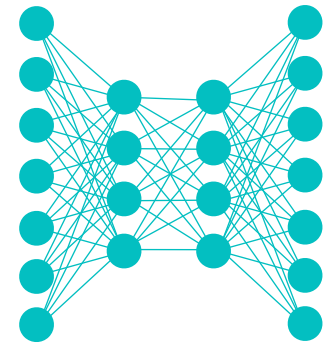Lecture Notes for

# Neural Networks
# and Machine Learning

Fully Convolutional Learning II:
Object Detection

# Logistics and Agenda

- Logistics
  - Lab grading update

- Agenda
  - Lab Town Hall Review and Final Project (if needed)
  - Upsampling
  - Full Convolutional Architectures
    - Semantic Segmentation Basics (last time)
    - Object Detection (this time):
      - RCNN, YOLO
    - Instance Segmentation (next time, probably):
      - Mask-RCNN, YOLACT

# Town Hall Revisited

# Final Project

Me: Predicts the next word correctly

Maybe I am a Language Model

- I have biometric data for pilots…
- Perhaps something to use for the final project…

# Basics: Upsampling Layers

# Last Time



Pooling/ Strided Convolution — Unpooling/Transpose Convolution

Convolutional Encoder-Decoder

Input — RGB Image

Pooling Indices

Output — Segmentation

Conv + Batch Normalisation + ReLU
Pooling — Upsampling — Softmax

Class 1 Proba. — Class 2 Proba. — Class 3 Proba. — ... — Class N Proba.

Channel-wise Softmax

# Decoder Network



Convolutional Encoder-Decoder

Pooling Indices

Conv + Batch Normalisation + ReLU
Pooling     Upsampling     Softmax

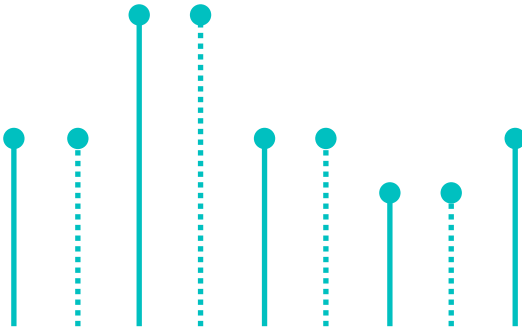Some researcher started calling this **deconvolution**.

If you use that term in this class, **you fail**.

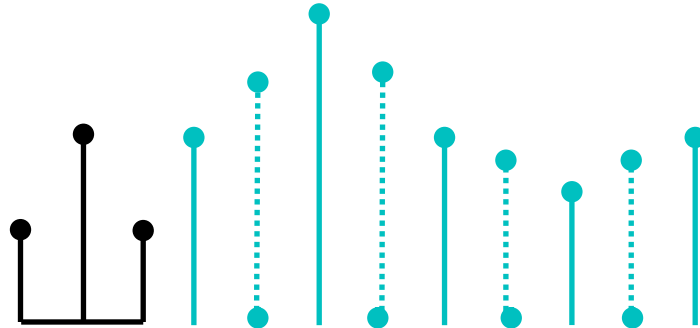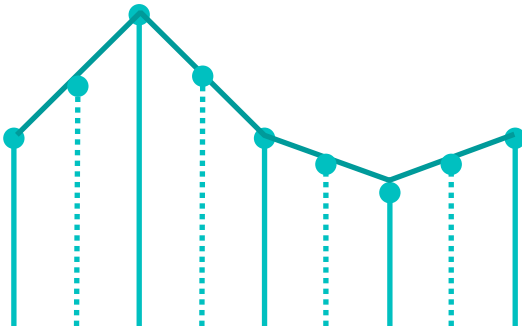This is upsampling and then convolution, but **now the interpolation filters are learned**!!

