**CMP717
Image Processing**

**Semantic Segmentation**

Erkut Erdem
Hacettepe University
Computer Vision Lab (HUCVL)

# Semantic Segmentation

- Joint recognition & segmentation
  - segmenting all the objects in a given image and identifying their visual categories
- aka scene parsing or image parsing

- Early studies aim at segmenting out a single object of a known category
  - Borenstein & Ullman, 2002, Liebe & Schiele, 2003, etc.

- More recent work depends on CNNs
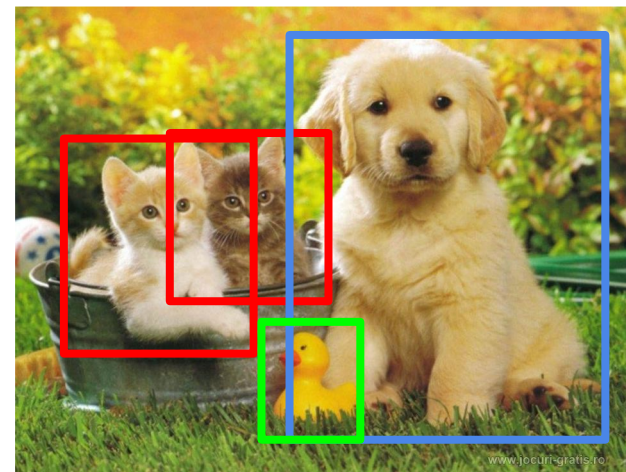  - Farabet et al., 2013, Pinheiro and Collobert, 2014, Long et al., 2015, Noh et al., 2015

# Computer Vision Tasks



| Classification | Classification + Localization | Object Detection | Semantic Segmentation |
|---|---|---|---|
| CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK |

Single object

Multiple objects

# Computer Vision Tasks

Classification | Classification + Localization | Object Detection | Semantic Segmentation
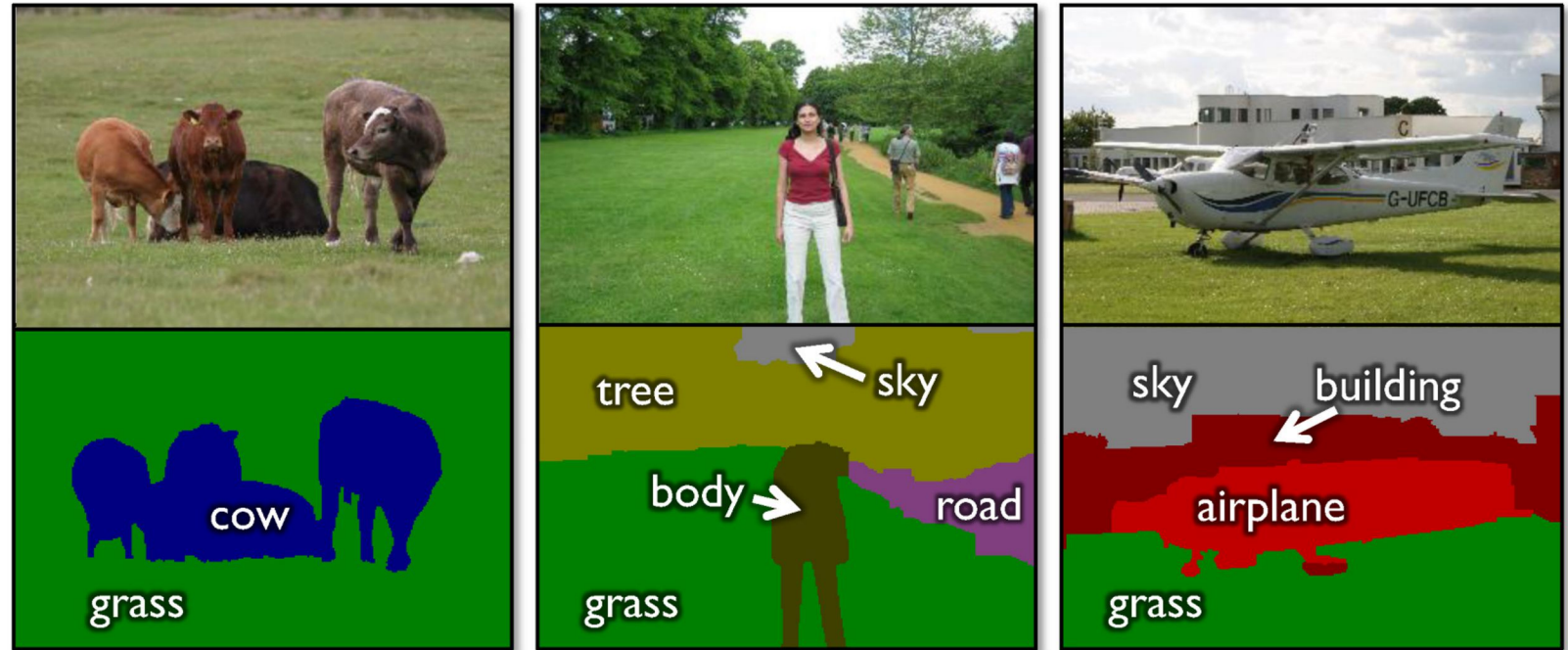


Today

F.-F. Li, A. Karpathy and J. Johnson

# Semantic Segmentation

Label every pixel in the image with a category label

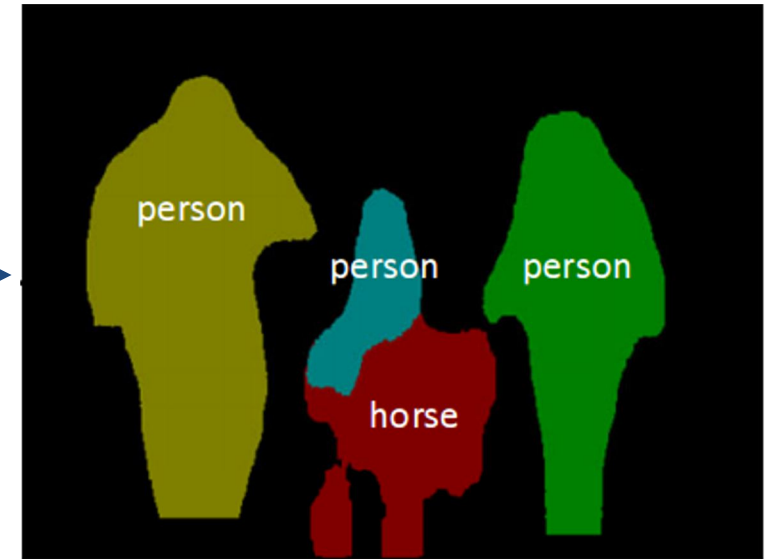Don't differentiate instances (cows)

Classic computer vision problem



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

F.-F. Li, A. Karpathy and J. Johnson

# Instance Segmentation

Detect instances, give category, label pixels

"simultaneous detection and segmentation" (SDS)

Lots of recent work (MS-COCO)

F.-F. Li, A. Karpathy and J. Johnson

# Early Studies of Semantic Segmentation

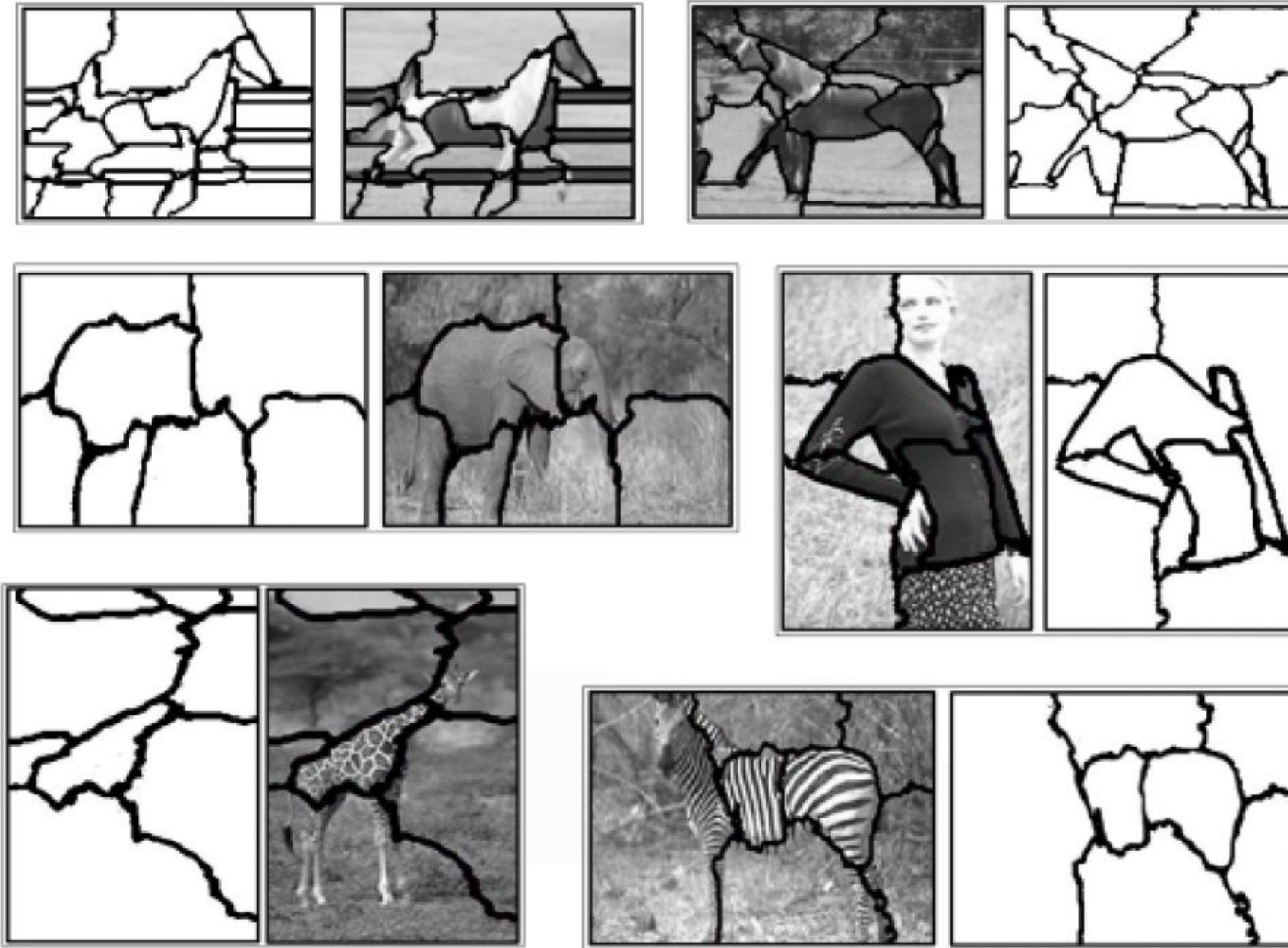- Given an image and object category, to segment the object



