# Lecture 23
# Deep Learning:
# Segmentation

## COS 429: Computer Vision
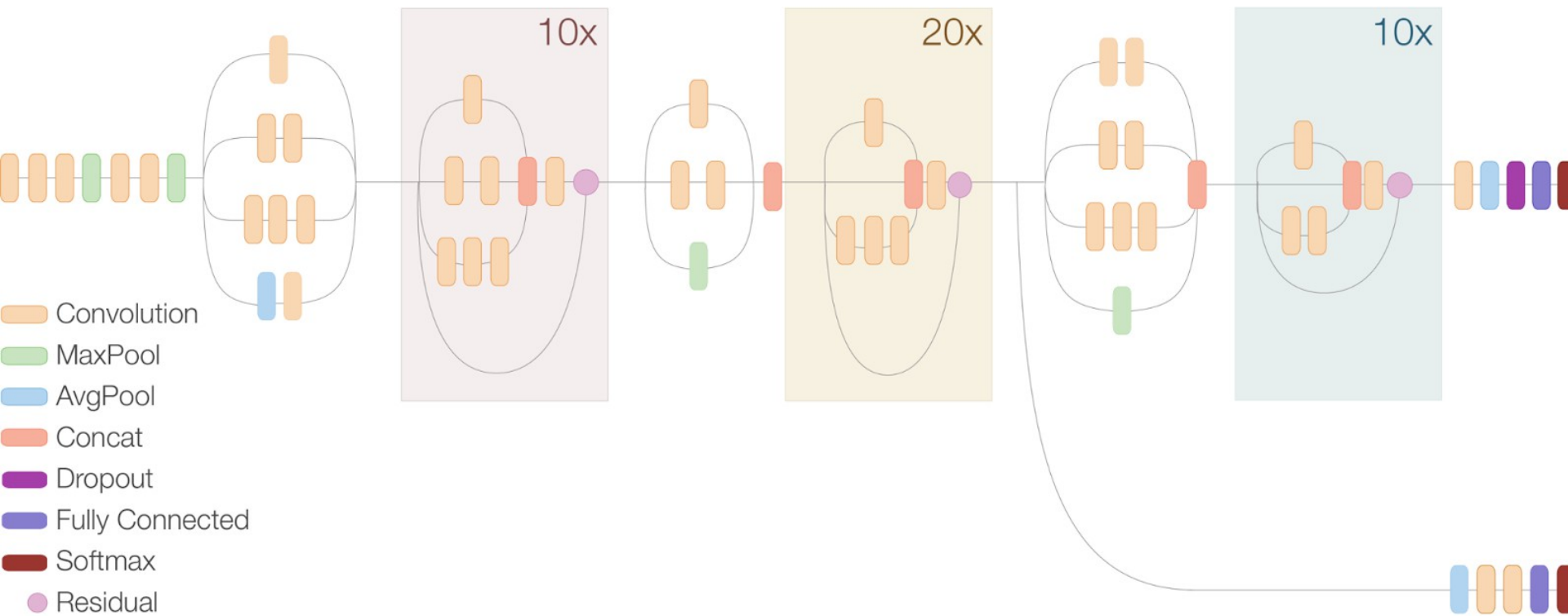
PRINCETON
UNIVERSITY

COS429 : 12.12.16 : Andras Ferencz

# Inception Resnet V2 Network



# Compressed View



10x

20x

10x

- ▭ Convolution
- ▭ MaxPool
- ▭ AvgPool
- ▭ Concat
- ▭ Dropout
- ▭ Fully Connected
- ▭ Softmax
- ● Residual

Slide Credit:

## ImageNet Classification Error (Top 5)



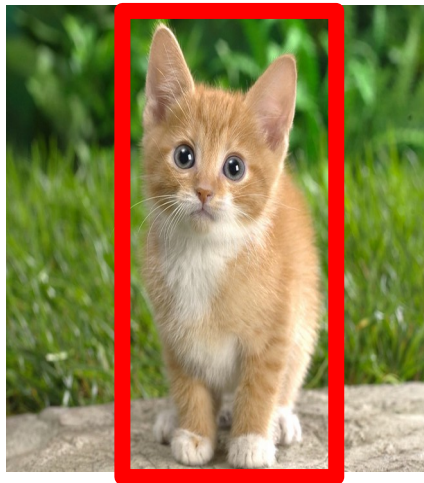Szegedy et al, Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning, arXiv 2016

3 : COS429 : L23 : 12.12.16 : Andras Ferencz            Slide Credit:

# Computer Vision Tasks

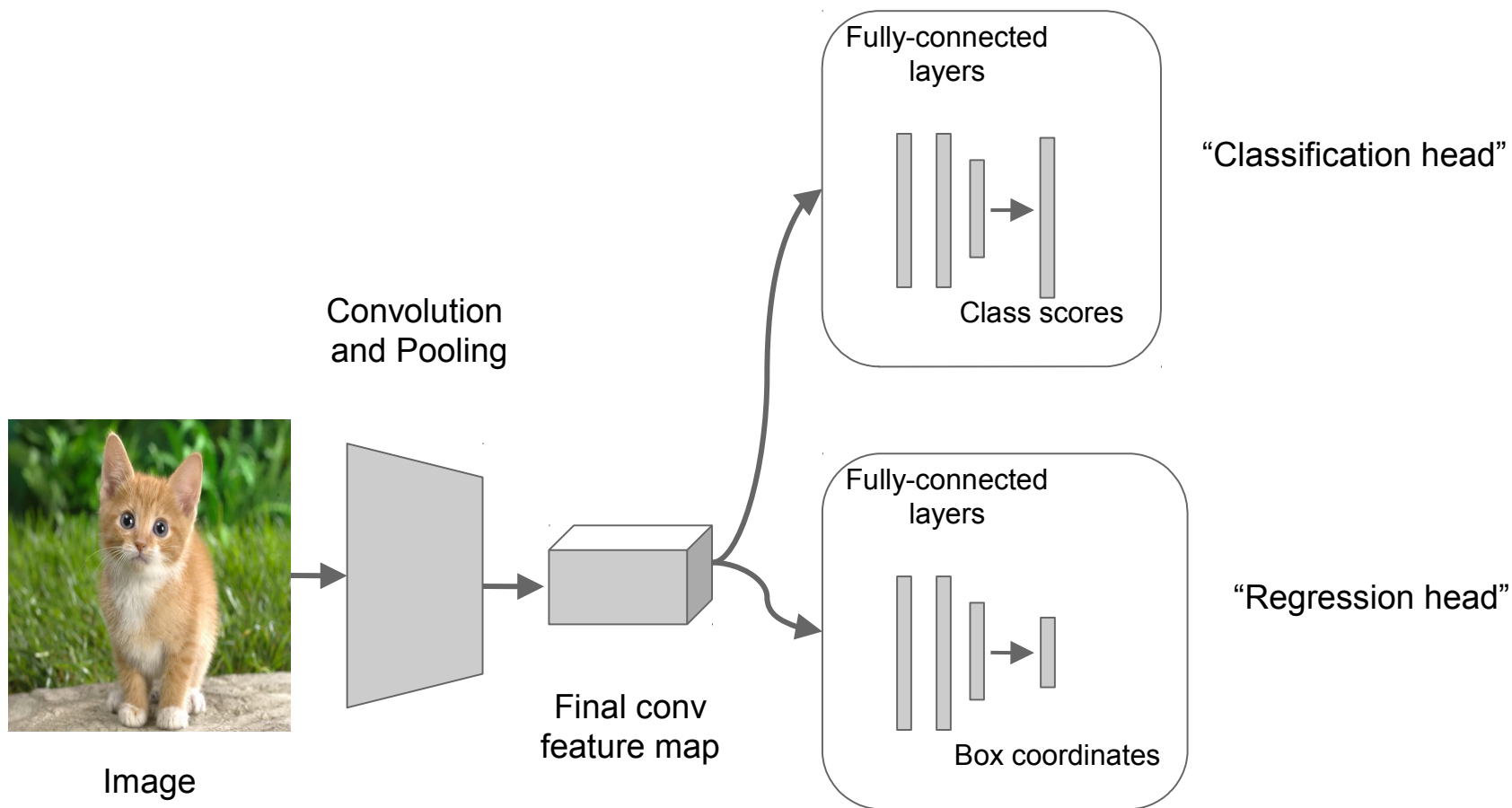| Classification | Classification + Localization | Object Detection | Instance Segmentation |
| --- | --- | --- | --- |