# Integer Upsampling via Interpolation
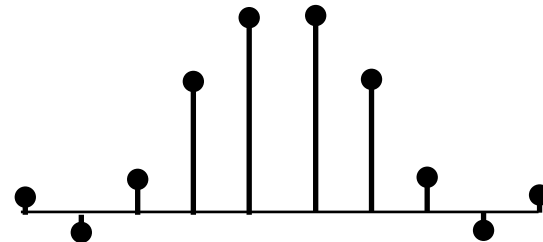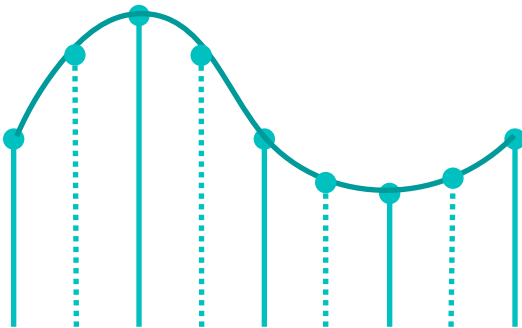
**Nearest Neighbor**

**Linear**

**Cubic**

All are equivalent to inserting zeros and applying convolutional filter

# Image Upsampling, Integer Factor

- Insert Zeros
- Convolve

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

| 1 | | 2 | | 3 | | 4 | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| 5 | | 6 | | 7 | | 8 | |
| | | | | | | | |
| 9 | | 10 | | 11 | | 12 | |
| | | | | | | | |
| 13 | | 14 | | 15 | | 16 | |
| | | | | | | | |

| 0.25 | 0.5 | 0.25 |
|---|---|---|
| 0.5 | 1 | 0.5 |
| 0.25 | 0.5 | 0.25 |

Bilinear Filtering



Bicubic Filter

# Learned Upsampling after Zero Insertion

- Learning the interpolation filter has some caveats:



four pixels + bias

two pixels + bias

one pixel + bias

four pixels + bias

four pixels + bias

3x3 not multiple of stride

4x4 multiple of stride

Bias needs to account for both when different numbers of pixels overlap with the kernel

Multiple of stride ensures that same number of active pixels overlap the kernel.

Stride = 2

# Image Upsampling, Integer Factor



**Nearest Neighbor**

UpSampling2D()

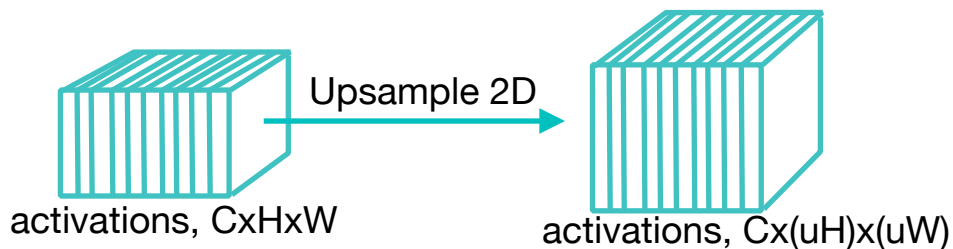**Bilinear**

UpSampling2D(interpolation='bilinear')

**Bicubic**

**Many Types of Upsampling, with varying computational cost:**

area, bicubic, gaussian, lanczos3, lanczos5, mitchellcubic



Upsample 2D

activations, CxHxW

activations, Cx(uH)x(uW)

# What about transpose convolution?

Convolution as Matrix Multiplication

| | | | | |
|---|---|---|---|---|
| $y$ | $x$ | $0$ | $0$ | $0$ |
| $z$ | $y$ | $x$ | $0$ | $0$ |
| $0$ | $z$ | $y$ | $x$ | $0$ |
| $0$ | $0$ | $z$ | $y$ | $x$ |
| $0$ | $0$ | $0$ | $z$ | $y$ |

x

| |
|---|
| $0$ |
| $a$ |
| $b$ |
| $c$ |
| $0$ |

=

| |
|---|
| $ax$ |
| $ay+bx$ |
| $az+by+cx$ |
| $bz+cy$ |
| $cz$ |

Transpose

| | | | | |
|---|---|---|---|---|
| $y$ | $z$ | $0$ | $0$ | $0$ |
| $x$ | $y$ | $z$ | $0$ | $0$ |
| $0$ | $x$ | $y$ | $z$ | $0$ |
| $0$ | $0$ | $x$ | $y$ | $z$ |
| $0$ | $0$ | $0$ | $x$ | $y$ |

x

| |
|---|
| $0$ |
| $a$ |
| $b$ |
| $c$ |
| $0$ |

=

| |
|---|
| $az$ |
| $ay+bz$ |
| $ax+by+cz$ |
| $bx+cy$ |
| $cx$ |

like convolving with "reversed coefficients"