Cow Image

Object Category Model → Segmentation

Segmented Cow

- Segmentation should (ideally) be
  - shaped like the object e.g. cow-like
  - obtained efficiently in an unsupervised manner
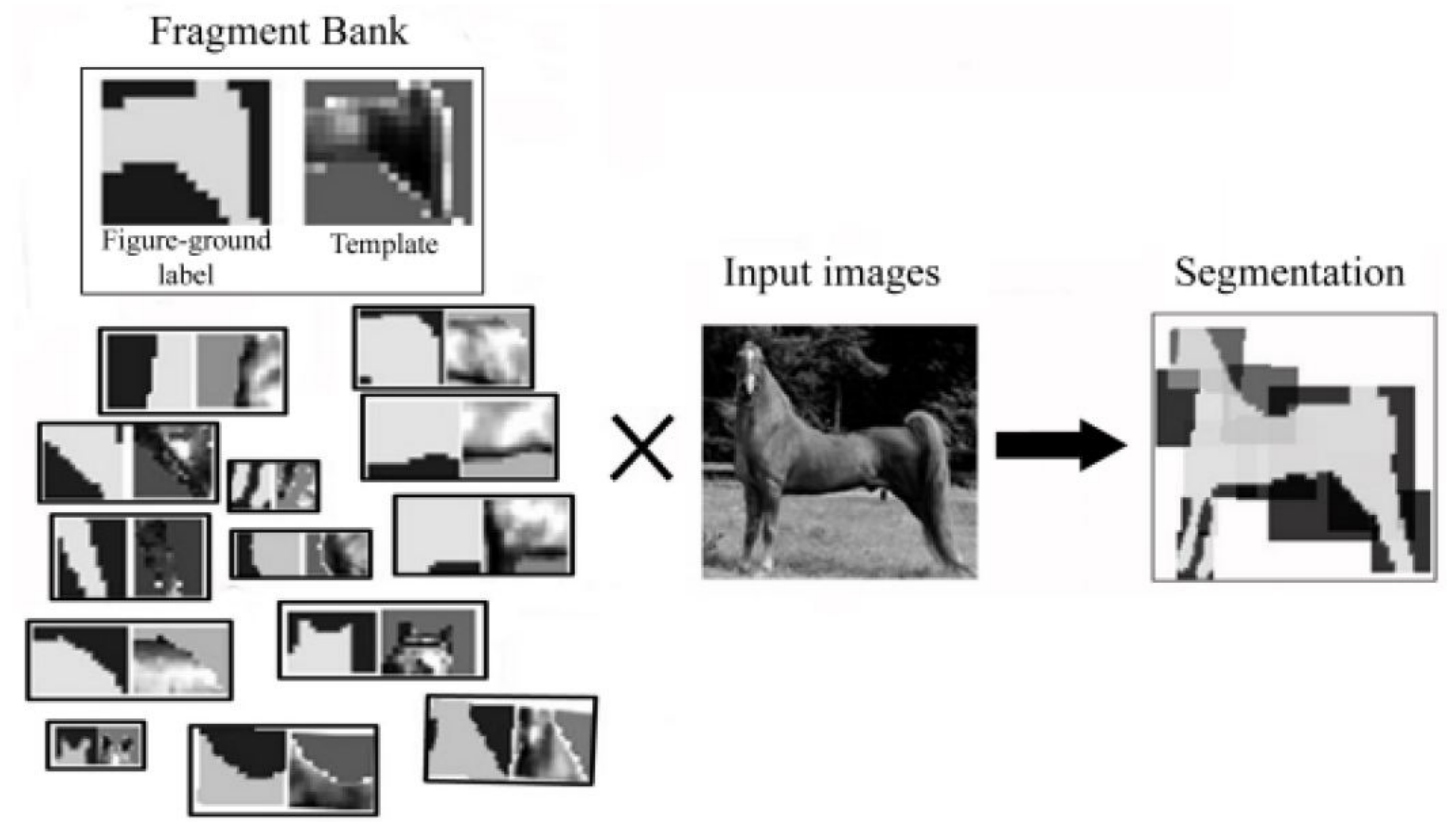  - able to handle self-occlusion

M. P. Kumar

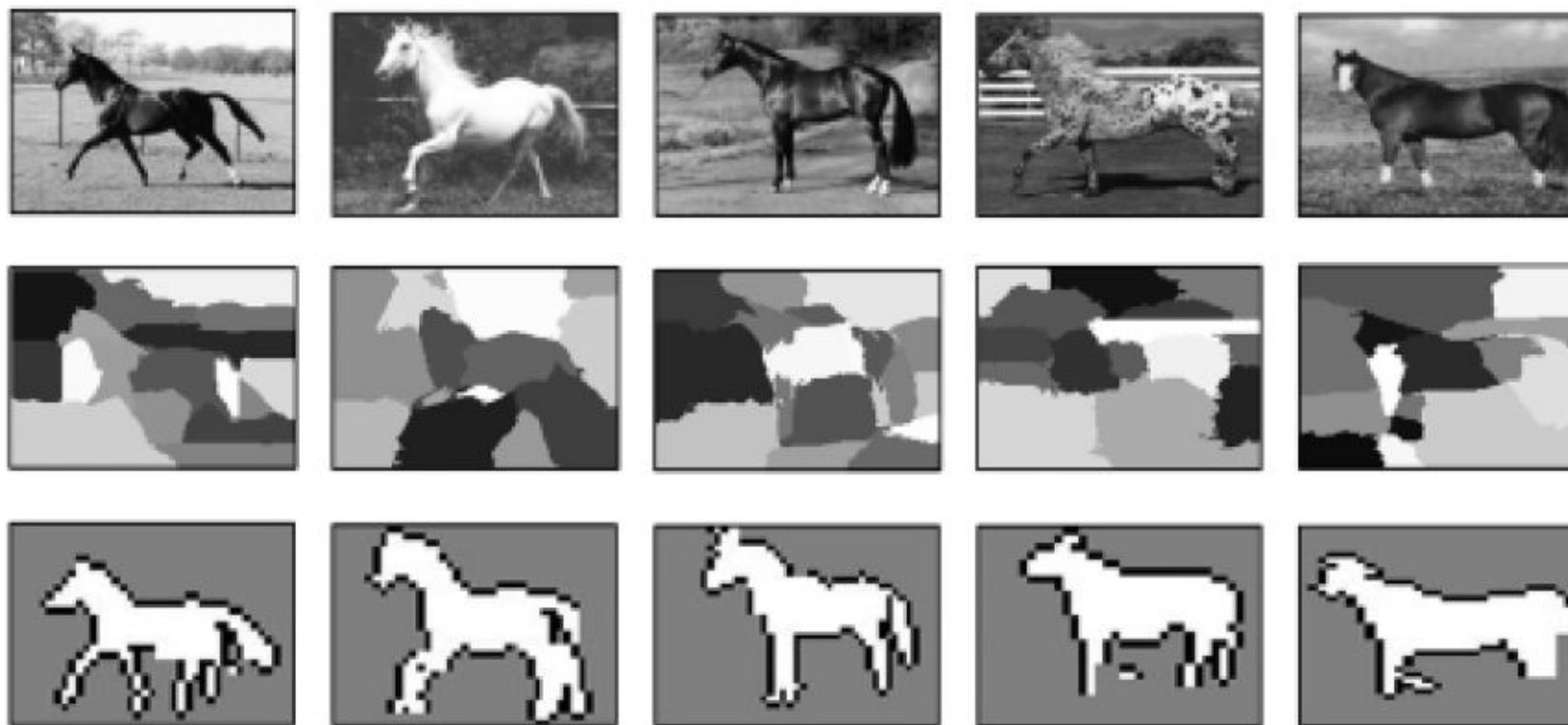# Examples of Bottom-up Segmentation

Using Normalized Cuts, Shi & Malik, 1997
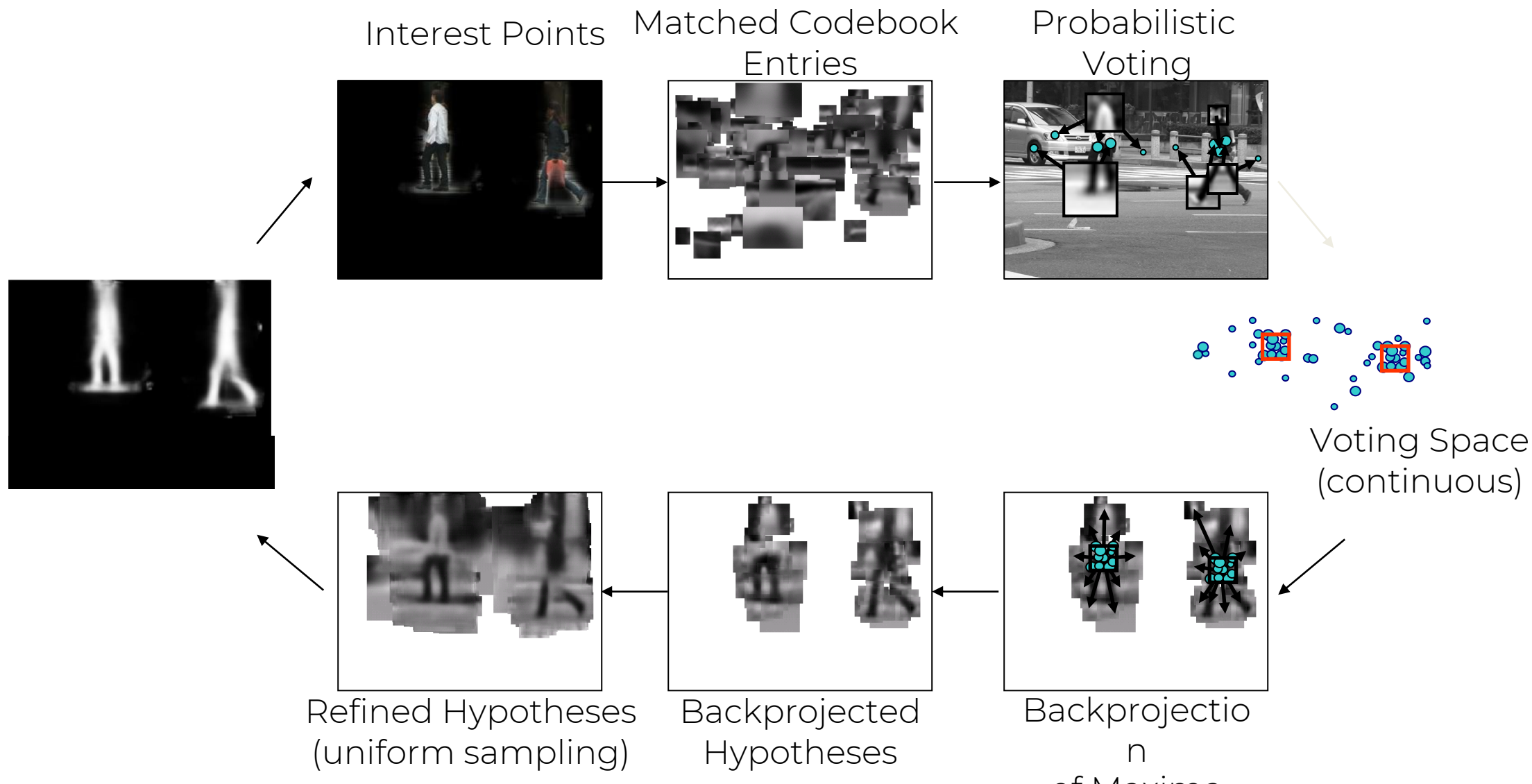
# Jigsaw approach: Borenstein and Ullman, 2002



Fragment Bank

Figure-ground label | Template

Input images × Segmentation

R. Fergus

# Jigsaw approach: Borenstein and Ullman, 2002

# Implicit Shape Model - Liebe and Schiele, 2003



Interest Points

Matched Codebook Entries

Probabilistic Voting

Voting Space (continuous)

Refined Hypotheses (uniform sampling)

