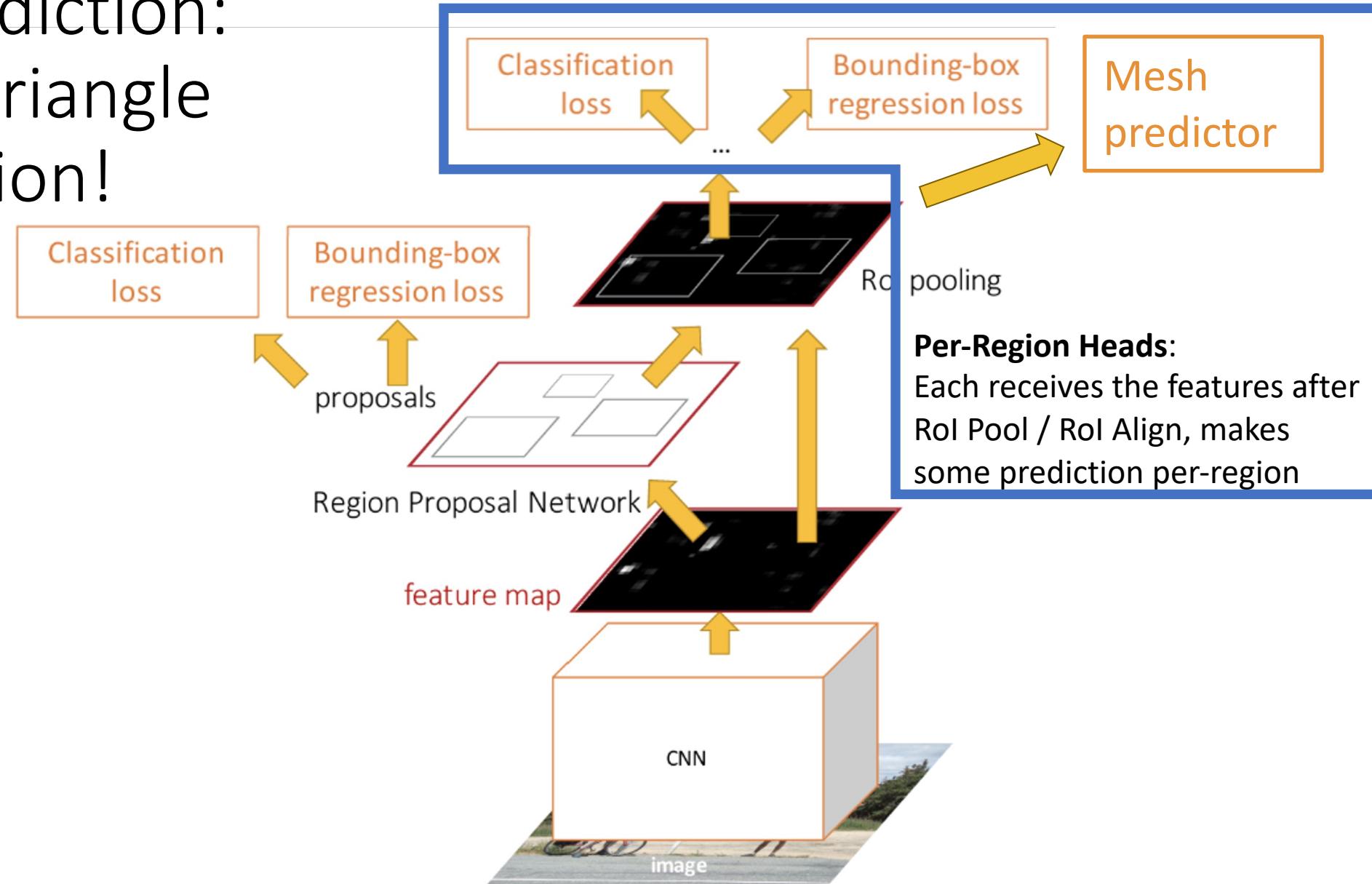


Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

Dense Captioning



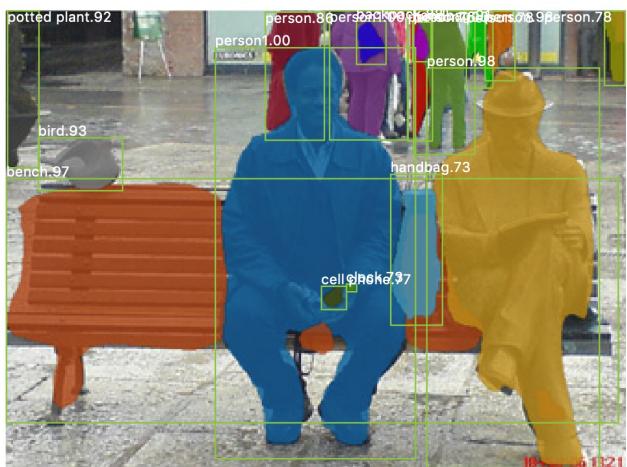
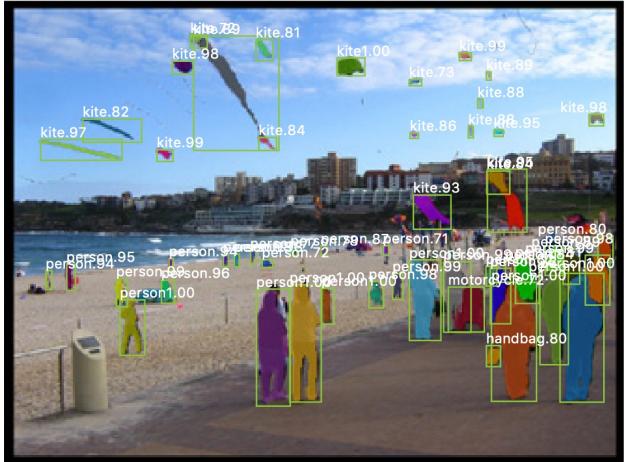
3D Shape Prediction: Predict a 3D triangle mesh per region!



3D Shape Prediction: Mask R-CNN + Mesh Head

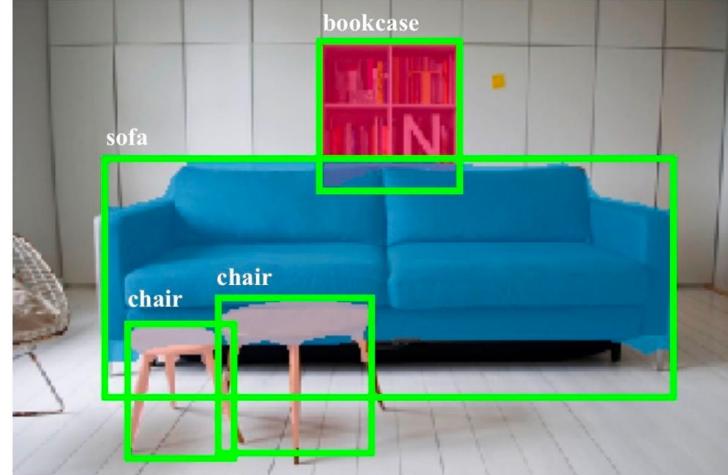
Mask R-CNN:

2D Image -> 2D shapes



Mesh R-CNN:

2D Image -> 3D shapes



More details
next time!

He, Gkioxari, Dollár, and
Girshick, "Mask R-CNN",
ICCV 2017

Gkioxari, Malik, and Johnson,