**Regular Convolution**

$b$

$a$

$c$

$z$ $y$ $x$

$ax$  $ay+bx$  $az+by+cx$  $bz+cy$  $cz$

**Transpose Convolution**

$b$

$a$

$c$

$x$ $y$ $z$

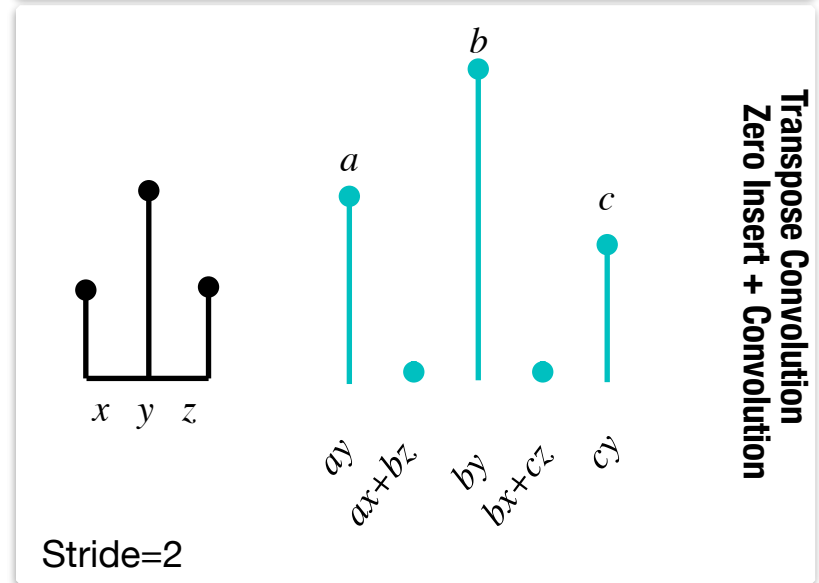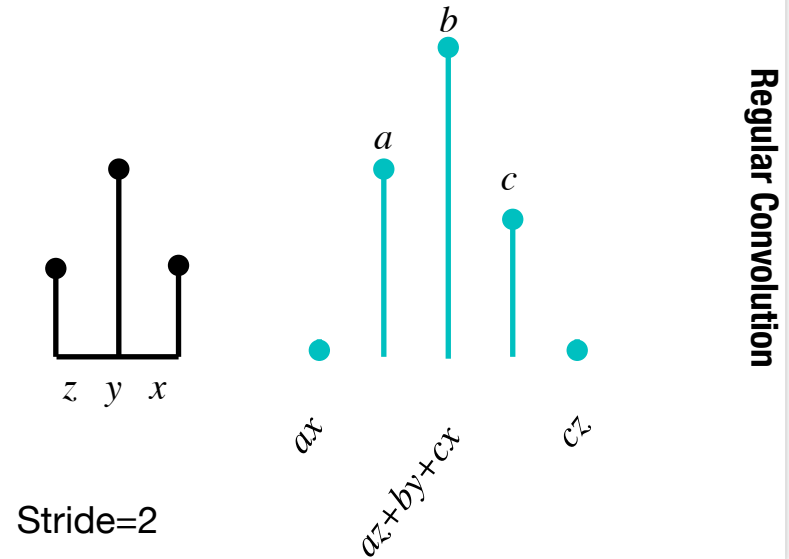$az$  $ay+bz$  $ax+by+cz$  $bx+cy$  $cx$