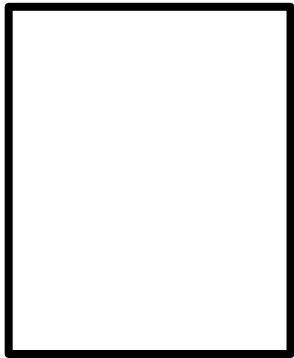CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK

Single object | Multiple objects

Slide Credit:

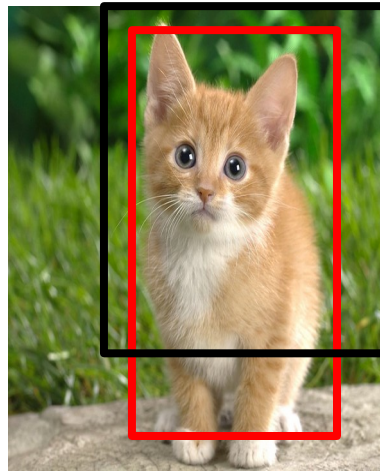# Simple Recipe for Classification + Localization

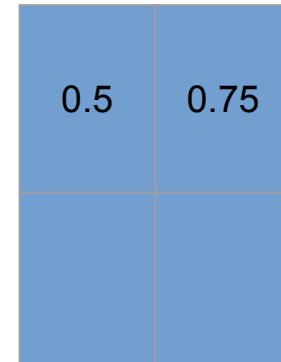**Step 2**: Attach new fully-connected "regression head" to the network

Slide Credit:

# Sliding Window: Overfeat

Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

Classification scores:
P(cat)

0.5    0.75

Slide Credit:

# Sliding Window: Overfeat



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

| 0.5 | 0.75 |
|-----|------|
| 0.6 | 0.8 |

Classification scores:
P(cat)

Slide Credit:

# Sliding Window: Overfeat

Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

| 0.5 | 0.75 |
|-----|------|
| 0.6 | 0.8  |

Classification scores:
P(cat)

Slide Credit:

# Sliding Window: Overfeat

Greedily merge boxes and scores (details in paper)



Network input:
3 x 221 x 221

Larger image:
3 x 257 x 257

0.8

Classification score:
P(cat)

Slide Credit:

# Sliding Window: Overfeat

In practice use many sliding window
locations and multiple scales

Window positions + score maps

Box regression outputs

Final Predictions

Sermanet et al, "Integrated Recognition, Localization and Detection using Convolutional Networks", ICLR 2014

Slide Credit:

# Efficient Sliding Window: Overfeat
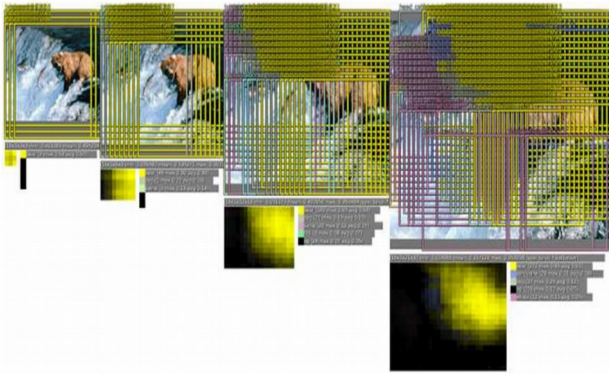
Efficient sliding window by converting fully-connected layers into convolutions



Convolution + pooling

Image:
3 x 221 x 221

Feature map:
1024 x 5 x 5

5 x 5 conv

5 x 5 conv

4096 x 1 x 1

1 x 1 conv

1024 x 1 x 1

1 x 1 conv

Class scores:
1000 x 1 x 1

1 x 1 conv

1 x 1 conv

4096 x 1 x 1

1024 x 1 x 1

Box coordinates:
(4 x 1000) x 1 x 1

11

Slide Credit:

# Efficient Sliding Window: Overfeat

**Training time:** Small image, 1 x 1 classifier output

**Test time:** Larger image, 2 x 2 classifier output, only extra compute at yellow regions
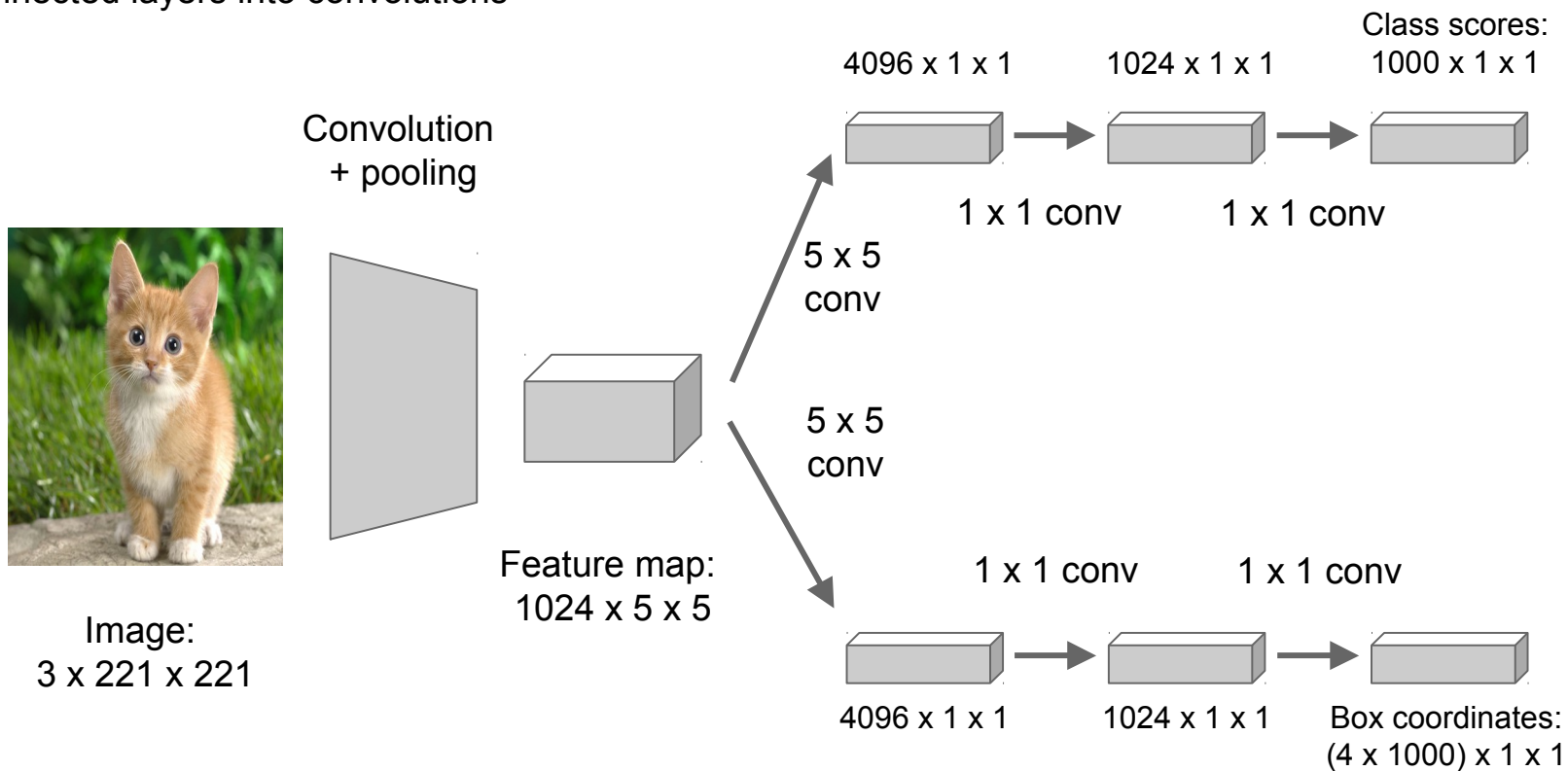


Sermanet et al, "Integrated Recognition, Localization and Detection using Convolutional Networks", ICLR 2014

Slide Credit:

# Computer Vision Tasks
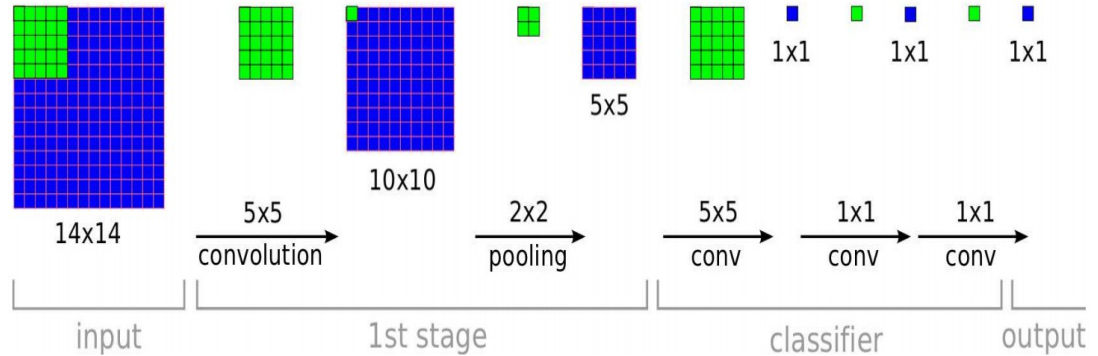
**Classification**    **Classification + Localization**    **Object Detection**    **Instance Segmentation**

# Region Proposals

- Find "blobby" image regions that are likely to contain objects
- "Class-agnostic" object detector
- Look for "blob-like" regions

Slide Credit:

# Region Proposals: Selective Search

Bottom-up segmentation, merging regions at multiple scales

Convert regions to boxes



Uijlings et al, "Selective Search for Object Recognition", IJCV 2013

15

Slide Credit:

# R-CNN
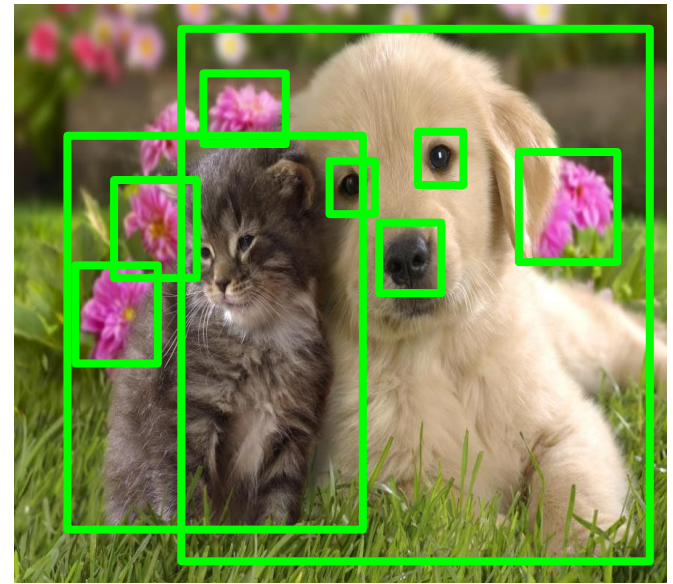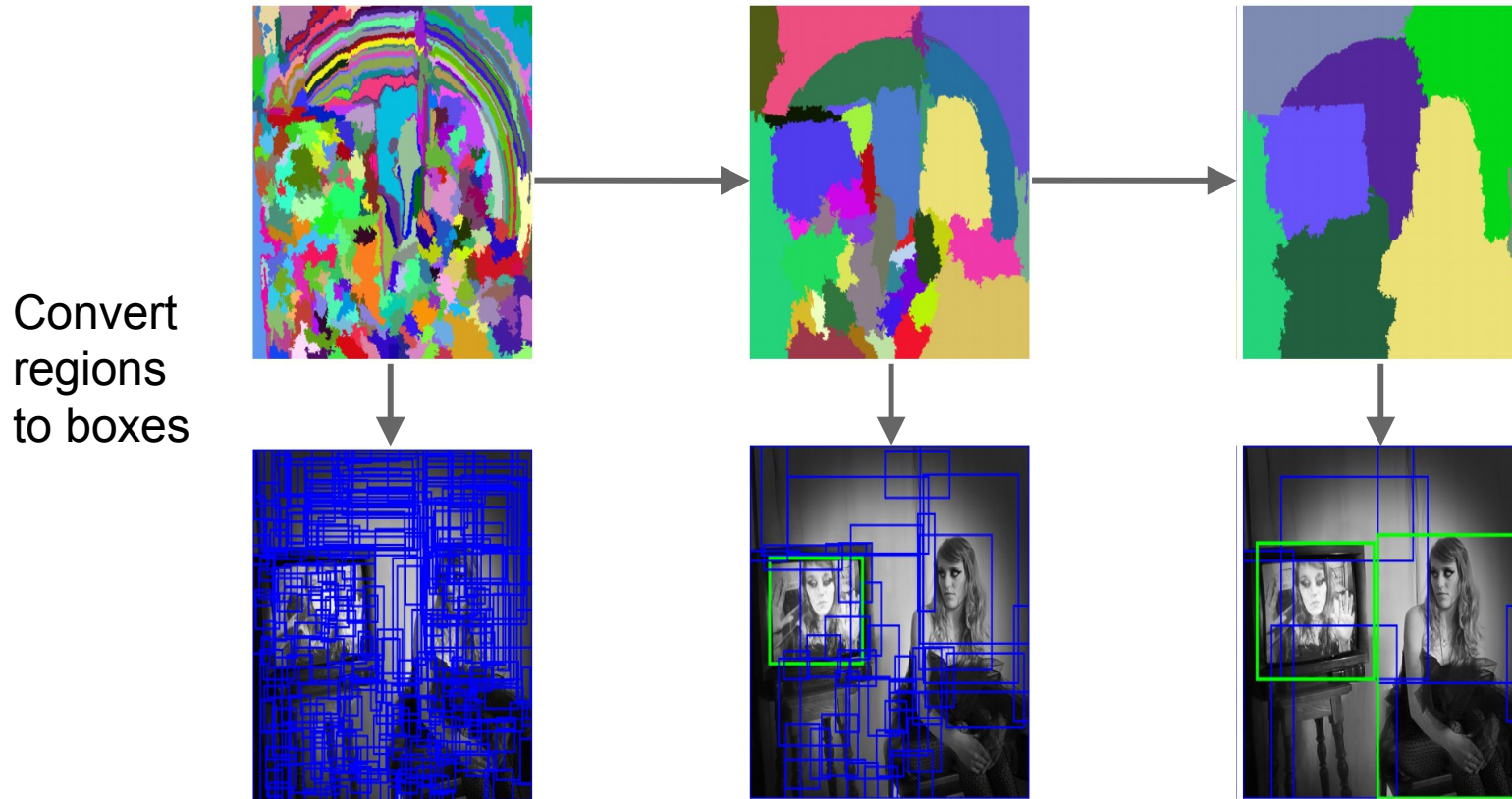


Apply bounding-box regressors

Classify regions with SVMs

Bbox reg | SVMs

Forward each region through ConvNet

Warped image regions

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

Input image

Girshick et al. CVPR14.

Post hoc component

Girschick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014

Slide credit: Ross Girschick

16

Slide Credit:

# Fast R-CNN



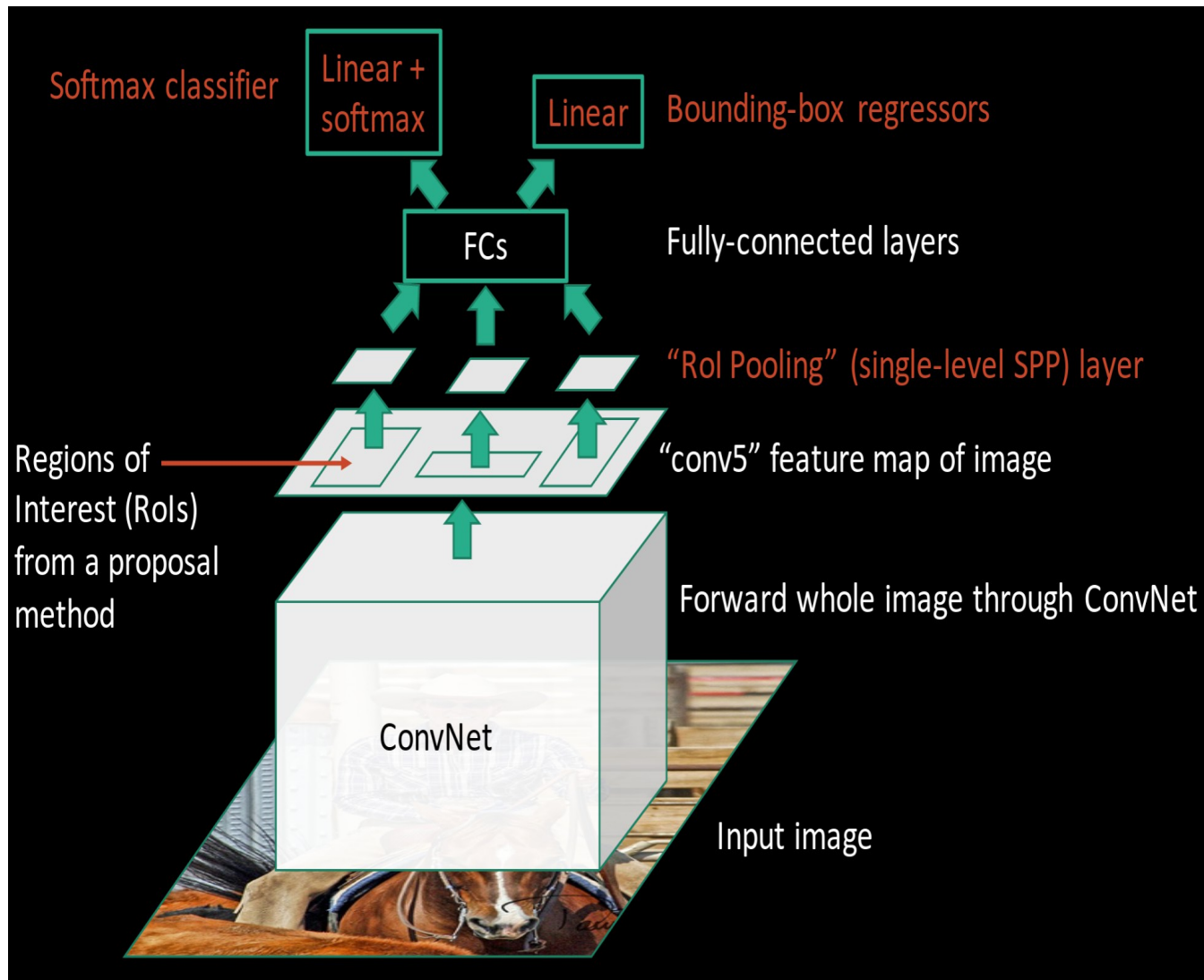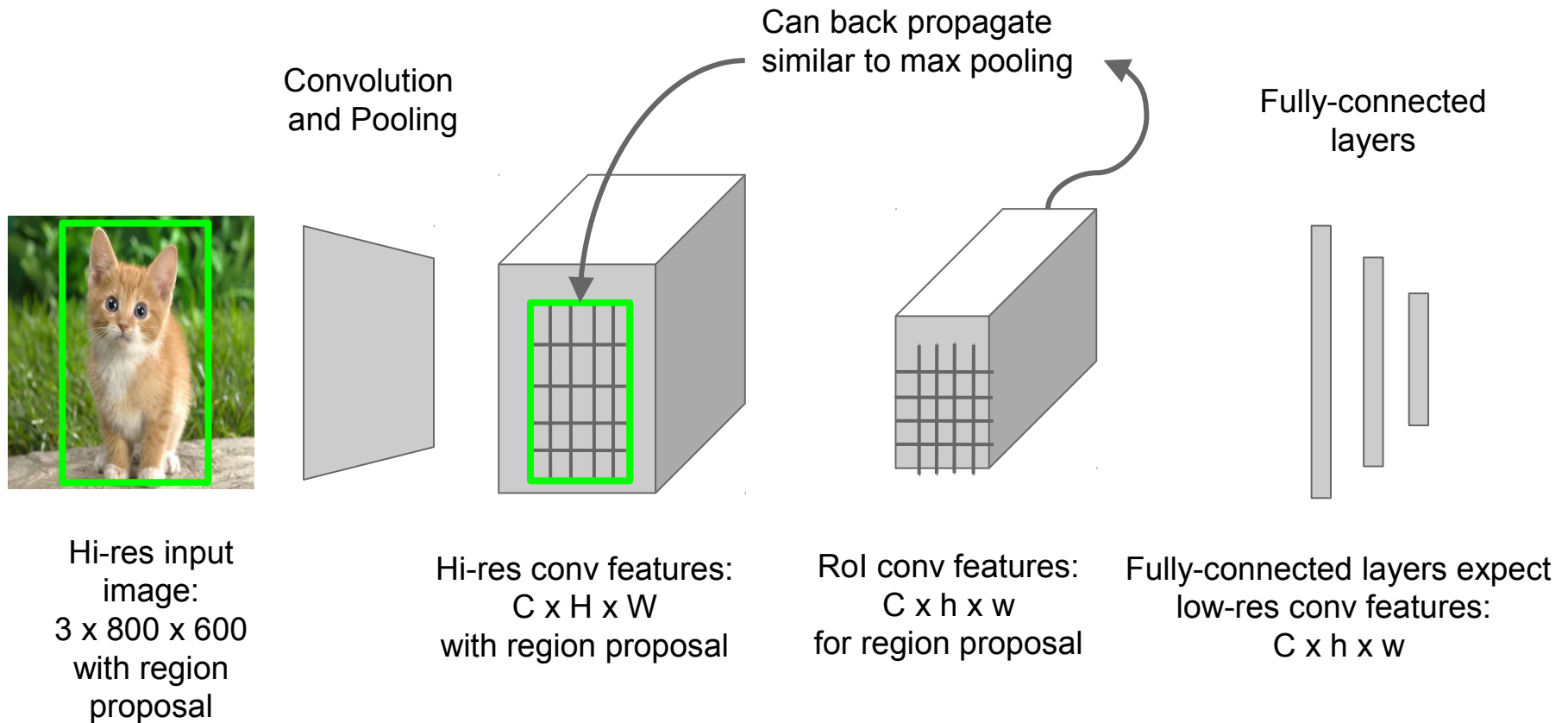**R-CNN Problems**:
Slow at test-time due to independent forward passes of the CNN

**Solution:**
Share computation of convolutional layers between proposals for an image

**R-CNN Problems**:
- Post-hoc training: CNN not updated in response to final classifiers and regressors
- Complex training pipeline

**Solution:**
Just train the whole system end-to-end all at once!

Slide Credit:

# Fast R-CNN: Region of Interest Pooling



Convolution and Pooling

Can back propagate similar to max pooling

Fully-connected layers

Hi-res input image:
3 x 800 x 600
with region proposal

Hi-res conv features:
C x H x W
with region proposal

RoI conv features:
C x h x w
for region proposal

Fully-connected layers expect low-res conv features:
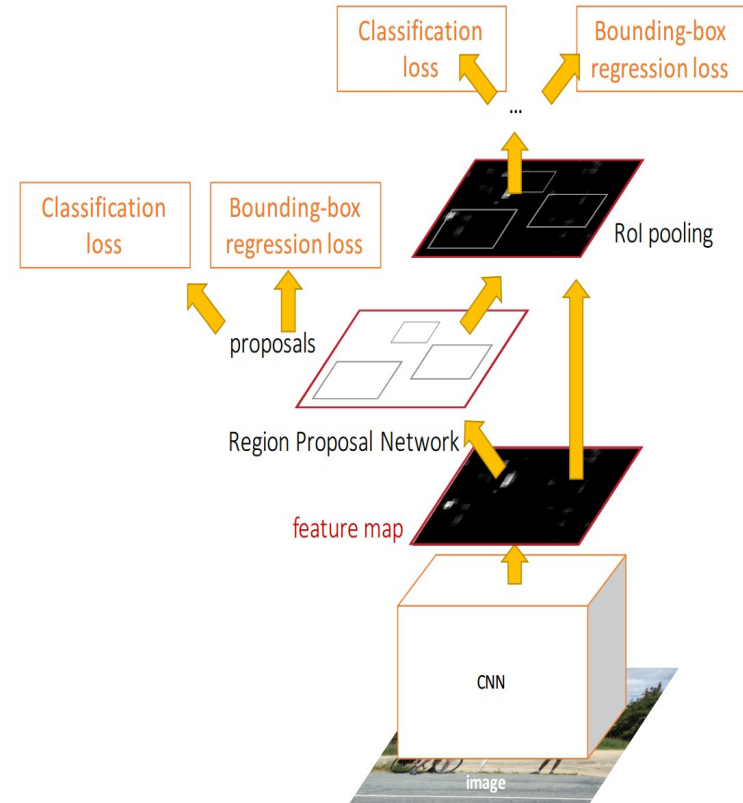C x h x w

