

Lecture 11: Segmentation and Pose Estimation

Thursday March 30, 2017



Announcements!

- HW #3 misunderstandings
 - Deadline for HW #3 re-try is next Thursday **April 6**
- Final Project milestones due next Tuesday **April 4**
- Vote for Final Day and Location

Python/Numpy of the Day

Decorators

How to write a decorator:

```
import time

def timeit(method):

    def timed(*args, **kw):
        ts = time.time()
        result = method(*args, **kw)
        te = time.time()

        print '%r (%r, %r) %2.2f sec' % \
              (method.__name__, args, kw, te-ts)
        return result

    return timed
```

How to use a decorator:

```
class Foo(object):

    @timeit
    def foo(self, a=2, b=3):
        time.sleep(0.2)

    @timeit
    def f1():
        time.sleep(1)
        print 'f1'
```

Grouping in vision

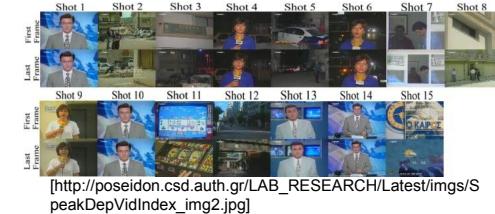
- Goals:
 - Gather features that belong together
 - Obtain an intermediate representation that compactly describes key image (video) parts
- Top down vs. bottom up segmentation
 - Top down: pixels belong together because they are from the same object
 - Bottom up: pixels belong together because they look similar
- Hard to measure success
 - What is interesting depends on the app.

Examples of grouping in vision



[Figure by J. Shi]

Determine image regions

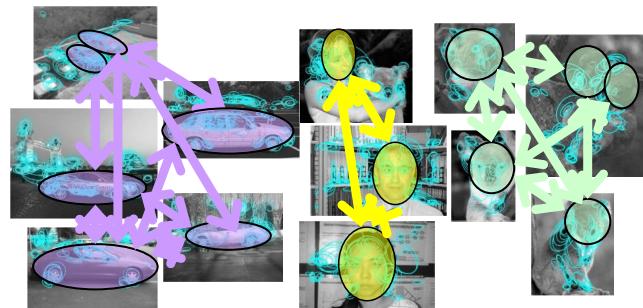


Group video frames into shots



[Figure by Wang & Suter]

Figure-ground

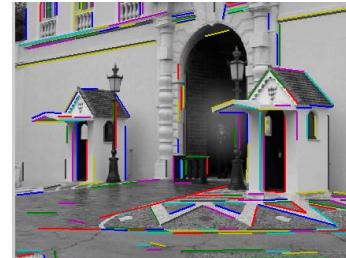


[Figure by Grauman & Darrell]

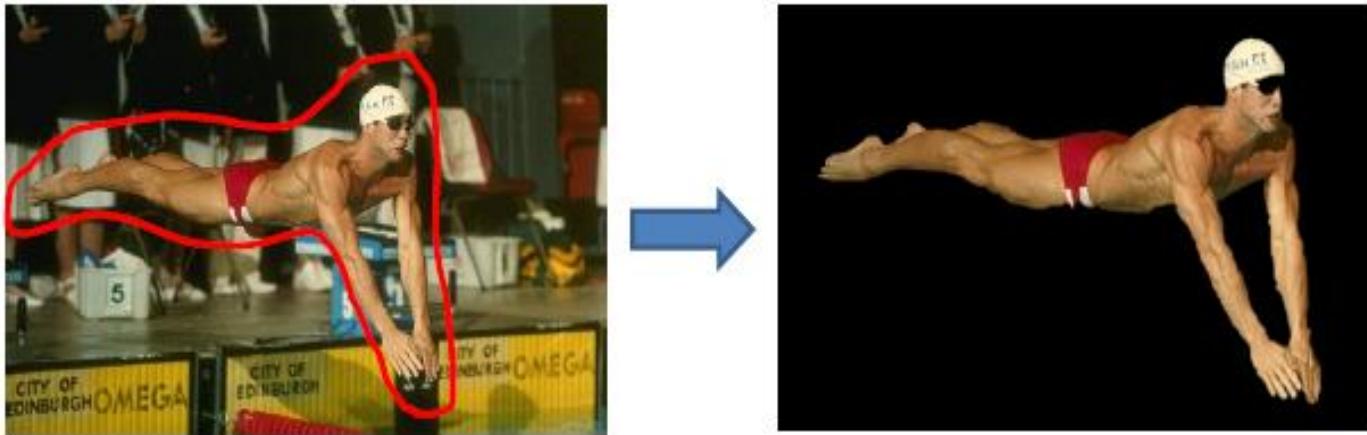
Object-level grouping

Edge and line detection

- Canny edge detector = smooth → derivative → thin → threshold → link
- Generalized Hough transform = points vote for shape parameters
- Straight line detector = canny + gradient orientations → orientation binning → linking → check for straightness



Segmentation



Computer Vision Tasks

Classification



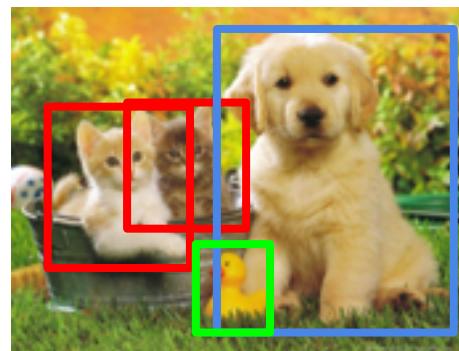
CAT

Classification + Localization



CAT

Object Detection



CAT, DOG, DUCK

Segmentation



CAT, DOG, DUCK

Single object

Multiple objects

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Computer Vision Tasks

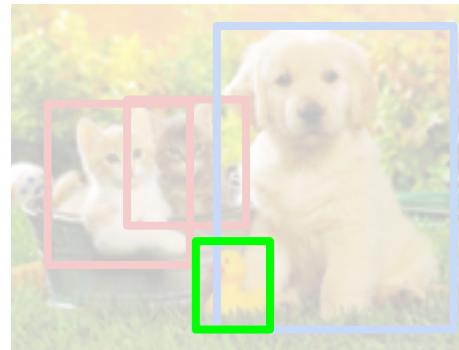
Classification



Classification
+ Localization



Object Detection

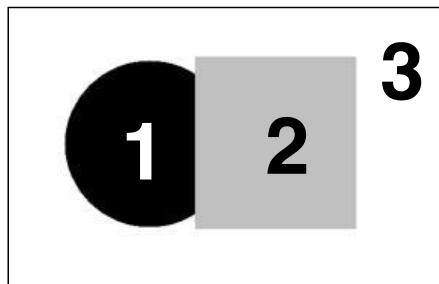


Segmentation

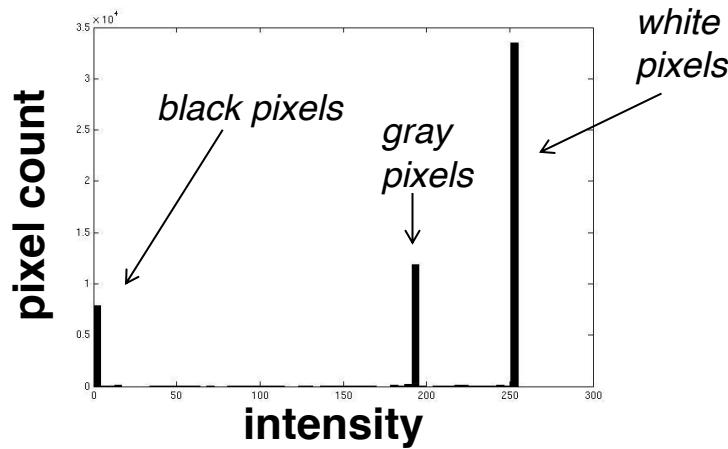


Today

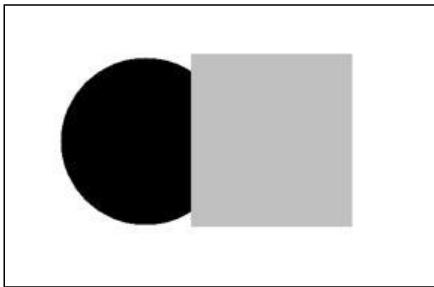
Image segmentation: toy example



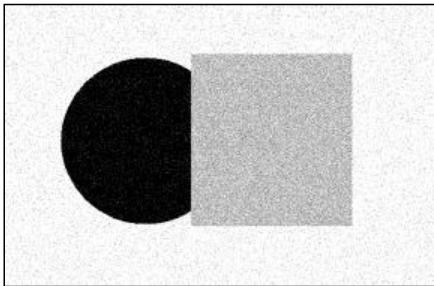
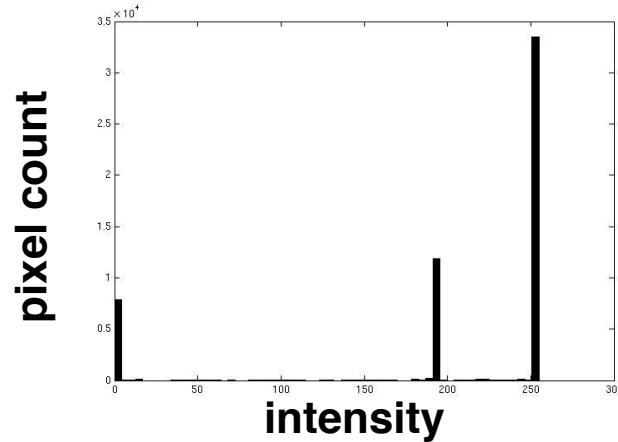
input image



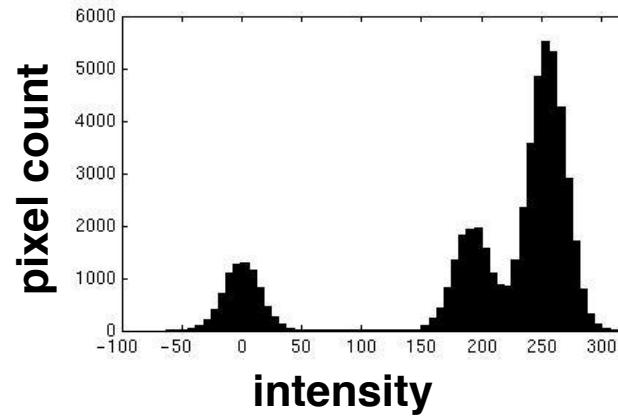
- These intensities define the three groups.
- We could label every pixel in the image according to which of these primary intensities it is.
 - i.e., *segment* the image based on the intensity feature.
- What if the image isn't quite so simple?

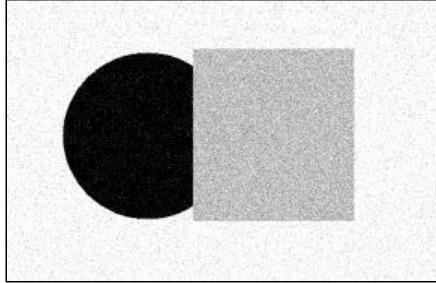


input image

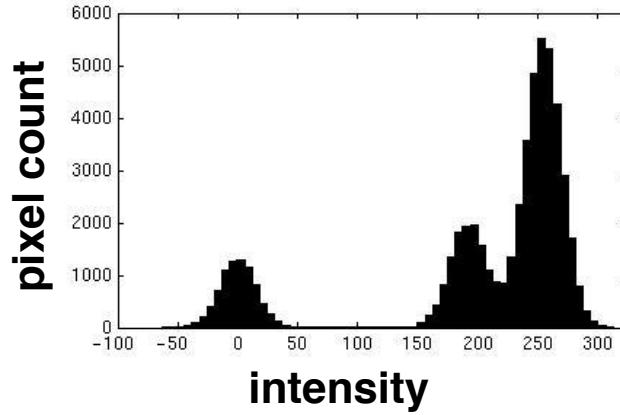


input image





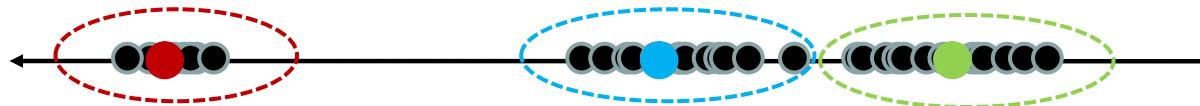
input image



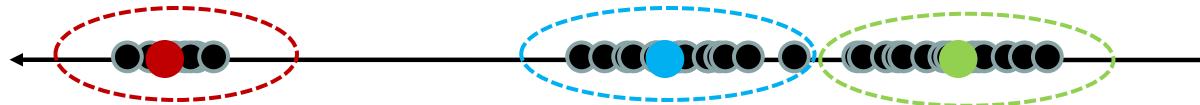
- Now how to determine the three main intensities that define our groups?
- We need to ***cluster***.

Clustering

- With this objective, it is a “chicken and egg” problem:
 - If we knew the **cluster centers**, we could allocate points to groups by assigning each to its closest center.

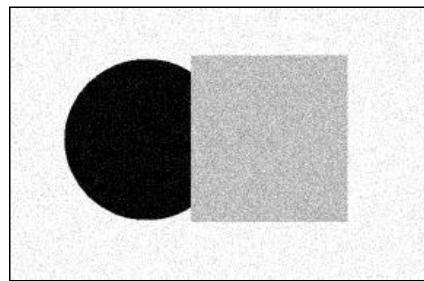


- If we knew the **group memberships**, we could get the centers by computing the mean per group.

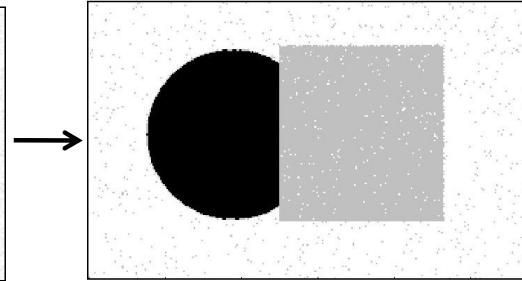


Smoothing out cluster assignments

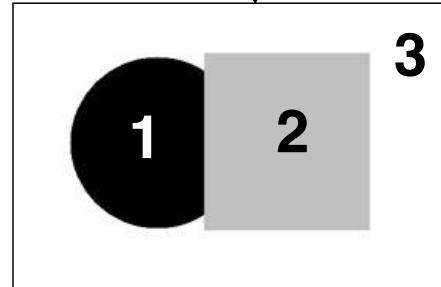
- Assigning a cluster label per pixel may yield outliers:



original

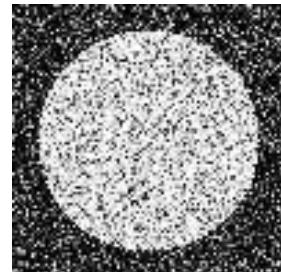


labeled by cluster center's
intensity



- How to ensure they are spatially smooth?

Solution



$P(\text{foreground} \mid \text{image})$

Encode dependencies between pixels

Normalizing constant

$$P(\mathbf{y}; \theta, \text{data}) = \frac{1}{Z} \prod_{i=1..N} f_1(y_i; \theta, \text{data}) \prod_{i, j \in \text{edges}} f_2(y_i, y_j; \theta, \text{data})$$

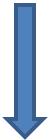
Labels to be predicted

Individual predictions

Pairwise predictions

Writing Likelihood as an “Energy”

$$P(\mathbf{y}; \theta, data) = \frac{1}{Z} \prod_{i=1..N} p_1(y_i; \theta, data) \prod_{i,j \in edges} p_2(y_i, y_j; \theta, data)$$

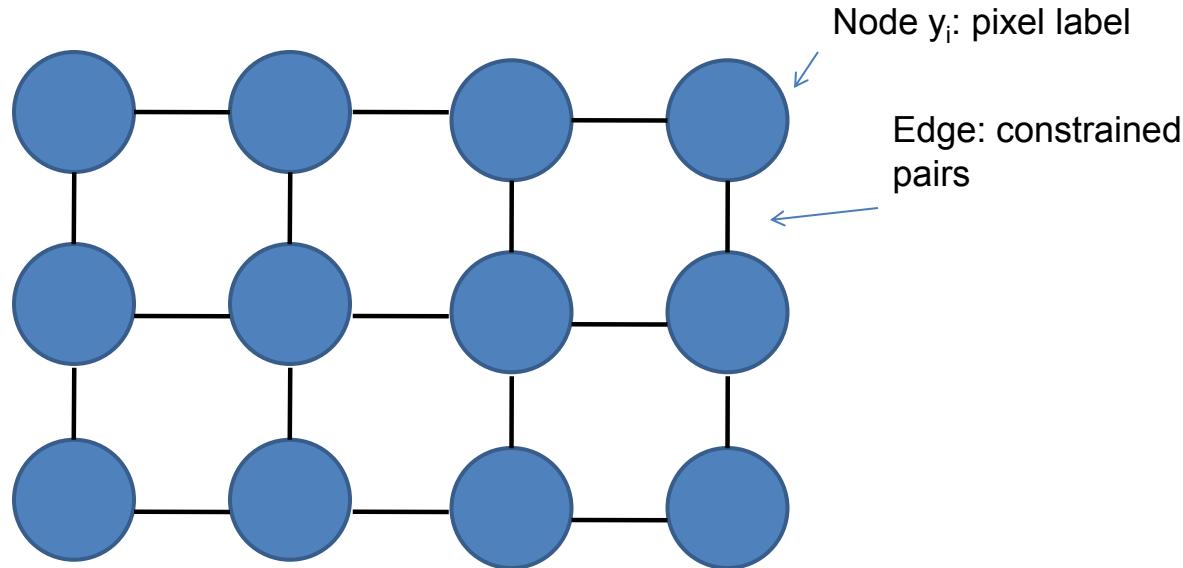


$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i,j \in edges} \psi_2(y_i, y_j; \theta, data)$$

“Cost” of assignment y_i

“Cost” of pairwise assignment y_i, y_j

Markov Random Fields



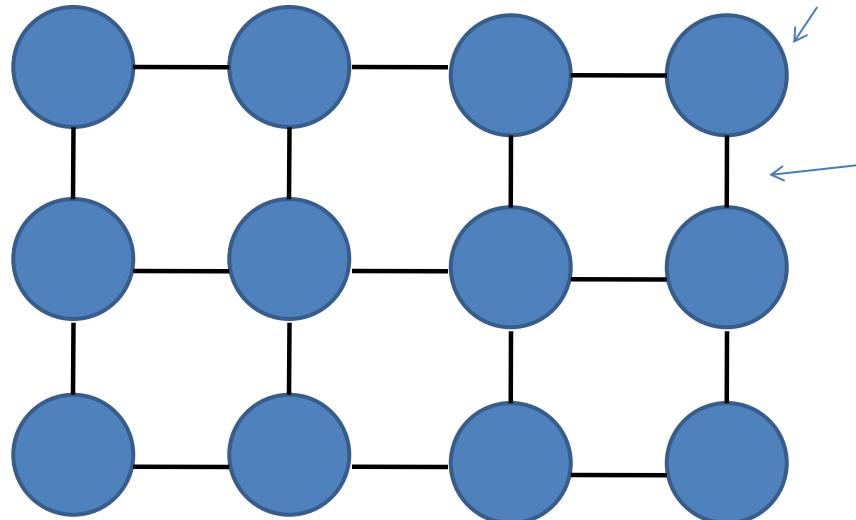
Cost to assign a label to
each pixel

Cost to assign a pair of labels to
connected pixels

$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i, j \in edges} \psi_2(y_i, y_j; \theta, data)$$

Markov Random Fields

- Example: “label smoothing” grid



Unary potential

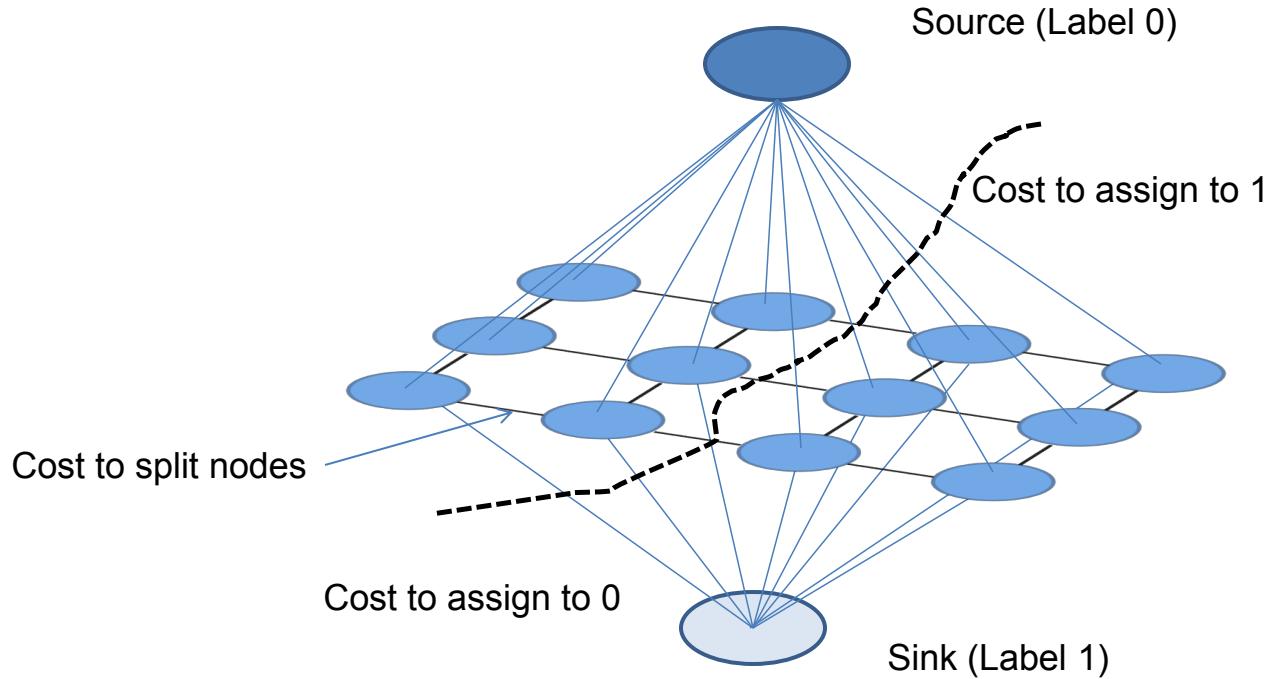
$$\begin{array}{l} 0: -\log P(y_i = 0 ; \text{data}) \\ 1: -\log P(y_i = 1 ; \text{data}) \end{array}$$

Pairwise Potential

0	1
0	K
1	K

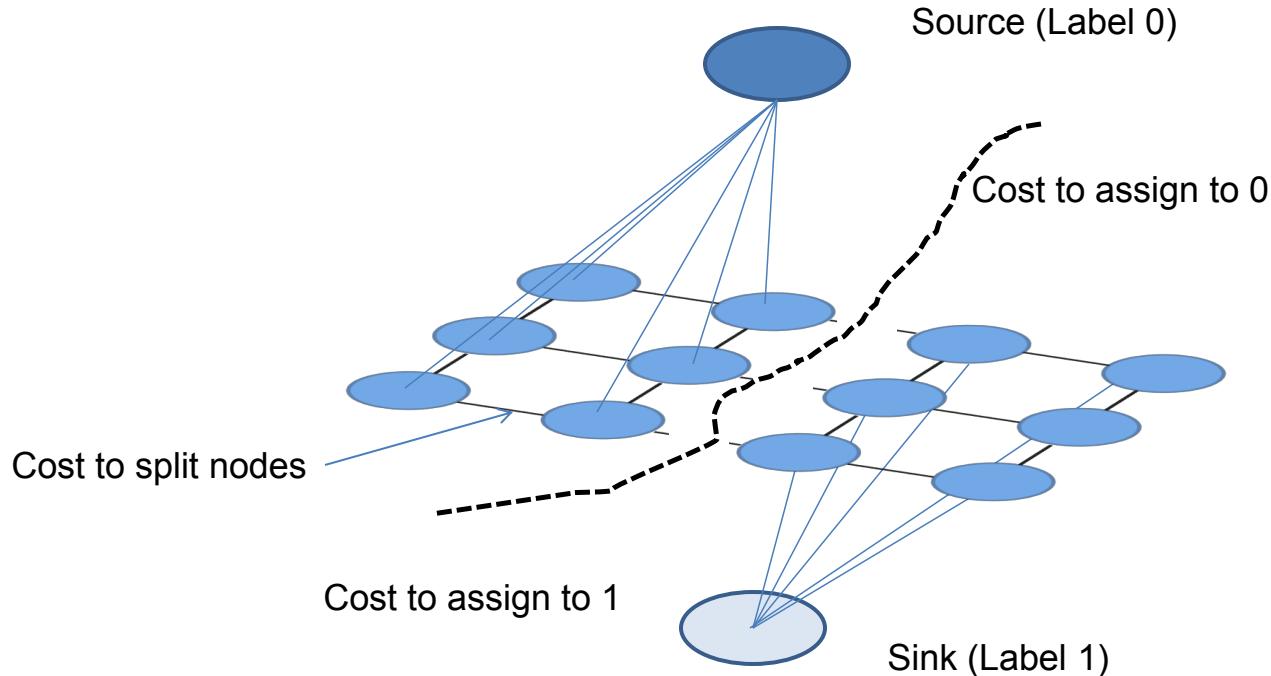
$$Energy(\mathbf{y}; \theta, \text{data}) = \sum_i \psi_1(y_i; \theta, \text{data}) + \sum_{i, j \in \text{edges}} \psi_2(y_i, y_j; \theta, \text{data})$$

Solving MRFs with graph cuts



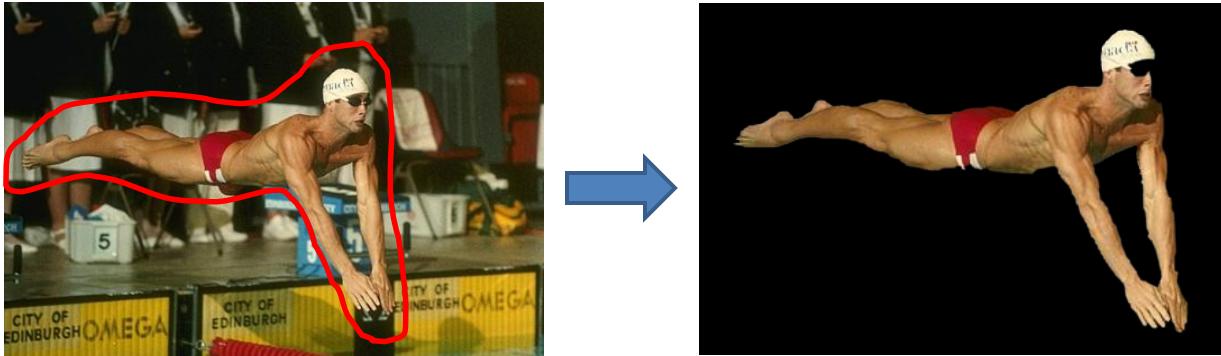
$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i, j \in edges} \psi_2(y_i, y_j; \theta, data)$$

Solving MRFs with graph cuts



$$Energy(\mathbf{y}; \theta, data) = \sum_i \psi_1(y_i; \theta, data) + \sum_{i, j \in edges} \psi_2(y_i, y_j; \theta, data)$$

GrabCut segmentation



User provides rough indication of foreground region.