# Transpose Convolution: Strides

Strided Convolution as Matrix Multiplication

$$
\begin{bmatrix}
y & x & 0 & 0 & 0 \\
0 & z & y & x & 0 \\
0 & 0 & 0 & z & y
\end{bmatrix}
\times
\begin{bmatrix}
0 \\ a \\ b \\ c \\ 0
\end{bmatrix}
=
\begin{bmatrix}
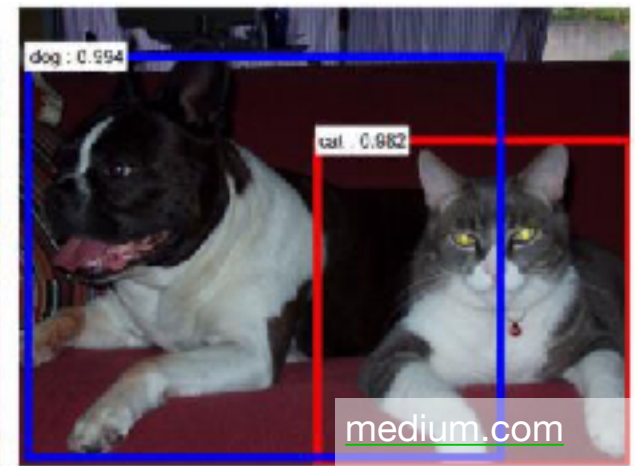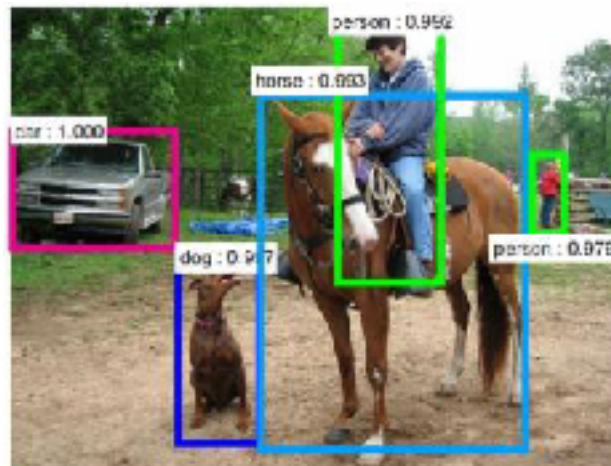ax \\ az+by+cx \\ cz
\end{bmatrix}
$$

Transpose

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

**Regular Convolution**

Stride=2

**Transpose Convolution
Zero Insert + Convolution**

Stride=2

# This time… Object Detection Methods

- Semantic segmentation has good mIoU values (up to 90%) but this is exaggerated by background recognition, many classes are <40%
- How to adapt these techniques to get bounding boxes, not semantic segmentations?
  - Could this be easier? More stable?
  - More consistent labeling?
  - Suitable for "higher risk" tracking applications?
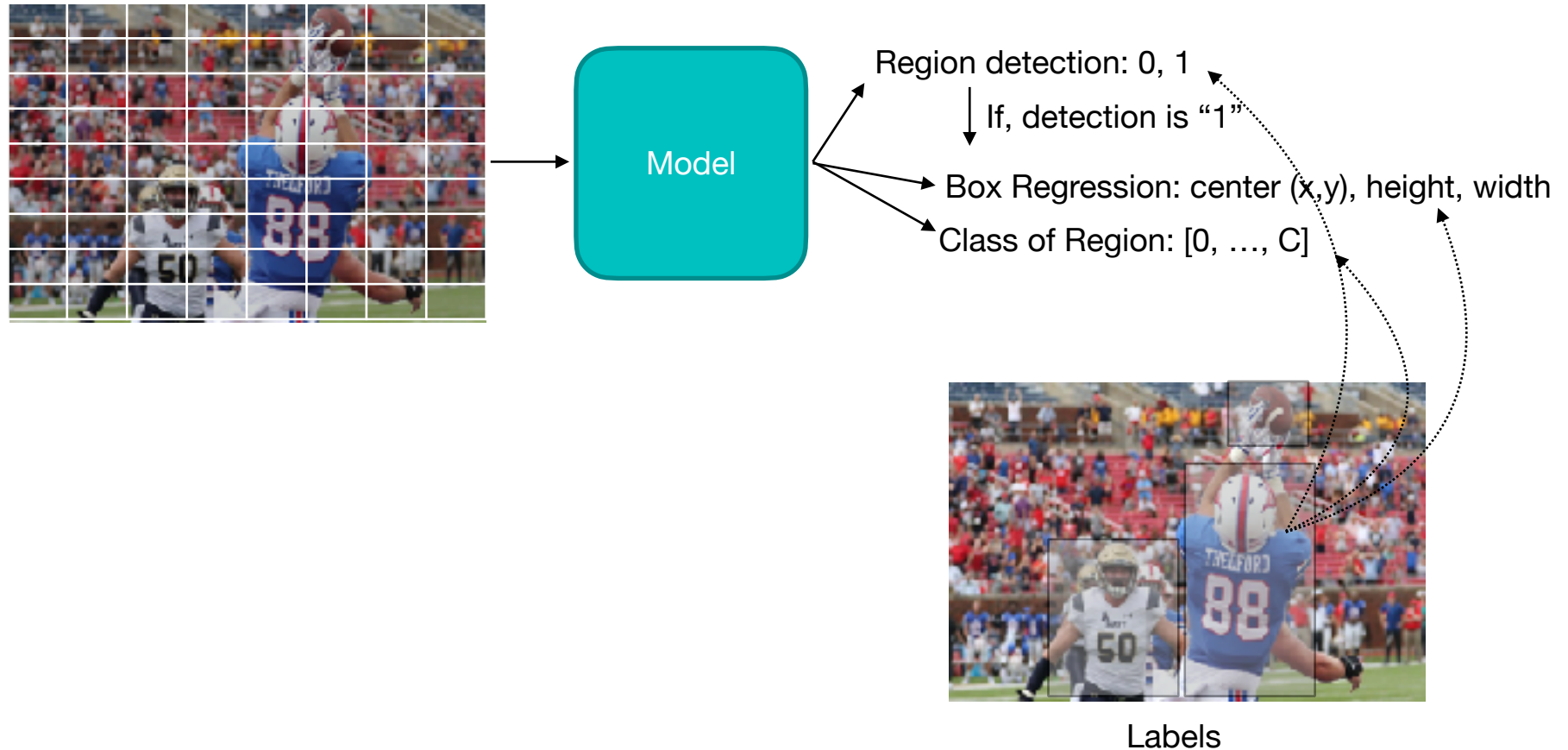