Backprojected Hypotheses

Backprojection of Maxima

R. Fergus

# Random Fields for segmentation

I = Image pixels (observed)
h = foreground/background labels (hidden) – one label per pixel
θ = Parameters

$$\underbrace{p(h\,|\,I,\theta)}_{\text{Posterior}}$$

# Random Fields for segmentation

I = Image pixels (observed)
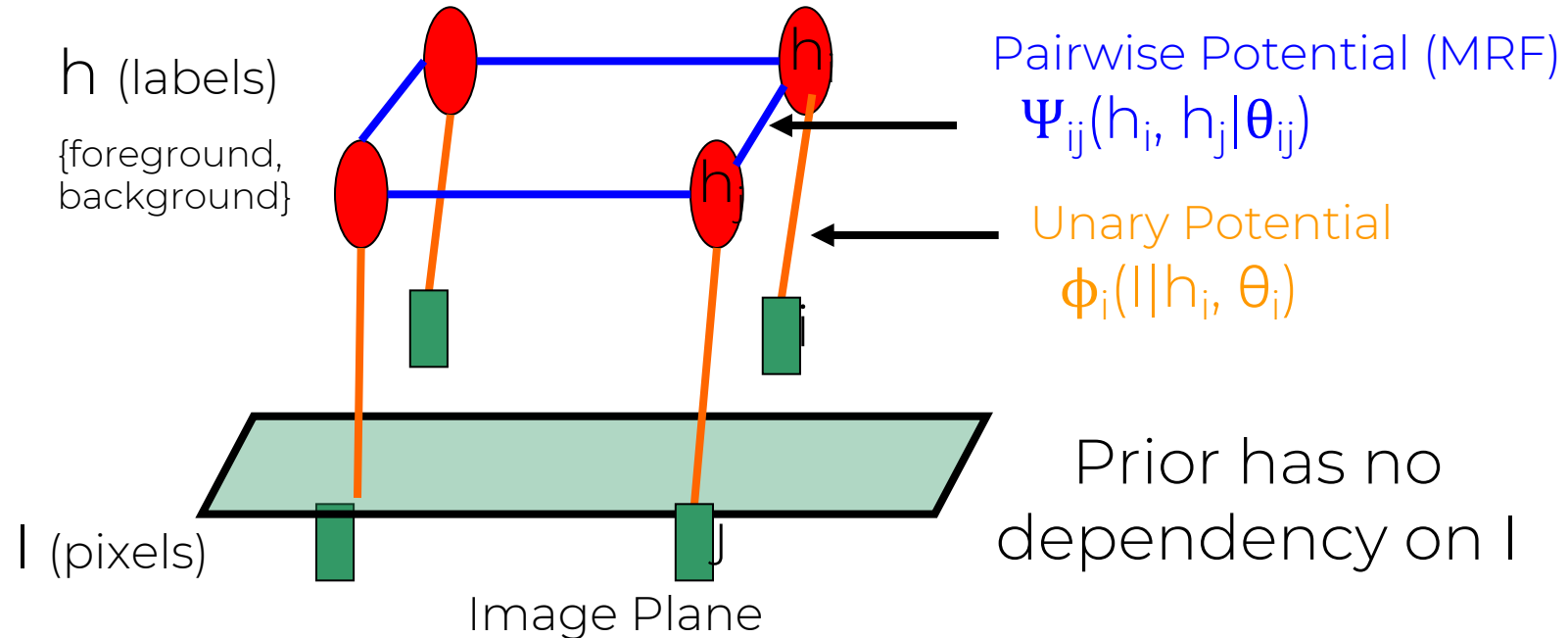h = foreground/background labels (hidden) – one label per pixel
θ = Parameters

$$\underbrace{p(h\,|\,I,\theta)}_{\text{Posterior}} \propto \underbrace{p(I,h\,|\,\theta)}_{\text{Joint}} = \underbrace{p(I\,|\,h,\theta)}_{\text{Likelihood}}\underbrace{p(h\,|\,\theta)}_{\text{Prior}}$$

1. Generative approach models joint
   → Markov random field (MRF)

2. Discriminative approach models posterior directly
   → Conditional random field (CRF)

# Generative Markov Random Field

$$p(h, I \mid \theta) = \boxed{p(I \mid h, \theta)}\boxed{p(h \mid \theta)}$$

$$= \frac{1}{Z(\theta)}\left[\underbrace{\prod_i \phi_i(I \mid h_i, \theta_i)}_{\text{Likelihood}}\underbrace{\prod_{ij} \psi_{ij}(h_i, h_j \mid \theta_{ij})}_{\text{MRF Prior}}\right]$$

h (labels)

{foreground, background}
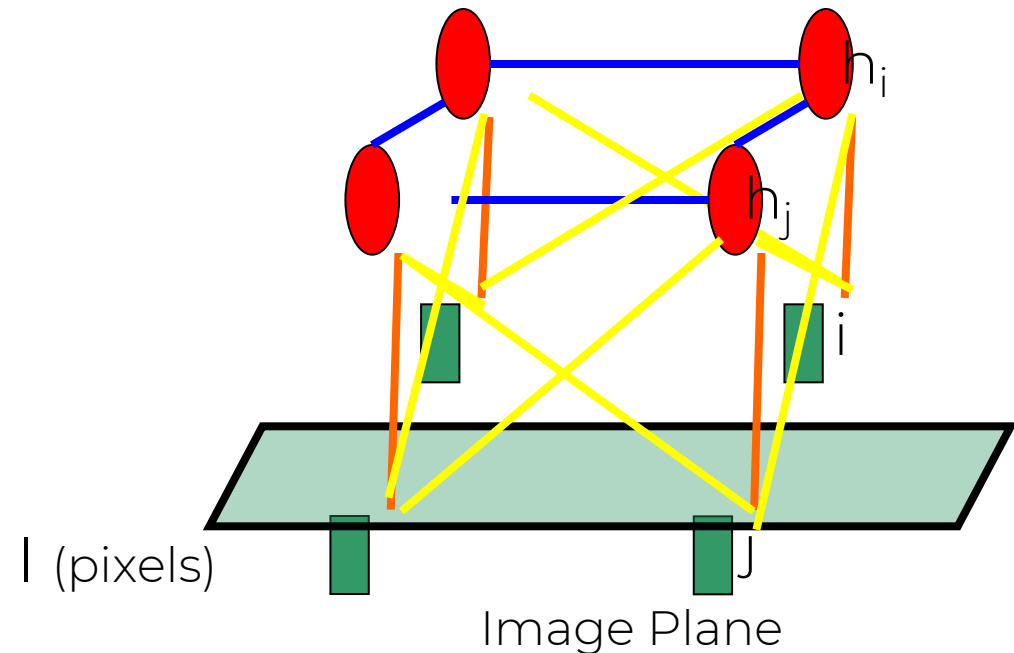
Pairwise Potential (MRF)
$\Psi_{ij}(h_i, h_j \mid \theta_{ij})$

Unary Potential
$\phi_i(I \mid h_i, \theta_i)$

Prior has no dependency on I

I (pixels)

Image Plane

# Conditional Random Field

Discriminative approach

Lafferty, McCallum and Pereira 2001

$$p(h \mid I, \theta) = \frac{1}{Z(I, \theta)} \left[ \prod_i \underbrace{\phi_i(h_i, I \mid \theta_i)}_{\text{Unary}} \prod_{ij} \underbrace{\psi_{ij}(h_i, h_j, I \mid \theta_{ij})}_{\text{Pairwise}} \right]$$

- Dependency on I allows introduction of pairwise terms that make use of image.

- For example, neighboring labels should be similar only if pixel colors are similar → Contrast term

e.g Kumar and Hebert 2003



I (pixels)

Image Plane

R. Fergus

# Conditional Random Fields for Segmentation

- Segmentation map x
- Image I

$$E(x; I) = \nu \sum_{i,j} w_{ij} |x(i) - x(j)| + \sum_{k} \lambda_k |x - x_{F_k, I}|$$
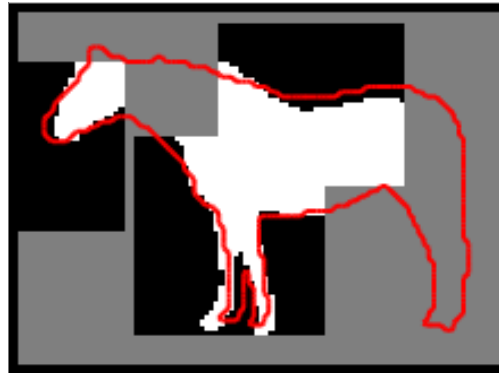
Low-level pairwise term

High-level local term

Pixel-wise similarity

# Levin & Weiss [ECCV 2006]

$$E(h; I) = \sum_i \lambda_i \left| h - h_{F_i, I} \right| + \sum_{ij} w(i, j) \left| h_i - h_j \right|$$

Consistency with fragments segmentation
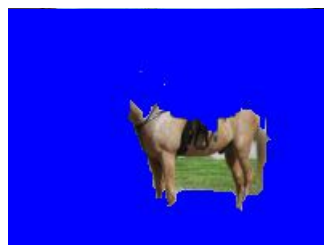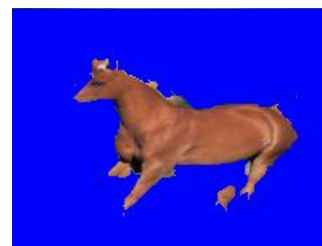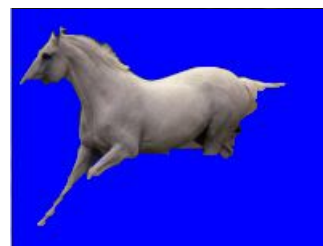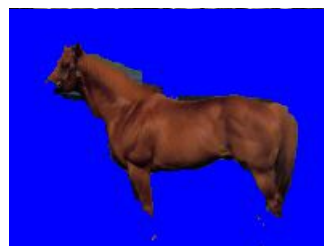
Segmentation alignment with image edges



Resulting min-cut segmentation

R. Fergus

# Levin & Weiss [ECCV 2006]

Levin & Weiss [ECCV 2006]

# Semantic Segmentation
## Joint Object recognition & segmentation
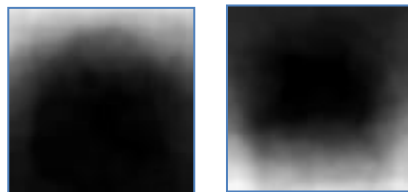


Sky

Building

Tree

Grass

$$E(x,\omega) = \sum_i \theta_i (\omega, x_i) + \sum_i \theta_i (x_i) + \sum_i \theta_i ( x_i) + \sum_{i,j} \theta_{ij} (x_i,x_j)$$

(color)   (location)   (class)   (edge aware Ising prior)
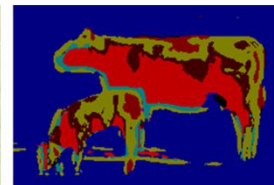
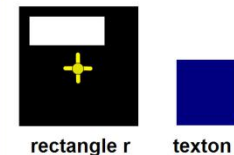$x_i \in \{1,...,K\}$ for K object classes

Location



sky          grass

Class (boosted textons)



(a) Input image    (b) Texton map    (c) Feature pair = (r,t)    (d) Superimposed rectangles

rectangle r    texton t

[TextonBoost; Shotton et al, '06]

C. Rother

# Semantic Segmentation
## Joint Object recognition & segmentation



(a)  (b) 69.6%  (c) 70.3%  (d) 72.2%

Class +
location

+
edges

+
color

[TextonBoost; Shotton et al, '06]

C. Rother

# Semantic Segmentation
## Joint Object recognition & segmentation

Good results …



[TextonBoost; Shotton et al, '06]

C. Rother

# Semantic Segmentation
## Joint Object recognition & segmentation

Failure cases...

C. Rother

# Deep Semantic Segmentation

# Deep Semantic Segmentation

Extract patch

# Deep Semantic Segmentation

Extract patch

Run through a CNN

CNN

# Deep Semantic Segmentation

Extract patch

Run through a CNN

Classify center pixel



CNN

COW

# Deep Semantic Segmentation

Extract patch

Run through a CNN

Classify center pixel



CNN

COW

# Deep Semantic Segmentation

Extract patch

Run through a CNN

Classify center pixel



CNN → COW

CNN

# Deep Semantic Segmentation



Extract patch → Run through a CNN → Classify center pixel

CNN → COW

CNN → GRASS

# Deep Semantic Segmentation



Extract patch

Run through a CNN

Classify center pixel

CNN

COW

Repeat for every pixel

cow

grass

F.-F. Li, A. Karpathy and J. Johnson

# Deep Semantic Segmentation

Extract patch

Run through a CNN

Classify center pixel



CNN → COW

Repeat for every pixel

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

# Deep Semantic Segmentation



Extract patch

Run through a CNN

Classify center pixel

CNN

COW

Repeat for every pixel

cow

grass

F.-F. Li, A. Karpathy and J. Johnson

# Semantic Segmentation Idea: Fully Convolutional

Run "fully convolutional"
network to get all pixels at
once



Smaller output
due to pooling

# Semantic Segmentation Idea: Fully Convolutional

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



Input:
3 x H x W

Conv

Conv

Conv

Conv

argmax

Convolutions:
D x H x W

Scores:
C x H x W

Predictions:
H x W

F.-F. Li and J. Johnson

# Semantic Segmentation Idea: Fully Convolutional

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



Input:
3 x H x W

Convolutions:
D x H x W

Scores:
C x H x W

Predictions:
H x W

Problem: convolutions at original image resolution will be very expensive ...