Goal: Automatically provide a pixel-level segmentation.

What is easy or hard about these cases for graphcut-based segmentation?



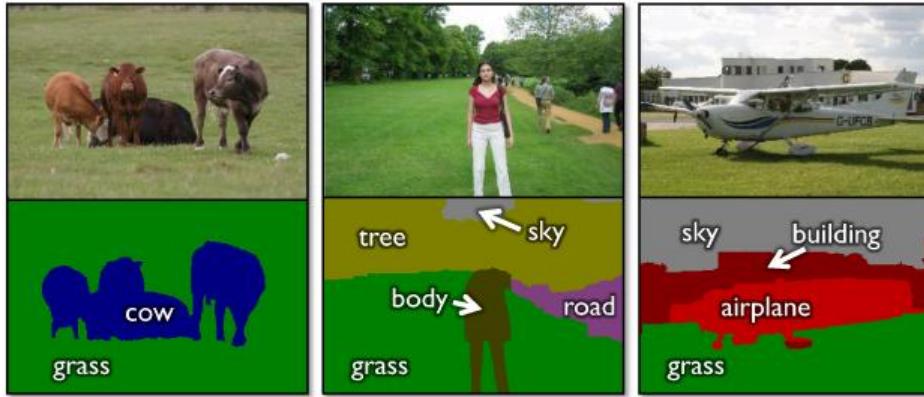
Easier examples



More difficult Examples

	Camouflage & Low Contrast	Fine structure	Harder Case
Initial Rectangle			
Initial Result			

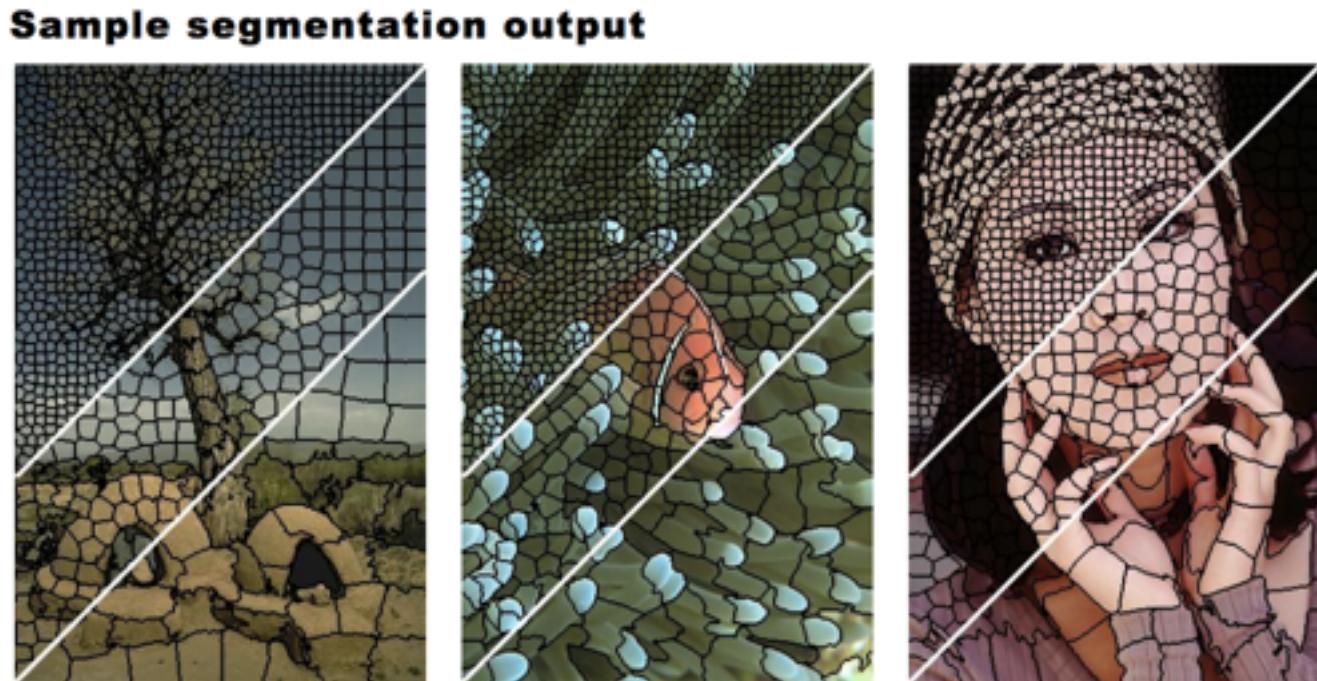
Using graph cuts for recognition



object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Unsupervised Segmentation

- Classic Example:
Superpixels
- Cluster pixels with their neighbors
- Break clusters when large gradient occurs between neighboring pixels
- Figure from Achanta et al.
— “SLIC Superpixels Compared to State-of-the-art Superpixel Methods,”
May 2012



Further reading and resources

- Graph cuts
 - <http://www.cs.cornell.edu/~rdz/graphcuts.html>
 - Classic paper: [What Energy Functions can be Minimized via Graph Cuts?](#)
(Kolmogorov and Zabih, ECCV '02/PAMI '04)
- Belief propagation

Yedidia, J.S.; Freeman, W.T.; Weiss, Y., "Understanding Belief Propagation and Its Generalizations", Technical Report, 2001:
<http://www.merl.com/publications/TR2001-022/>
- Normalized cuts and image segmentation (Shi and Malik)
<http://www.cs.berkeley.edu/~malik/papers/SM-ncut.pdf>
- N-cut implementation
<http://www.seas.upenn.edu/~timothée/software/ncut/ncut.html>

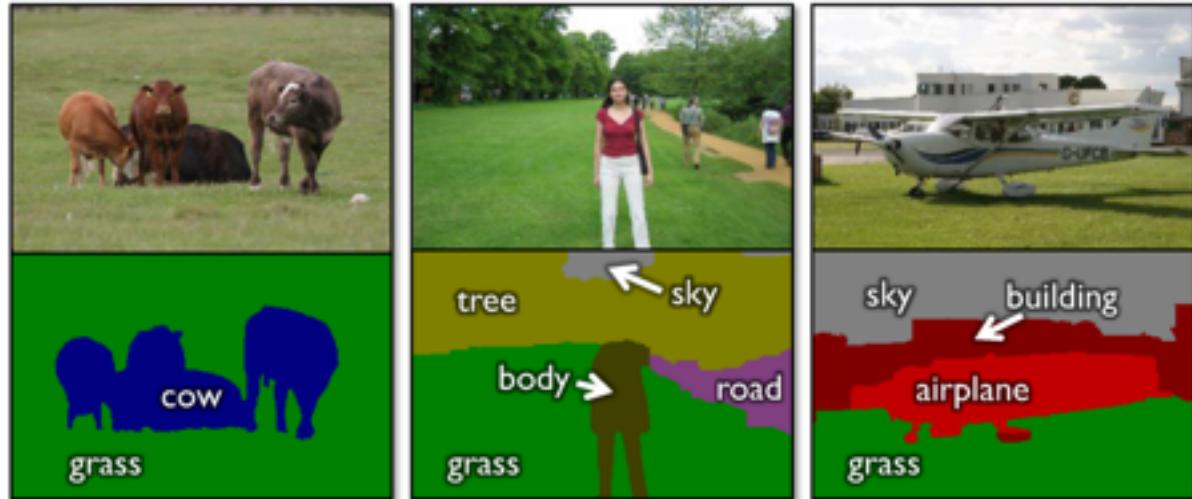




Semantic Segmentation

Label every pixel!

Don't differentiate instances (cows)



object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Figure credit: Shotton et al, "TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context", IJCV 2007

Instance Segmentation

Detect instances,
give category, label
pixels

“simultaneous
detection and
segmentation” (SDS)

Lots of recent work
(MS-COCO
Challenges)

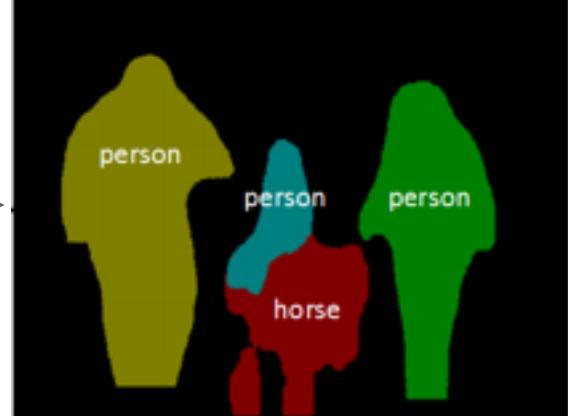


Figure credit: Dai et al, “Instance-aware Semantic Segmentation via Multi-task Network Cascades”, arXiv 2015

Semantic Segmentation

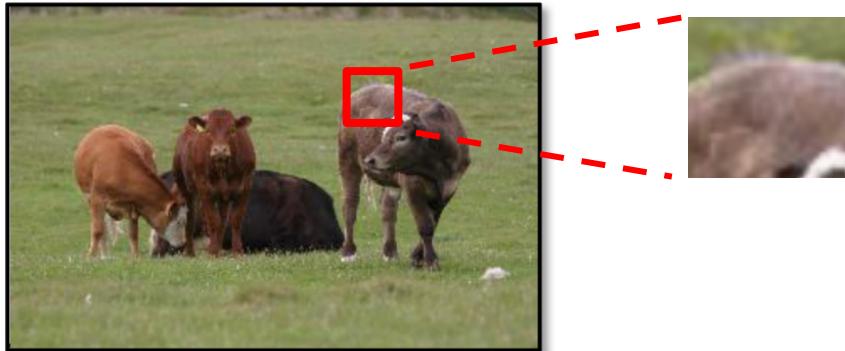
Semantic Segmentation



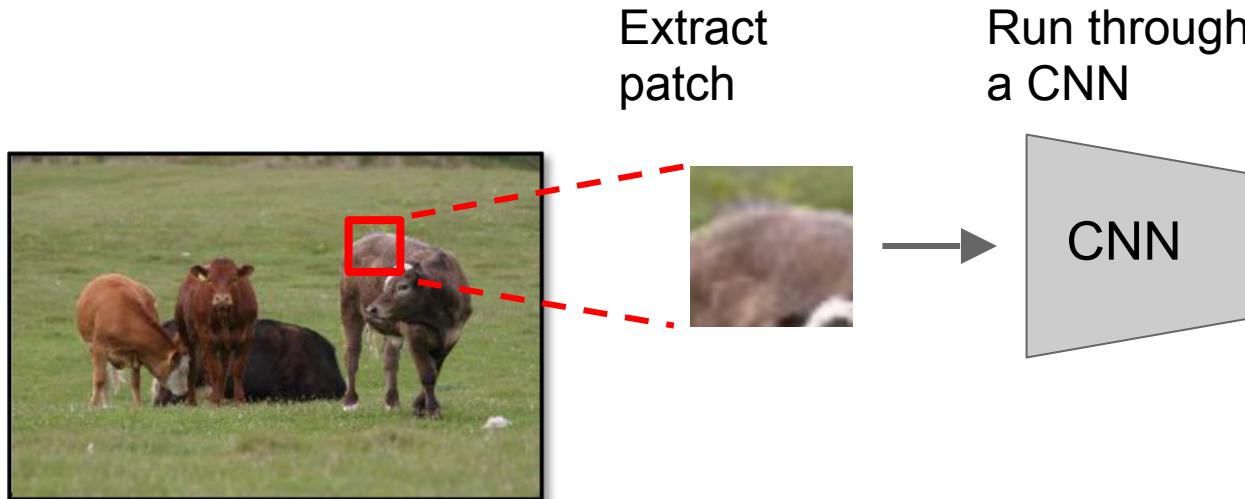
* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Semantic Segmentation