medium.com

# Object Detection with RCNN



# Research!

A history in naming one network five different times
with five different papers
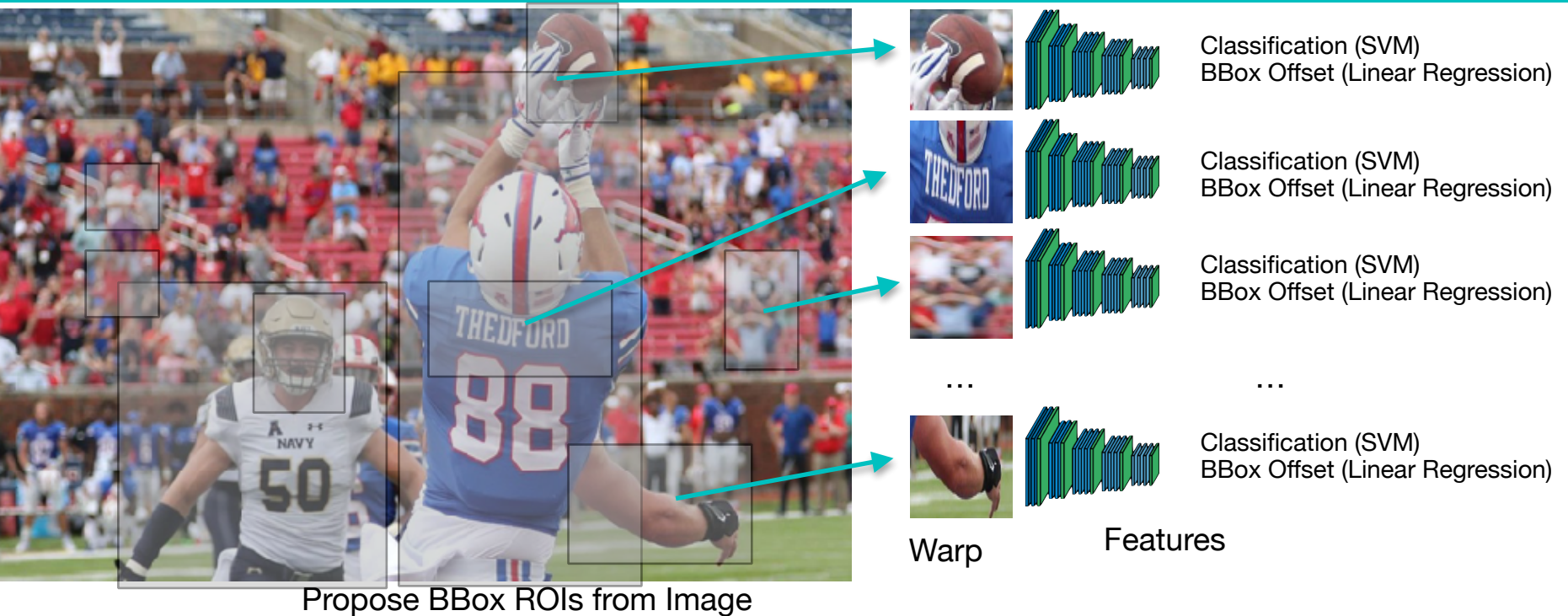each time changing one thing about the architecture

# General Structure



Model

Region detection: 0, 1

If, detection is "1"

Box Regression: center (x,y), height, width

Class of Region: [0, ..., C]

Labels

# 2014: R-CNN



Classification (SVM)
BBox Offset (Linear Regression)

Classification (SVM)
BBox Offset (Linear Regression)

Classification (SVM)
BBox Offset (Linear Regression)

…                    …

Classification (SVM)
BBox Offset (Linear Regression)

Warp            Features

Propose BBox ROIs from Image

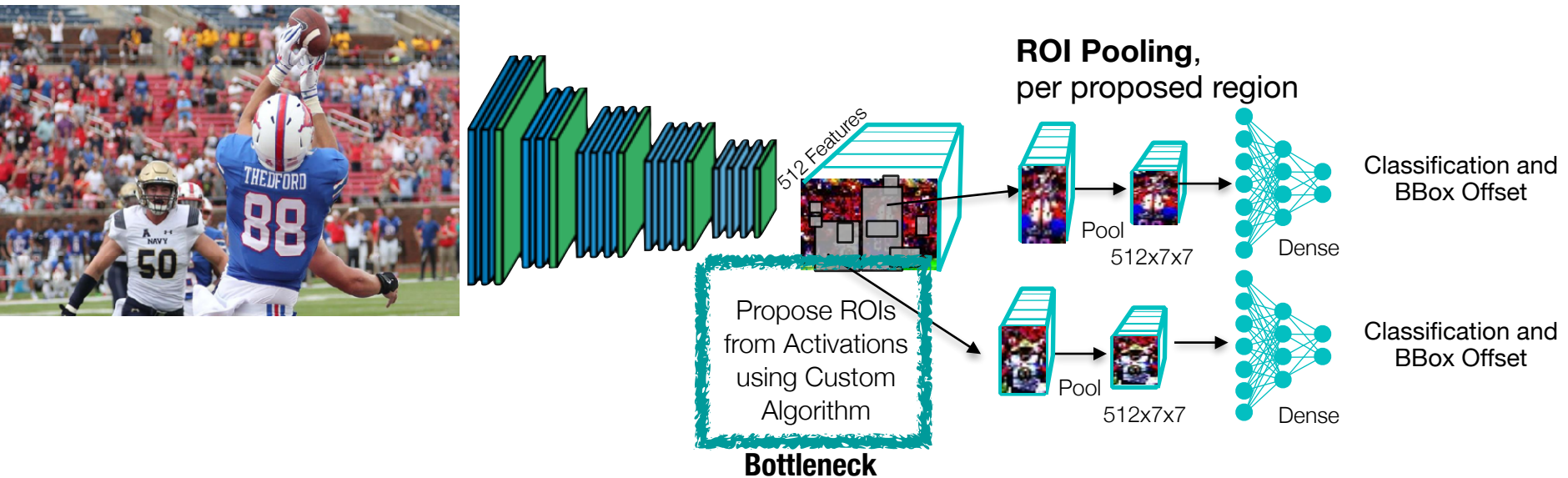- Too Slow to Be Useful

- SVM and BBox Regression Trained Separately

- Fine Tuned Existing ConvNet (for Warped Images)