F.-F. Li and J. Johnson

# Semantic Segmentation Idea: Fully Convolutional

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



Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

# Semantic Segmentation Idea: Fully Convolutional

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



Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

Downsampling:
Pooling, strided convolution

Upsampling: ???

F.-F. Li and J. Johnson

# Semantic Segmentation: Upsampling



forward/inference

backward/learning

pixelwise prediction

segmentation g.t.

96  256  384  384  256  4096  4096  21

21

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

F.-F. Li, A. Karpathy and J. Johnson

# Semantic Segmentation: Upsampling



forward/inference

backward/learning

pixelwise prediction

segmentation g.t.

96

256

384

384

256

4096

4096

21

21

Learnable upsampling!

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

F.-F. Li, A. Karpathy and J. Johnson

# Semantic Segmentation: Upsampling



image   conv1   pool1   conv2   pool2   conv3   pool3   conv4   pool4   conv5   pool5   conv6-7   32x upsampled prediction (FCN-32s)

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

# Semantic Segmentation: Upsampling



image  conv1  pool1  conv2  pool2  conv3  pool3  conv4  pool4  conv5  pool5  conv6-7

32x upsampled prediction (FCN-32s)

"skip connections"

2x conv7
pool4

16x upsampled prediction (FCN-16s)

4x conv7
2x pool4
pool3

8x upsampled prediction (FCN-8s)

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

F.-F. Li, A. Karpathy and J. Johnson

# Semantic Segmentation: Upsampling



Skip connections = Better results

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

# Learnable Upsampling: "Deconvolution"

Typical 3 x 3 convolution, stride 1 pad 1

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

Typical 3 x 3 convolution, stride 1 pad 1



Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

Typical 3 x 3 convolution, stride 1 pad 1



Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

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

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling: "Deconvolution"

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



Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling: "Deconvolution"

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



Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling: "Deconvolution"

3 x 3 deconvolution, stride 2 pad 1

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

3 x 3 deconvolution, stride 2 pad 1

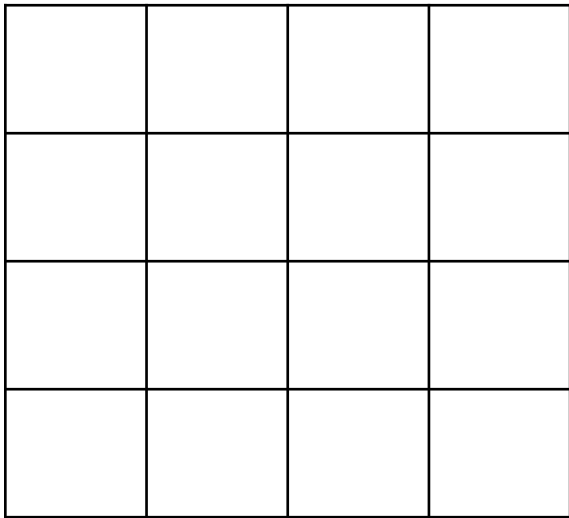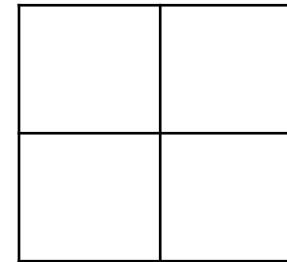Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

3 x 3 deconvolution, stride 2 pad 1

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

Sum where output overlaps

3 x 3 deconvolution, stride 2 pad 1

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

3 x 3 deconvolution, stride 2 pad 1

Sum where output overlaps

Same as backward pass for normal convolution!

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: "Deconvolution"

3 x 3 deconvolution, stride 2 pad 1

Sum where output overlaps

Same as backward pass for normal convolution!

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

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

**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

# Semantic Segmentation: Upsampling



Normal VGG

"Upside down" VGG

6 days of training on Titan X...

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

F.-F. Li, A. Karpathy and J. Johnson

# 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

# Instance Segmentation

Similar to R-CNN,
but with segments

Proposal
Generation

External
Segment
proposals



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

# Instance Segmentation

Similar to R-CNN, but with segments



Proposal Generation → External Segment proposals → Feature Extraction → Box CNN

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

# Instance Segmentation

Similar to R-CNN,
but with segments



Proposal
Generation

External
Segment
proposals

Feature
Extraction

Box
CNN

Region
CNN

Mask out background
with mean image

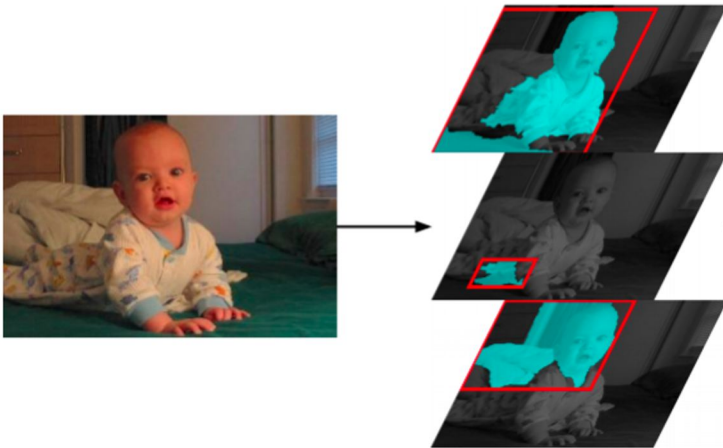Hariharan et al, "Simultaneous Detection and Segmentation", ECCV 2014

# Instance Segmentation



Similar to R-CNN,
but with segments

Proposal Generation — External Segment proposals — Feature Extraction — Box CNN / Region CNN — Region Classification

Person?
+1.8

Mask out background
with mean image

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

F.-F. Li, A. Karpathy and J. Johnson

# Instance Segmentation

Similar to R-CNN, but with segments



Proposal Generation — External Segment proposals — Feature Extraction — Region Classification — Region Refinement

Box CNN

Region CNN

Person? +1.8

Mask out background with mean image

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

F.-F. Li, A. Karpathy and J. Johnson

# Instance Segmentation: Hypercolumns



Region Classification

Region Refinement

Person?
+1.8

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

# Instance Segmentation: Hypercolumns

F.-F. Li, A. Karpathy and J. Johnson

# 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

F.-F. Li, A. Karpathy and J. Johnson

# Instance Segmentation: Cascades

Similar to
Faster R-CNN



CONVs

conv feature map

Won COCO 2015
challenge
(with ResNet)

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

# Instance Segmentation: Cascades

Similar to
Faster R-CNN

Region proposal network
(RPN)

box instances (RoIs)
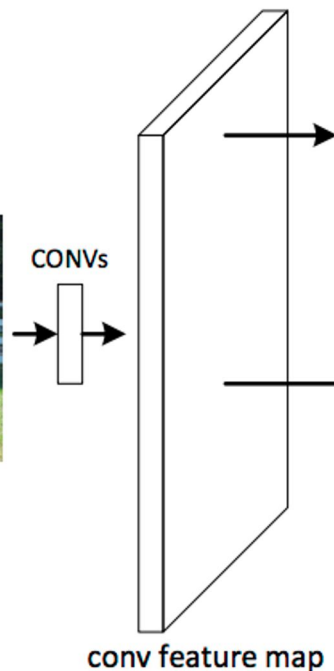
CONVs

CONVs

CONVs

conv feature map

Won COCO 2015
challenge
(with ResNet)

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

F.-F. Li, A. Karpathy and J. Johnson
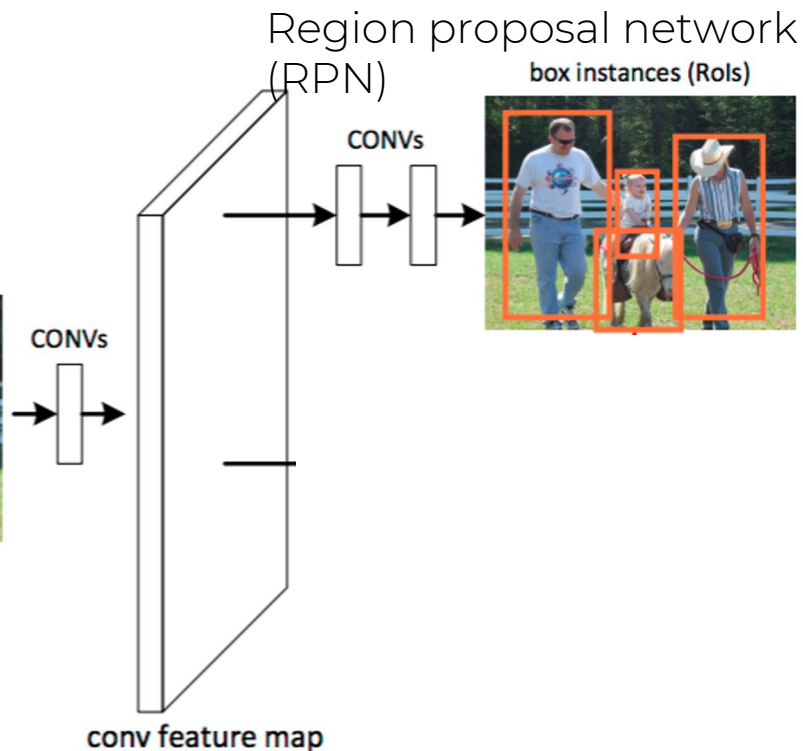
# Instance Segmentation: Cascades

Similar to
Faster R-CNN

Region proposal network
(RPN)



Reshape boxes to
fixed size,
figure / ground
logistic regression

Won COCO 2015
challenge
(with ResNet)

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

# Instance Segmentation: Cascades

Similar to
Faster R-CNN



Region proposal network
(RPN)

Reshape boxes to
fixed size,
figure / ground
logistic regression

Mask out
background, predict
object class

Won COCO 2015
challenge
(with ResNet)

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

F.-F. Li, A. Karpathy and J. Johnson

# Instance Segmentation: Cascades

Similar to
Faster R-CNN

Region proposal network
(RPN)

Reshape boxes to
fixed size,
figure / ground
logistic regression

Learn entire model
end-to-end!