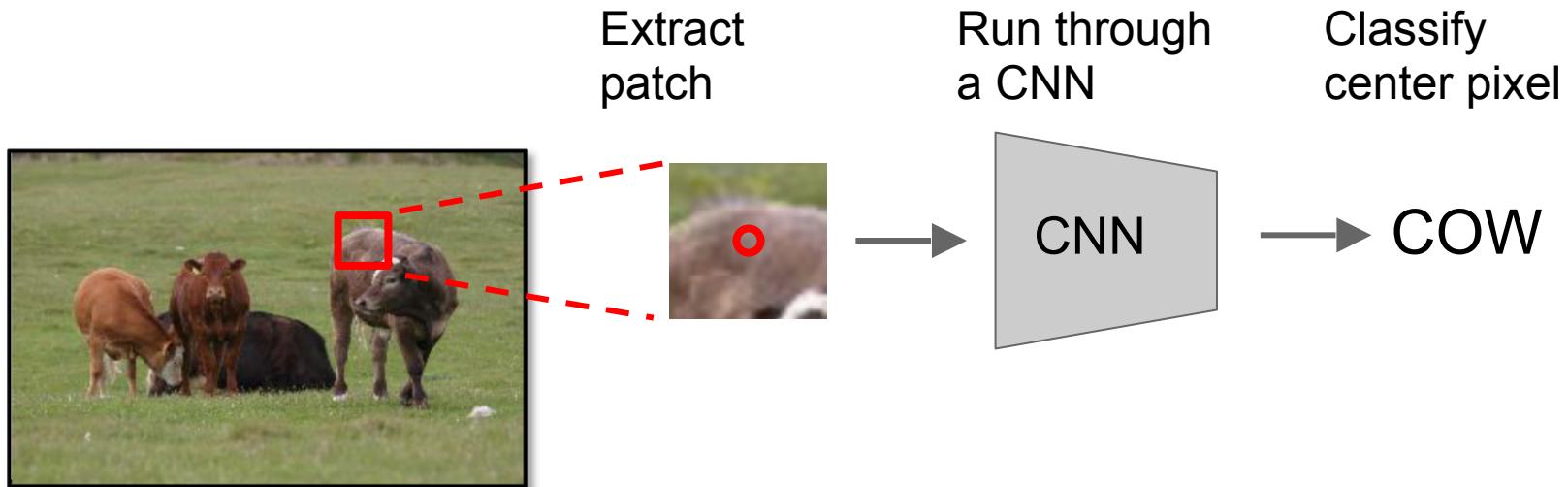
Extract
patch



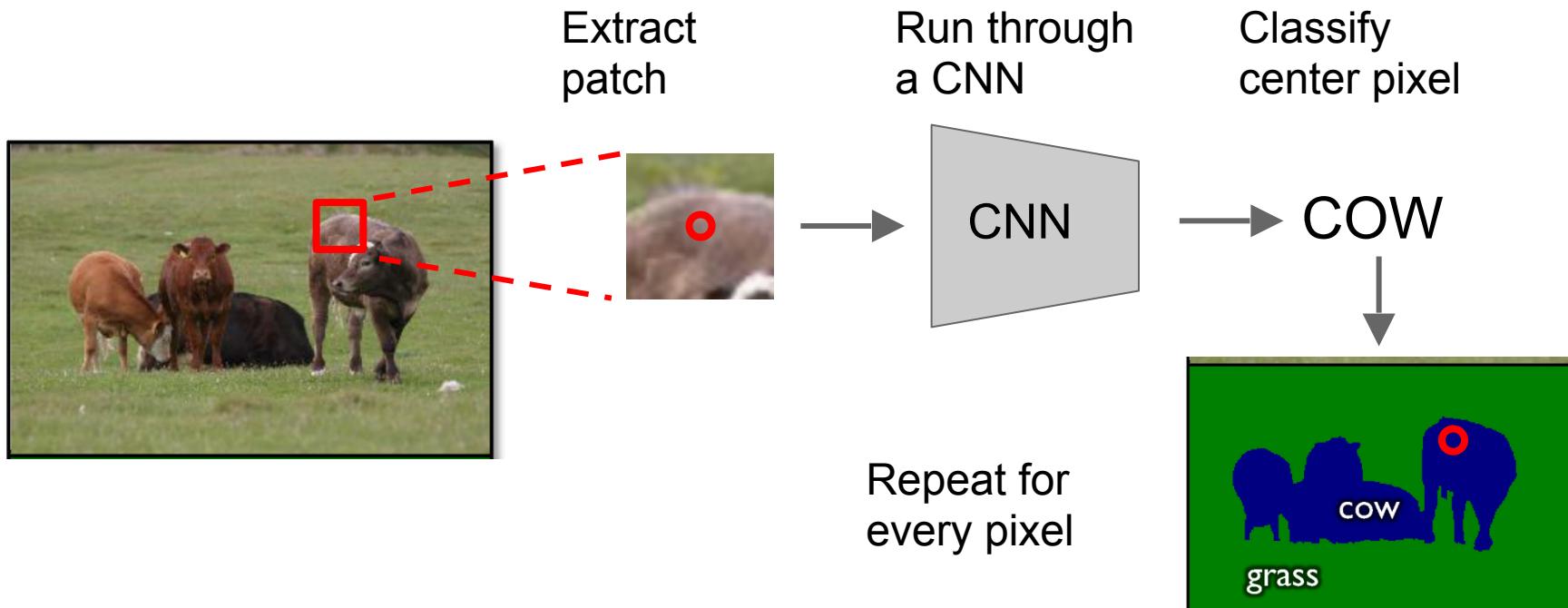
Semantic Segmentation



Semantic Segmentation



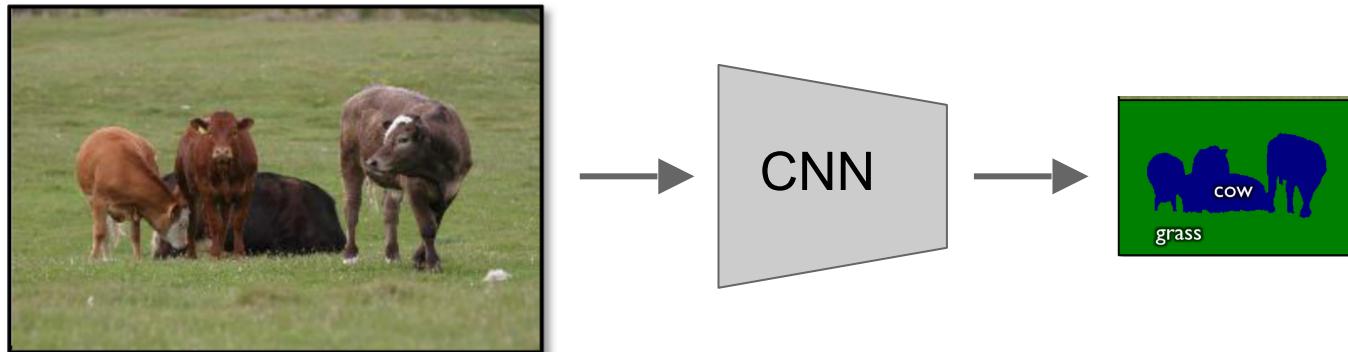
Semantic Segmentation



* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Semantic Segmentation

Run “fully convolutional” network
to get all pixels at once



Smaller output
due to pooling

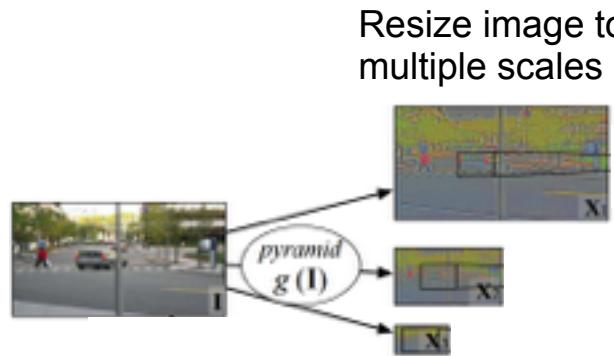
Semantic Segmentation: Multi-Scale



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

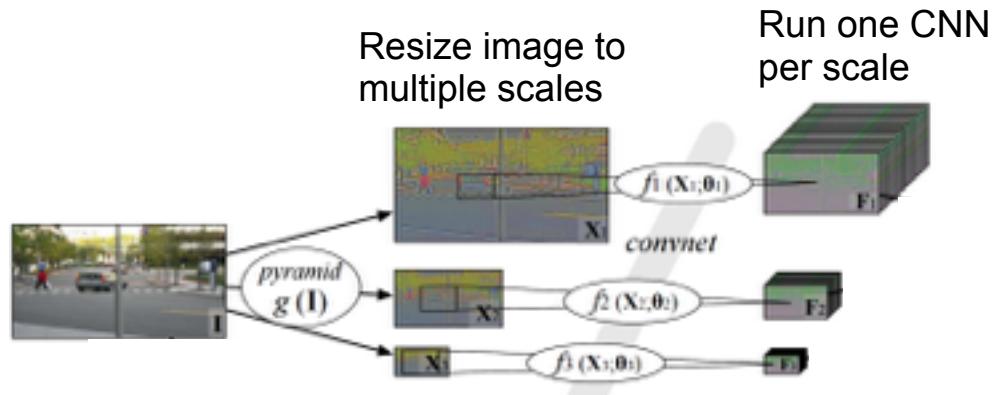
Semantic Segmentation: Multi-Scale



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

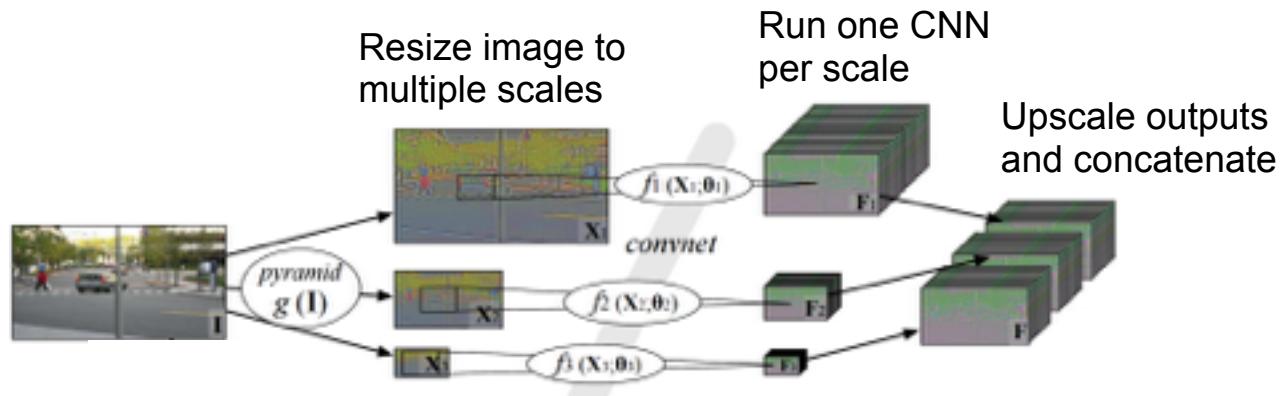
Semantic Segmentation: Multi-Scale



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

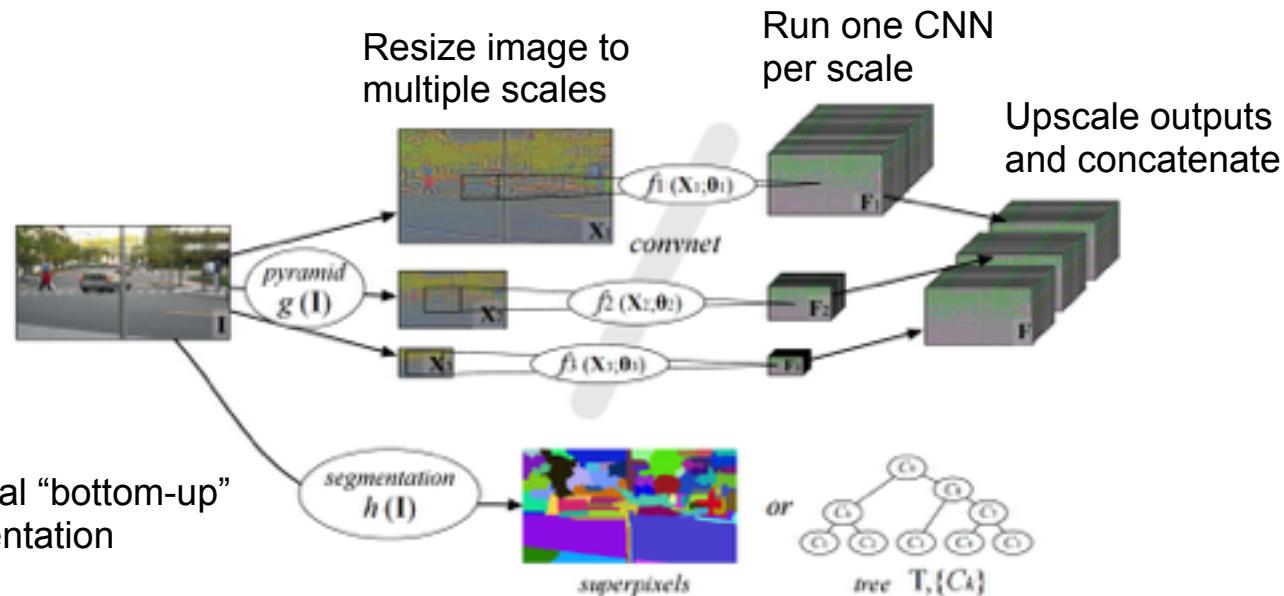
Semantic Segmentation: Multi-Scale



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

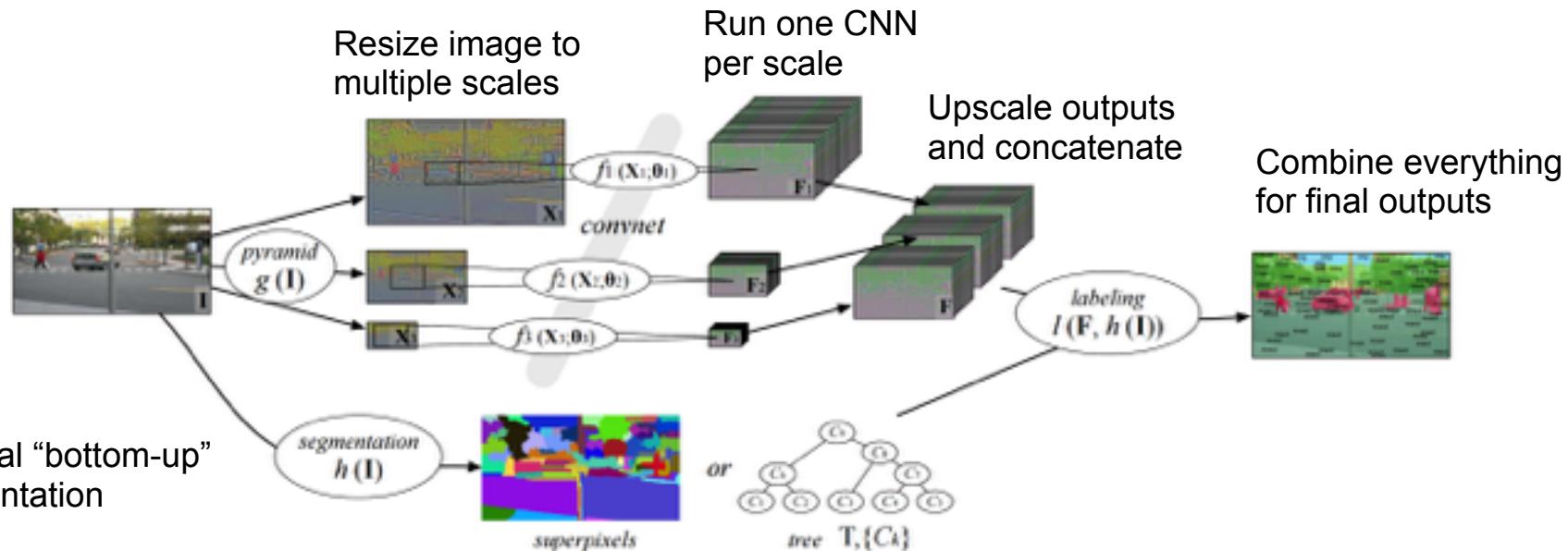
Semantic Segmentation: Multi-Scale



Farabet et al, “Learning Hierarchical Features for Scene Labeling,” TPAMI 2013

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

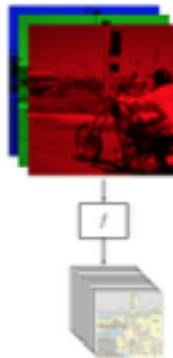
Semantic Segmentation: Multi-Scale



Farabet et al, “Learning Hierarchical Features for Scene Labeling,” TPAMI 2013

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

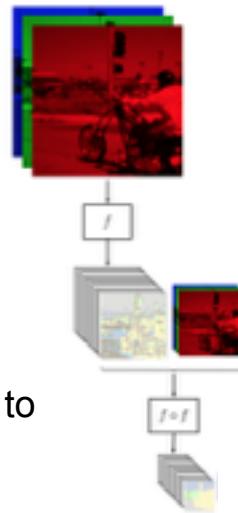
Semantic Segmentation: Refinement



Apply CNN once
to get labels

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation: Refinement

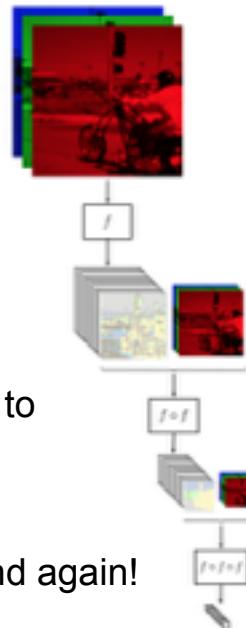


Apply CNN once
to get labels

Apply AGAIN to
refine labels

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation: Refinement



Apply CNN once
to get labels

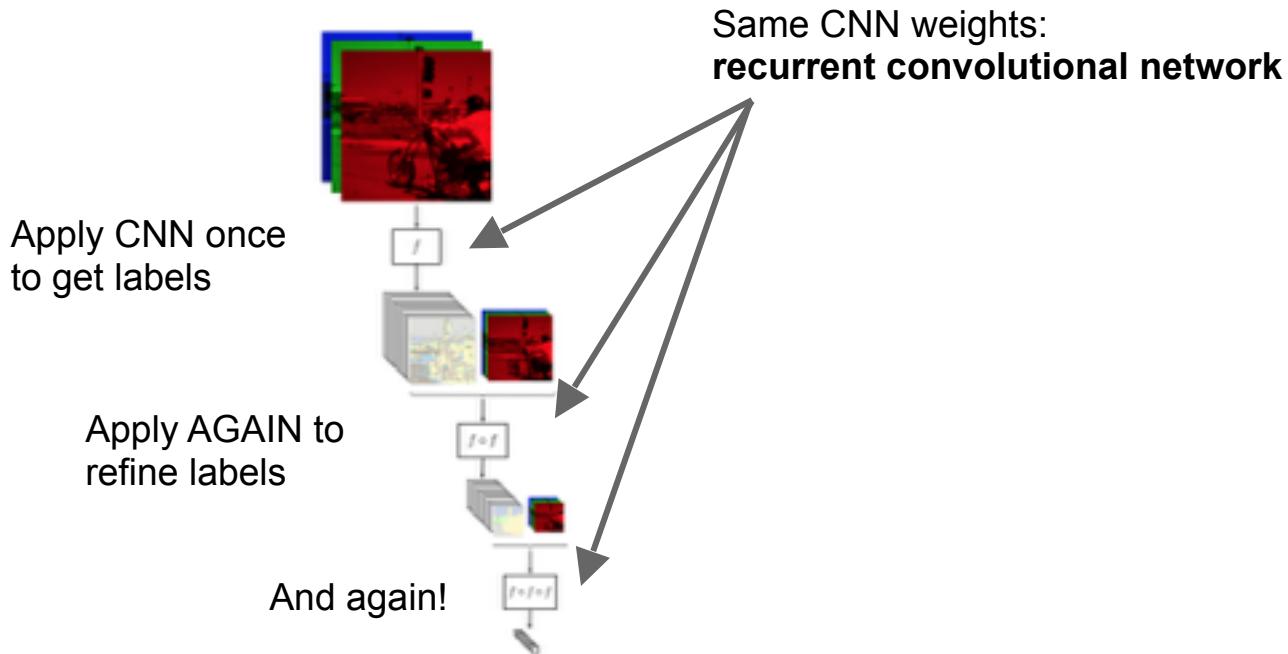
Apply AGAIN to
refine labels

And again!

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Semantic Segmentation: Refinement

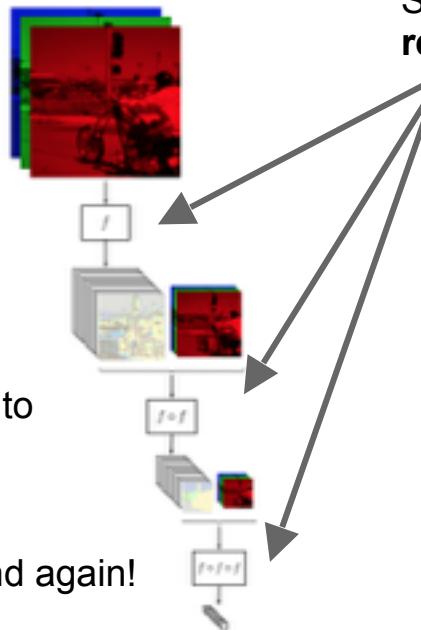


Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Semantic Segmentation: Refinement

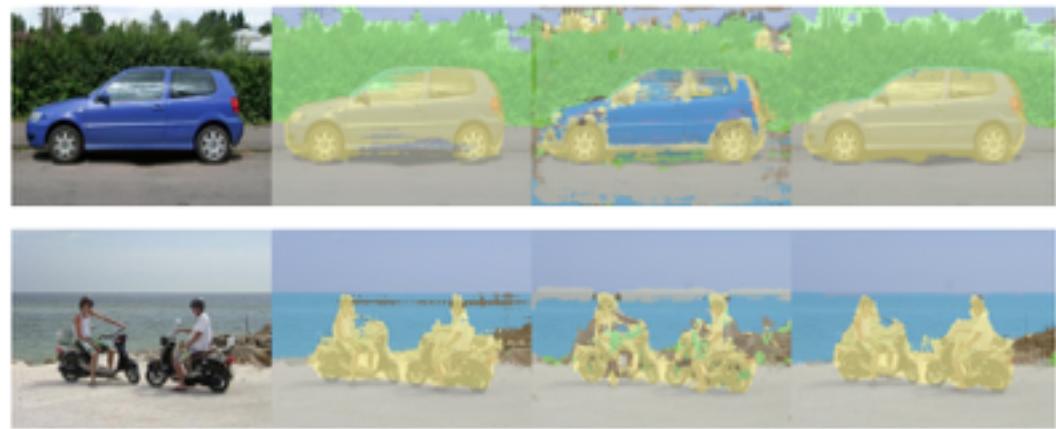
Apply CNN once
to get labels



Same CNN weights:
recurrent convolutional network

Apply AGAIN to
refine labels

And again!

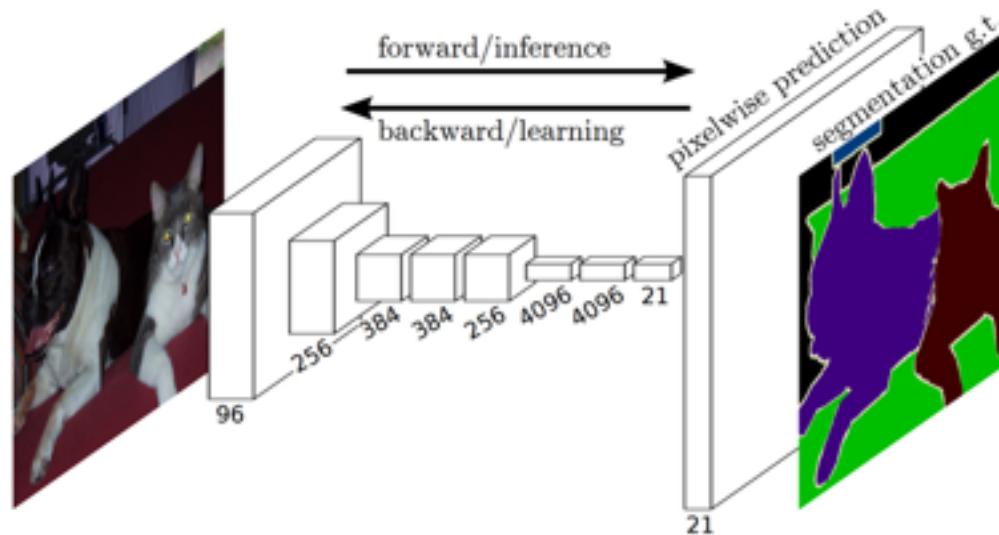


More iterations improve results

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

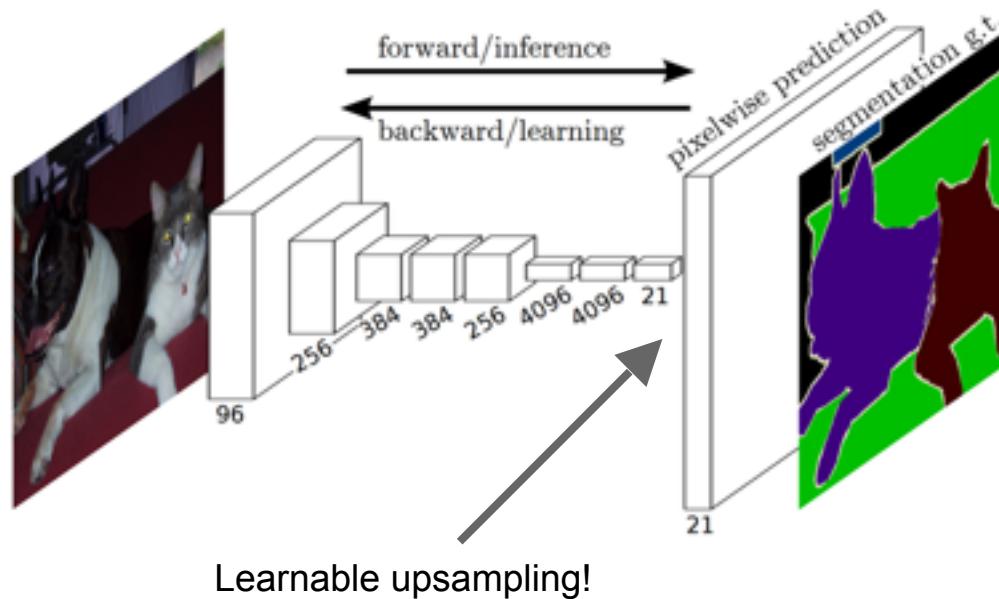
Semantic Segmentation: Upsampling



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

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

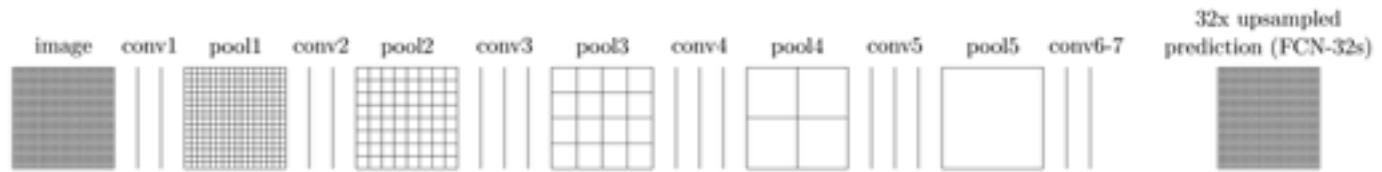
Semantic Segmentation: Upsampling



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

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

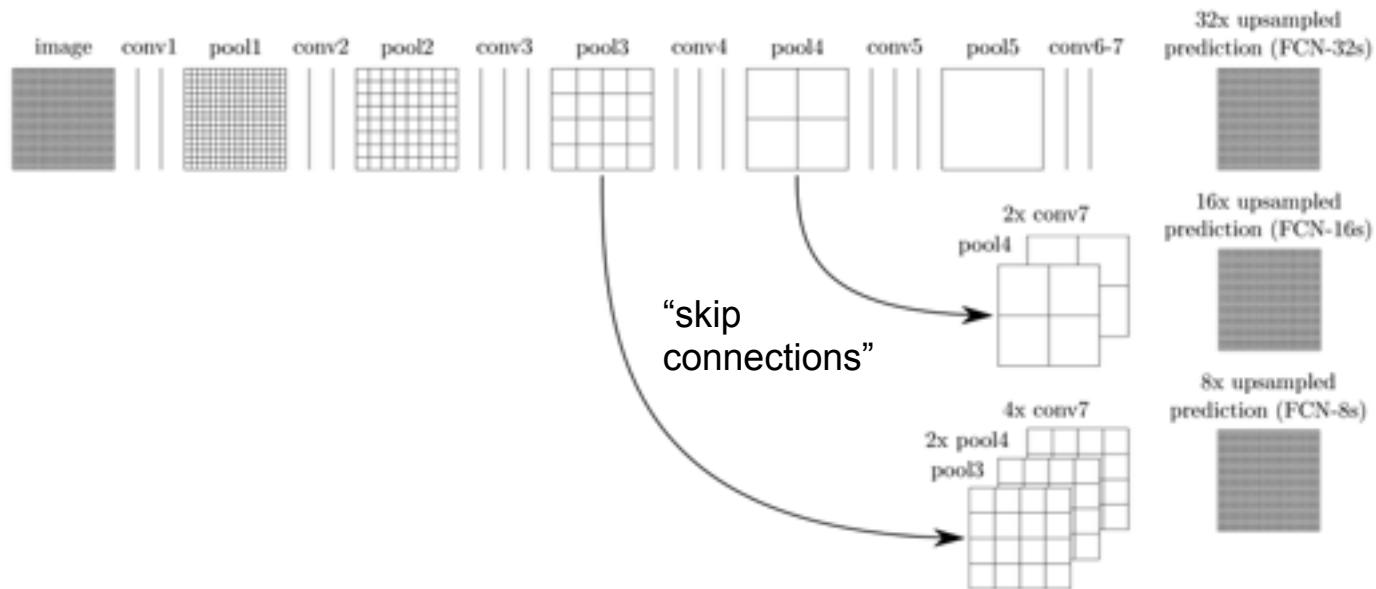
Semantic Segmentation: Upsampling



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

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

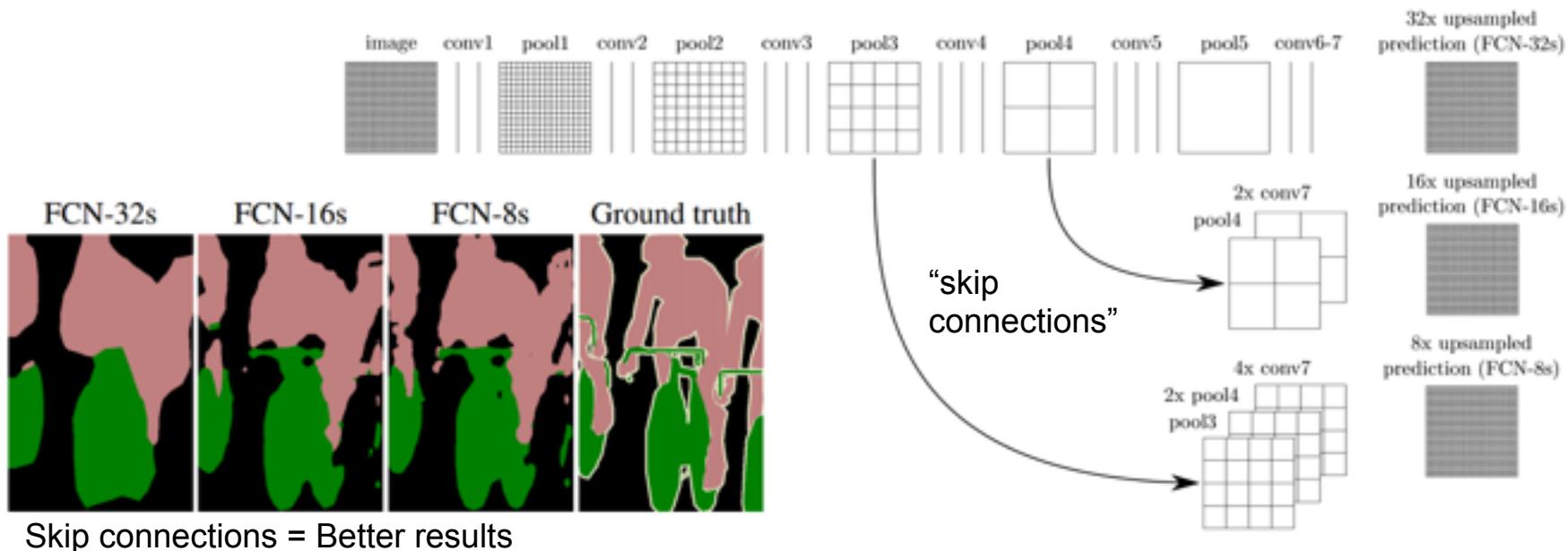
Semantic Segmentation: Upsampling



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

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Semantic Segmentation: Upsampling