- ~50 Seconds per Image when Deployed
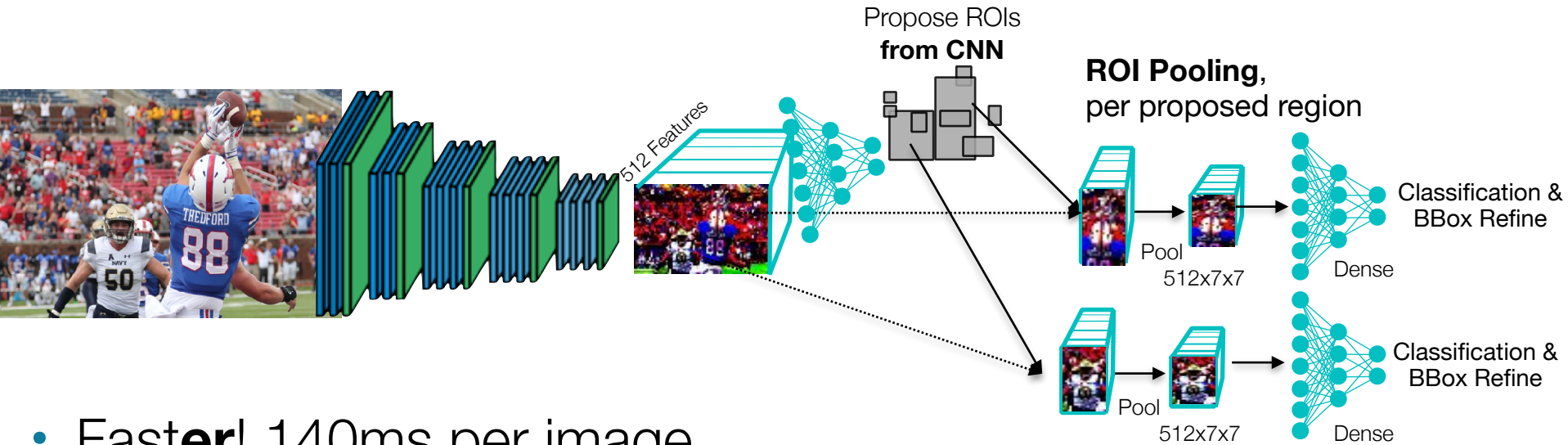
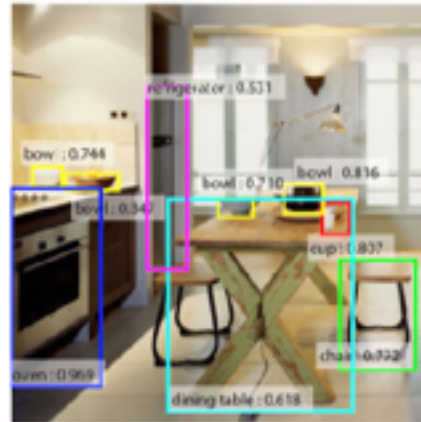Girshick et al. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, 2014

# 2015: Fast R-CNN



- Fast! 2.3 seconds per image (not ~50)
- But still not real time…

Girshick et al. Fast R-CNN, 2015, February

# 2015: Faster R-CNN



Propose ROIs **from CNN**

**ROI Pooling**, per proposed region

512 Features

Pool 512x7x7

Dense

Classification & BBox Refine

Pool 512x7x7

Dense

Classification & BBox Refine

- Fast**er**! 140ms per image (7 FPS)
- Highly Accurate

Ren et al. Faster R-CNN, 2015, November                                                                                    19

Lecture Notes for CS8321 Neural Networks and Machine Learning     |     Professor Eric C. Larson     |

ROI Pooling,
per proposed region

512 Features

Propose ROIs
from CNN

Pool
512x7x7

Classification &
BBox Refine

Dense

Pool
512x7x7

Classification &
BBox Refine

Dense

512 Features

Regress coordinates
of box proposal
to ground truth

Ground Truth Boxes

Actual, Labeled
BBox

Predicted

- Pre-train image classifier backbone (like VGG)
- Train region proposal network
  - Sliding window (1x1 conv) across CNN activations from backbone classifier
  - Regress multiple bounding boxes (usually $k$ proposals)
  - Regress "object-ness" of each box
- Train Fast R-CNN on generated ROI proposals
- Fix weights of classifier pipeline, fine tune RPN
- Fix RPN and fine tune classifier
- Rinse, repeat fine tuning

$$l_{box} = \sum_i \hat{p}_i \left[ (x - \hat{x}_i)^2 + (y - \hat{y}_i)^2 + (\log w - \log \hat{w}_i)^2 + (\log h - \log \hat{h}_i)^2 \right]$$

$$l_{class} = \sum_c CE(c, \hat{c})$$

$$l_{obj} = \sum_i CE(p_i, \hat{p}_i)$$

Ren et al. Faster R-CNN, 2015, November