18

Slide Credit:

# Faster R-CNN: Training

In the paper: Ugly pipeline
- Use alternating optimization to train RPN, then Fast R-CNN with RPN proposals, etc.
- More complex than it has to be

Since publication: Joint training!
One network, four losses
- RPN classification (anchor good / bad)
- RPN regression (anchor -> proposal)
- Fast R-CNN classification (over classes)
- Fast R-CNN regression (proposal -> box)



Slide credit: Ross Girschick

# Faster R-CNN: Results

| | R-CNN | Fast R-CNN | Faster R-CNN |
|---|---|---|---|
| Test time per image (with proposals) | 50 seconds | 2 seconds | **0.2 seconds** |
| (Speedup) | 1x | 25x | **250x** |
| mAP (VOC 2007) | 66.0 | **66.9** | **66.9** |

Slide Credit:
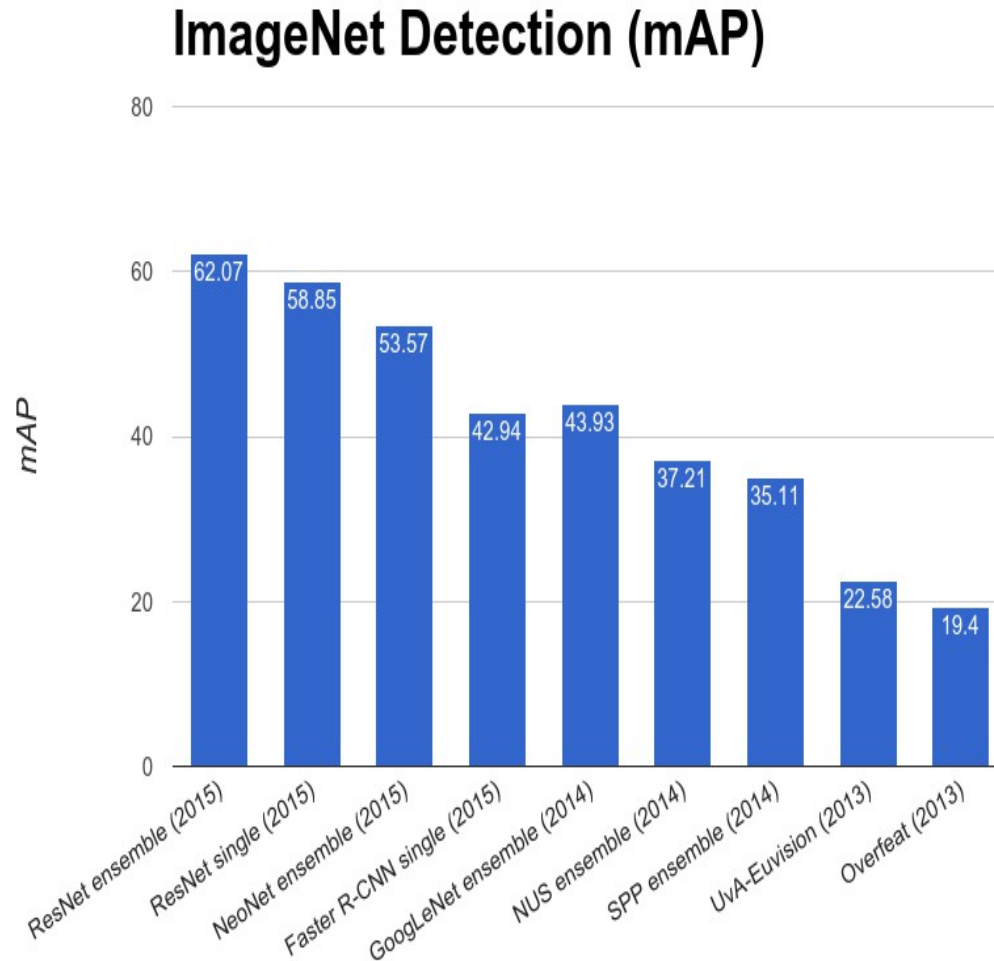
# Object Detection State-of-the-art:
## ResNet 101 + Faster R-CNN + some extras

| training data | COCO train | | COCO trainval | |
|---|---|---|---|---|
| test data | COCO val | | COCO test-dev | |
| mAP | @.5 | @[.5, .95] | @.5 | @[.5, .95] |
| baseline Faster R-CNN (VGG-16) | 41.5 | 21.2 | | |
| baseline Faster R-CNN (ResNet-101) | 48.4 | 27.2 | | |
| +box refinement | 49.9 | 29.9 | | |
| +context | 51.1 | 30.0 | 53.3 | 32.2 |
| +multi-scale testing | 53.8 | 32.5 | **55.7** | **34.9** |
| ensemble | | | **59.0** | **37.4** |

He et. al, "Deep Residual Learning for Image Recognition", arXiv 2015

Slide Credit:

# ImageNet Detection 2013 - 2015



## ImageNet Detection (mAP)

Bar chart of mAP values:
- ResNet ensemble (2015): 62.07
- ResNet single (2015): 58.85
- NeoNet ensemble (2015): 53.57
- Faster R-CNN single (2015): 42.94
- GoogLeNet ensemble (2014): 43.93
- NUS ensemble (2014): 37.21
- SPP ensemble (2014): 35.11
- UvA-Euvision (2013): 22.58
- Overfeat (2013): 19.4

Slide Credit:

# YOLO: You Only Look Once

## Detection as Regression

Divide image into S x S grid

Within each grid cell predict:
    B Boxes: 4 coordinates + confidence
    Class scores: C numbers

Regression from image to
7 x 7 x (5 * B + C) tensor

Direct prediction using a CNN

Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", arXiv 2015

# YOLO: You Only Look Once
## Detection as Regression

Faster than Faster R-CNN, but not as good

Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", arXiv 2015

| Real-Time Detectors | Train | mAP | FPS |
|---|---|---|---|
| 100Hz DPM [30] | 2007 | 16.0 | 100 |
| 30Hz DPM [30] | 2007 | 26.1 | 30 |
| Fast YOLO | 2007+2012 | 52.7 | **155** |
| YOLO | 2007+2012 | **63.4** | 45 |
| **Less Than Real-Time** | | | |
| Fastest DPM [37] | 2007 | 30.4 | 15 |
| R-CNN Minus R [20] | 2007 | 53.5 | 6 |
| Fast R-CNN [14] | 2007+2012 | 70.0 | 0.5 |
| Faster R-CNN VGG-16[27] | 2007+2012 | 73.2 | 7 |
| Faster R-CNN ZF [27] | 2007+2012 | 62.1 | 18 |

Slide Credit:

# Computer Vision Tasks



| Classification | Classification + Localization | Object Detection | Segmentation |
|---|---|---|---|
| CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK |
| Single object | | Multiple objects | |

Slide Credit:

# Today

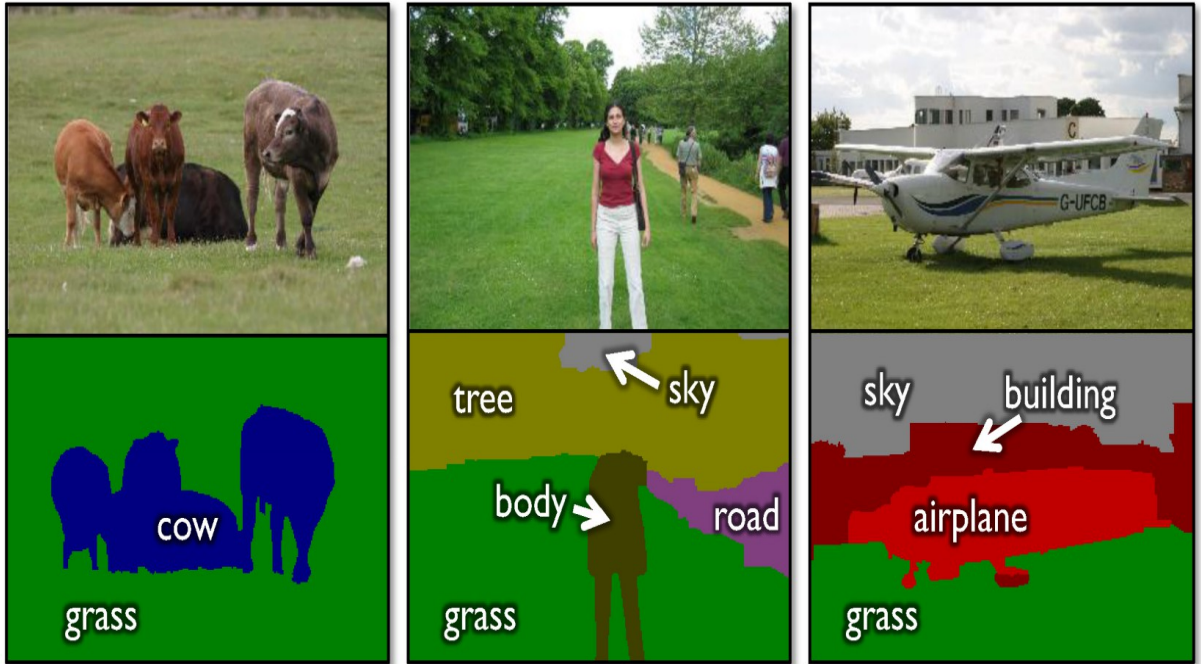**Classification**   **Classification + Localization**   **Object Detection**   **Segmentation**