Mask out
background, predict
object class

Won COCO 2015
challenge
(with ResNet)

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

F.-F. Li, A. Karpathy and J. Johnson

# Instance Segmentation: Cascades



Predictions          Ground truth

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

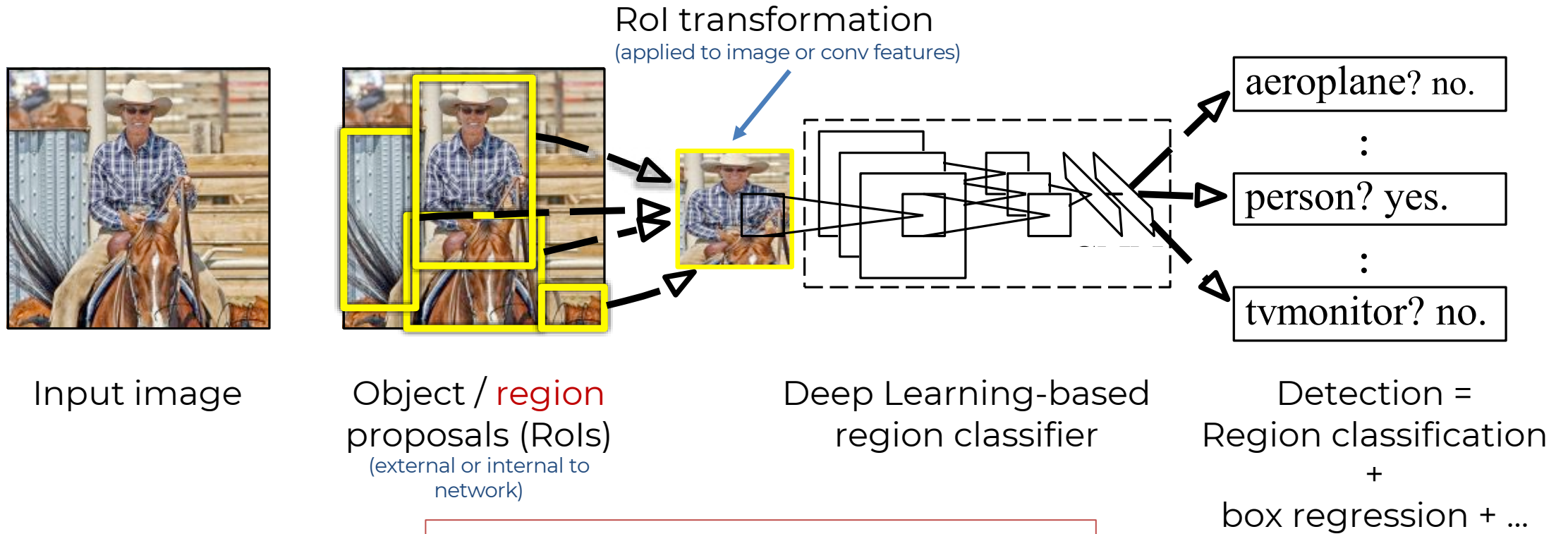F.-F. Li, A. Karpathy and J. Johnson

# Mask R-CNN: Model Overview

R-CNN-style detection system

1. Backbone architecture
2. Feature Pyramid Network (FPN)
3. Region Proposal Network (RPN)
4. Region of interest feature alignment (RoIAlign)
5. Multi-task network head
   - Box classifier
   - Box regressor
   - Mask predictor
   - Keypoint predictor

Modular composition of many recent ideas
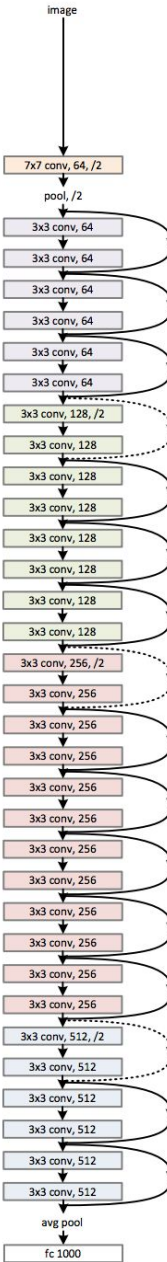
R. Girshick

# 0. R-CNN-style Approach to Object Detection



RoI transformation
(applied to image or conv features)

Input image

Object / region proposals (RoIs)
(external or internal to network)

Deep Learning-based region classifier

Detection =
Region classification
+
box regression + …

aeroplane? no.
:
person? yes.
:
tvmonitor? no.

General formula for
Region-based CNN models

Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014
Uijlings, van de Sande, Gevers, Smeulders. Selective Search for Object Recognition. IJCV 2013

R. Girshick

# 1. Backbone ConvNet
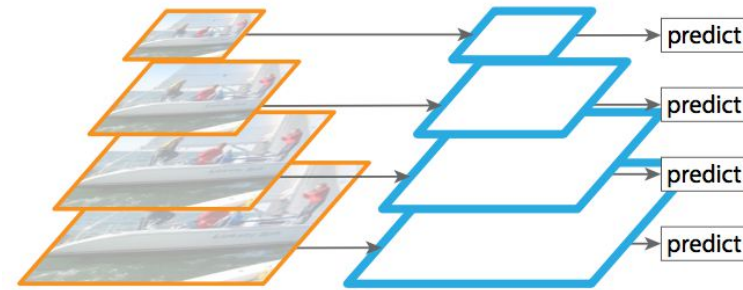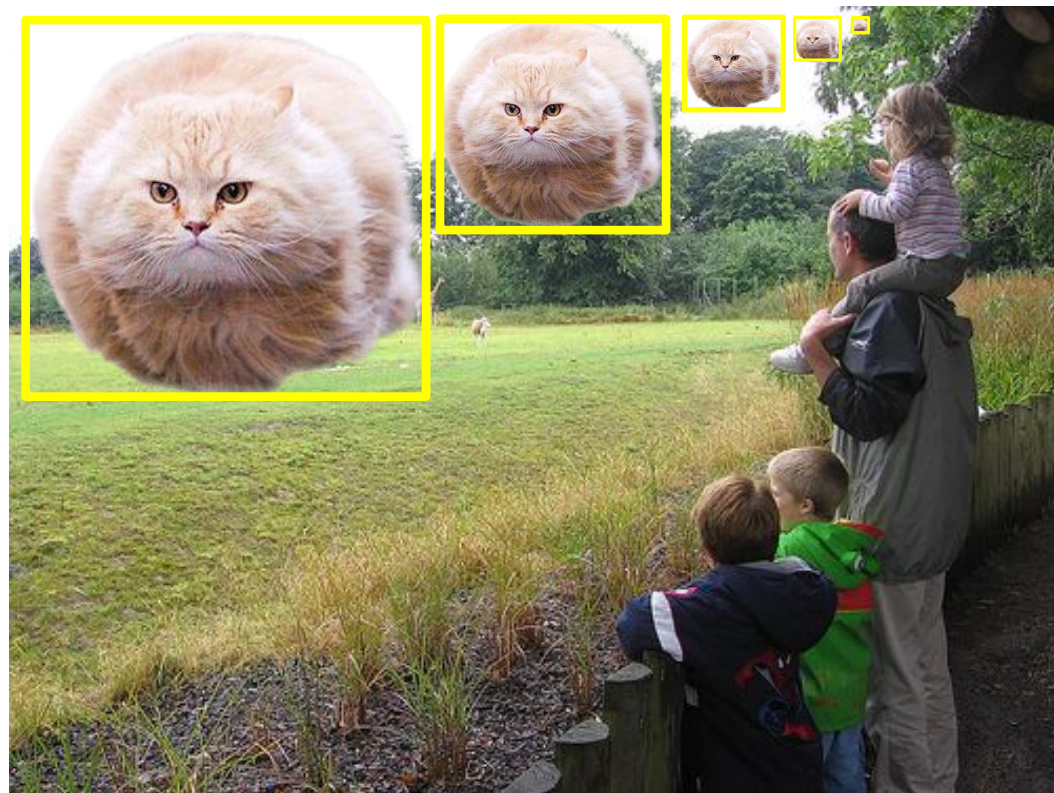
Use any standard ConvNet as a "backbone architecture"

- AlexNet, VGG, ResNet, Inception, Inception-ResNet, ResNeXt, …

- Use "same" padding everywhere (preserves integer scales)

- Prefer fully convolutional networks (ignoring cls head; see next slide)

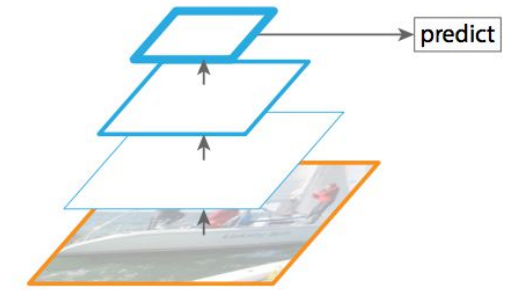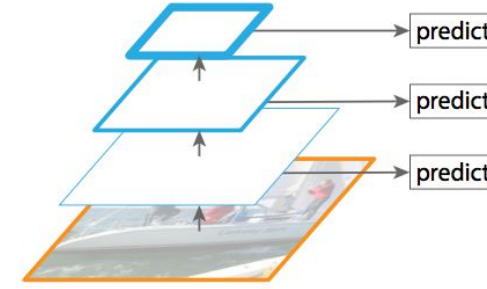- Pre-train on ImageNet classification (or similar)



R. Girshick

# 2. Scale Invariant Detection



(a) Featurized image pyramid
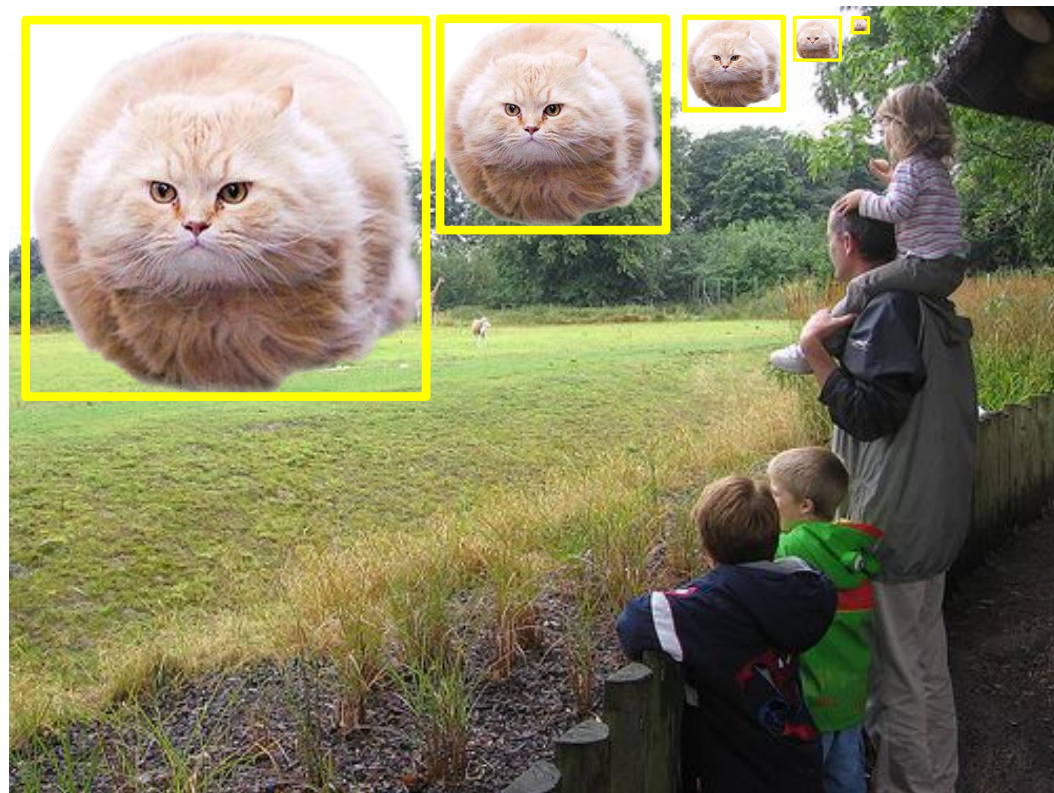SLOW!
(Viola & Jones, HOG, DPM, multiscale Fast R-CNN, ...)

(b) Single feature map
Do nothing; give up; Fast
(Fast R-CNN, YOLO, ...)

(c) Pyramidal feature hierarchy
Fast!
(≈ SSD, ...)

Popular strategies for detecting objects over large scale changes
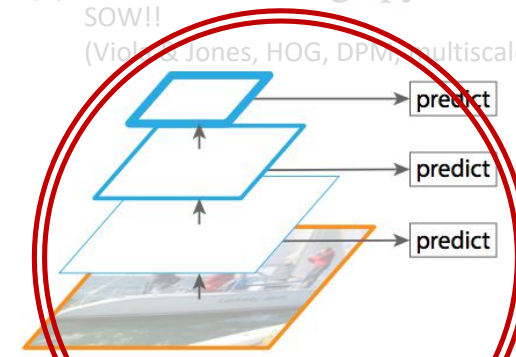
R. Girshick

# 2. Scale Invariant Detection



(a) Featurized image pyramid
SOW!!
(Viola & Jones, HOG, DPM, multiscale Fast R-CNN, ...)

(b) Single feature map
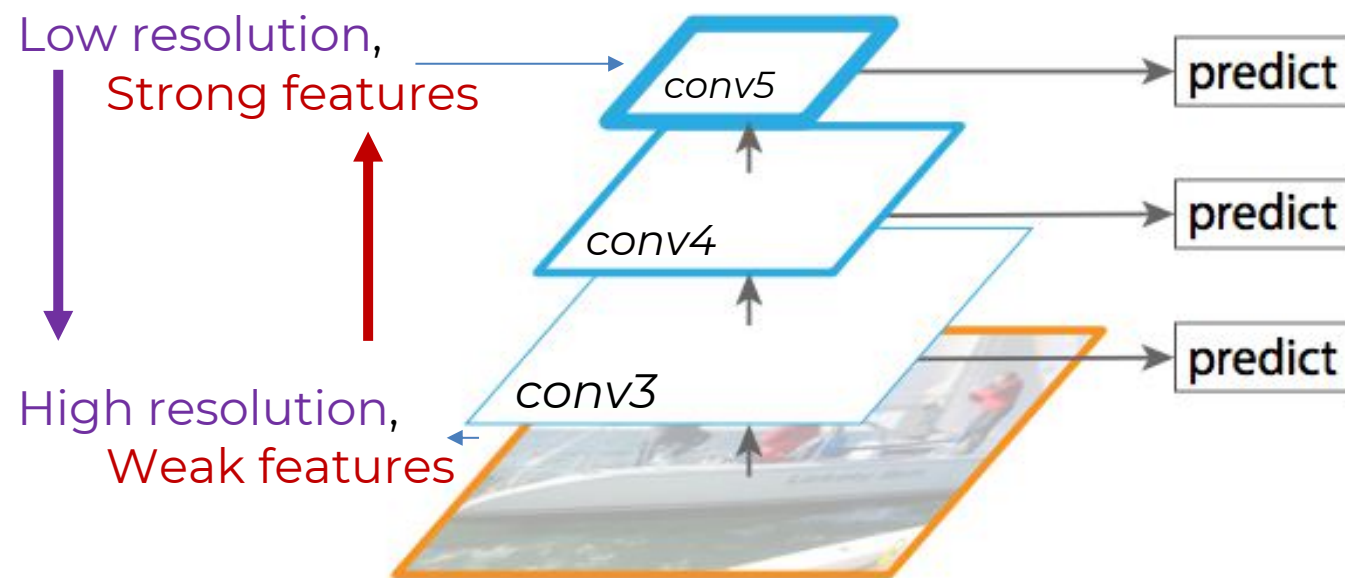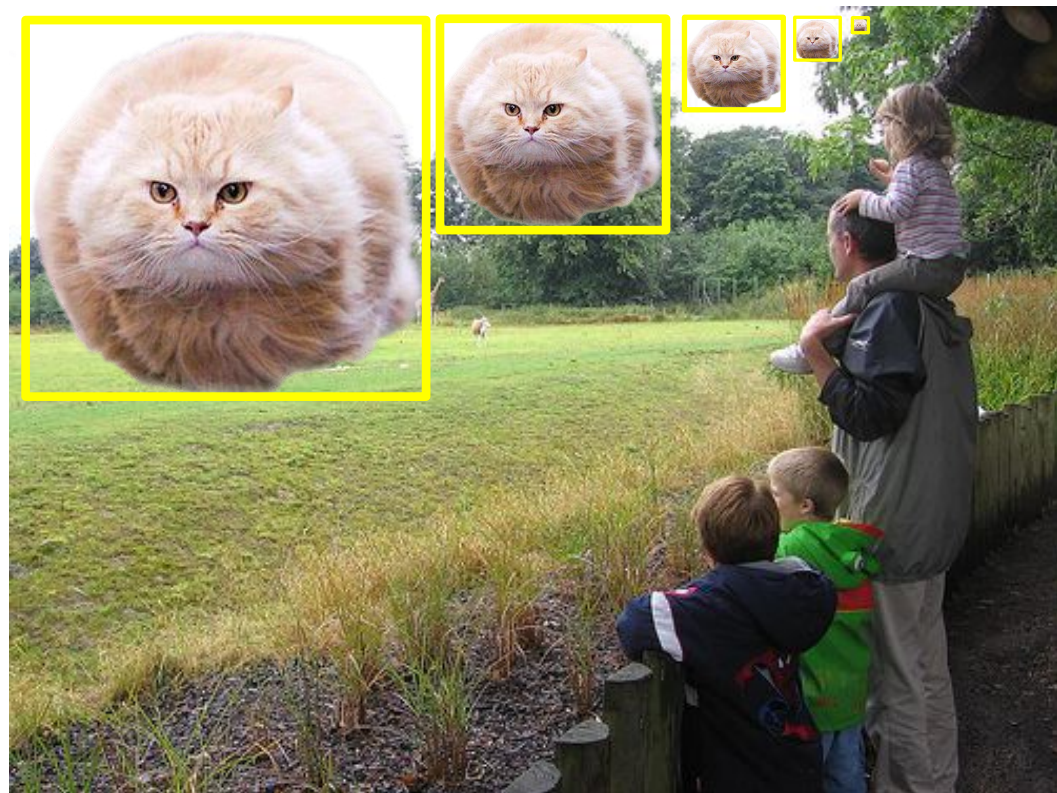Do nothing; give up; Fast!
(Fast R-CNN, YOLO, ...)

predict
predict
predict

(c) Pyramidal feature hierarchy
Fast!
(≈ SSD, ...)

Let's examine this one

R. Girshick

# Compromise Feature Quality, but Fast ("Free")



Low resolution, Strong features

High resolution, Weak features
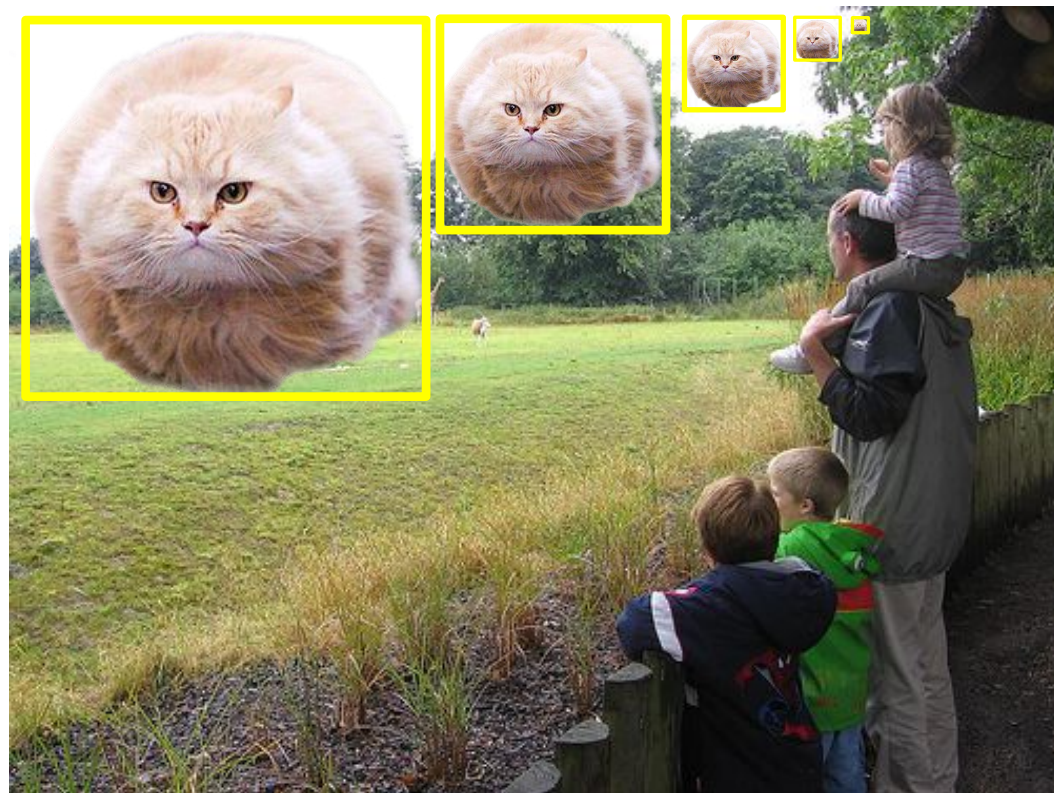
conv5

conv4

conv3

predict

predict

predict

The "native" in-network feature pyramid poses an inherent tradeoff

R. Girshick

# 2. … + Feature Pyramid Network (FPN)

[Optional, but recommended]



Lin et al. Feature Pyramid Networks for Object Detection. In CVPR 2017.

R. Girshick

FPN

Backbone ConvNet

1/16 image resolution
Res4 stage output

1/8 image resolution
Res3 stage output

1/4 image resolution
Res2 stage output

FPN weights use random initialization

predict

predict

predict

(Illustrated with 3 levels for simplicity; typically use 5 or 6)

2x up

1x1 conv

+

Input

R. Girshick

Lin et al. Feature Pyramid Networks for Object Detection. In CVPR 2017.

# No Compromise on Feature Quality, still Fast



Low resolution, Strong features

High resolution, **Strong** features

FPN pyramid

predict
predict
predict

2x up

1x1 conv

Lin et al. Feature Pyramid Networks for Object Detection. In CVPR 2017.
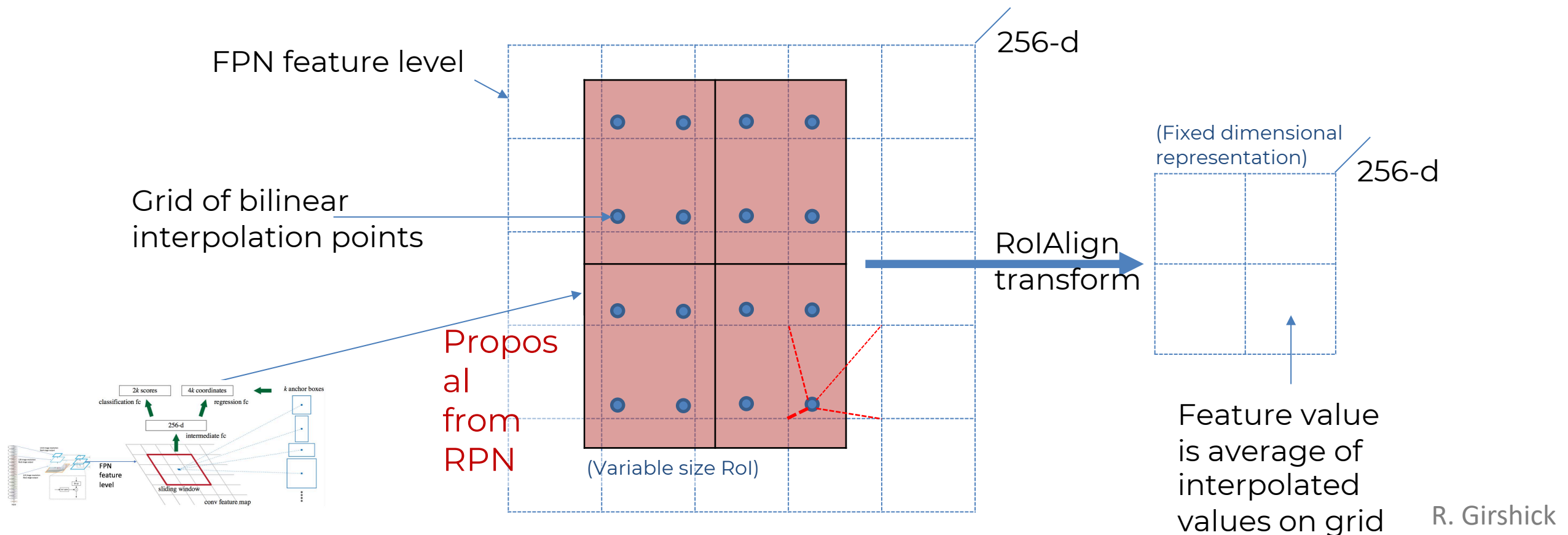
R. Girshick

# 3. … + Region Proposal Network (RPN)

Proposals = sliding window object/not-object classifier
+ box regression *inside the same network*



Anchors are prototypical object boxes

Ren et al. Faster R-CNN: Region Proposal Networks for Real-Time Object Detection. In NIPS 2015

R. Girshick

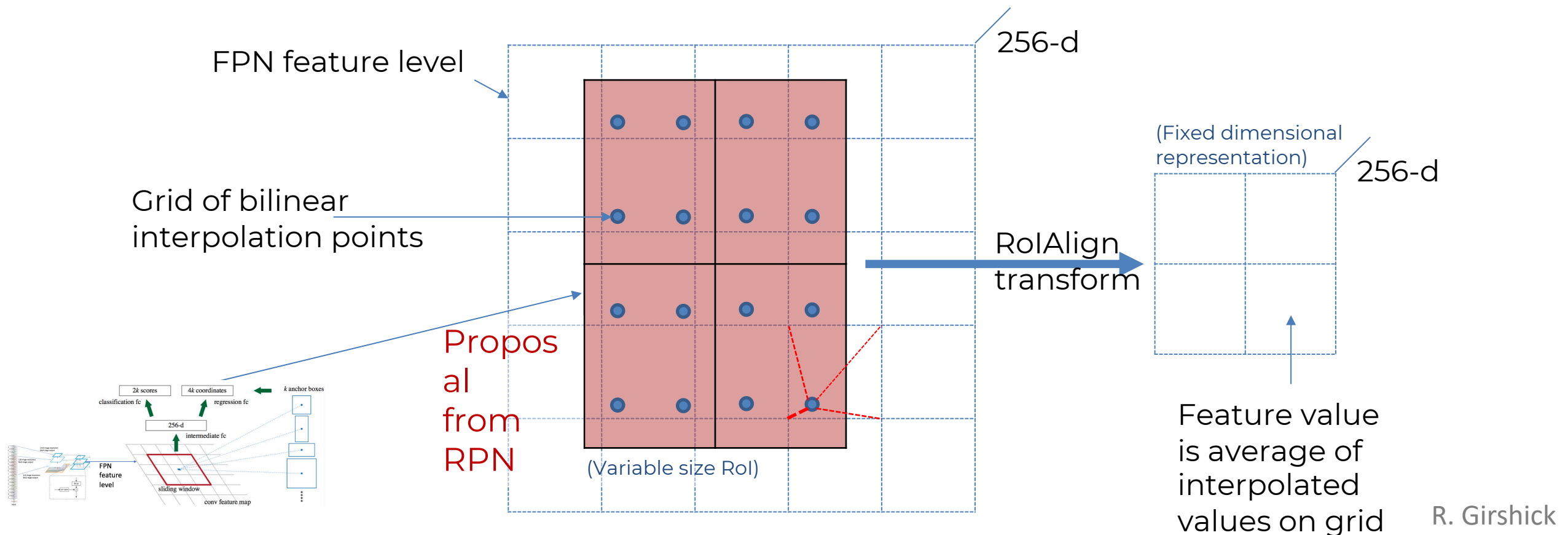# 4. ... + RoIAlign Transform (on each Proposal)

Smoothly normalize features and predictions into coordinate frame free of scale and aspect ratio



FPN feature level

256-d

Grid of bilinear interpolation points

(Fixed dimensional representation)

256-d

RoIAlign transform

Proposal from RPN

(Variable size RoI)

Feature value is average of interpolated values on grid

R. Girshick

# 4. … + RoIAlign Transform (on each Proposal)

Key: No coordinate quantization
(cf. RoIPool in Fast R-CNN, etc.)



FPN feature level

256-d

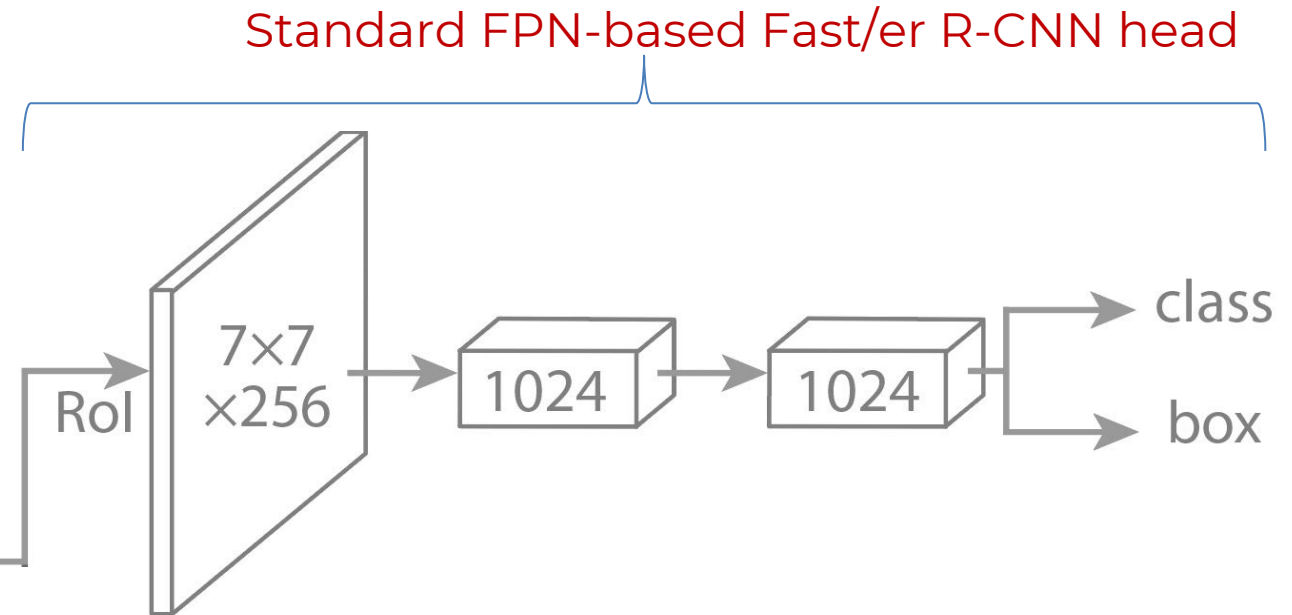Grid of bilinear interpolation points

(Fixed dimensional representation)

256-d

RoIAlign transform

Proposal from RPN

(Variable size RoI)

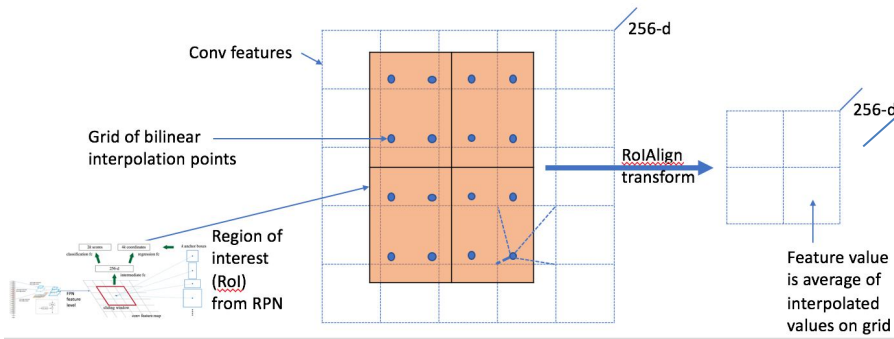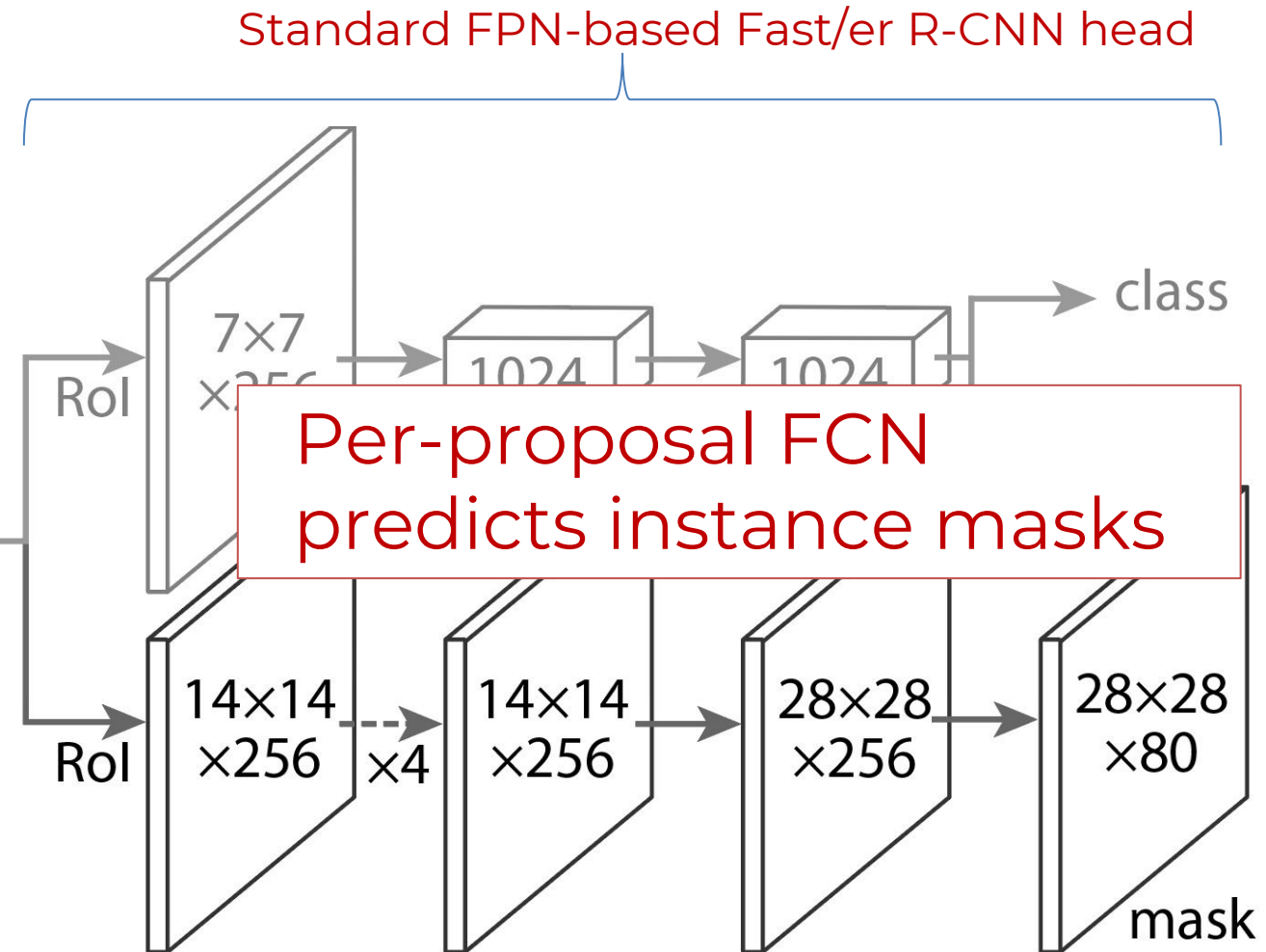Feature value is average of interpolated values on grid

R. Girshick

# 5. … + Task-specific Heads (on each Proposal)

Task specific heads for …
- Bounding box regression
- Object classification
- Instance mask prediction
- Human keypoint prediction



Standard FPN-based Fast/er R-CNN head

RoIAlign transformed features

RoI    7×7 ×256 → 1024 → 1024 → class / box

256-d

Conv features

Grid of bilinear interpolation points

RoIAlign transform

256-d

Region of interest (RoI) from RPN

Feature value is average of interpolated values on grid

R. Girshick

# 5. … + Task-specific Heads (on each Proposal)

Task specific heads for …
- Bounding box regression
- Object classification
- Instance mask prediction
- Human keypoint prediction

Standard FPN-based Fast/er R-CNN head

RoIAlign transformed features



Conv features

256-d

Grid of bilinear interpolation points

256-d

RoIAlign transform

Region of interest (RoI) from RPN

Feature value is average of interpolated values on grid

RoI

7×7 ×256

1024

1024

class

Per-proposal FCN predicts instance masks

RoI

14×14 ×256

×4

14×14 ×256

28×28 ×256

28×28 ×80

mask

# Mask R-CNN: Training

Same as "image centric" Fast/er R-CNN training
- Use precomputed proposals for faster experimentation
- Use joint / end-to-end training for sharing features

But with training targets for masks

# Example Mask Training Targets

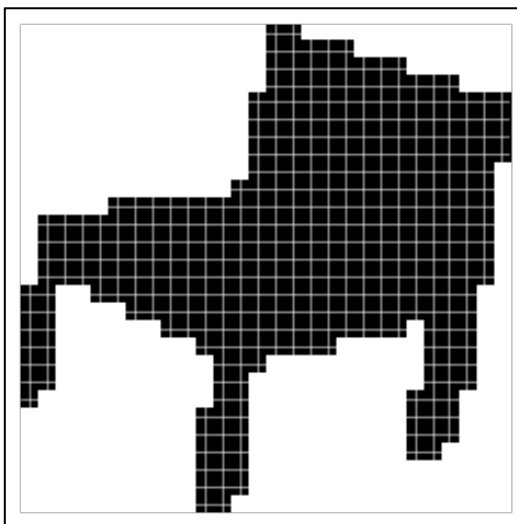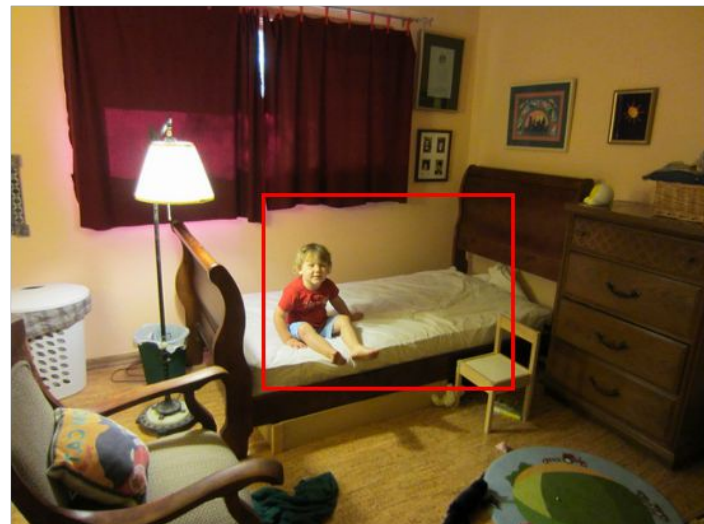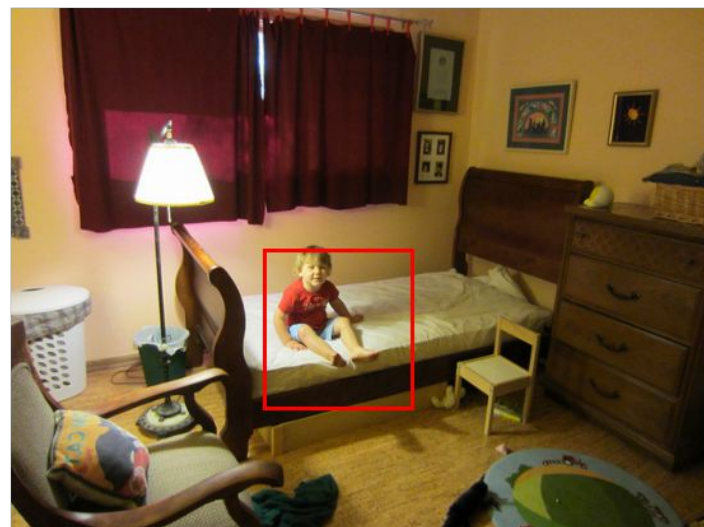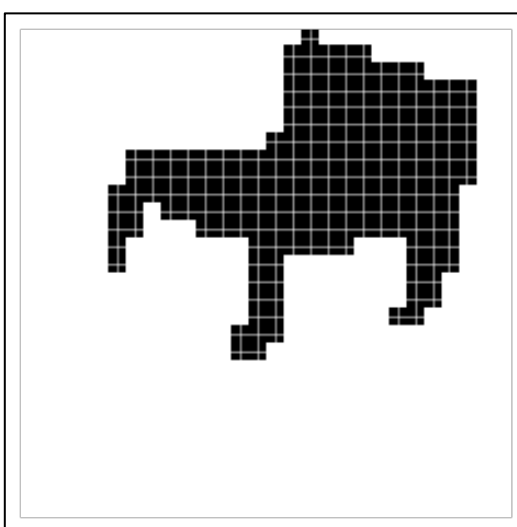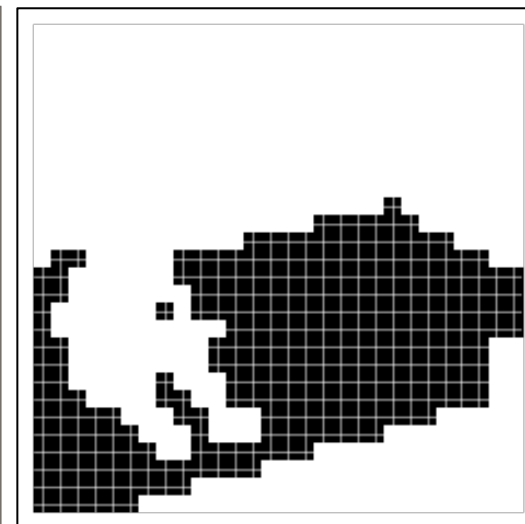| Image with training proposal | 28x28 mask target | Image with training proposal | 28x28 mask target |

# Mask R-CNN: Inference
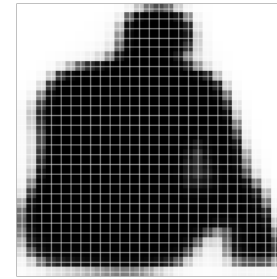
1. Perform Faster R-CNN inference
    - Generate proposals (RPN)
    - Score the proposals
    - Regress from proposals to refined detection boxes
    - Apply NMS and take the top $K$ (= 100, e.g.)

2. Run RoIAlign and mask head on top-$K$ refined, post-NMS boxes
    - Fast (only compute masks for top-$K$ detections)
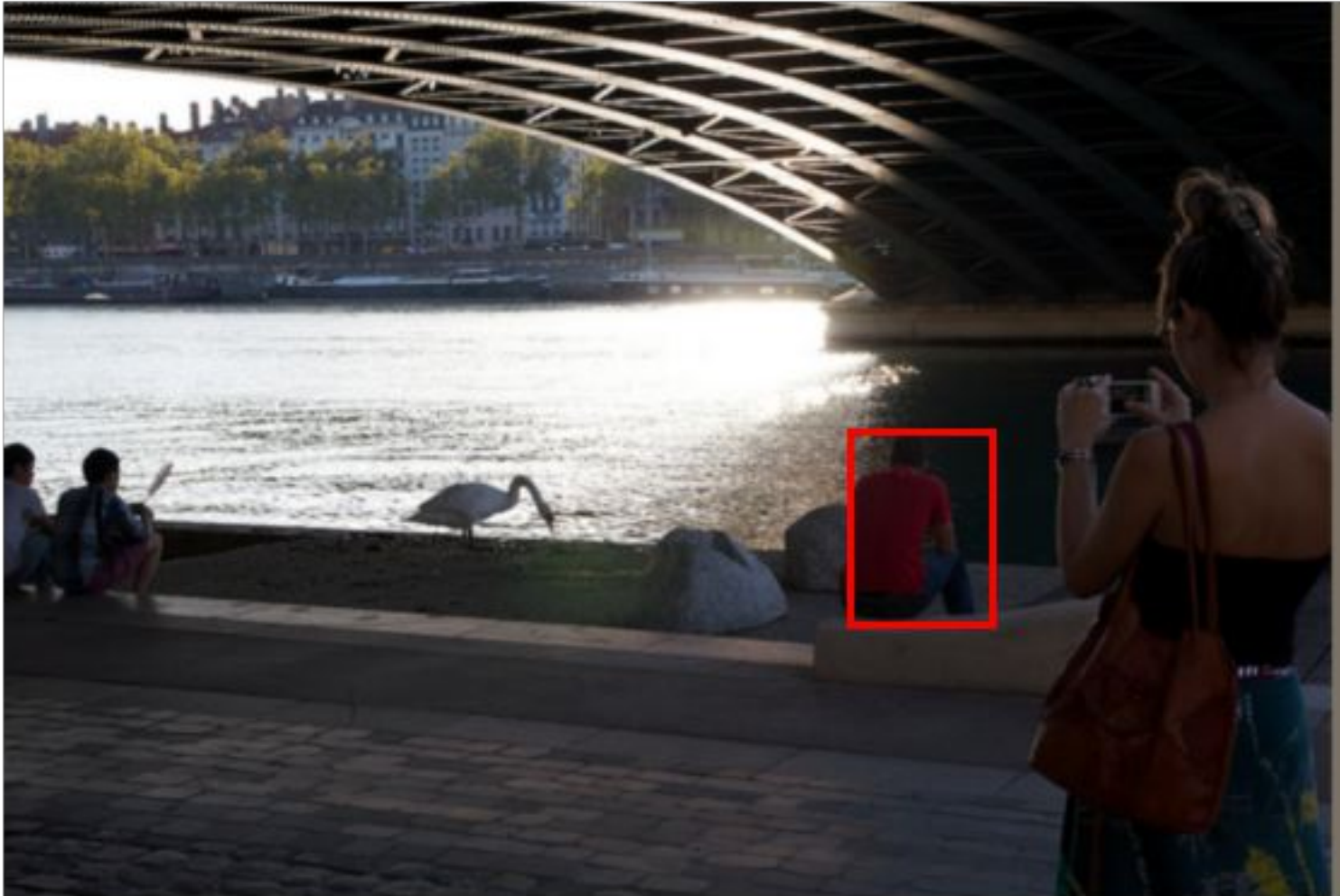    - Improves accuracy (uses *refined* detection boxes, not proposals)

# Mask Prediction



Validation image with box detection shown in red

28x28 soft prediction from Mask R-CNN
(enlarged)



Soft prediction resampled to image coordinates
(bilinear and bicubic interpolation work equally well)



Final prediction (threshold at 0.5)
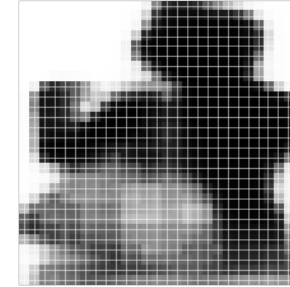


R. Girshick

# Mask Prediction

28x28 soft prediction



Resized soft prediction    Final mask



Validation image with box detection shown in red

R. Girshick

# Mask Prediction



Validation image with box detection shown in red
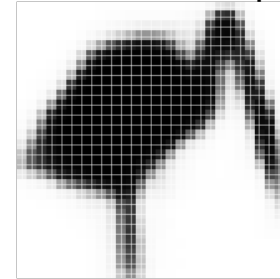
28x28 soft prediction

Resized Soft prediction

Final mask

R. Girshick

R. Girshick

R. Girshick

# Segmentation Overview

- 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