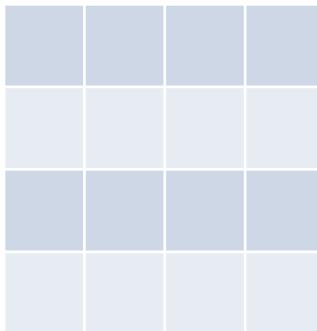


Long, Shelhamer, and Darrell, “Fully Convolutional Networks for Semantic Segmentation”, CVPR 2015

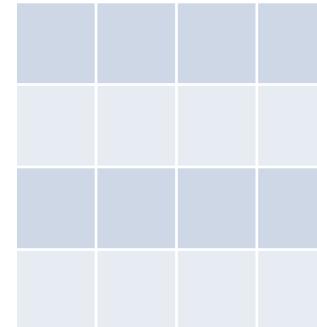
* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 1 pad 1



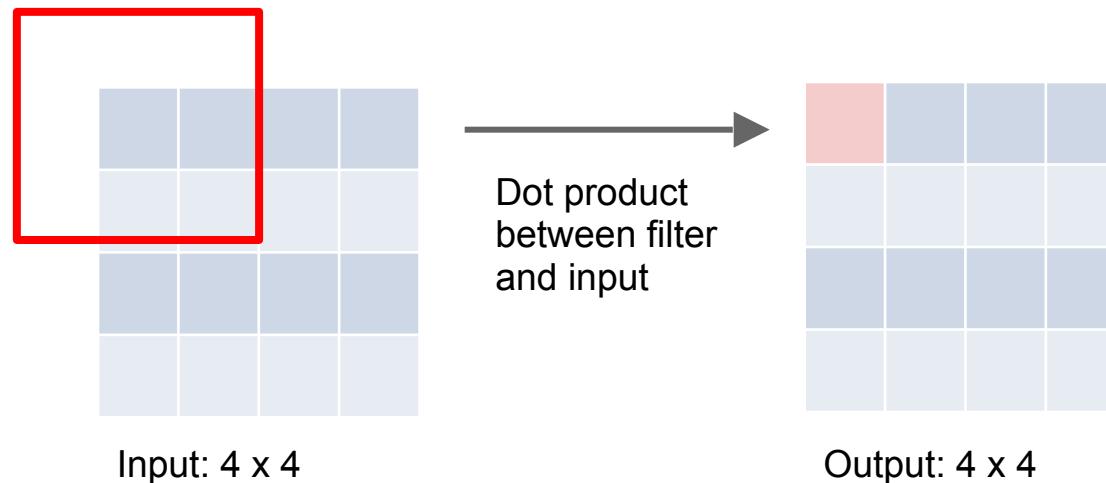
Input: 4×4



Output: 4×4

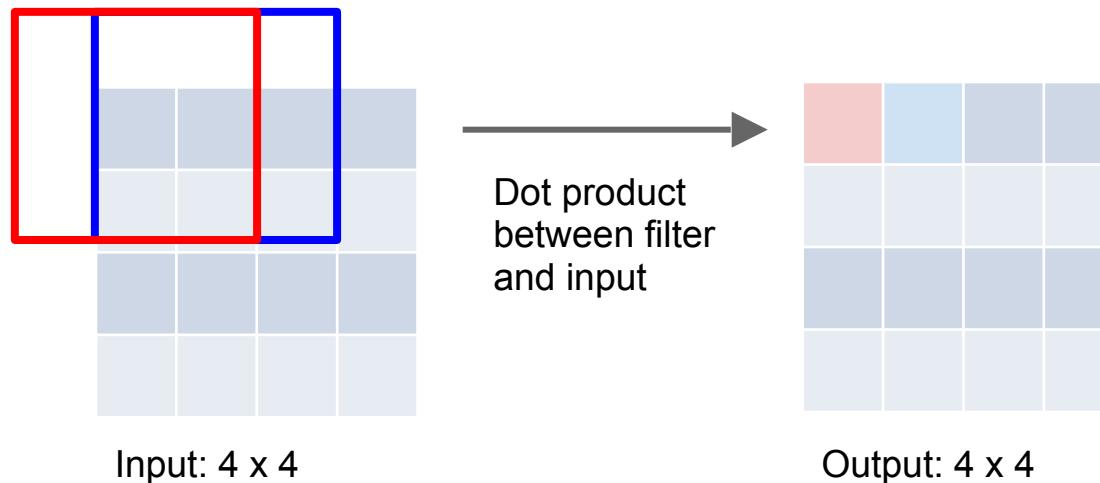
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 1 pad 1



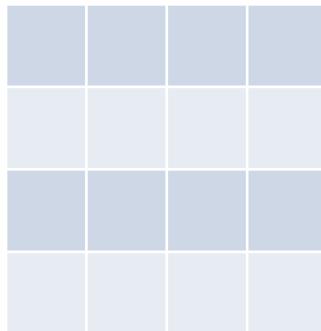
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 1 pad 1

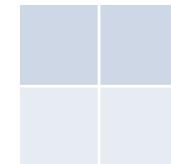


Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, **stride 2** pad 1



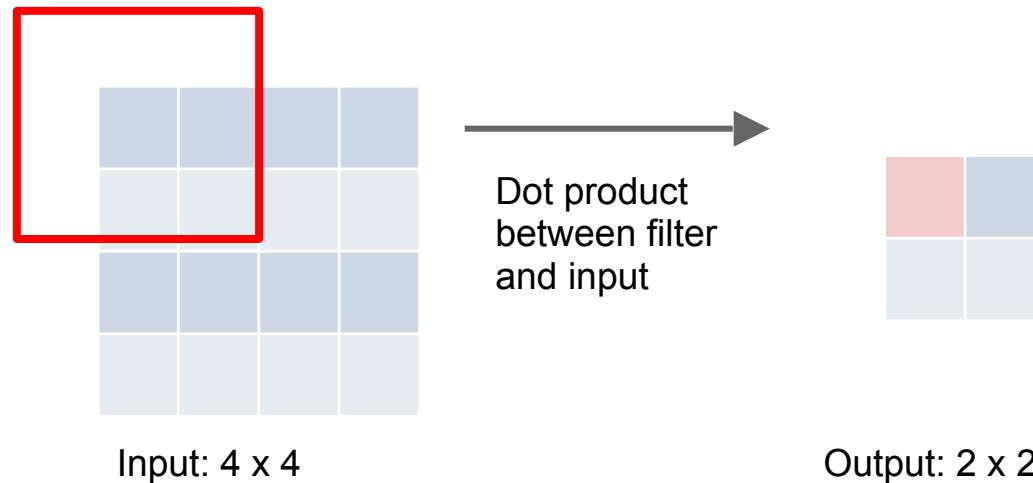
Input: 4×4



Output: 2×2

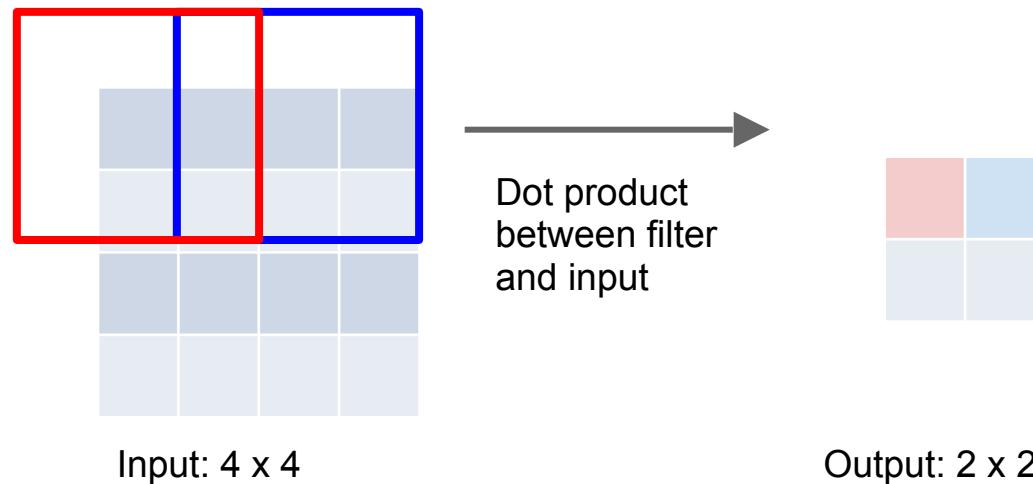
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 2 pad 1



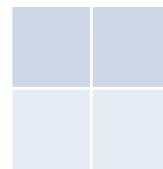
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 2 pad 1

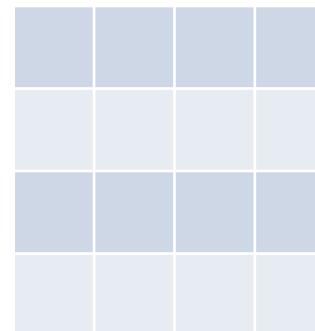


Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1



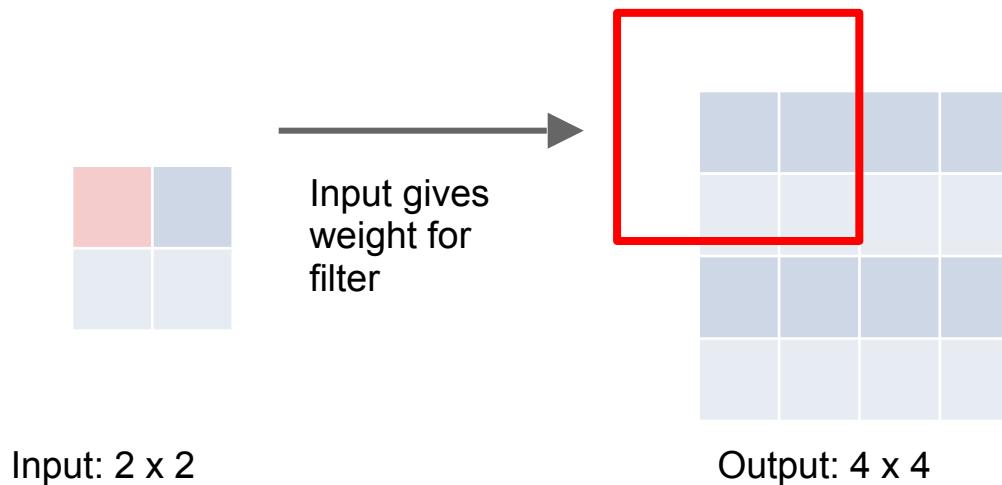
Input: 2 x 2



Output: 4 x 4

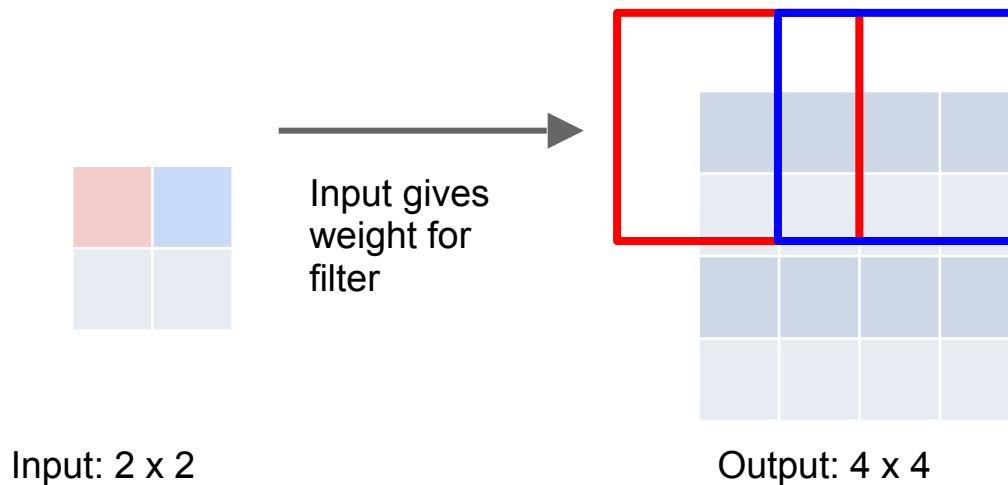
Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1

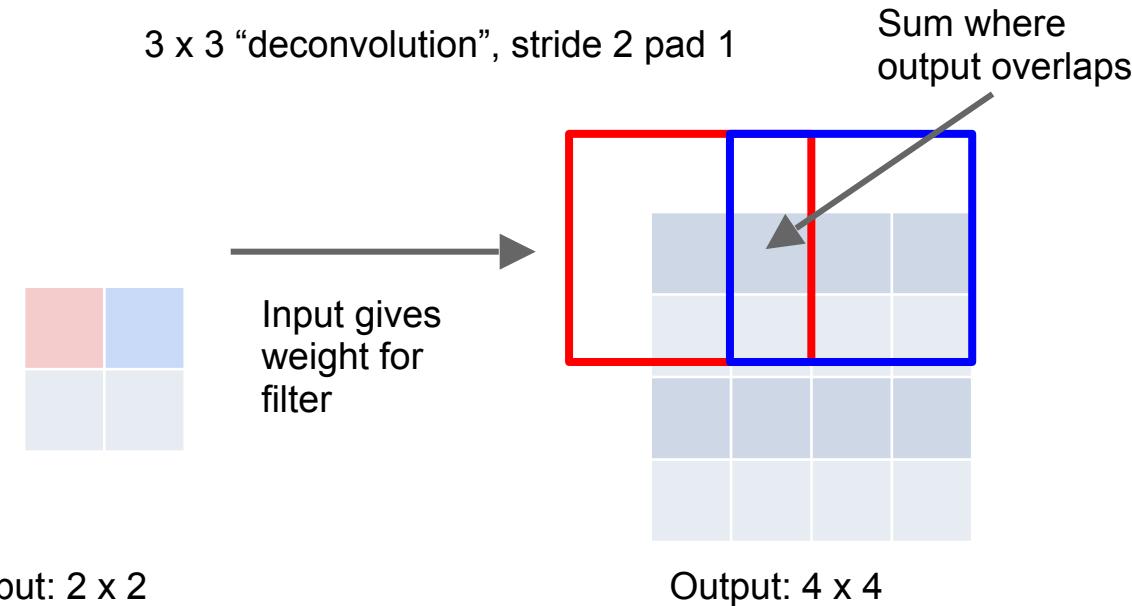


Learnable Upsampling: “Deconvolution”

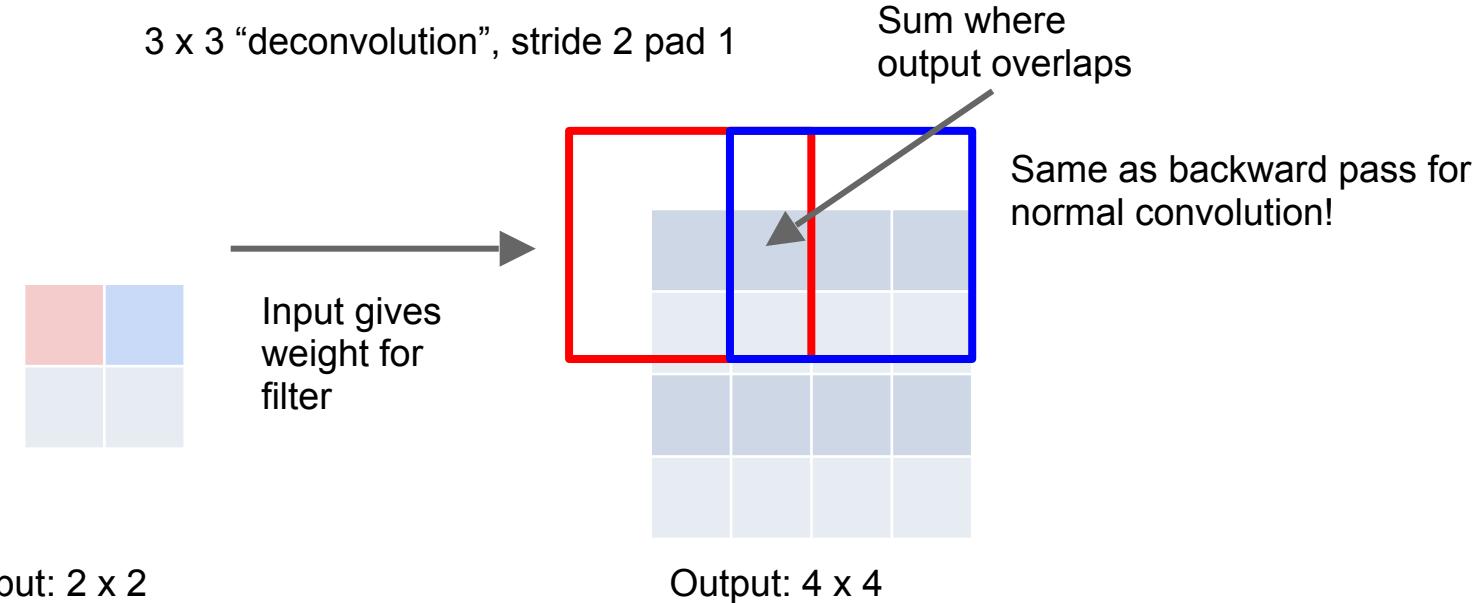
3 x 3 “deconvolution”, stride 2 pad 1



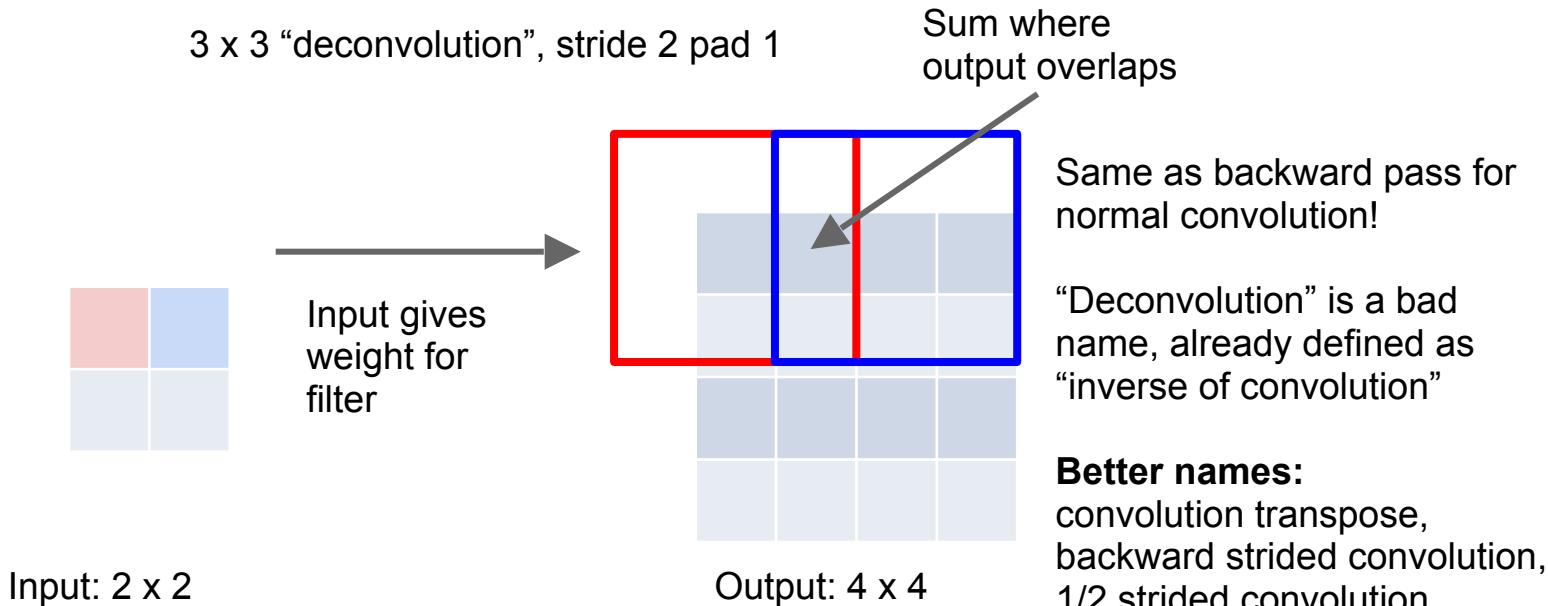
Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”

¹It is more proper to say “convolutional transpose operation” rather than “deconvolutional” operation. Hence, we will be using the term “convolutional transpose” from now.

Im et al, “Generating images with recurrent adversarial networks”, arXiv 2016

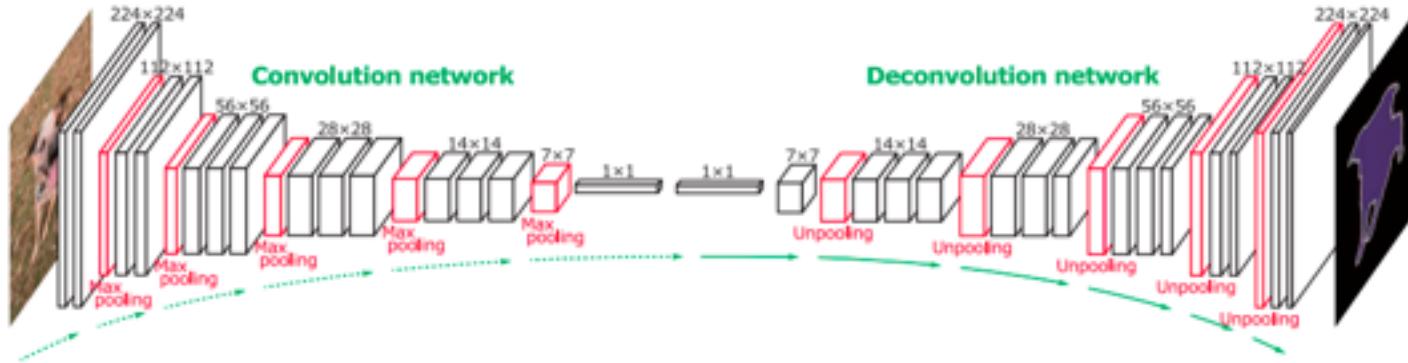
A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions)

Radford et al, “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks”, ICLR 2016

“Deconvolution” is a bad name, already defined as “inverse of convolution”

Better names:
convolution transpose,
backward strided convolution,
1/2 strided convolution,
upconvolution

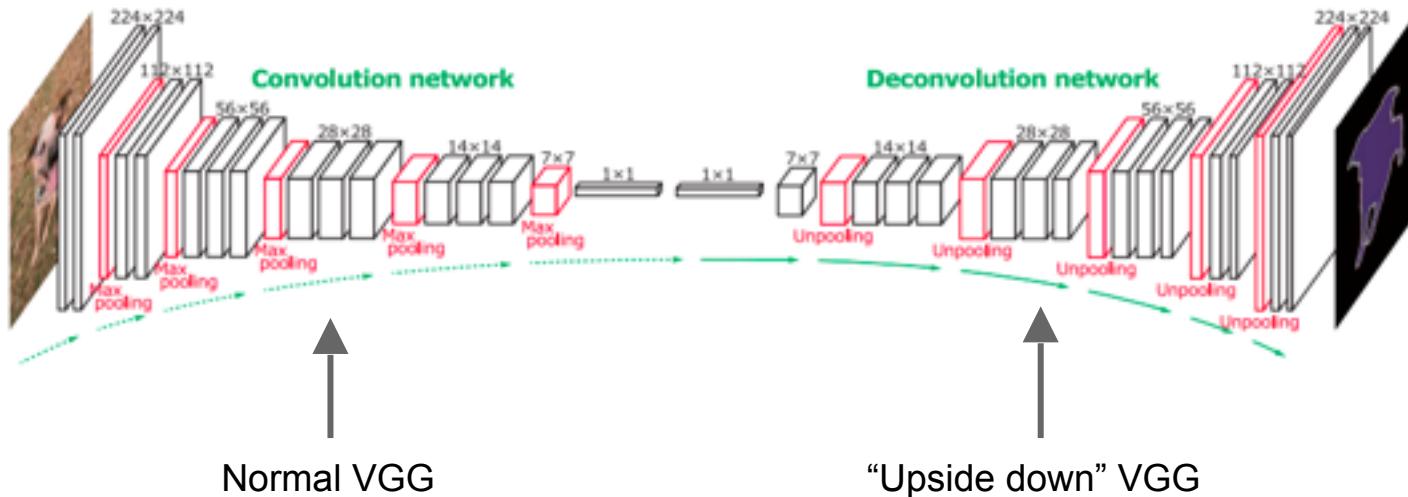
Semantic Segmentation: Upsampling



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

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Semantic Segmentation: Upsampling



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

6 days of training on Titan X...

Instance Segmentation

Instance Segmentation

Detect instances,
give category, label
pixels

“simultaneous
detection and
segmentation” (SDS)

Lots of recent work
(MS-COCO)

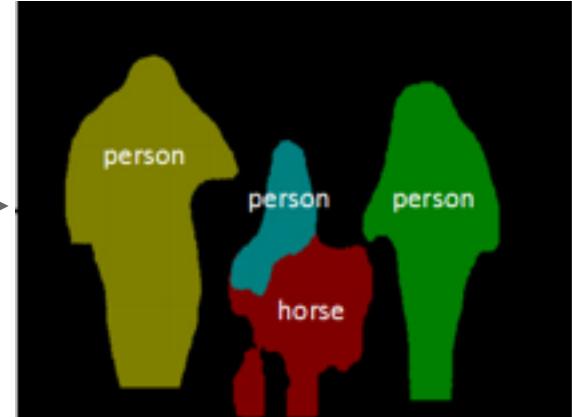


Figure credit: Dai et al, “Instance-aware Semantic Segmentation via Multi-task Network Cascades”, arXiv 2015

Instance Segmentation

Similar to R-CNN, but
with segments



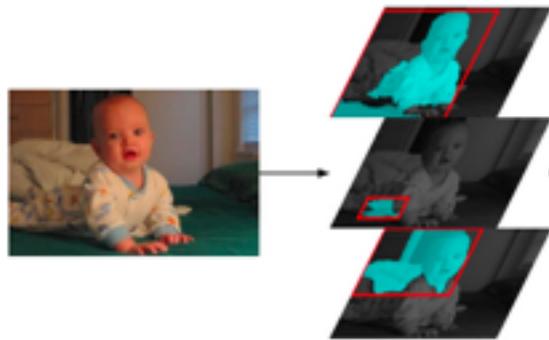
Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Instance Segmentation

Proposal Generation
External Segment proposals

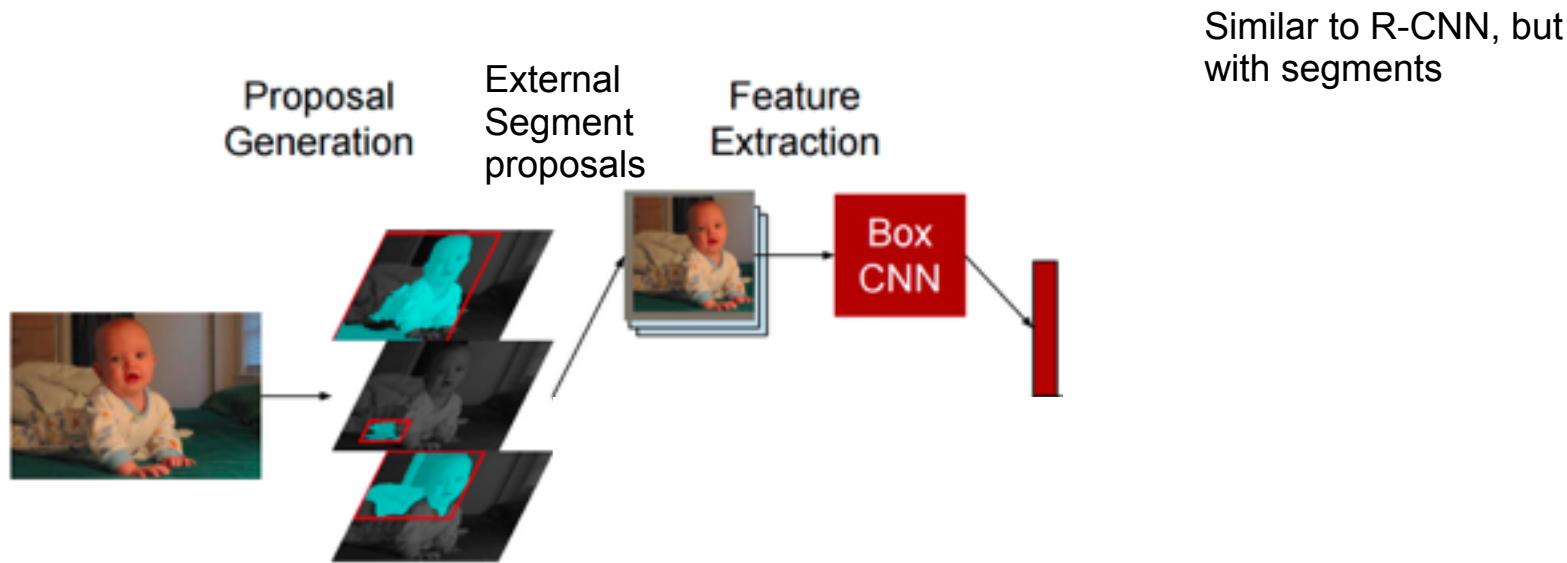
Similar to R-CNN, but
with segments



Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

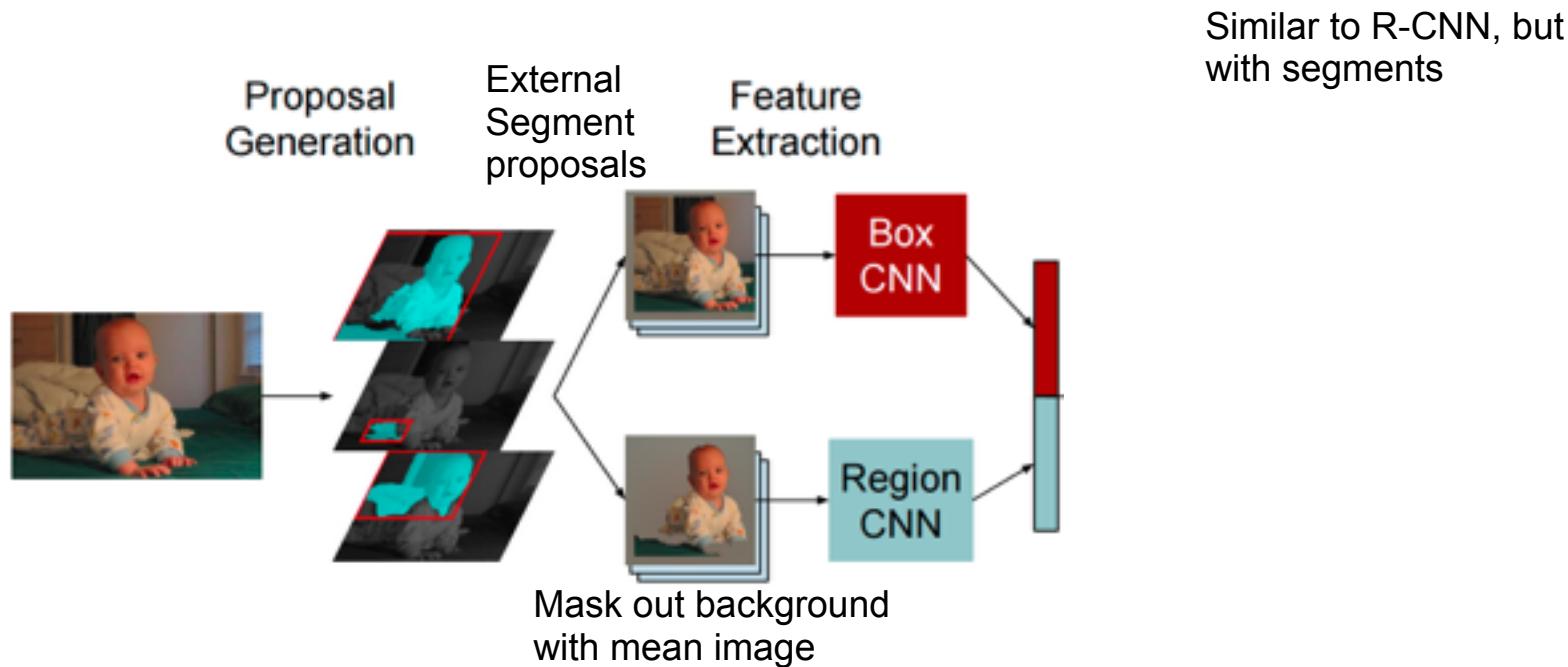
Instance Segmentation



Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

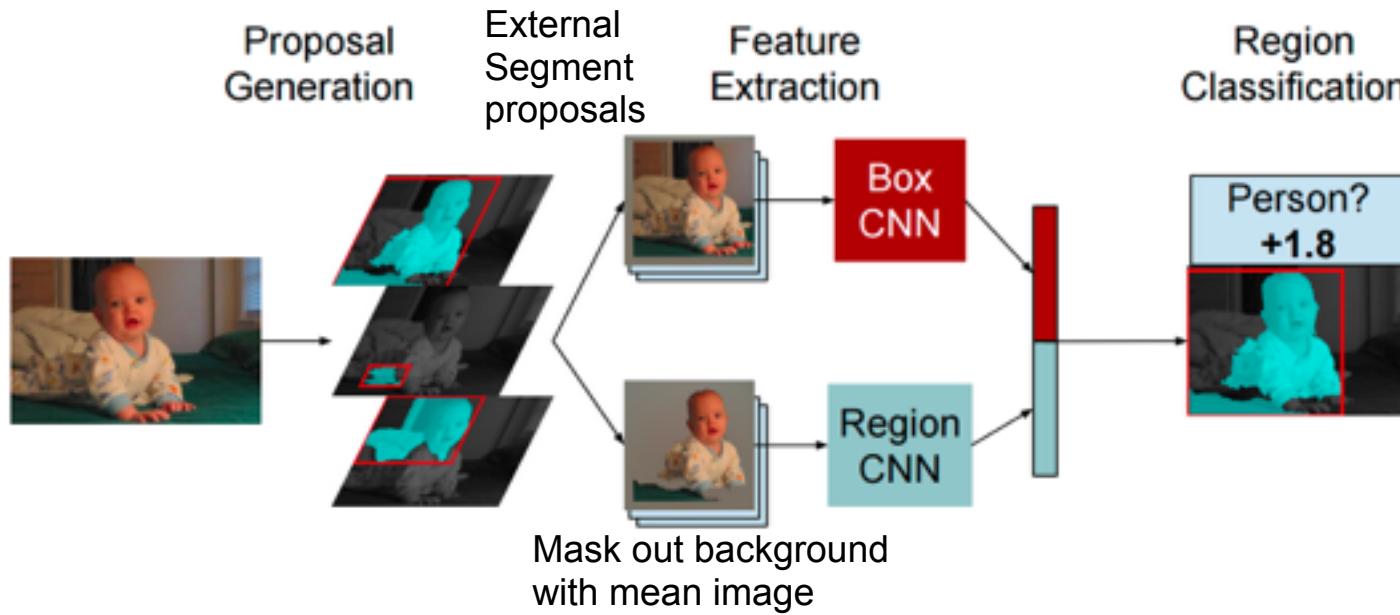
Instance Segmentation



Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

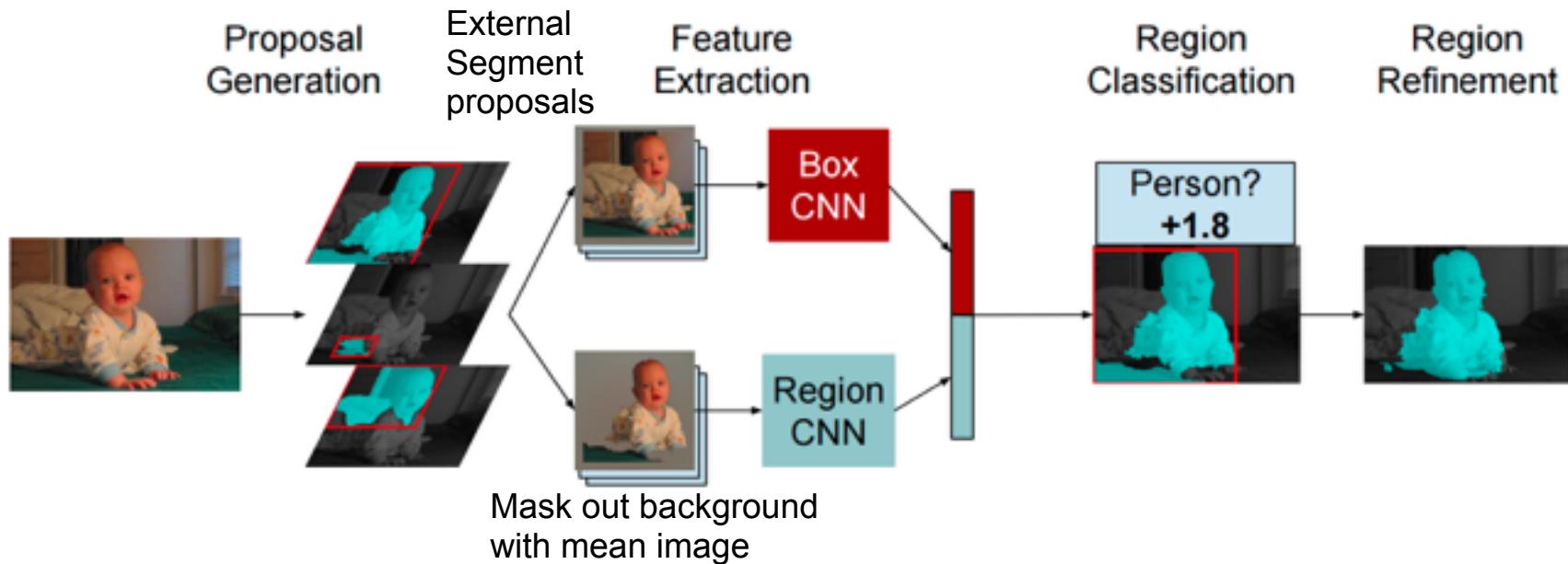
Instance Segmentation



Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

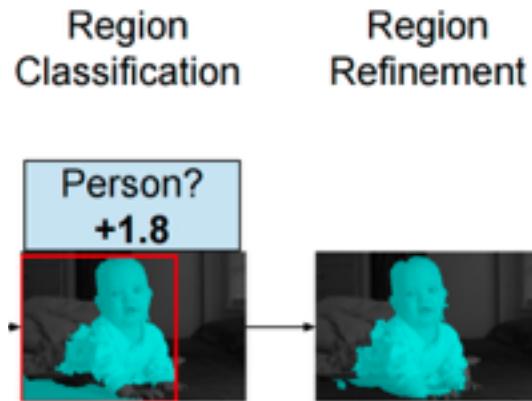
Instance Segmentation



Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Hypercolumns



Hariharan et al, "Hypercolumns for Object Segmentation and Fine-grained Localization", CVPR 2015

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades

Similar to
Faster R-CNN



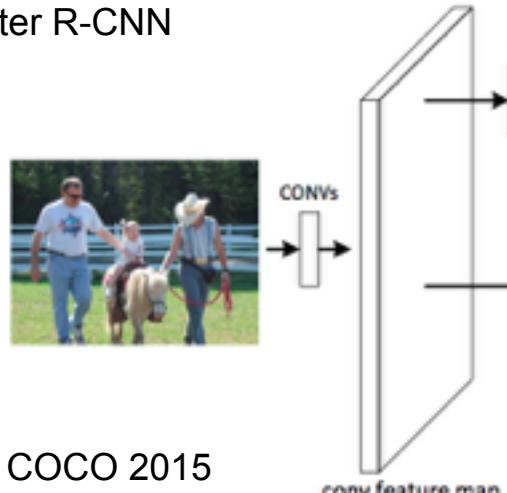
Won COCO 2015
challenge
(with ResNet)

Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades

Similar to
Faster R-CNN



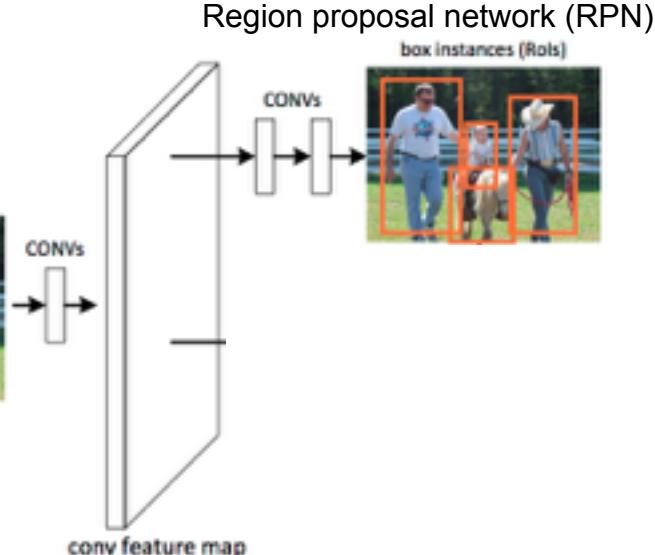
Won COCO 2015
challenge
(with ResNet)

Dai et al, “Instance-aware Semantic Segmentation via Multi-task Network Cascades”, arXiv 2015

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades

Similar to
Faster R-CNN



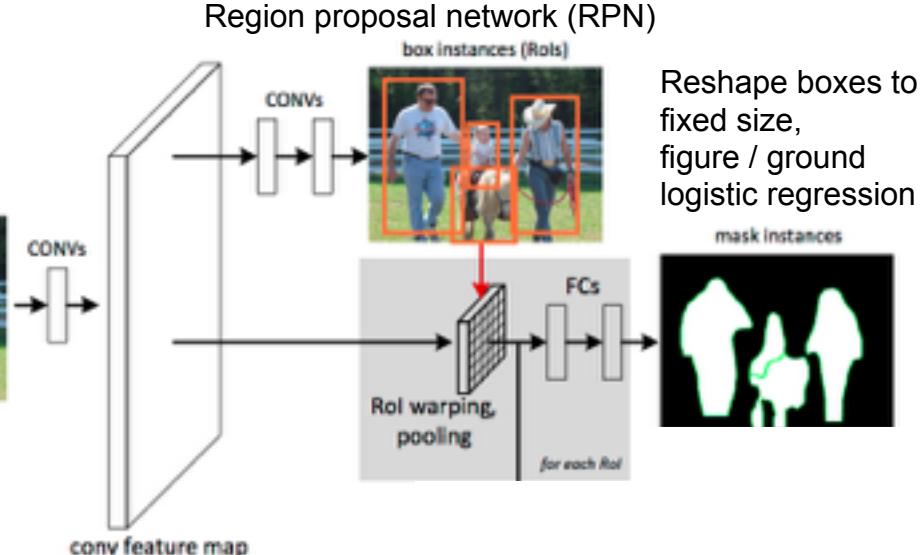
Won COCO 2015
challenge
(with ResNet)

Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades

Similar to
Faster R-CNN



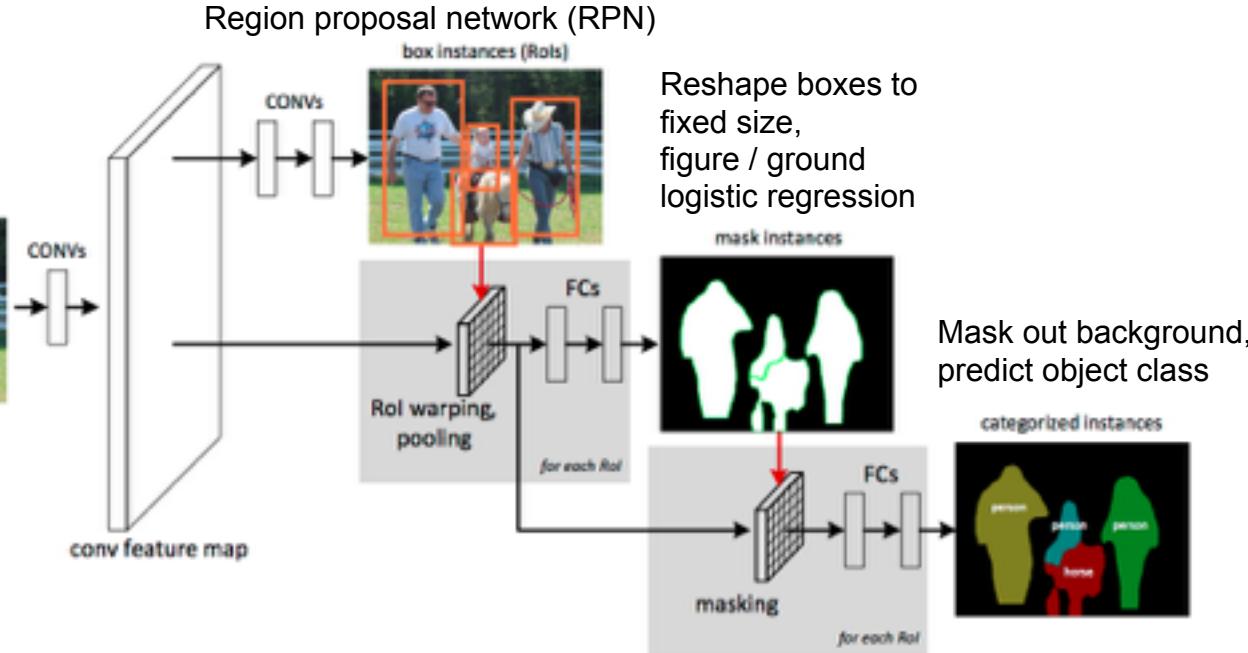
Won COCO 2015
challenge
(with ResNet)

Dai et al, “Instance-aware Semantic Segmentation via Multi-task Network Cascades”, arXiv 2015

* Original slides borrowed from Andrej Karpathy
and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades

Similar to
Faster R-CNN



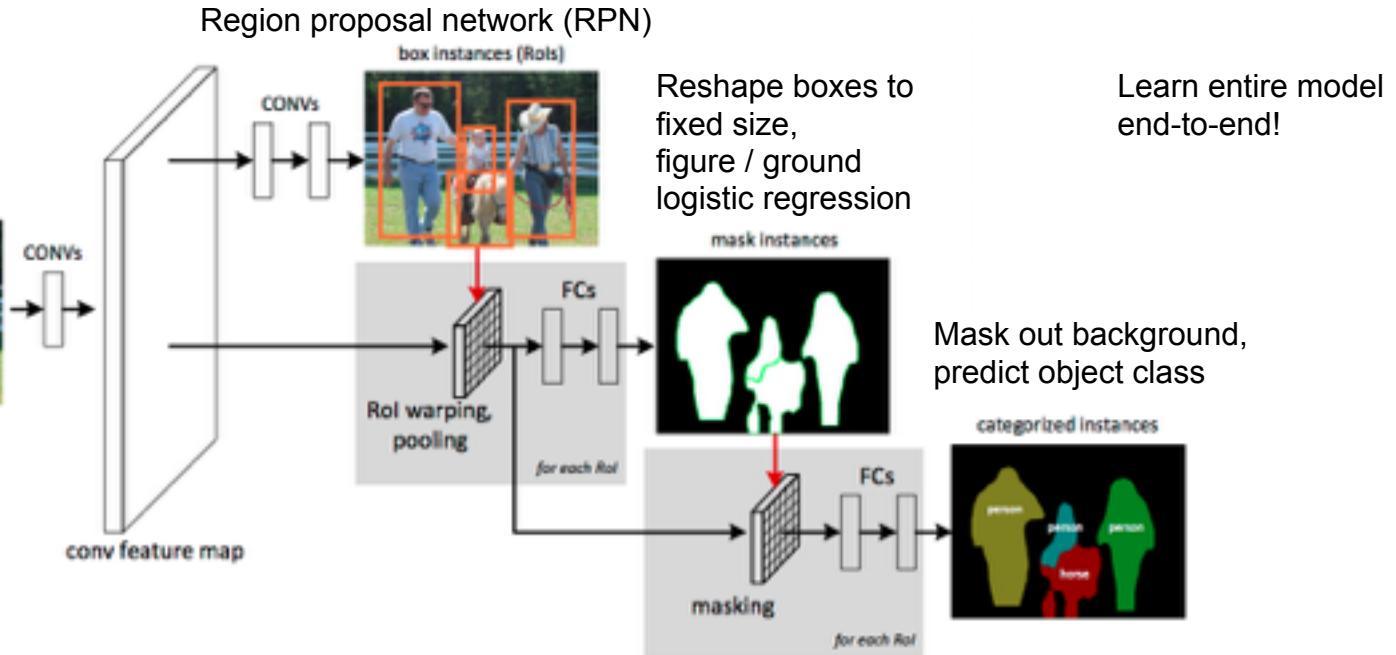
Won COCO 2015 challenge
(with ResNet)

Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades

Similar to
Faster R-CNN

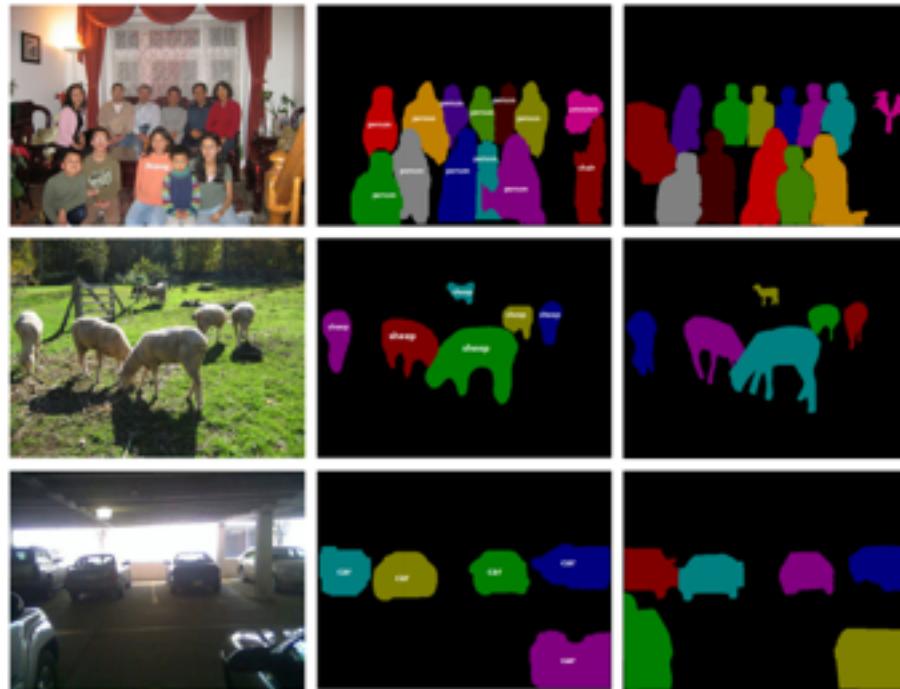


Won COCO 2015 challenge
(with ResNet)

Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Instance Segmentation: Cascades



Predictions

Ground truth

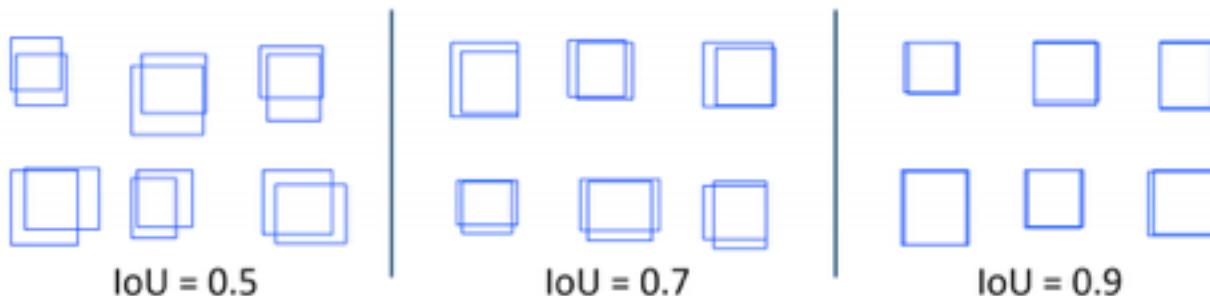
Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

* Original slides borrowed from Andrej Karpathy and Li Fei-Fei, Stanford cs231n

Evaluation metrics

Detection Score: Average Precision (mAP is mean AP over all object categories)

- AP is averaged over multiple IoU values between 0.5 and 0.95 (and categories, size)
- More comprehensive comparison metric than the traditional AP at Intersection over Union (IoU) threshold of 0.5.



Evaluation Metrics

Detection Score: AP

- AP is averaged over multiple IoU values between 0.5 and 0.95
- AP is averaged over groups of objects or over object size in an image
- More comprehensive metric than the traditional AP at Intersection over Union (IoU) threshold of 0.5.

A < 32x32



A > 96x96

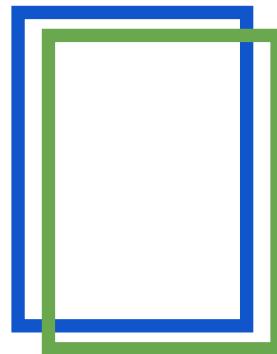


32x32 < A < 96x96

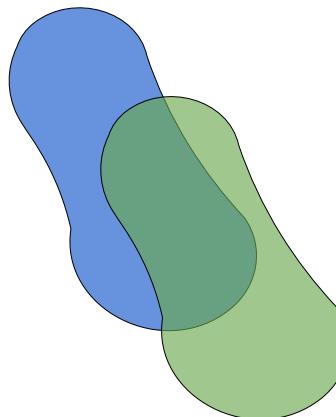


Evaluating Detection or Segmentation

To calculate AP we need:



Bounding Box IoU



Mask IoU

Segmentation Overview

- Graph Cut

- Classic algorithm to quickly get foreground from background segmentation
- Not trained to consider object shape prior
- Can fail with complicated backgrounds

- Unsupervised Segmentation

- Superpixels
- Open area of research in ML

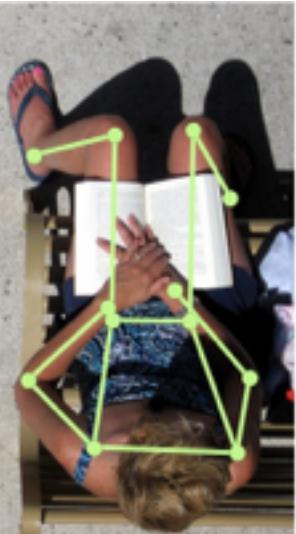
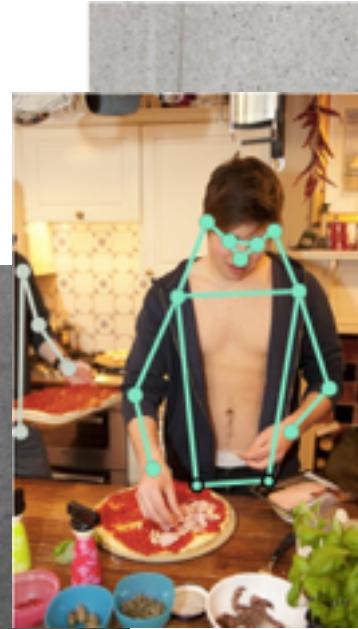
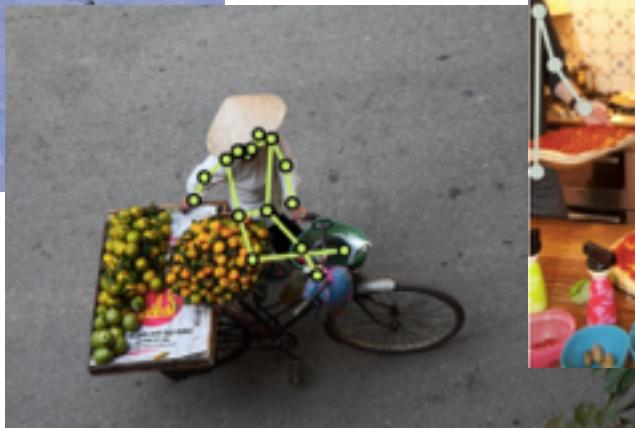
- Semantic segmentation

- Classify all pixels
- Fully convolutional models, downsample then upsample
- Learnable upsampling: fractionally strided convolution
- Skip connections can help

- Instance Segmentation

- Detect instance, generate mask
- Similar pipelines to object detection

Pose Estimation (aka Keypoints)



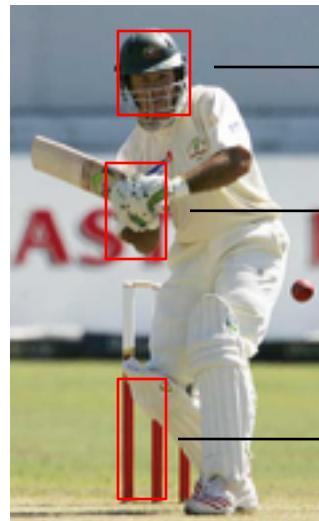
Classic Problem: Activity Recognition



What is this person doing?

Human pose estimation & Object detection

Human pose estimation is challenging.



Difficult part appearance

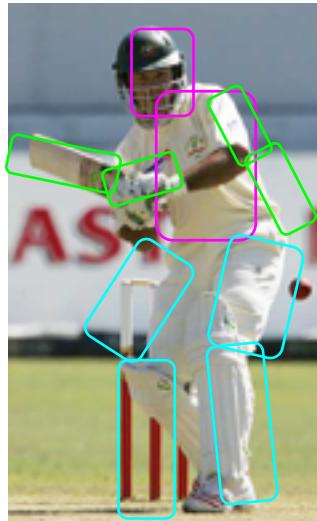
Self-occlusion

Image region looks like a body part

- Felzenszwalb & Huttenlocher, 2005
- Ren et al, 2005
- Ramanan, 2006
- Ferrari et al, 2008
- Yang & Mori, 2008
- Andriluka et al, 2009
- Eichner & Ferrari, 2009

Human pose estimation & Object detection

Human pose estimation is challenging.

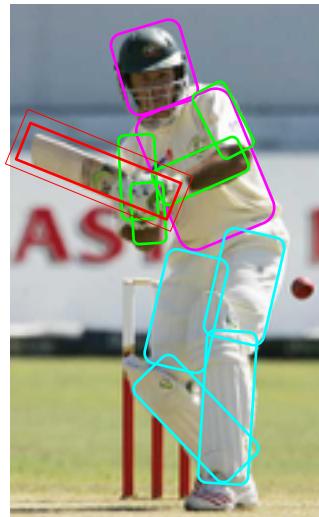


- Felzenszwalb & Huttenlocher, 2005
- Ren et al, 2005
- Ramanan, 2006
- Ferrari et al, 2008
- Yang & Mori, 2008
- Andriluka et al, 2009
- Eichner & Ferrari, 2009

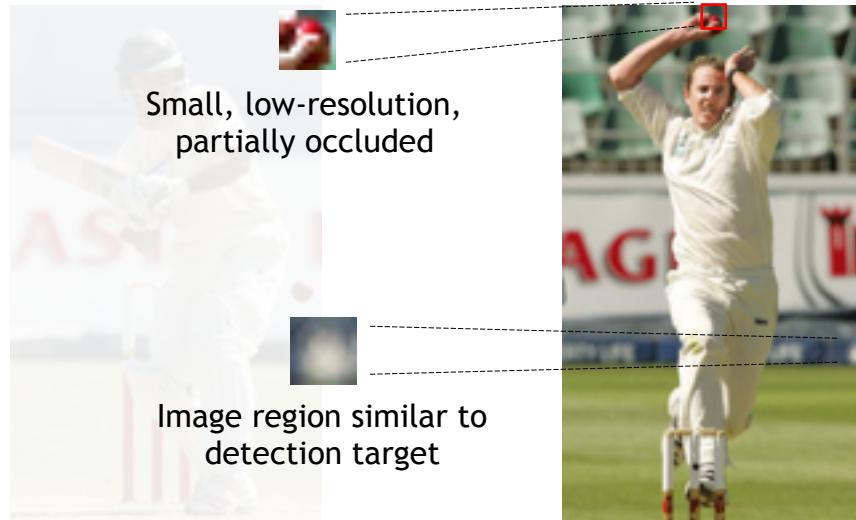
Human pose estimation & Object detection

Facilitate

Given the
object is
detected.

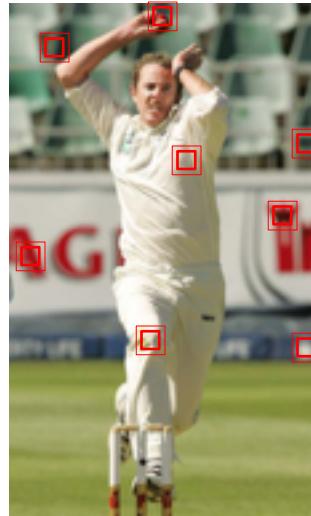


Human pose estimation & Object detection



- Viola & Jones, 2001
- Lampert et al, 2008
- Divvala et al, 2009
- Vedaldi et al, 2009

Human pose estimation & Object detection

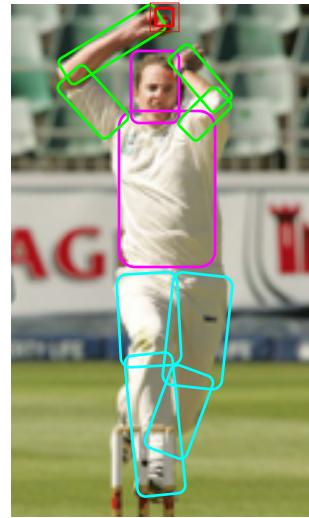


Object
detection is
challenging

- Viola & Jones, 2001
- Lampert et al, 2008
- Divvala et al, 2009
- Vedaldi et al, 2009

Human pose estimation & Object detection

Facilitate

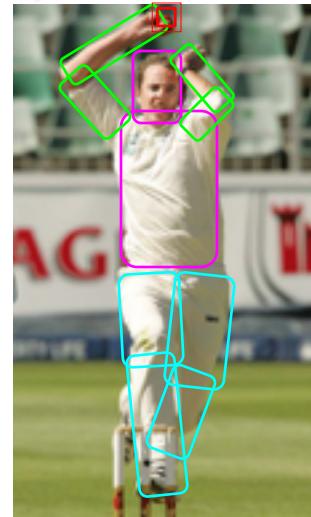
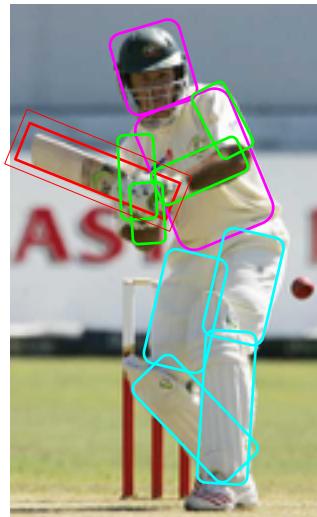


Given the
pose is
estimated.

Human pose estimation & Object detection



Mutual Context



Multi-Person Pose Estimation using Part Affinity Fields

Zhe Cao, Shih-En Wei, Tomas Simon, Yaser Sheikh
Carnegie Mellon University

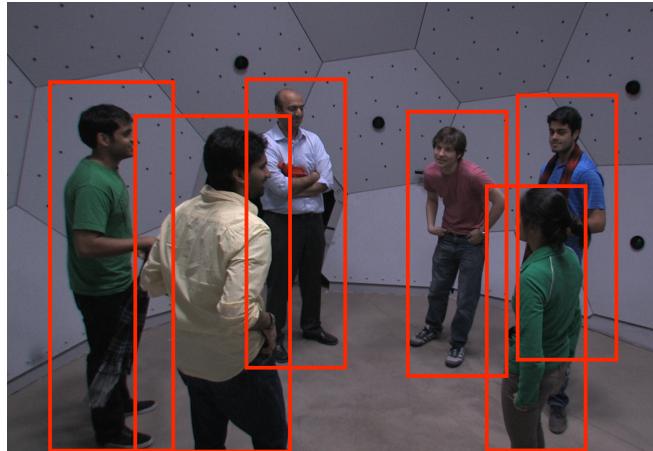


Top-down Approach: Person Detection + Pose Estimation



Top-down

Top-down Approach: Person Detection + Pose Estimation



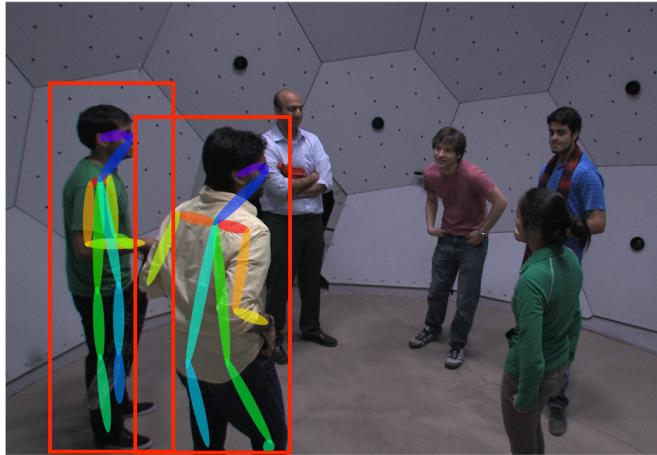
Top-down

Top-down Approach: Person Detection + Pose Estimation



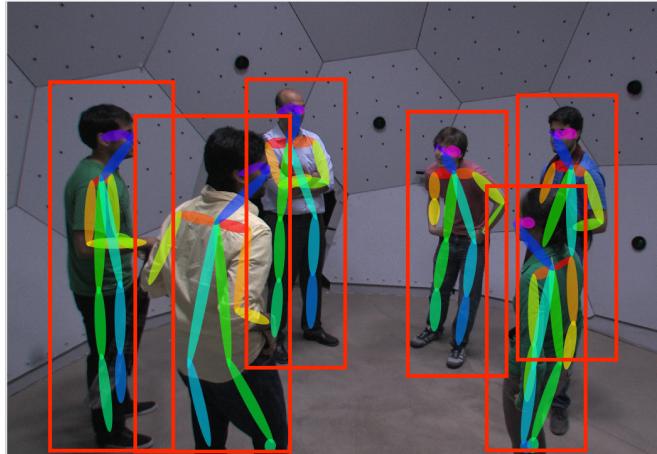
Top-down

Top-down Approach: Person Detection + Pose Estimation



Top-down

Top-down Approach: Person Detection + Pose Estimation



Top-down

Our Method: Parts Detection + Parts Association



Top-down

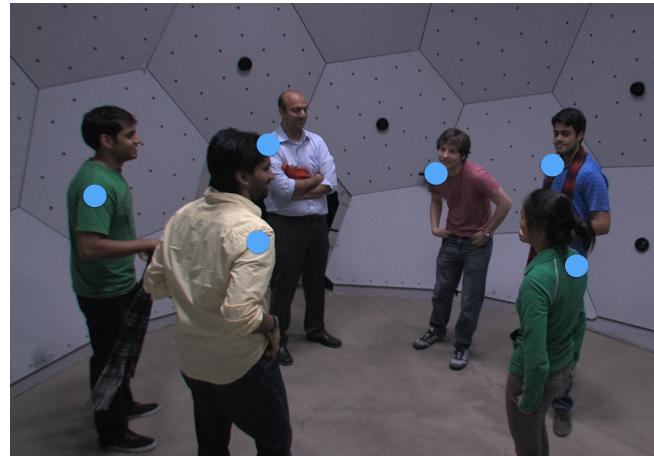


Ours

Our Method: Parts Detection + Parts Association

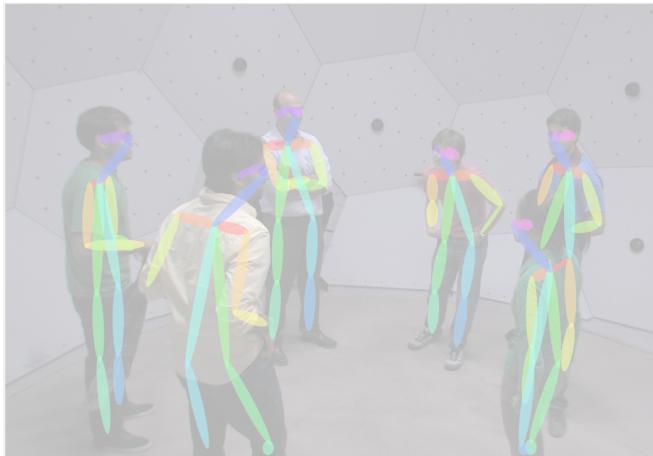


Top-down



Ours

Our Method: Parts Detection + Parts Association

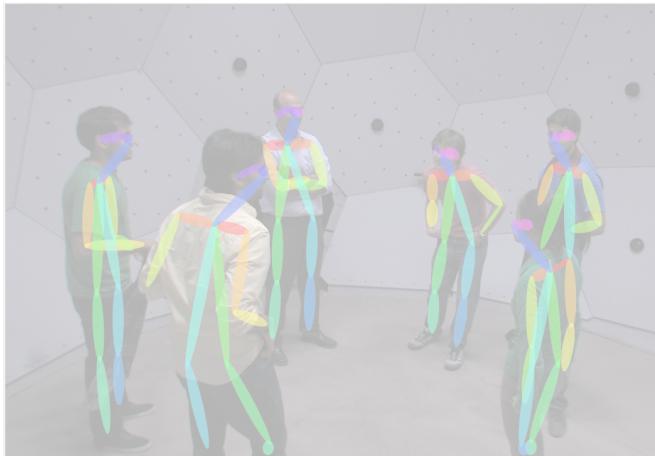


Top-down

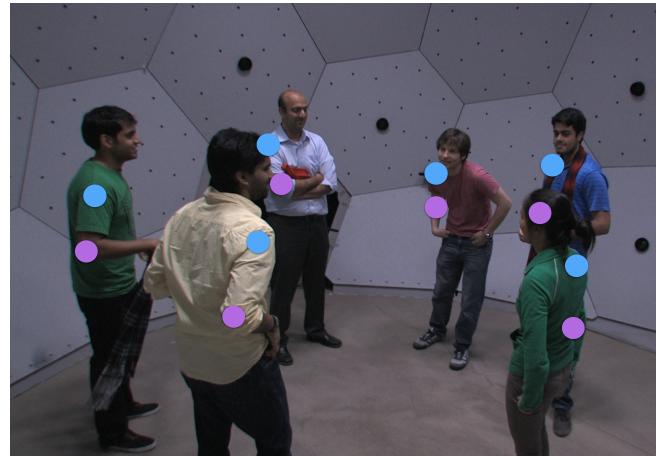


Ours

Our Method: Parts Detection + Parts Association



Top-down

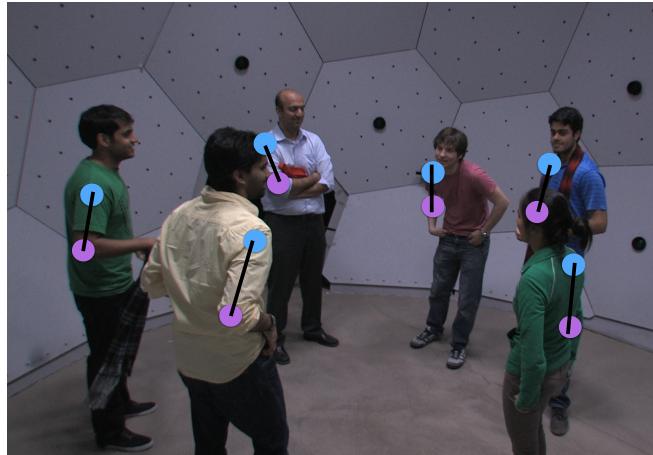


Ours

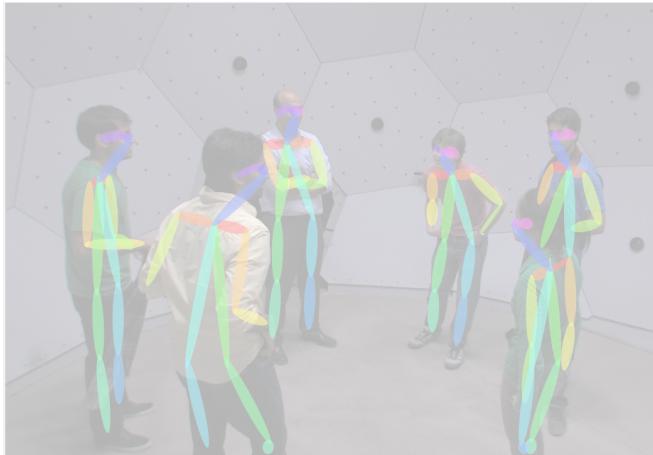
Our Method: Parts Detection + Parts Association



Top-down



Our Method: Parts Detection + Parts Association

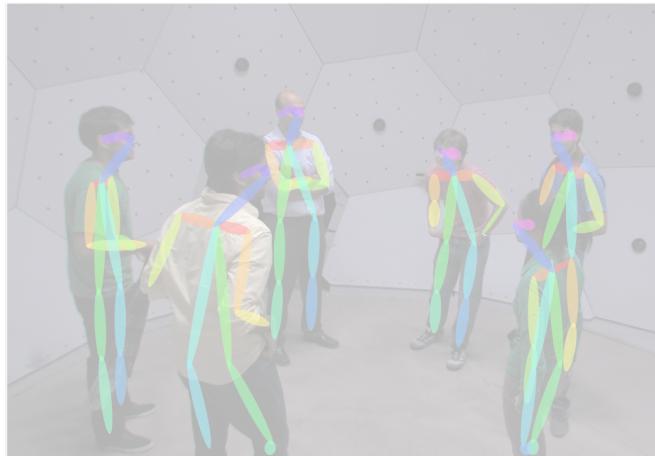


Top-down

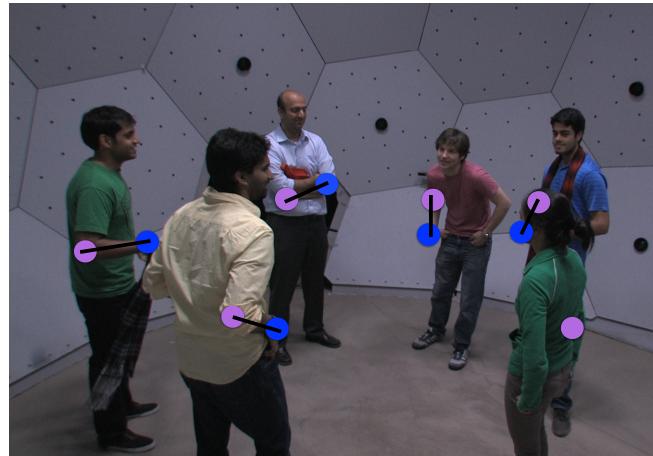


Part Affinity Fields

Our Method: Parts Detection + Parts Association



Top-down



Our Method: Parts Detection + Parts Association

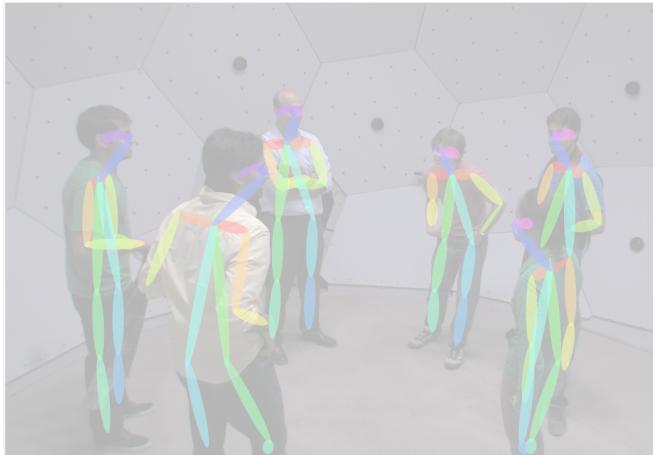


Top-down



Part Affinity Fields

Our Method: Parts Detection + Parts Association

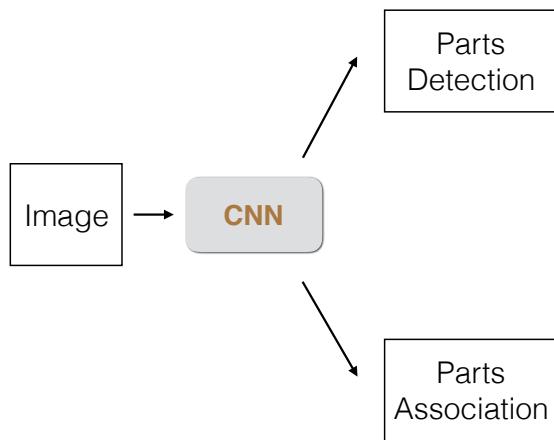


Top-down

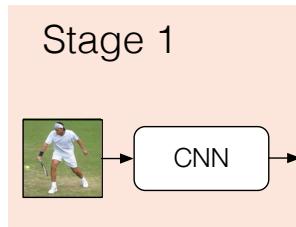


Ours

Novelty: Jointly Learning Parts Detection and Parts Association



Sequential Prediction with Learned Spatial Context



Right shoulder



Right wrist

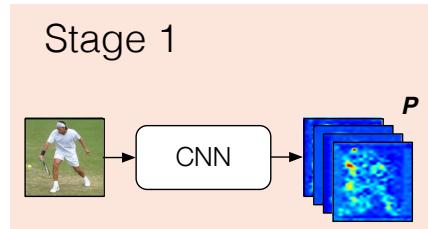


Right knee

:

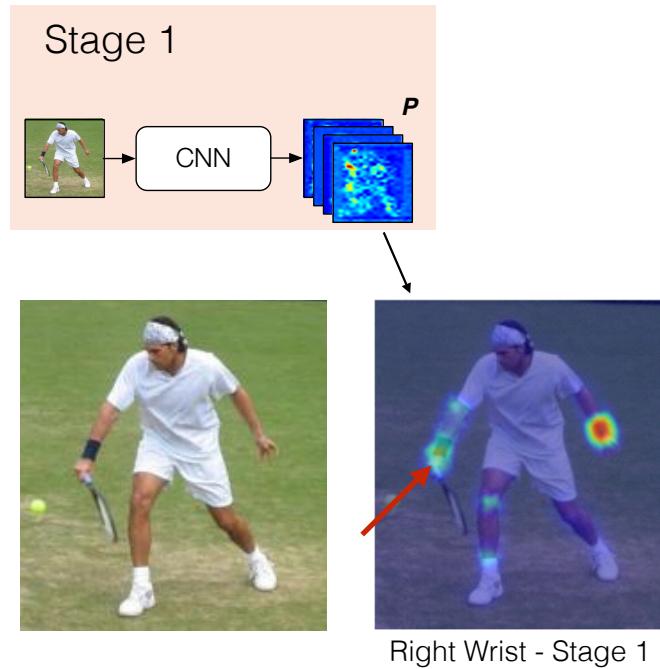
Convolutional Pose Machines, Wei, Ramakrishna, Kanade, Sheikh, CVPR 2016

Sequential Prediction with Learned Spatial Context

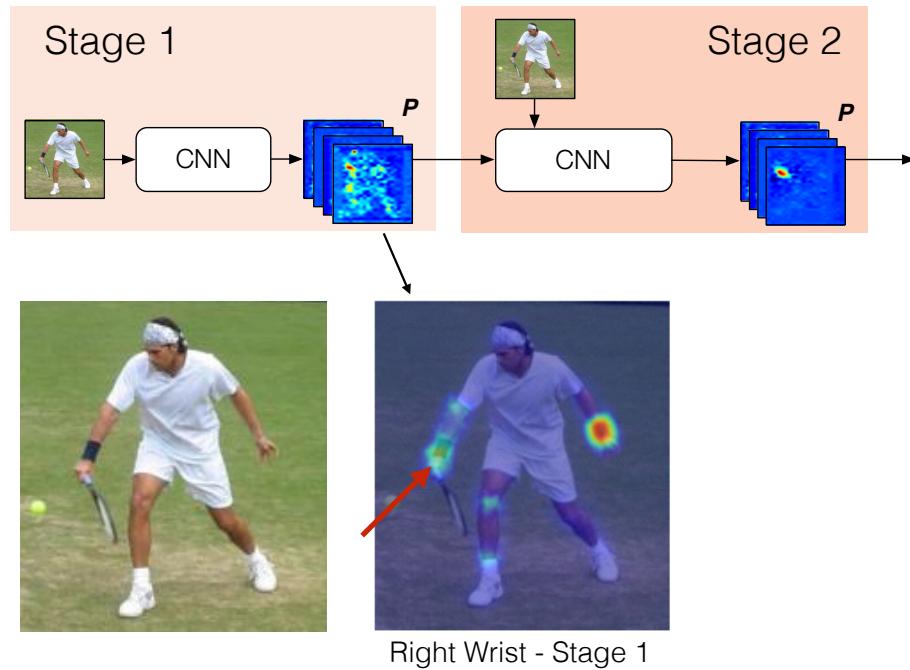


Convolutional Pose Machines, Wei, Ramakrishna, Kanade, Sheikh, CVPR 2016

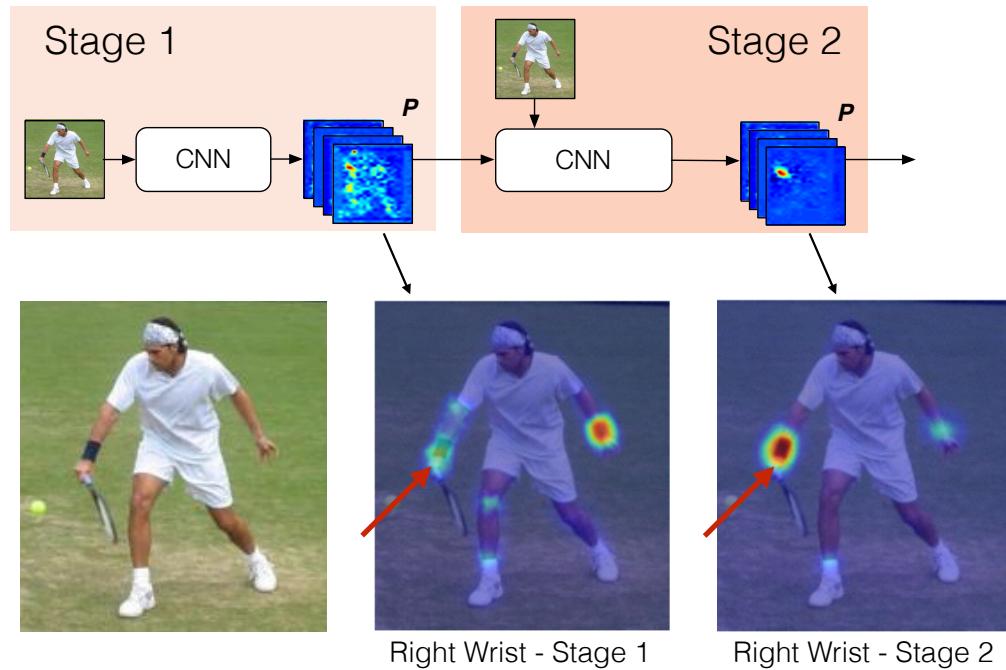
Sequential Prediction with Learned Spatial Context



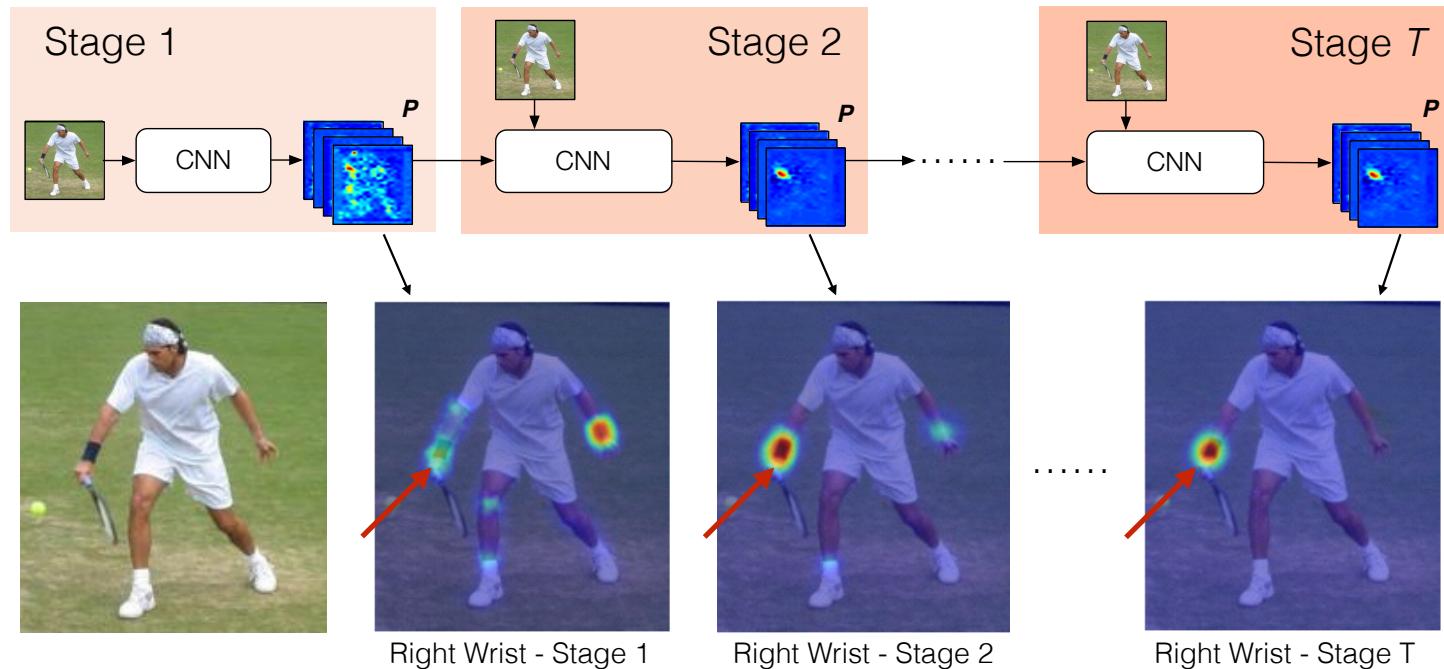
Sequential Prediction with Learned Spatial Context



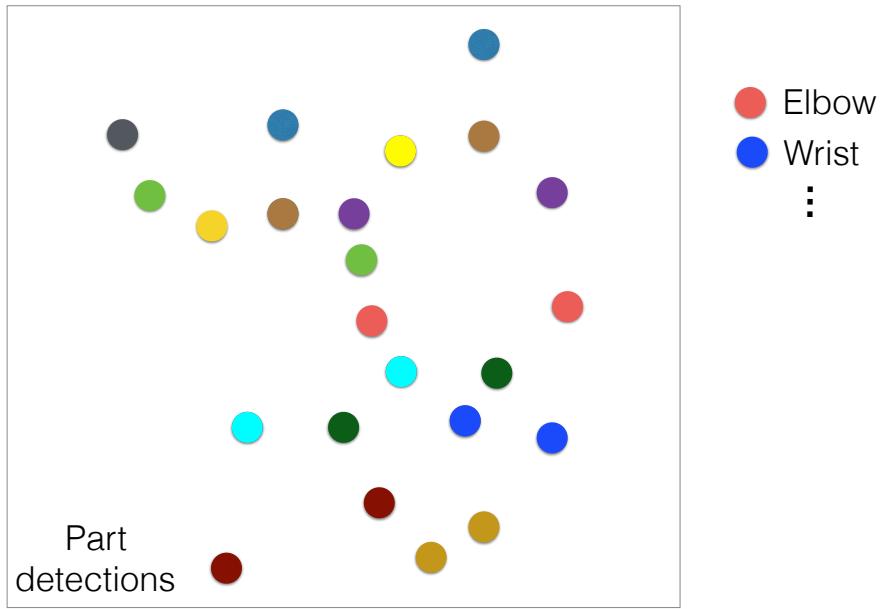
Sequential Prediction with Learned Spatial Context



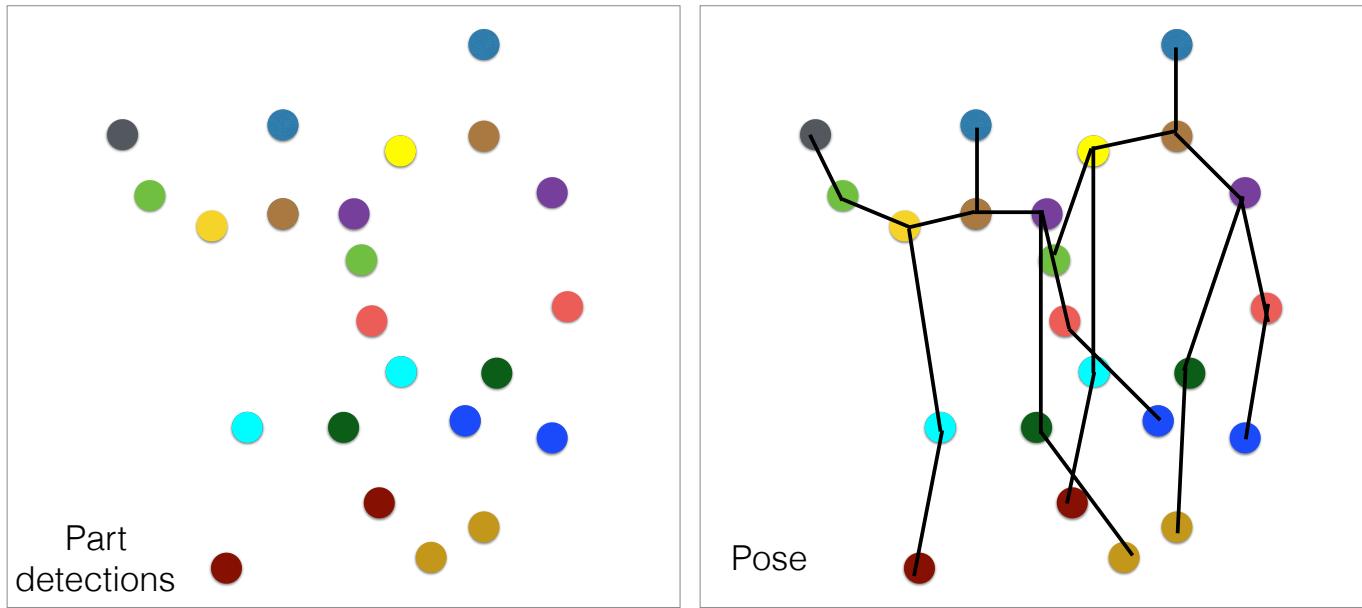
Sequential Prediction with Learned Spatial Context



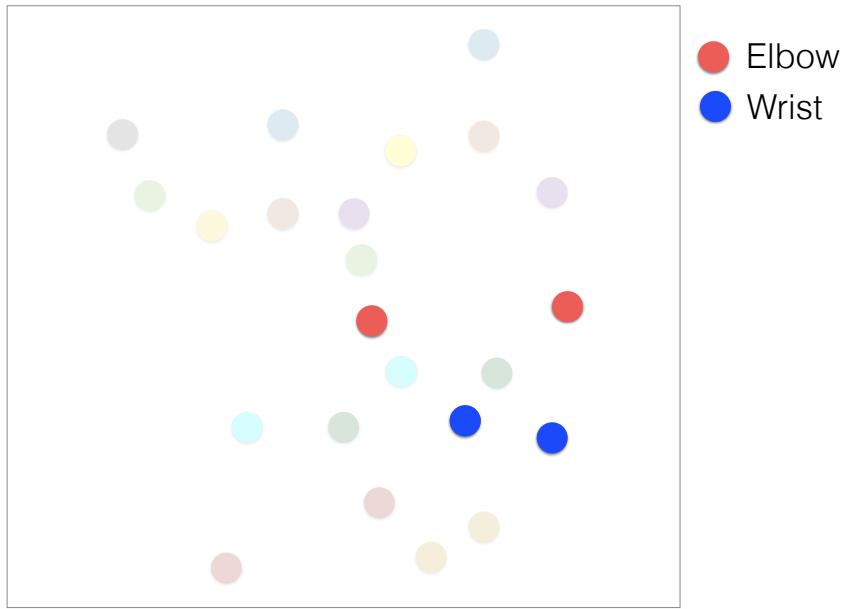
Part-Person Association for Multi-Person Pose Estimation



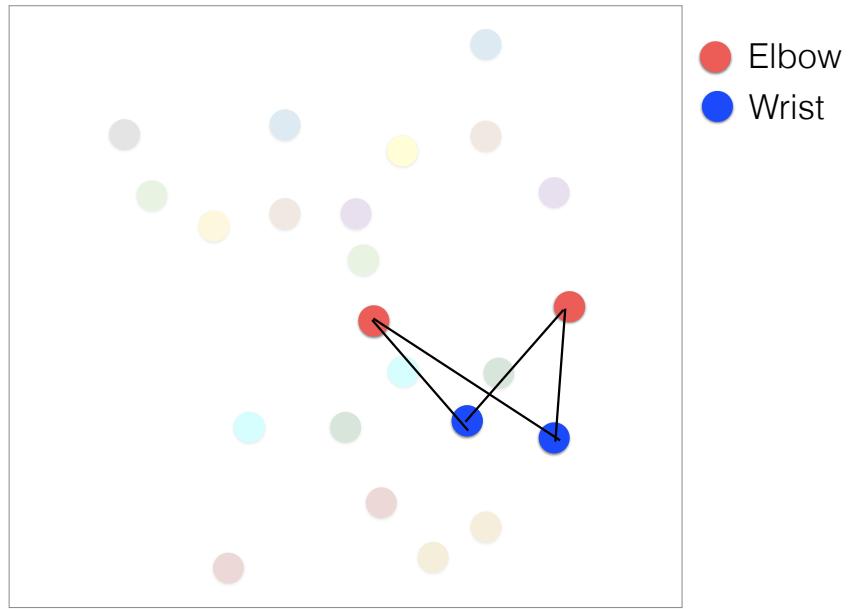
Part-Person Association for Multi-Person Pose Estimation



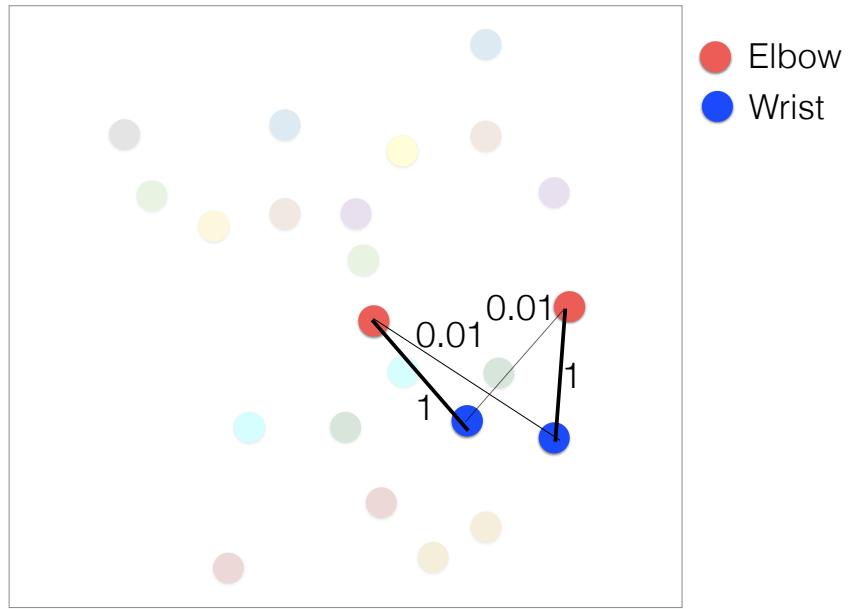
Part-to-Part Association for Multi-Person Pose Estimation



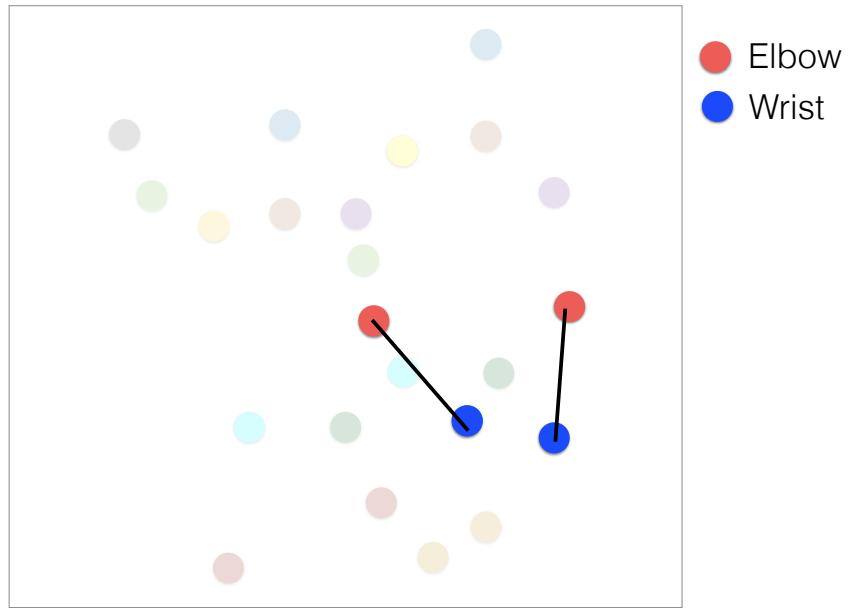
Part Affinity Score Guides the Connection



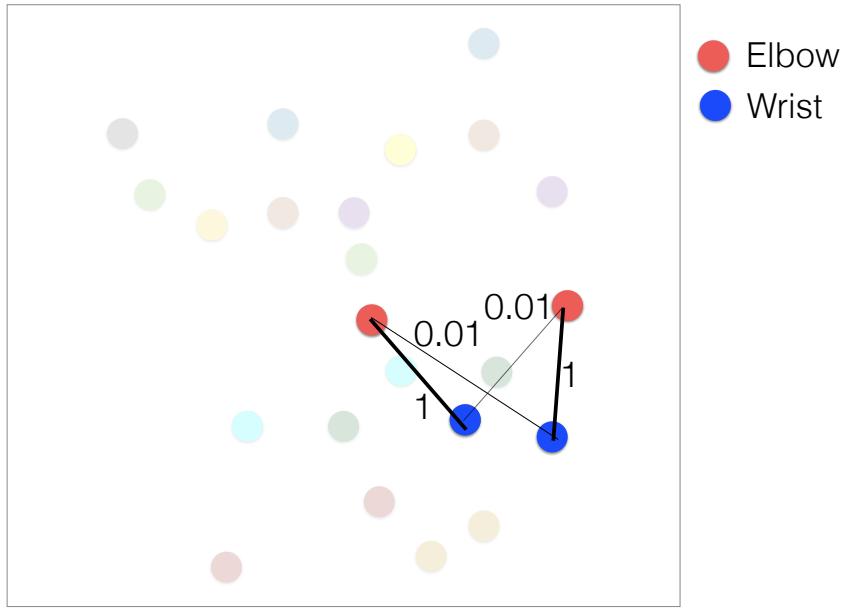
Part Affinity Score Guides the Connection



Part Affinity Score Guides the Connection



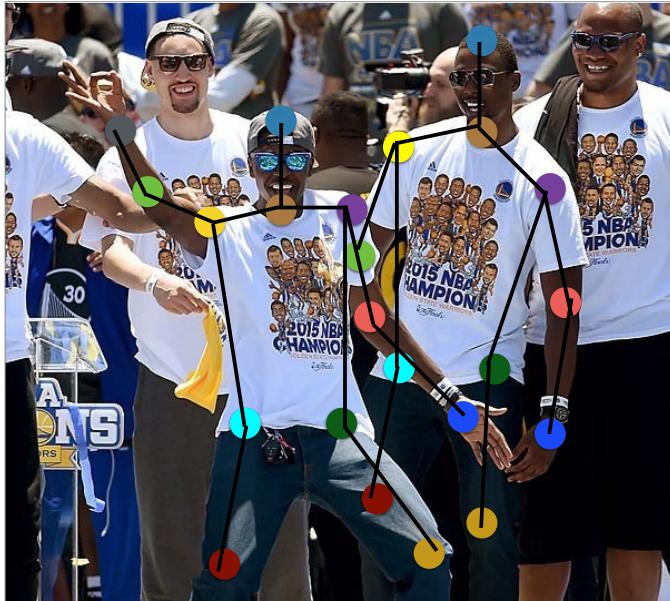
How to Obtain the Part Affinity Score



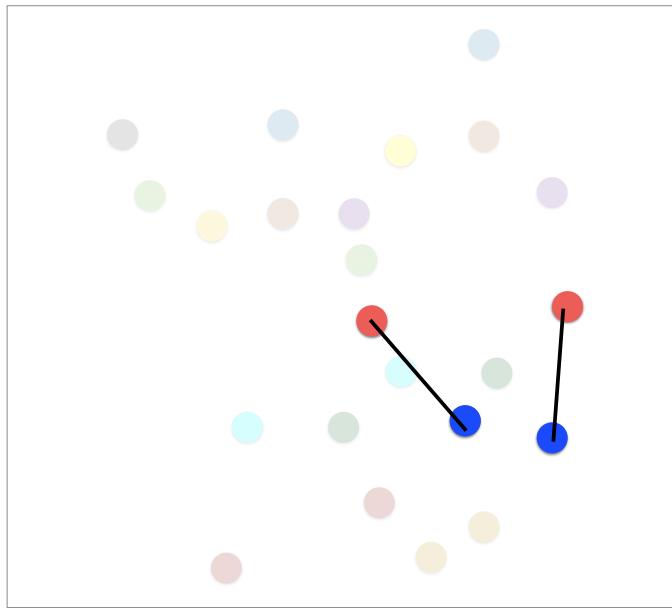
Part Affinity Score is Dependent on Visual Appearance



Part Affinity Score is Dependent on Visual Appearance

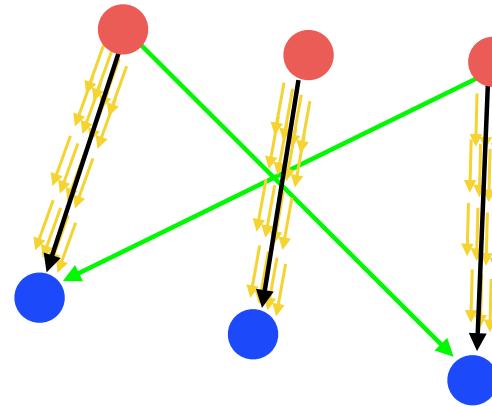


Key Idea: Encode the Part Affinity Score on the Image Plane



Part Affinity Fields
encode **direction** and **position**

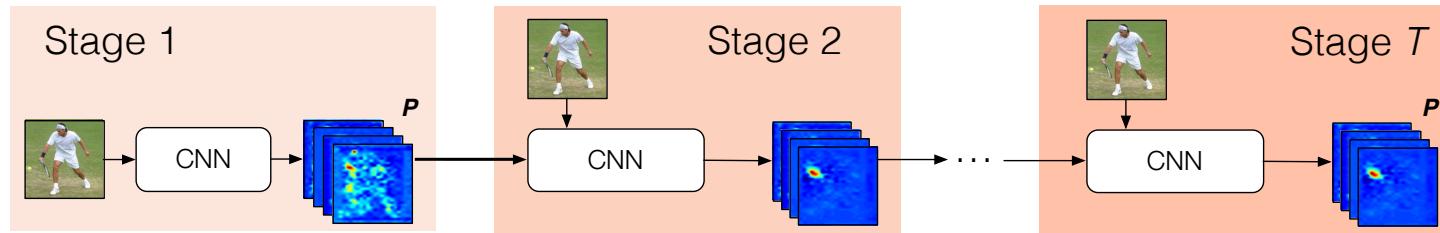
Part Affinity Fields Avoid Spatial Ambiguity



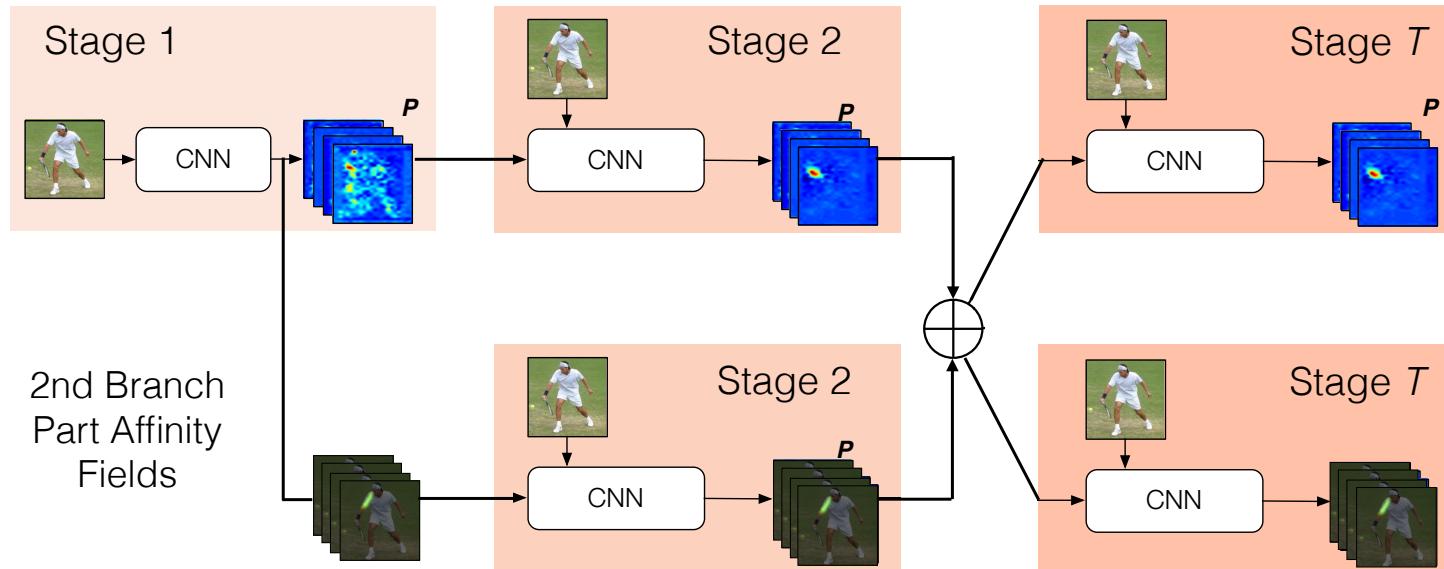
● Elbow
● Wrist

→ Correct Connection
→ Wrong Connection

Jointly Learning Parts Detection and Parts Association



Jointly Learning Parts Detection and Parts Association

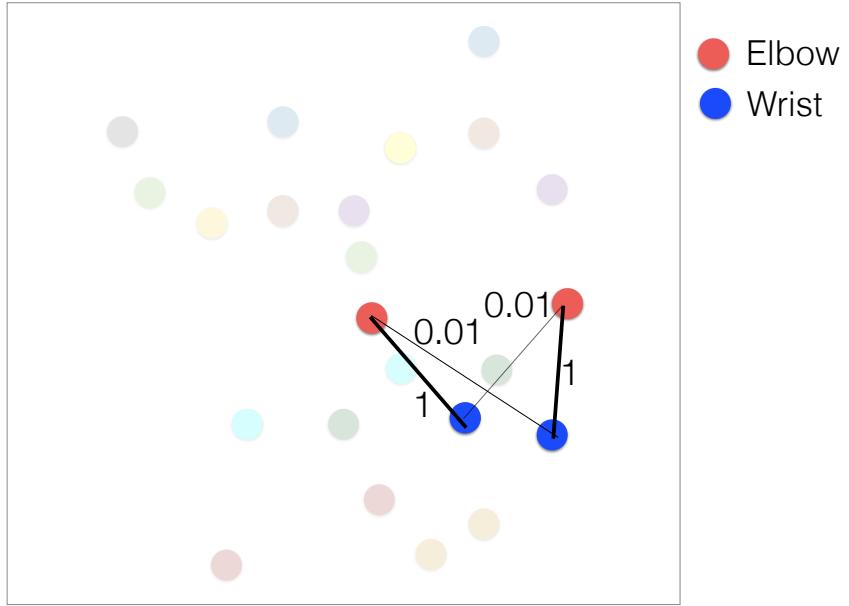




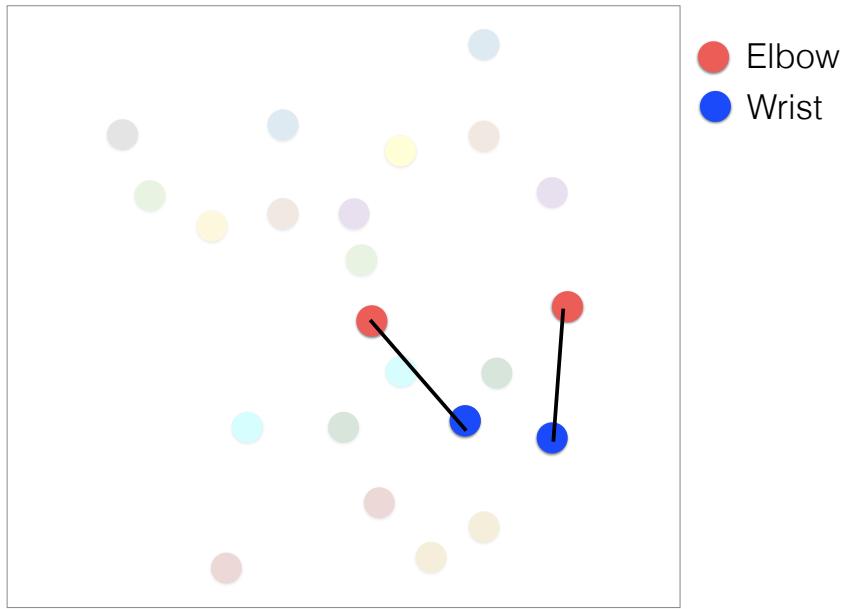
Bkg



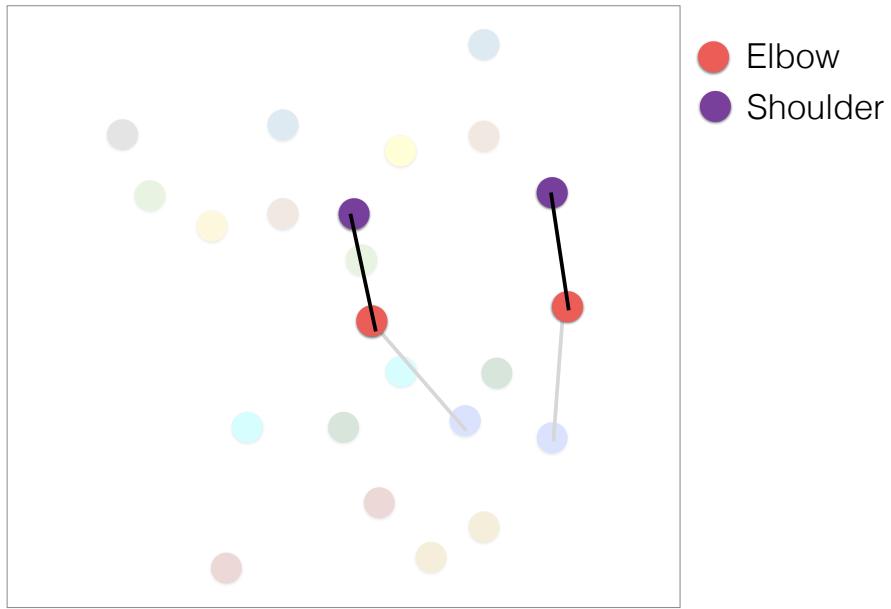
Greedy Algorithm for Body Parts Association



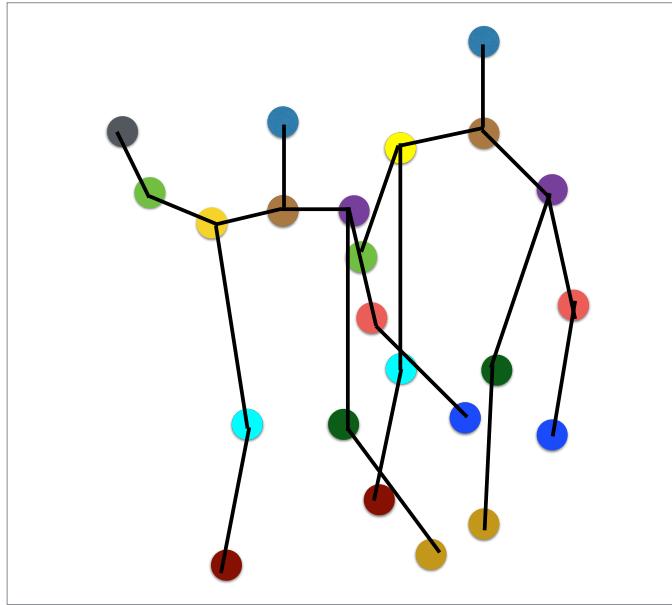
Greedy Algorithm for Body Parts Association

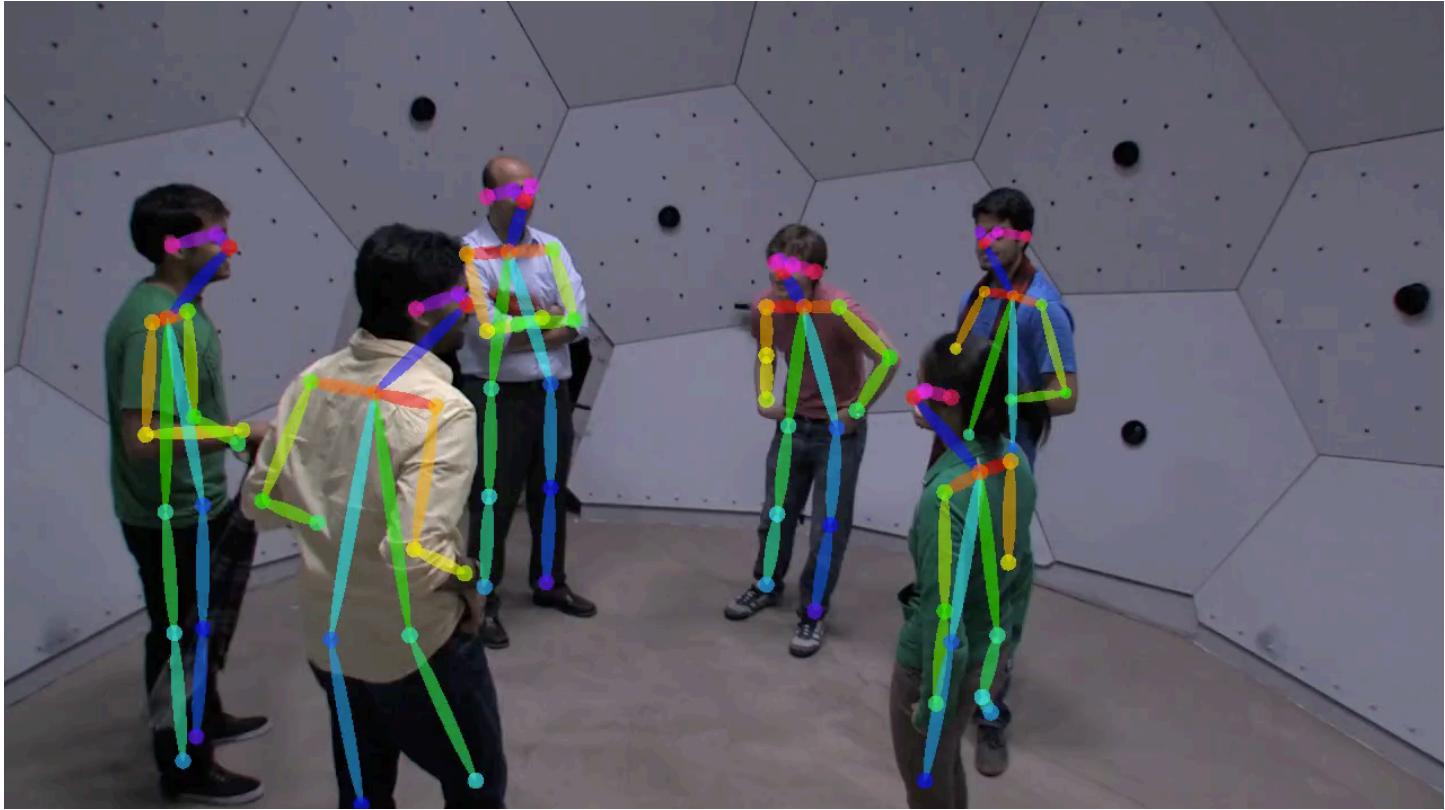


Greedy Algorithm for Body Parts Association



Greedy Algorithm for Body Parts Association





Results on COCO Challenge Validation Set

Top-down

Method	AP on val
GT bbox + CPM [1]	63
SSD [2] + CPM [1]	53
Our Method	58.5
Ours + Refinement	61

[1] Convolutional Pose Machines [Wei et al. 2016]

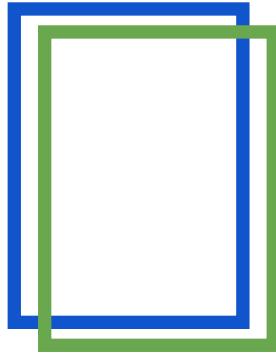
[2] SSD: Single Shot MultiBox Detector [Liu et al. 2015]



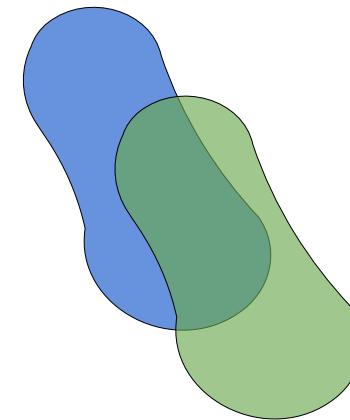
<https://youtu.be/pW6nZXeWIGM>

Evaluation: Keypoints

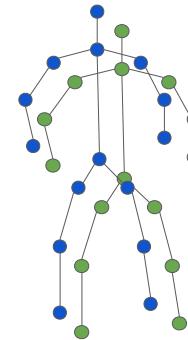
To calculate AP we need:



Bounding Box IoU



Mask IoU



Object
Keypoint
Similarity

Object Keypoint Similarity -OKS