"Mesh R-CNN", ICCV 2019

Toward a foundation model for image segmentation

Segment Anything (SAM)

100-200 masks



200-300 masks



Summary: Many Computer Vision Tasks!

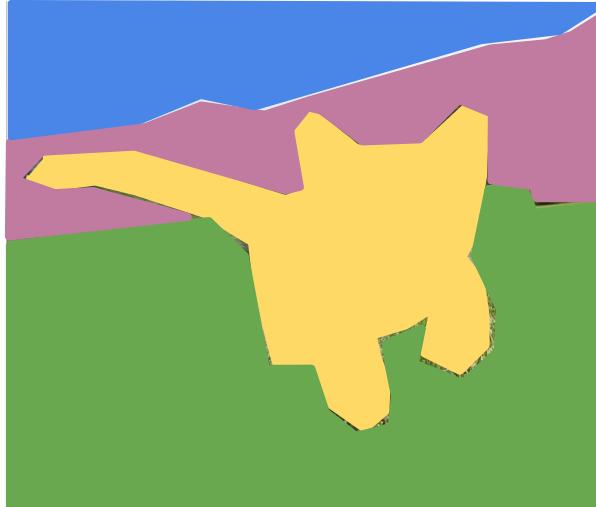
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Objects

Instance Segmentation



DOG, DOG, CAT

[This image is CC0 public domain](#)

Reading and Practice

Textbook D2L:

Chapter 14.8.4 Mask R-CNN

Chapter 14.9 Semantic Segmentation and the Dataset

Chapter 14.10 Transposed Convolution

Chapter 14.11 Fully Convolutional Networks