Today

Slide Credit:

# Semantic Segmentation

Label every pixel!

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

Slide Credit:

# Instance Segmentation

Detect instances, give category, label pixels

"simultaneous detection and segmentation" (SDS)

Lots of recent work (MS-COCO)



person
person  person
horse

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

Slide Credit:

# Semantic Segmentation

Extract
patch

Slide Credit:

# Semantic Segmentation

Extract patch

Run through a CNN



CNN

Slide Credit:

# Semantic Segmentation

Extract patch      Run through a CNN      Classify center pixel



CNN → COW

Slide Credit:

# Semantic Segmentation



Extract patch

Run through a CNN

Classify center pixel

CNN → COW

Repeat for every pixel

cow

grass

Slide Credit:

# Semantic Segmentation

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



CNN

cow

grass

Smaller output
due to pooling

Slide Credit:

# Semantic Segmentation: Multi-Scale



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

Slide Credit:

# Semantic Segmentation: Multi-Scale

Resize image to
multiple scales



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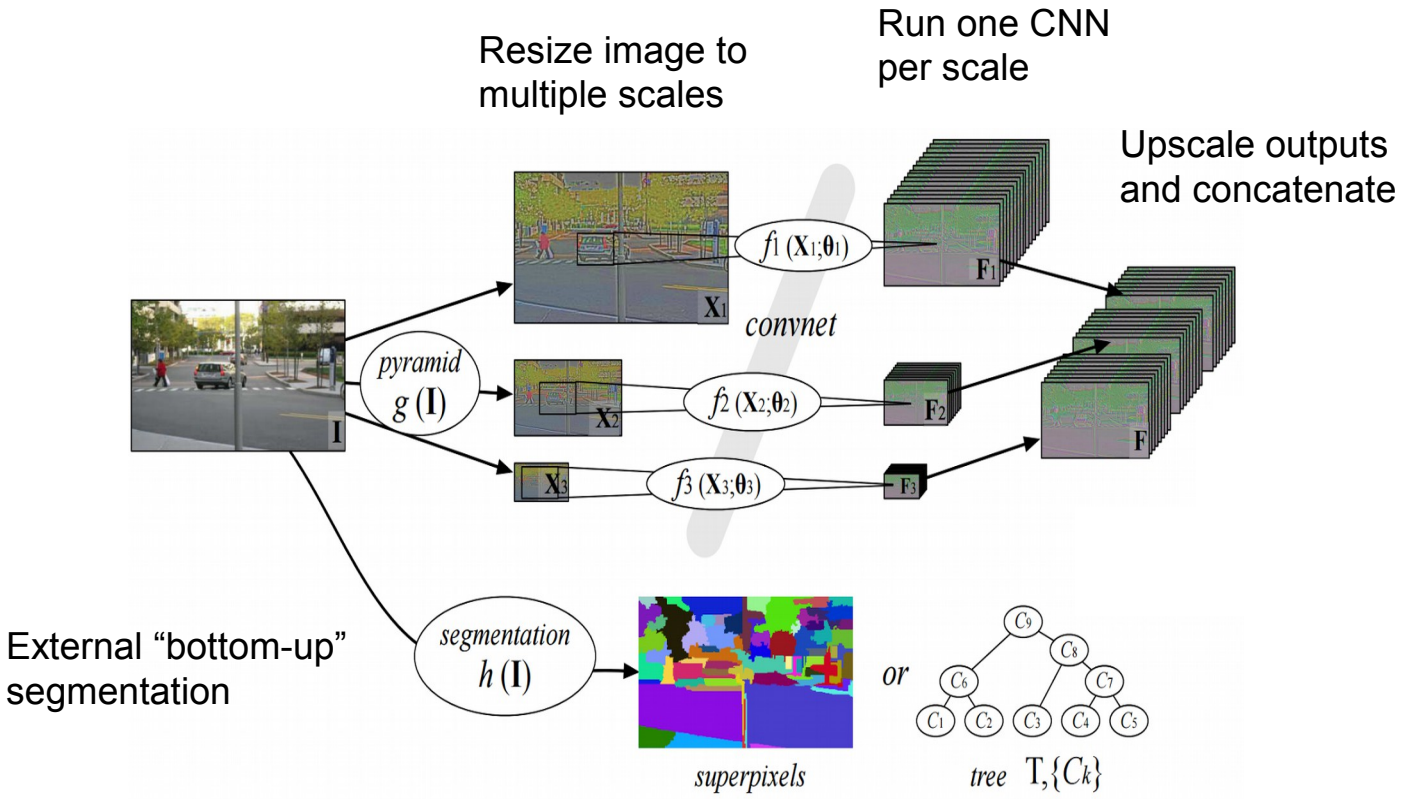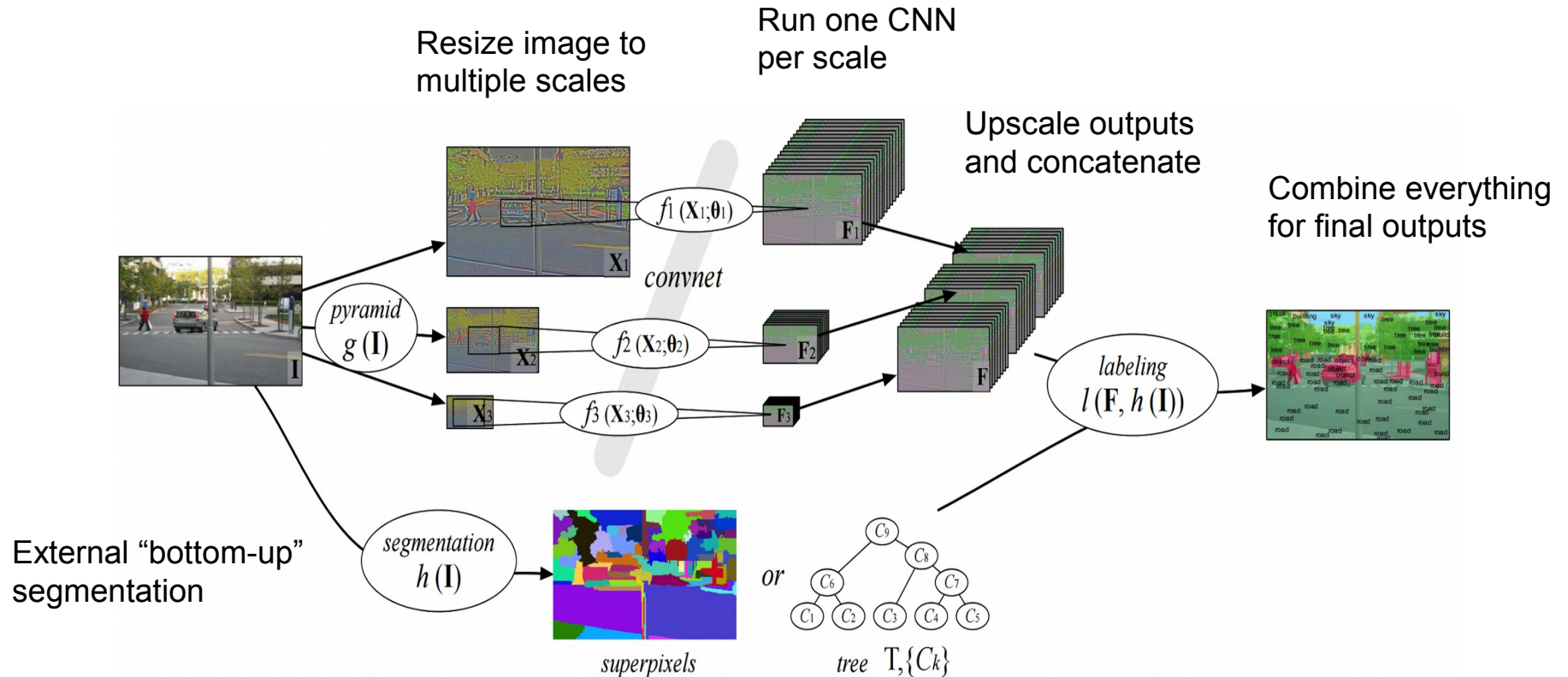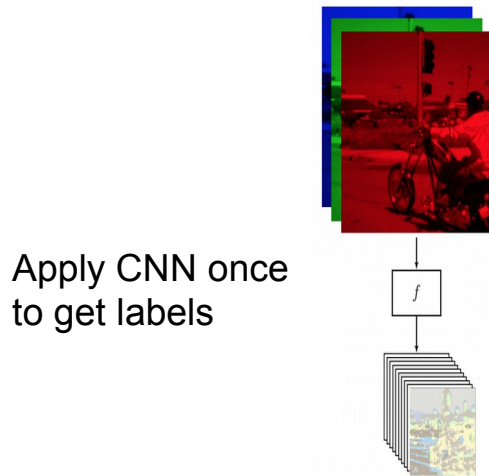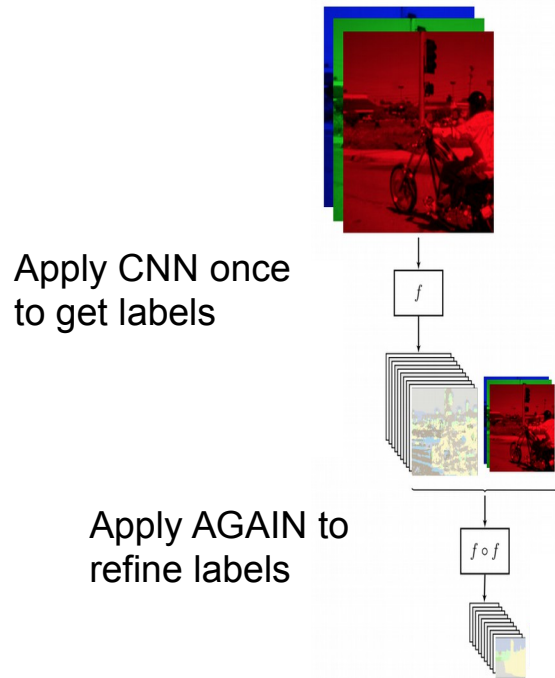
Slide Credit:

# Semantic Segmentation: Multi-Scale

Resize image to multiple scales

Run one CNN per scale



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

Slide Credit:

# Semantic Segmentation: Multi-Scale



Resize image to multiple scales

Run one CNN per scale

Upscale outputs and concatenate

$$pyramid\ g\,(\mathbf{I})$$

$$f_1\,(\mathbf{X}_1;\boldsymbol{\theta}_1)$$

$$f_2\,(\mathbf{X}_2;\boldsymbol{\theta}_2)$$

$$f_3\,(\mathbf{X}_3;\boldsymbol{\theta}_3)$$

convnet

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

Slide Credit:

# Semantic Segmentation: Multi-Scale

Resize image to multiple scales

Run one CNN per scale

Upscale outputs and concatenate

$f_1(X_1; \theta_1)$

$X_1$

*convnet*

$F_1$

*pyramid* $g(I)$

$X_2$

$f_2(X_2; \theta_2)$

$F_2$

$F$

$X_3$

$f_3(X_3; \theta_3)$

$F_3$

$I$

External "bottom-up" segmentation

*segmentation* $h(I)$

*superpixels*

*or*

$C_9$
$C_8$
$C_6$ $C_7$
$C_1$ $C_2$ $C_3$ $C_4$ $C_5$

*tree* $T, \{C_k\}$

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

# Semantic Segmentation: Multi-Scale

Resize image to
multiple scales

Run one CNN
per scale

Upscale outputs
and concatenate

Combine everything
for final outputs

$f_1(X_1; \theta_1)$

$X_1$    *convnet*    $F_1$

*pyramid*
$g(I)$

$X_2$    $f_2(X_2; \theta_2)$    $F_2$

$I$

$X_3$    $f_3(X_3; \theta_3)$    $F_3$

$F$

*labeling*
$l(F, h(I))$

External "bottom-up"
segmentation

*segmentation*
$h(I)$

*superpixels*

*or*

$C_9$, $C_8$, $C_6$, $C_7$, $C_1$, $C_2$, $C_3$, $C_4$, $C_5$

*tree* $T, \{C_k\}$

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

# Semantic Segmentation: Refinement

Apply CNN once
to get labels



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

Slide Credit:

# 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

Slide Credit:

# Semantic Segmentation: Refinement

Same CNN weights:
**recurrent convolutional network**

Apply CNN once
to get labels

Apply AGAIN to
refine labels

And again!



More iterations improve results

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

42

Slide Credit:

# Semantic Segmentation: Upsampling



Learnable upsampling!

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

Slide Credit:

# 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

Slide Credit:

# Semantic Segmentation: Upsampling



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

Slide Credit:

# 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

Slide Credit:

# 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

Slide Credit:

# 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

Slide Credit:

# Learnable Upsampling: "Deconvolution"

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



Input: 4 x 4

Output: 2 x 2

Slide Credit:

# 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

Slide Credit:

# 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

Slide Credit:

# Learnable Upsampling: "Deconvolution"

3 x 3 "deconvolution", stride 2 pad 1



Input: 2 x 2                              Output: 4 x 4

Slide Credit:

# Learnable Upsampling: "Deconvolution"

3 x 3 "deconvolution", stride 2 pad 1



Input gives
weight for
filter

Input: 2 x 2

Output: 4 x 4

54

Slide Credit:

# Learnable Upsampling: "Deconvolution"

3 x 3 "deconvolution", stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

Same as backward pass for normal convolution!

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

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

Input: 2 x 2

Output: 4 x 4

Slide Credit:

# Semantic Segmentation: Upsampling



Normal VGG      "Upside down" VGG

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

6 days of training on Titan X…

56

# 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

Slide Credit:

# Instance Segmentation

Similar to R-CNN, but
with segments



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

Slide Credit:

# Instance Segmentation

Similar to R-CNN, but with segments

Proposal Generation    External Segment proposals



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

Slide Credit:

# Instance Segmentation



Proposal Generation

External Segment proposals

Feature Extraction

Similar to R-CNN

Box CNN

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

60

Slide Credit:

# Instance Segmentation

Mask out background with mean image

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

Slide Credit:

# Instance Segmentation

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

Box CNN

Region CNN

Person? +1.8

Mask out background with mean image

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

Slide Credit:

# Instance Segmentation

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

Box CNN

Region CNN

Mask out background with mean image

Person? +1.8

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

63

Slide Credit:

# 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

Slide Credit:

# Instance Segmentation: Cascades

Similar to
Faster R-CNN

Region proposal network (RPN)



Won COCO 2015
challenge
(with ResNet)

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

Slide Credit:

# 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

# Instance Segmentation: Cascades



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

**Predictions**          **Ground truth**

Slide Credit:

# 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

Slide Credit:

# Quick overview of Other Topics

Slide Credit:

# Recurrent Neural Networks (RNN)
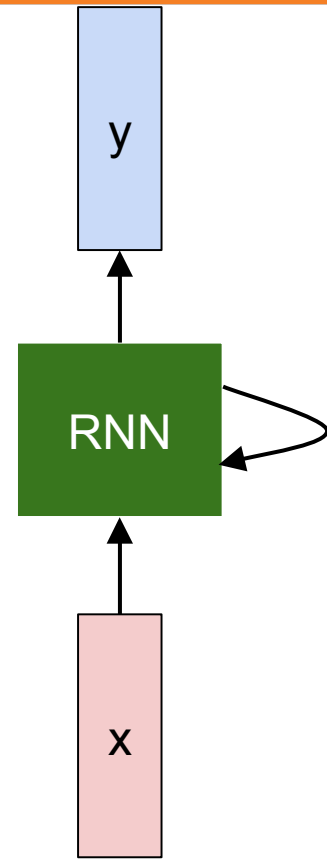


one to one     one to many     many to one     many to many     many to many

**Vanilla Neural Networks**

# Recurrent Neural Networks (RNN)



| one to one | one to many | many to one | many to many | many to many |

e.g. **Image Captioning**
image -> sequence of words

71

Slide Credit:

# Recurrent Neural Networks (RNN)



e.g. **Sentiment Classification**
sequence of words -> sentiment

Slide Credit:

# Recurrent Neural Networks (RNN)



one to one | one to many | many to one | many to many | many to many

e.g. **Machine Translation**
seq of words -> seq of words

# Recurrent Neural Networks (RNN)



| one to one | one to many | many to one | many to many | many to many |

e.g. **Video classification on frame level**

74

Slide Credit:

y

RNN

x

# Character RNN during training

```
tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt   h ne etie h,hregtrs nigtike,aoaenns lng
```

train more

```
"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."
```

train more

```
Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.
```

train more

```
"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.
```

Slide Credit:

PANDARUS:

Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

VIOLA:

Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.

Slide Credit:

## Generated C code

```c
static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem_info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)
    seq_puts(s, "policy ");
}
```

# Searching for interpretable cells



quote detection cell

79 : COS429 : L23 : 12.12.16 : Andras Ferencz                    Slide Credit:

# Sequential Processing of fixed inputs



Multiple Object Recognition with
Visual Attention, Ba et al.

Slide Credit:

# Sequential Processing of fixed outputs



DRAW: A Recurrent
Neural Network For
Image Generation,
Gregor et al.

Slide Credit:

# Image Captioning



Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Slide Credit:

# Recurrent Neural Network



**Convolutional Neural Network**

Slide Credit:

# Soft Attention for Captioning



Distribution over
L locations

Distribution
over vocab

CNN

Image:
H x W x 3

Features:
L x D

Weighted
combination
of features

Weighted
features: D

First word

a1    a2    d1

h0    h1    h2

z1    y1    z2    y2

Xu et al, "Show, Attend and Tell: Neural
Image Caption Generation with Visual
Attention", ICML 2015

# Soft Attention for Captioning



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Slide Credit:

# Soft Attention for Captioning



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Slide Credit:

# Spatial Transformer Networks



Can we make this function differentiable?

Input image:
H x W x 3

Cropped and
rescaled image:
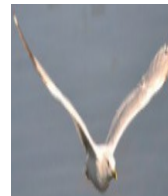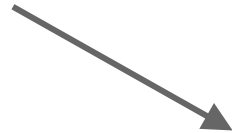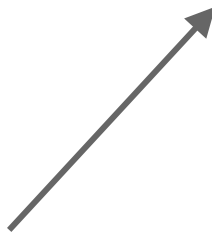X x Y x 3

Box Coordinates:
(xc, yc, w, h)

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

Slide Credit:

# Spatial Transformer Networks

Can we make this function differentiable?

Input image:
H x W x 3

Box Coordinates:
(xc, yc, w, h)

Cropped and rescaled image:
X x Y x 3

**Idea**: Function mapping *pixel coordinates* (xt, yt) of output to *pixel coordinates* (xs, ys) of input

Network attends to input by predicting $\theta$

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

$\mathcal{T}_\theta(G)$

Repeat for all pixels in *output* to get a **sampling grid**

Then use **bilinear interpolation** to compute output

$U$       $V$

Jaderberg et al, "Spatial Transformer Networks", NIPS 2015

Slide Credit:

# Spatial Transformer Networks

**Grid generator** uses $\theta$ to compute sampling grid

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$
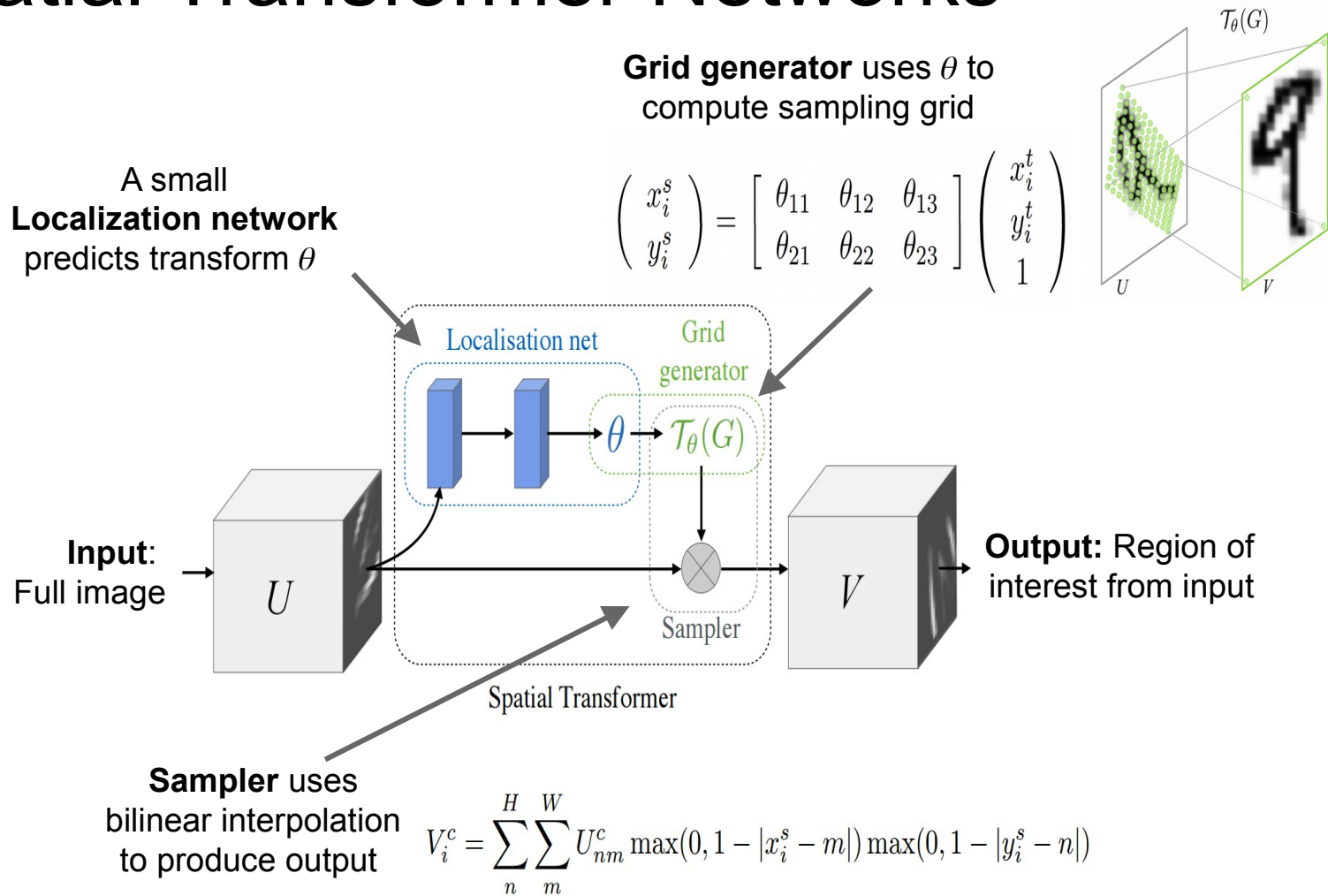
$\mathcal{T}_\theta(G)$

A small **Localization network** predicts transform $\theta$

Localisation net

Grid generator

$\theta \rightarrow \mathcal{T}_\theta(G)$

**Input**: Full image

$U$

Sampler

$V$

**Output:** Region of interest from input

Spatial Transformer

**Sampler** uses bilinear interpolation to produce output

$$V_i^c = \sum_n^H \sum_m^W U_{nm}^c \max(0, 1 - |x_i^s - m|) \max(0, 1 - |y_i^s - n|)$$

Slide Credit:

# Spatial Transformer Networks

Insert spatial transformers into a classification network and it learns to attend and transform the input

Differentiable "attention / transformation" module

Slide Credit: