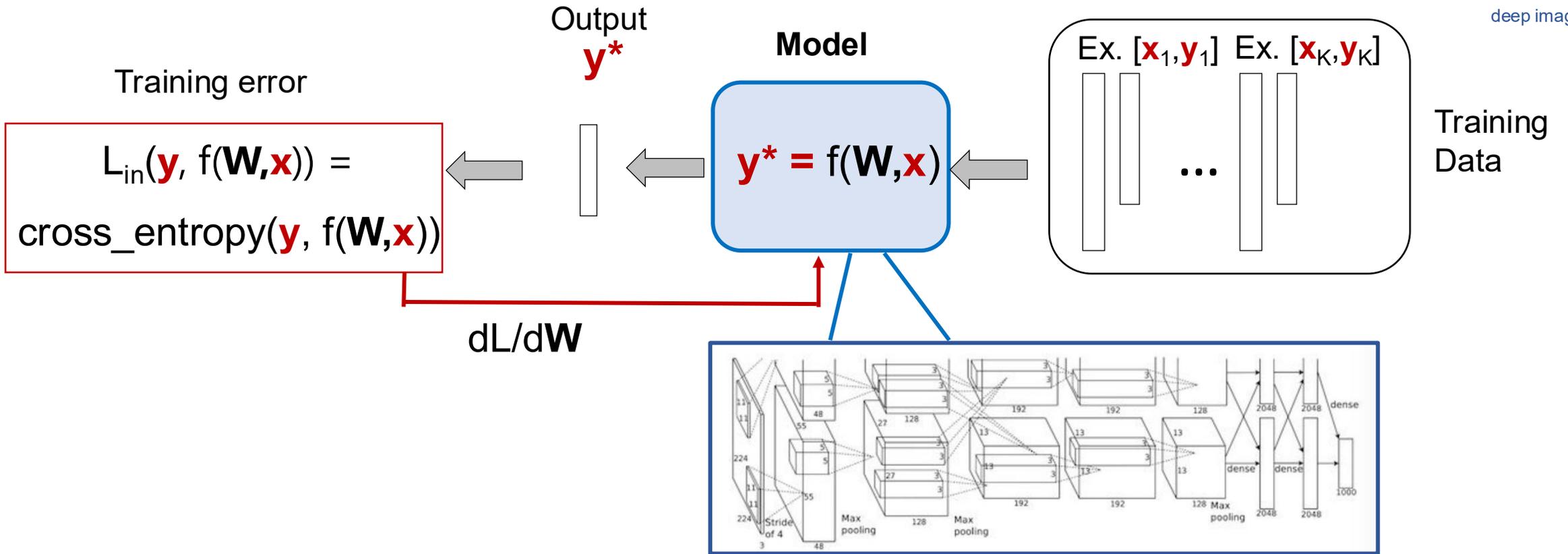


# Lecture 14: Object detection and Segmentation

Machine Learning and Imaging

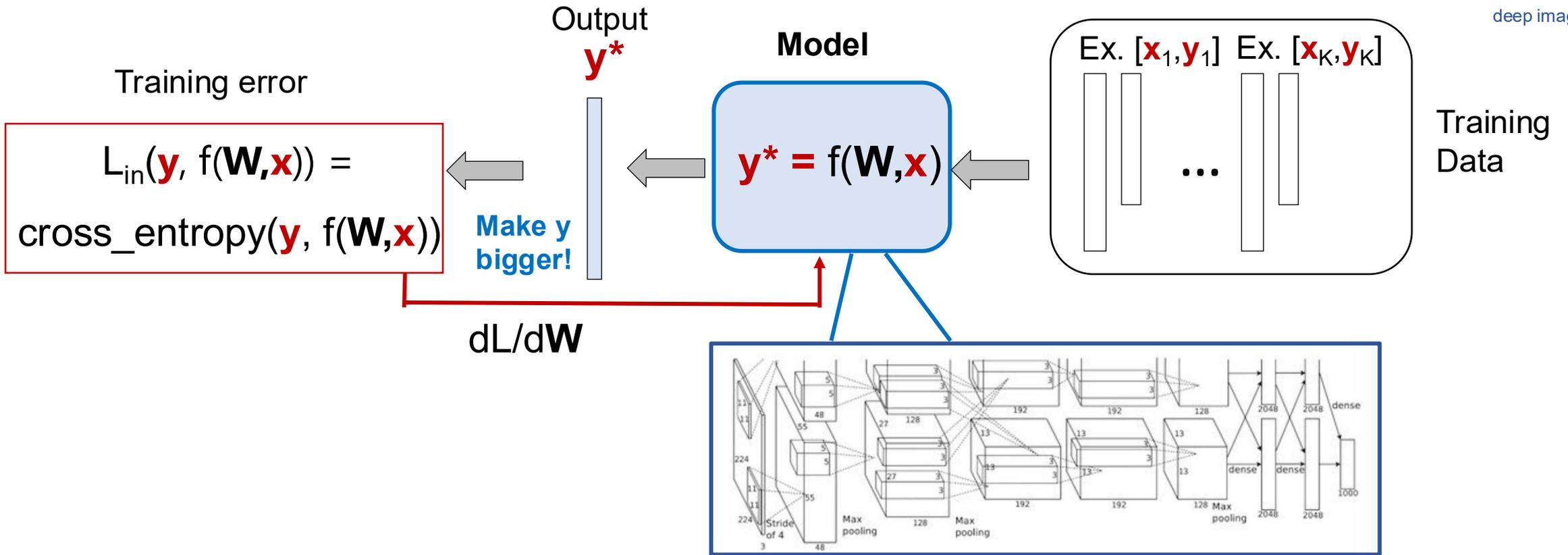
BME 548L  
Roarke Horstmeyer



Dimensional analysis for classification:

Input  $\mathbf{x}$ :  $\sim \mathbb{R}^{1000}$

Output  $\mathbf{y}^*$ :  $\sim \mathbb{R}^2 - \mathbb{R}^{10}$



Dimensional analysis for classification:

Input  $\mathbf{x}$ :  $\sim \mathbb{R}^{1000}$

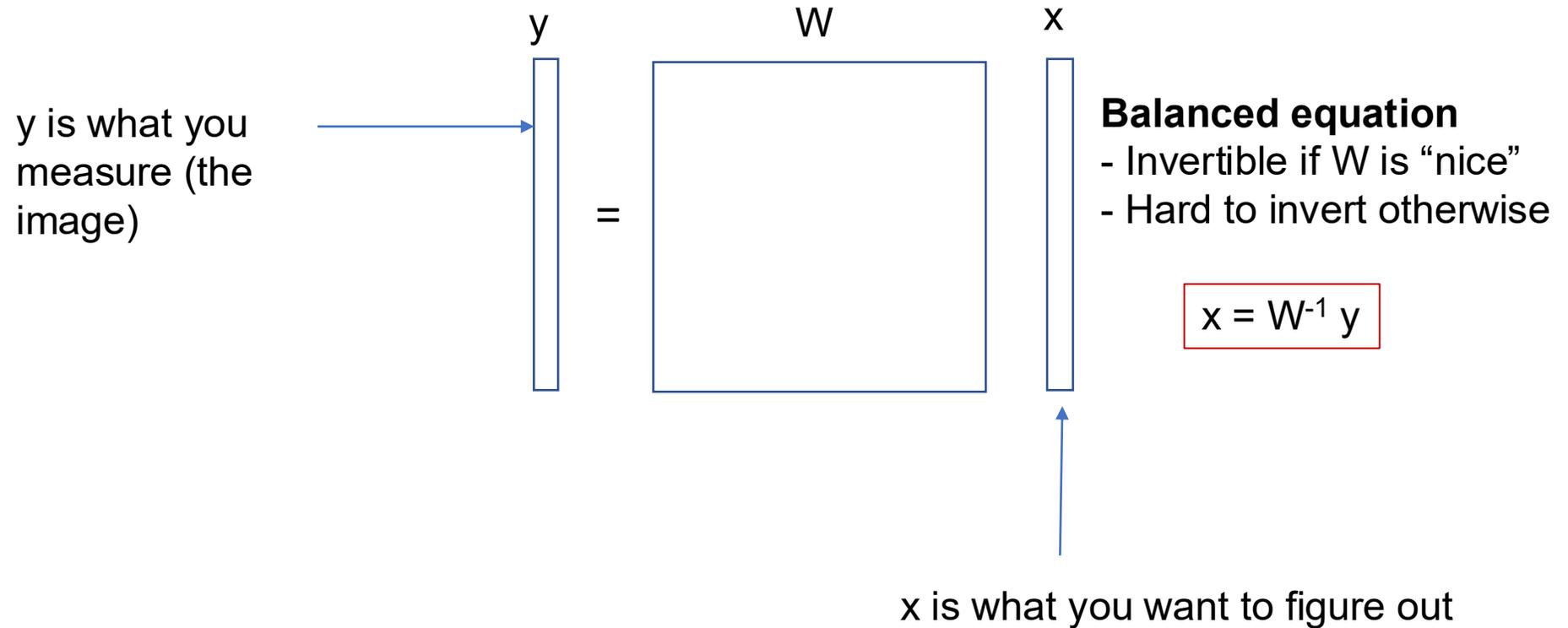
Output  $\mathbf{y}^*$ :  $\sim \mathbb{R}^2 - \mathbb{R}^{10}$



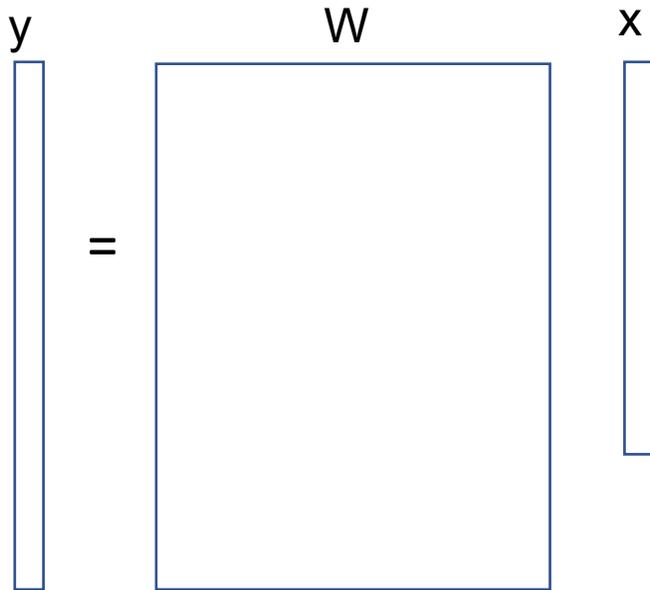
This class – let's make  $\mathbf{y}^*$  bigger!

- Object detection
- Segmentation
- Creating 3D volumes
- Better resolution

# Over-determined, under-determined and balanced inverse equations



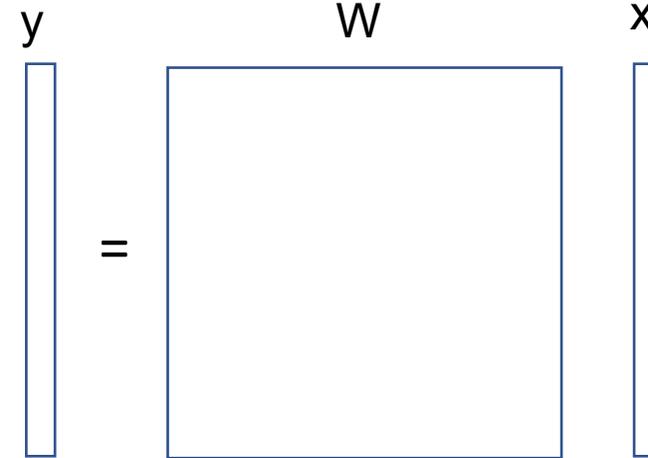
# Over-determined, under-determined and balanced inverse equations



## Over-determined equation

- Unique solution can exist
- If not, it's easy to get close
- Good place – more measurements than unknowns

$$x = W^+ y$$

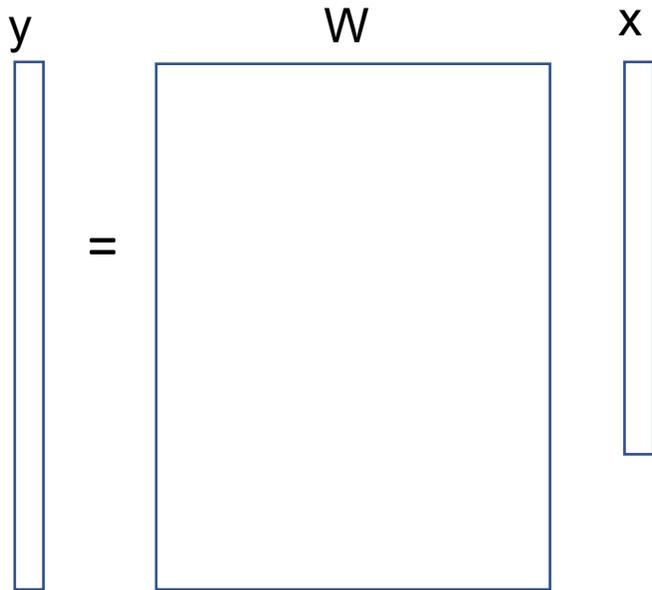


## Balanced equation

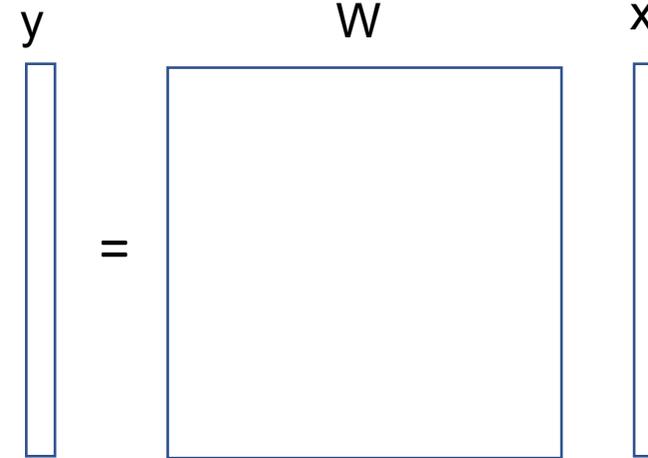
- Invertible if  $W$  is “nice”
- Hard to invert otherwise

$$x = W^{-1} y$$

# Over-determined, under-determined and balanced inverse equations

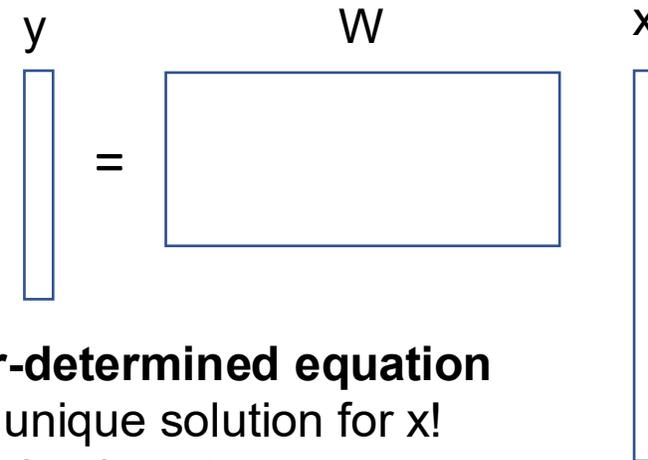


- Over-determined equation**
- Unique solution can exist
  - If not, it's easy to get close
  - Good place – more measurements than unknowns



- Balanced equation**
- Invertible if  $W$  is “nice”
  - Hard to invert otherwise

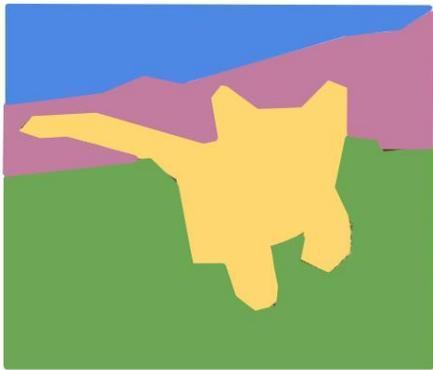
$$x = W^{-1} y$$



- Under-determined equation**
- No unique solution for  $x$ !
  - Hard to invert
  - Not a good place to be

# Other Computer Vision Tasks

## Semantic Segmentation

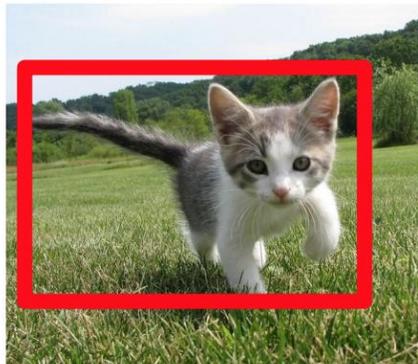


GRASS, CAT,  
TREE, SKY

No objects, just pixels

**Balanced equation**

## Classification + Localization

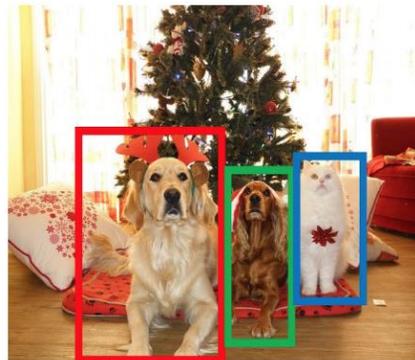


CAT

Single Object

**Over-determined**

## Object Detection



DOG, DOG, CAT

Multiple Object

**Over-determined**

## Instance Segmentation



DOG, DOG, CAT

**Over-determined**

## Super-resolution

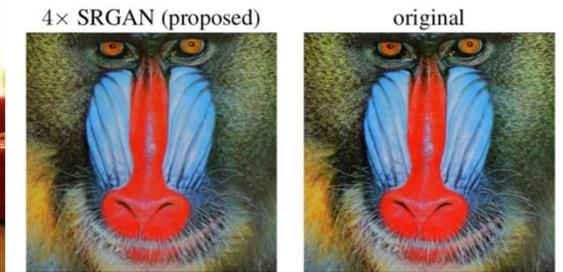
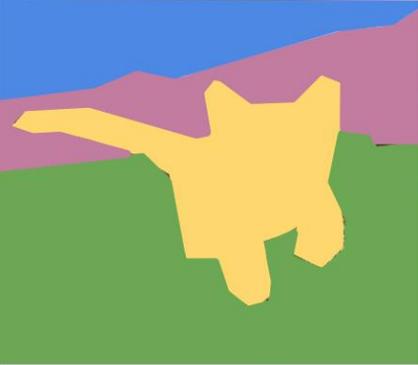
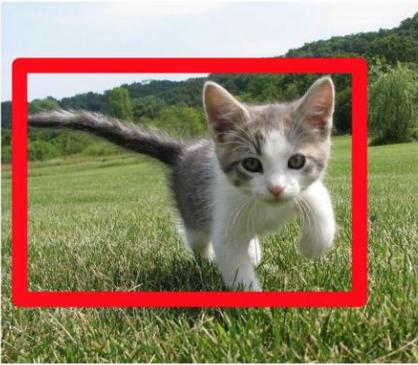
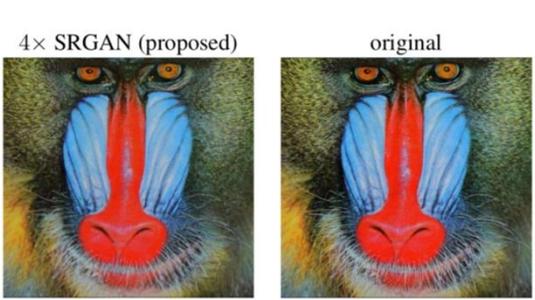


Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4x upscaling]

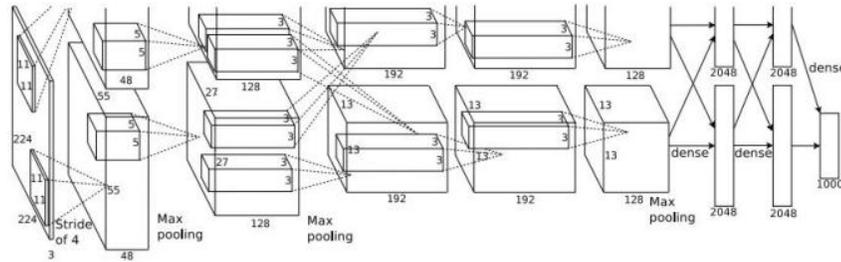
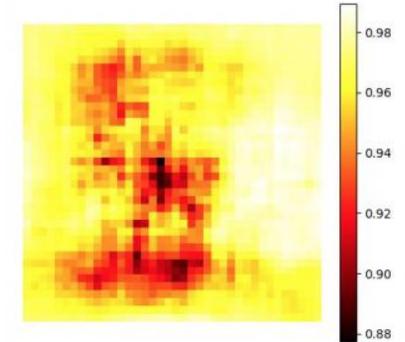
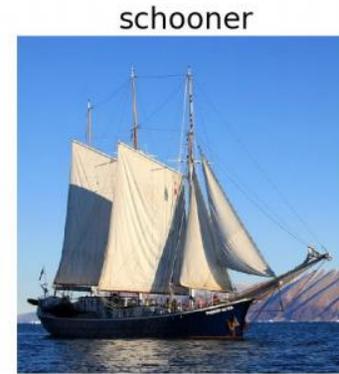
**Under-determined**

# Other Computer Vision Tasks

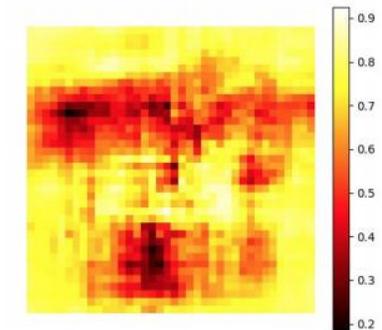
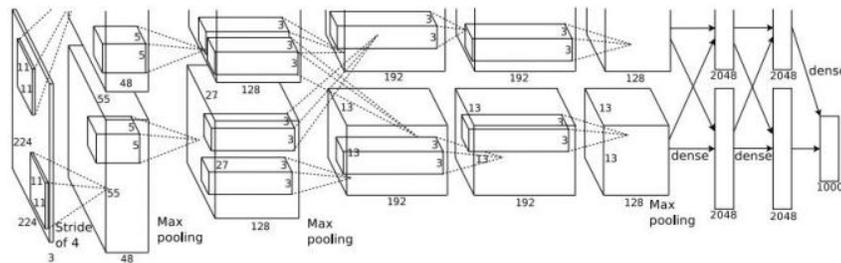
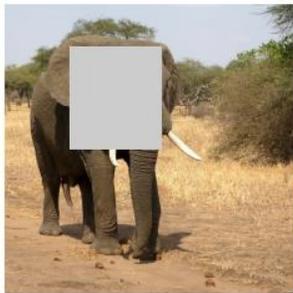
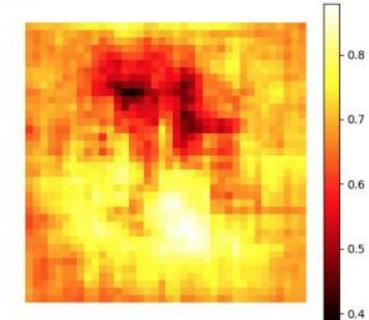
Semantic Segmentation	Classification + Localization	Object Detection	Instance Segmentation	Super-resolution
 <p>GRASS, CAT, TREE, SKY</p> <p>No objects, just pixels</p>	 <p>CAT</p> <p>Single Object</p>	 <p>DOG, DOG, CAT</p>	 <p>DOG, DOG, CAT</p>	 <p>4x SRGAN (proposed)      original</p> <p>Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4x upscaling]</p>
<p>Balanced equation</p>	<p>Over-determined</p>	<p>Over-determined</p>	<p>Over-determined</p>	<p>Under-determined</p>

This image is CC0 public domain

# Approach #1: Sliding window + occlusion map (last lecture)



African elephant, *Loxodonta africana*

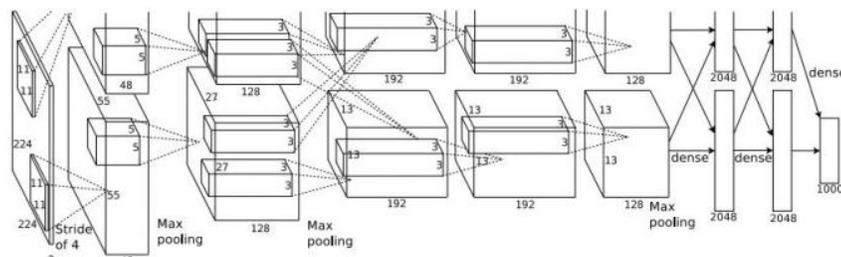
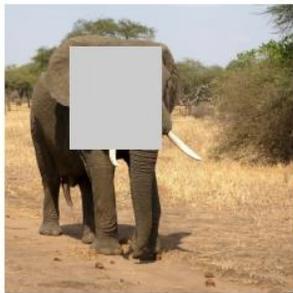
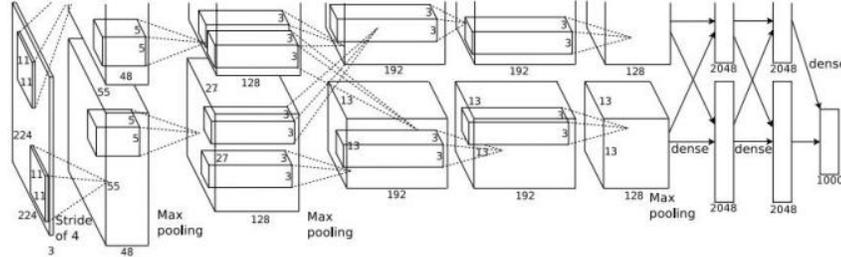
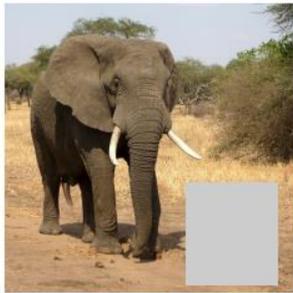
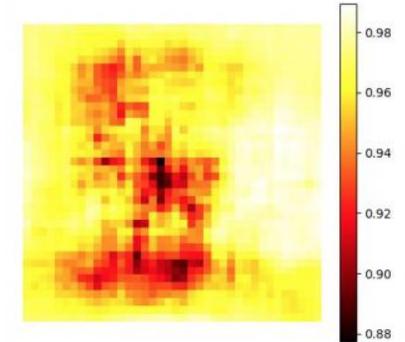
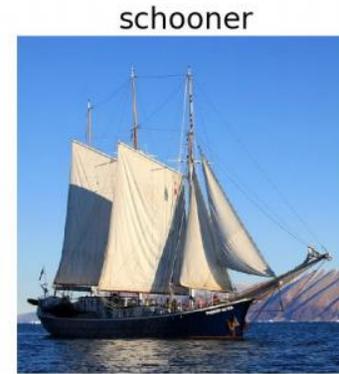


Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

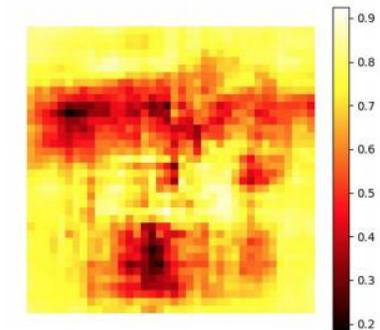
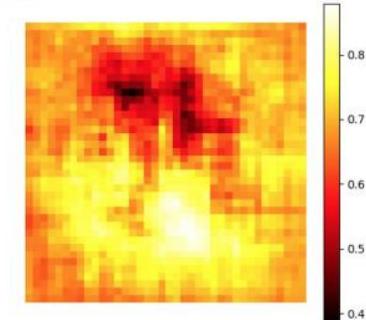
Boat image is [CC0 public domain](#)  
 Elephant image is [CC0 public domain](#)  
 Go-Karts image is [CC0 public domain](#)

# Approach #1: Sliding window + occlusion map (last lecture)

Problem: Inefficient – not sharing information between different sliding window positions (even w/ lots of overlap)



African elephant, *Loxodonta africana*



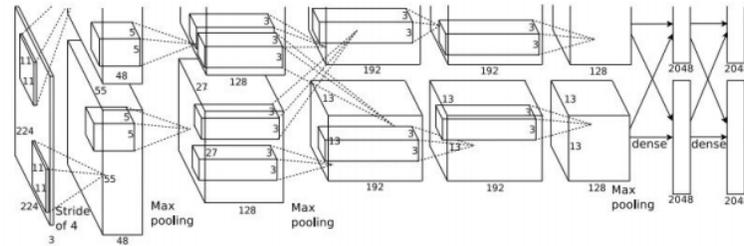
Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

Boat image is [CC0 public domain](#)  
 Elephant image is [CC0 public domain](#)  
 Go-Karts image is [CC0 public domain](#)

# Classification + Localization



This image is CC0 public domain



**Fully Connected:**  
4096 to 1000

**Class Scores**  
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

**Vector:** 4096  
**Fully Connected:** 4096 to 4

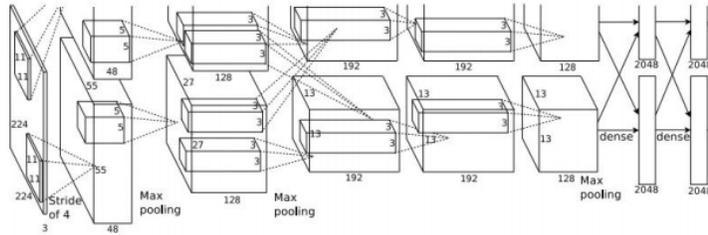
**Box Coordinates**  
(x, y, w, h)

Treat localization as a regression problem!

# Classification + Localization



This image is CC0 public domain



Fully Connected:  
4096 to 1000

**Class Scores**

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

Correct label:  
Cat

**Softmax Loss**

Vector: 4096  
Fully Connected:  
4096 to 4

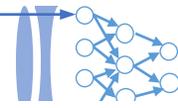
**Box Coordinates**  
(x, y, w, h)

**L2 Loss**

Correct box:  
(x', y', w', h')

Treat localization as a regression problem!



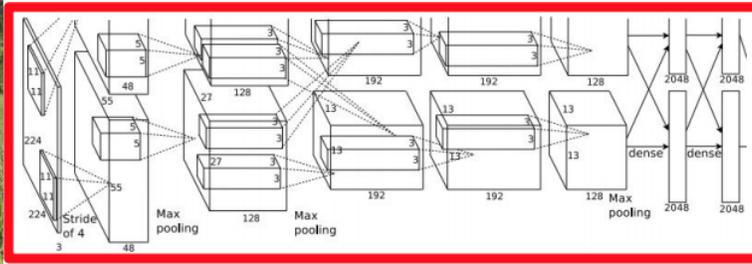


ing

# Classification + Localization



This image is CC0 public domain



Often pretrained on ImageNet (Transfer learning)

Vector: 4096

Fully Connected: 4096 to 1000

## Class Scores

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

Correct label: Cat

Softmax Loss

+

Loss

Fully Connected: 4096 to 4

Box Coordinates (x, y, w, h)

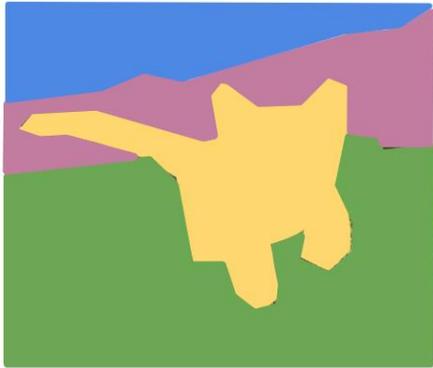
L2 Loss

Correct box: (x', y', w', h')

Treat localization as a regression problem!

# Other Computer Vision Tasks

## Semantic Segmentation

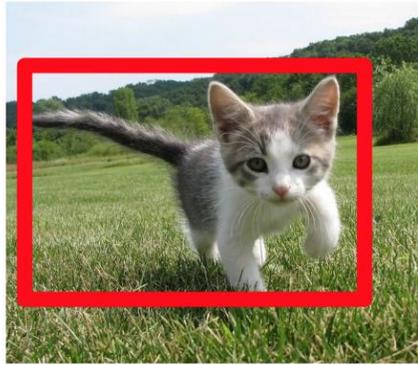


GRASS, CAT,  
TREE, SKY

No objects, just pixels

**Balanced equation**

## Classification + Localization

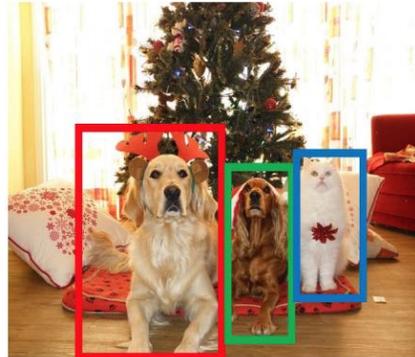


CAT

Single Object

**Over-determined**

## Object Detection

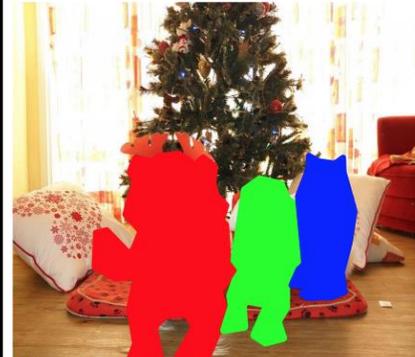


DOG, DOG, CAT

Multiple Object

**Over-determined**

## Instance Segmentation



DOG, DOG, CAT

**Over-determined**

## Super-resolution

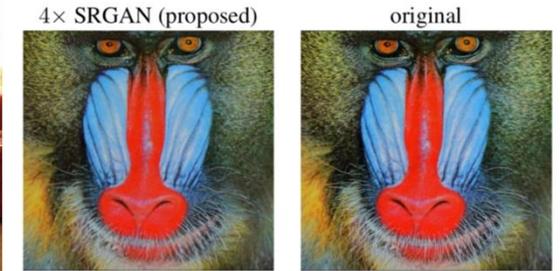
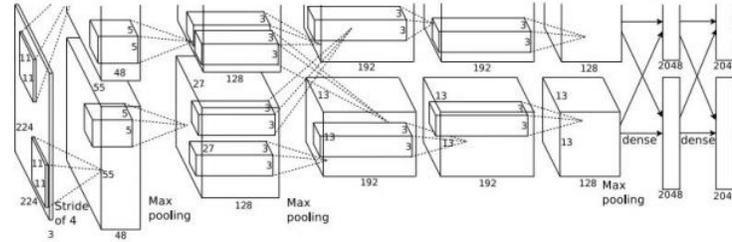


Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4x upscaling]

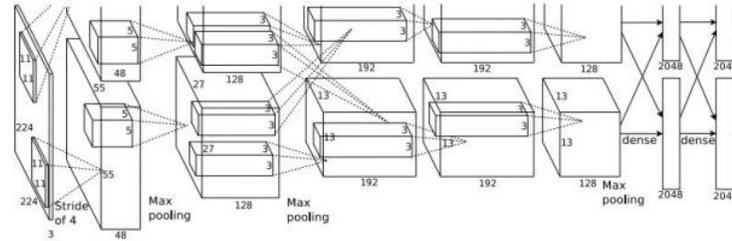
**Under-determined**

# Object Detection as Regression?

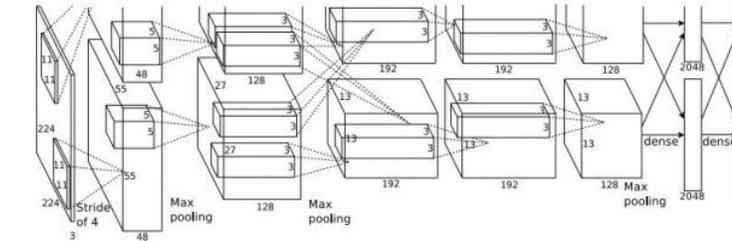
Each image needs a different number of outputs!



CAT:  $(x, y, w, h)$       4 numbers



DOG:  $(x, y, w, h)$   
 DOG:  $(x, y, w, h)$       16 numbers  
 CAT:  $(x, y, w, h)$

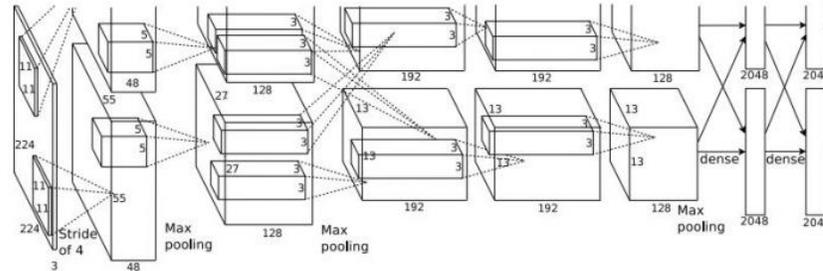


DUCK:  $(x, y, w, h)$       Many numbers!  
 DUCK:  $(x, y, w, h)$       numbers!

....

# Object Detection as Classification: Sliding Window

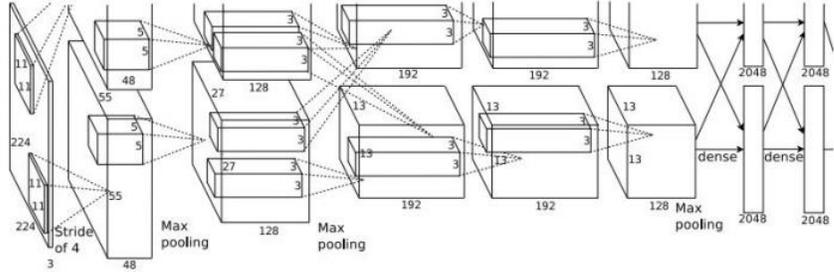
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? NO  
Background? YES

# Object Detection as Classification: Sliding Window

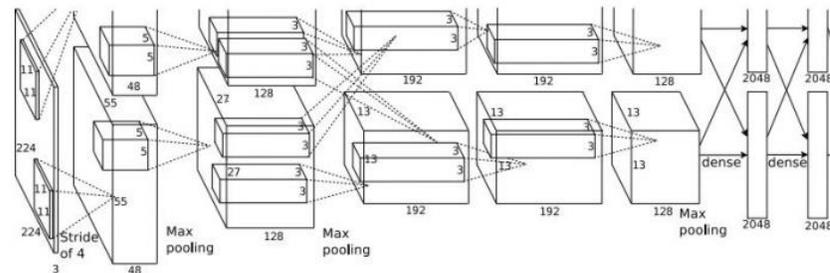
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
 Cat? NO  
 Background? NO

# Object Detection as Classification: Sliding Window

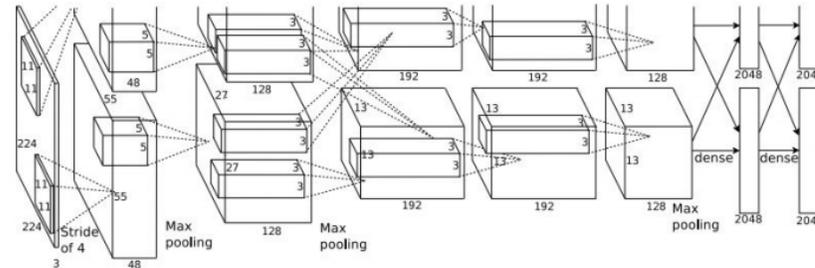
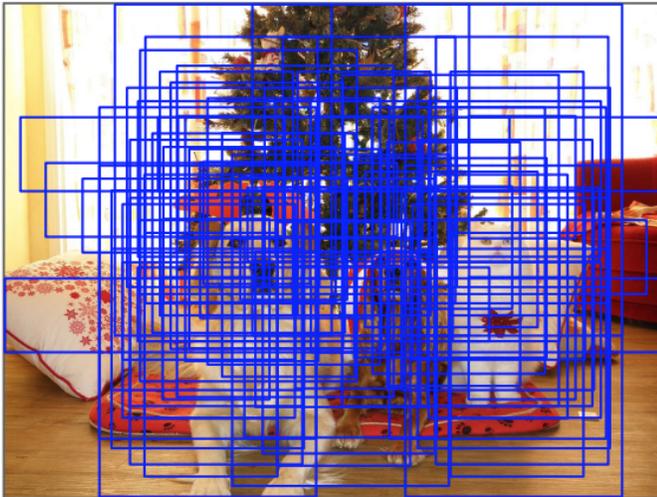
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? YES  
Background? NO

# Object Detection as Classification: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? YES  
Background? NO

**Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!**

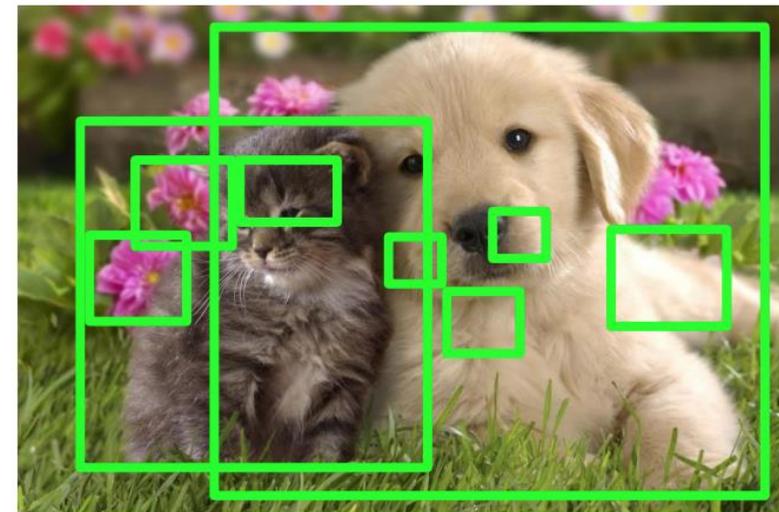
# Solution: First apply a fixed ROI scheme to pull out “blobs” of interest



(Image source: van de Sande et al. ICCV'11)

# Region Proposals / Selective Search

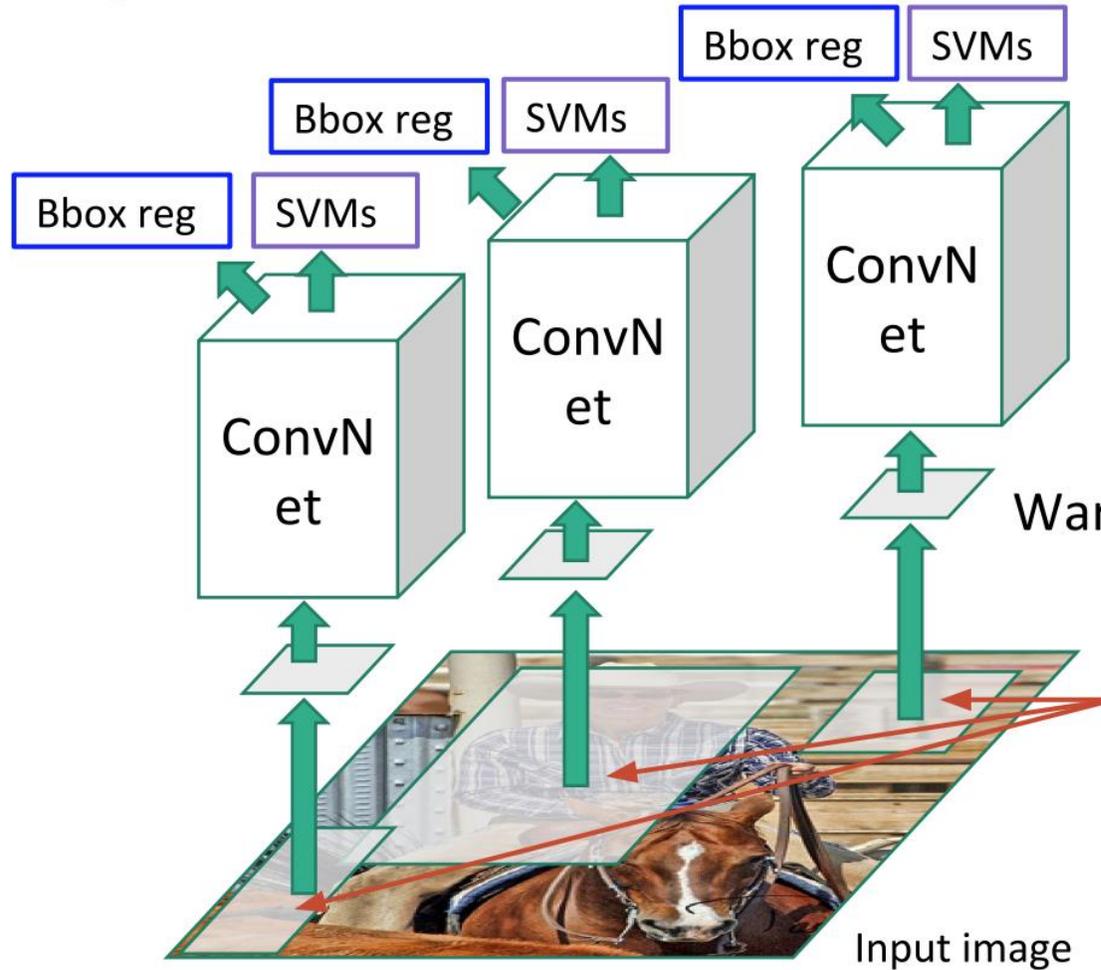
- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



Note: Training dataset has marked boxes, so don't necessarily need to do selective search for training, just evaluation/testing

Alexe et al, "Measuring the objectness of image windows", TPAMI 2012  
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013  
Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014  
Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

# R-CNN



Linear Regression for bounding box offsets

Classify regions with SVMs

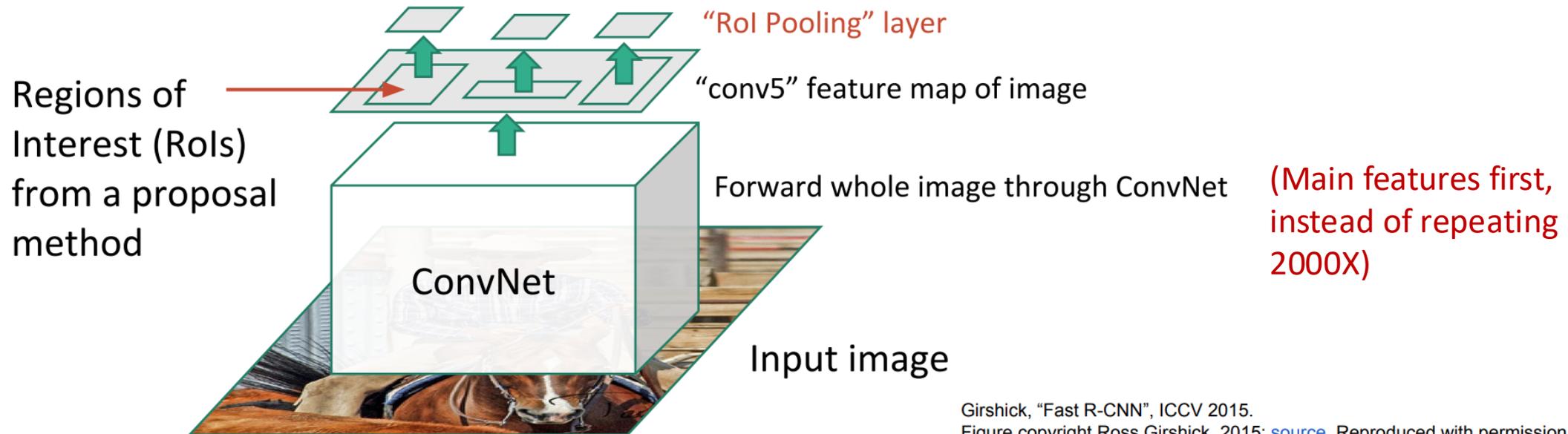
Forward each region through ConvNet

Warped image regions

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

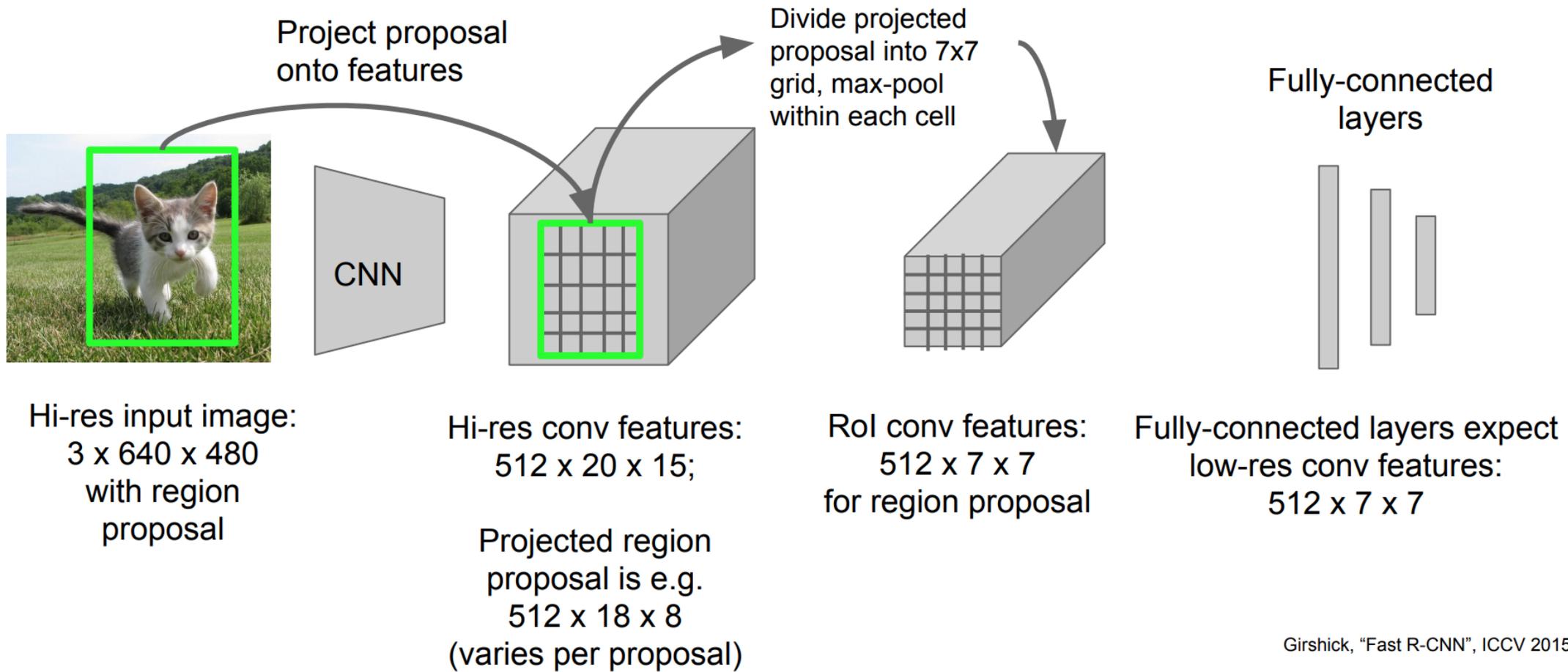
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



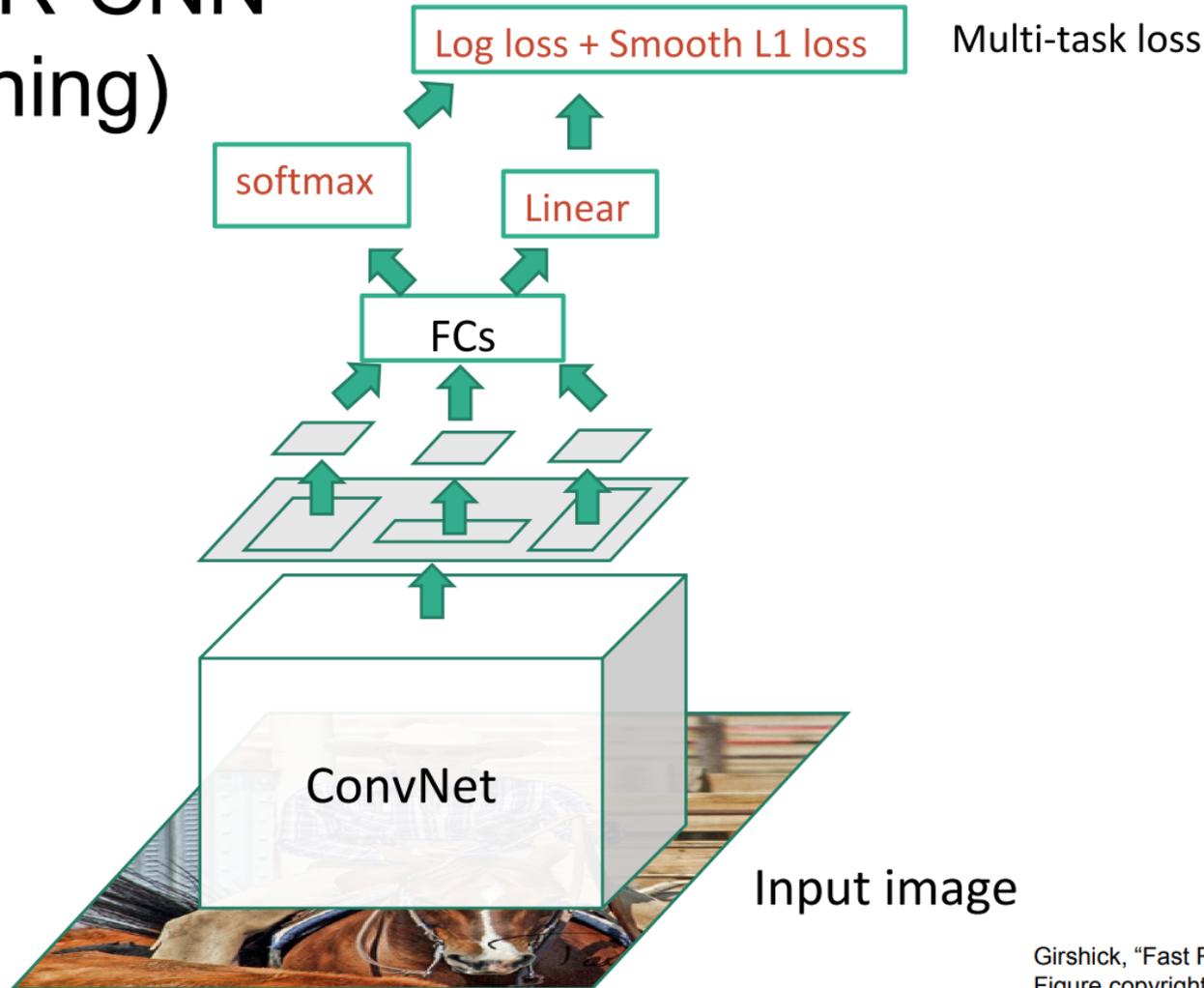
Girshick, "Fast R-CNN", ICCV 2015.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN: RoI Pooling



Girshick, "Fast R-CNN", ICCV 2015.

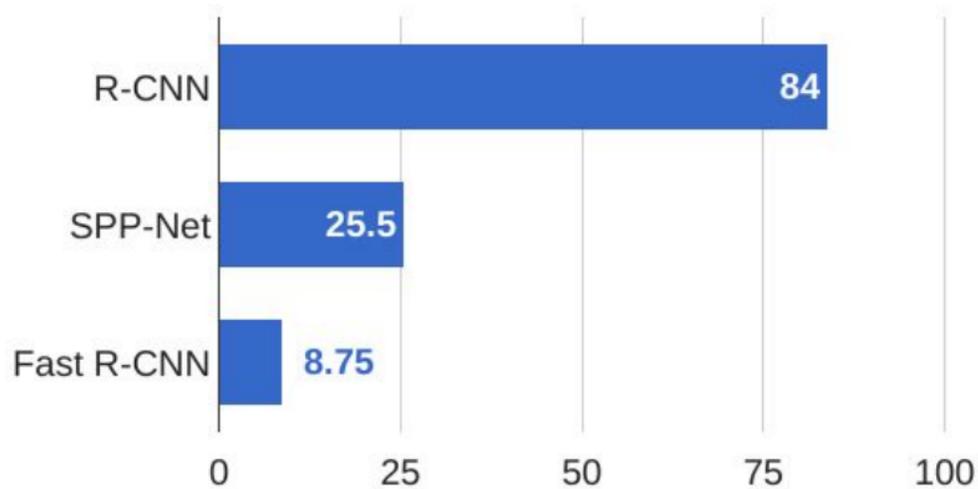
# Fast R-CNN (Training)



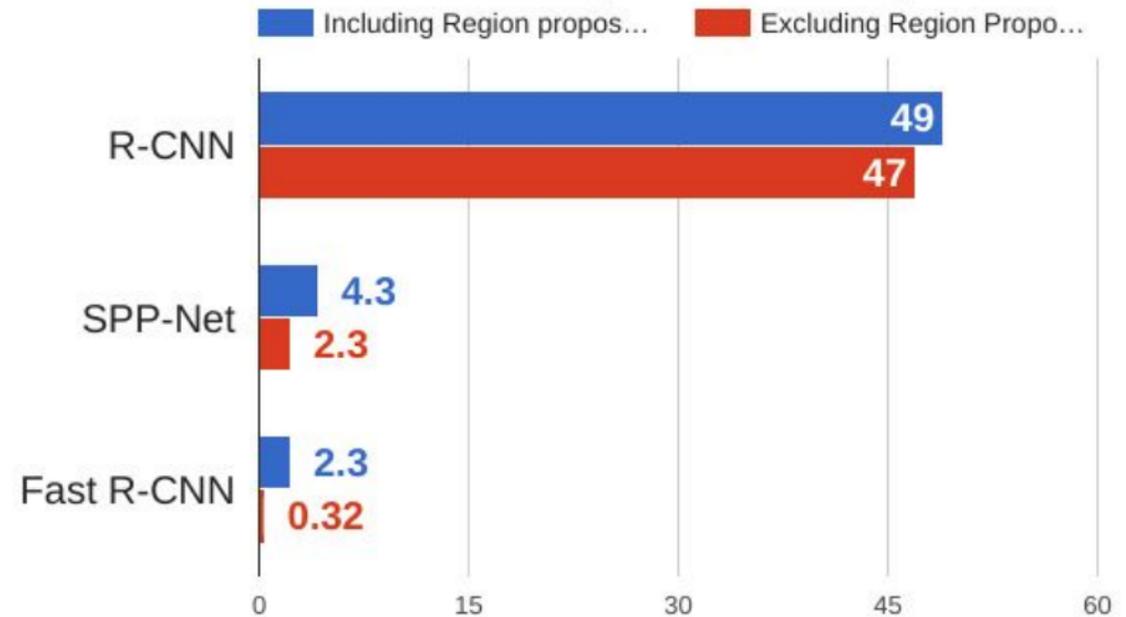
Girshick, "Fast R-CNN", ICCV 2015.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN vs SPP vs Fast R-CNN

### Training time (Hours)



### Test time (seconds)



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014  
 Girshick, "Fast R-CNN", ICCV 2015

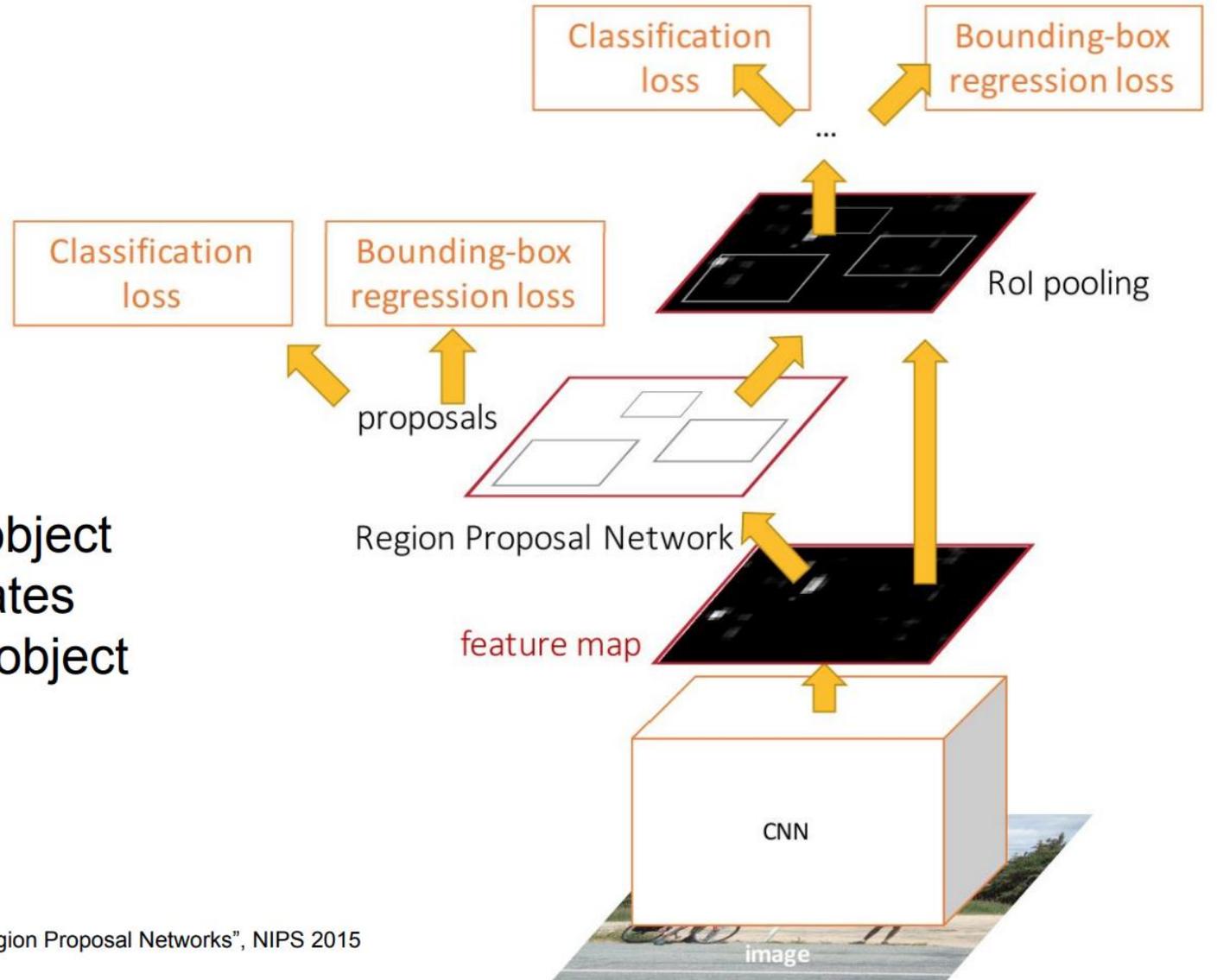
# Faster R-CNN:

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

Jointly train with 4 losses:

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates

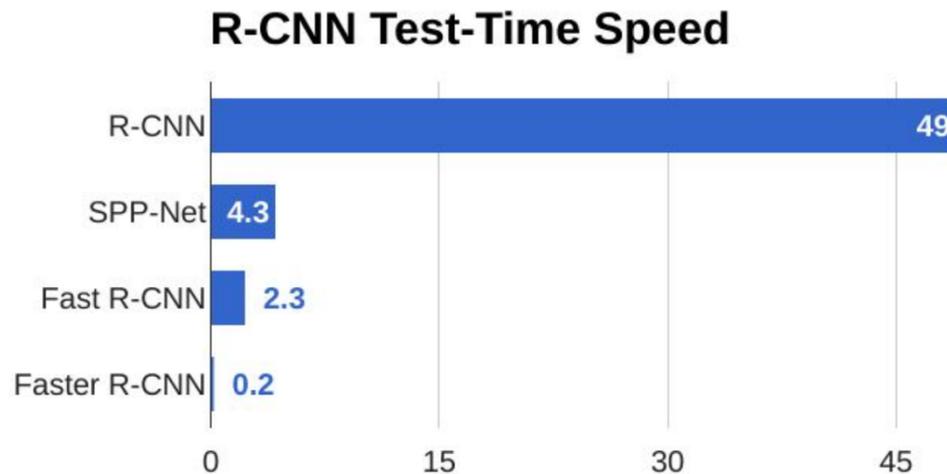


Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015  
 Figure copyright 2015, Ross Girshick; reproduced with permission

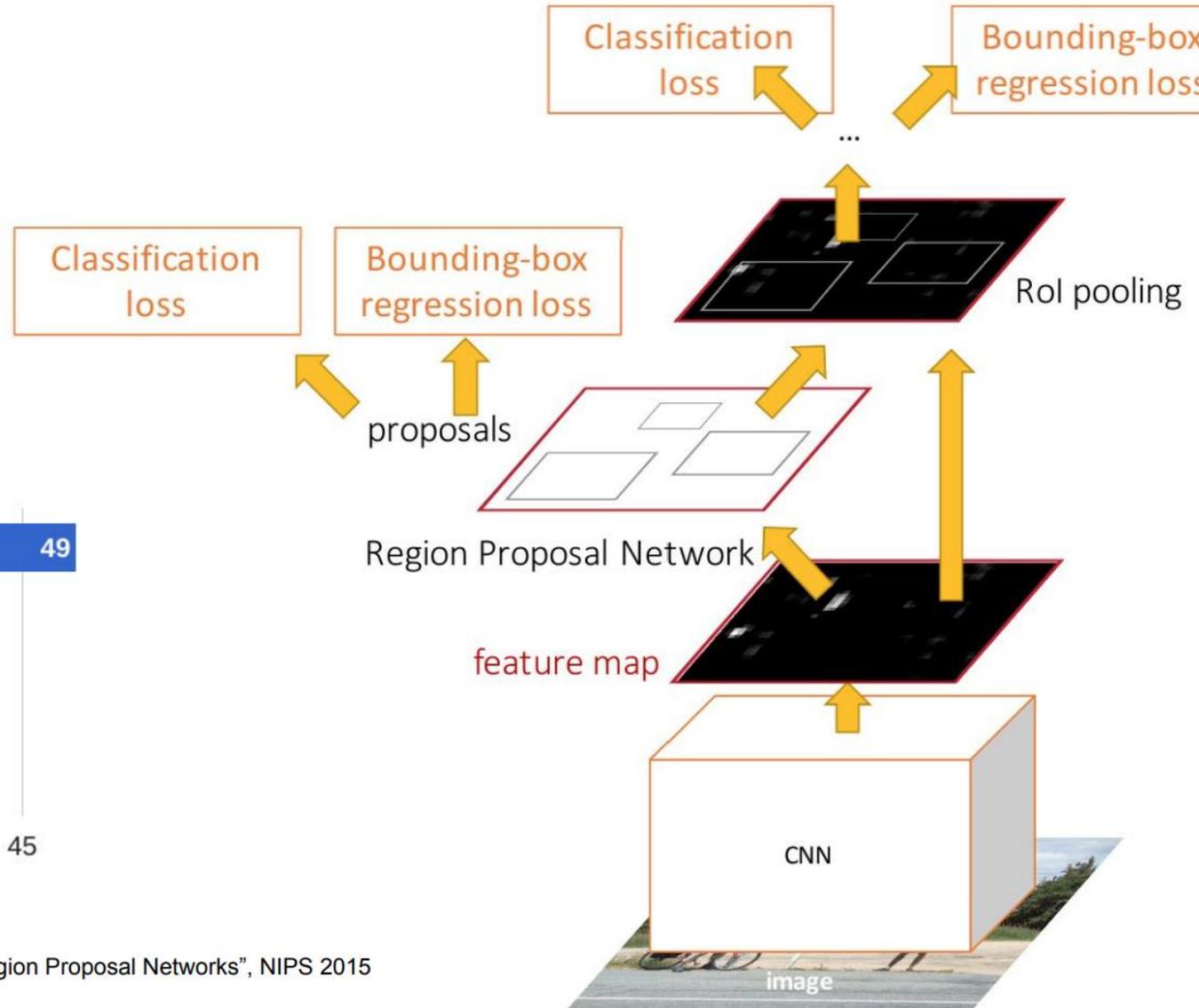
# Faster R-CNN:

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015  
 Figure copyright 2015, Ross Girshick; reproduced with permission

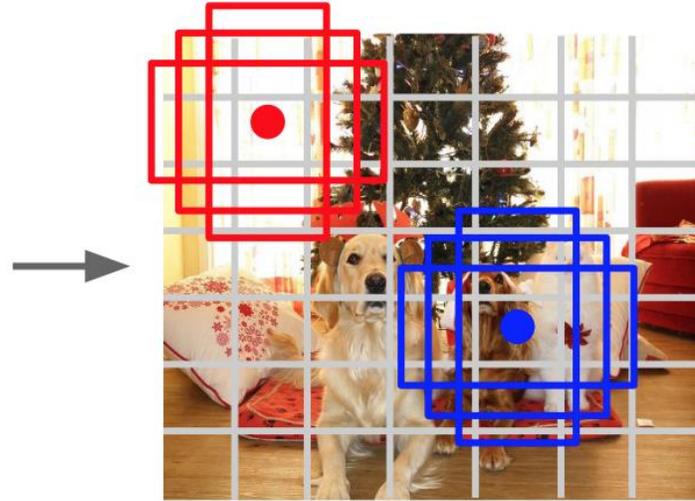


# Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network!



Input image  
3 x H x W



Divide image into grid  
7 x 7

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
(dx, dy, dh, dw, confidence)
- Predict scores for each of  $C$  classes (including background as a class)

Output:  
 $7 \times 7 \times (5 * B + C)$

Redmon et al, "You Only Look Once:  
Unified, Real-Time Object Detection", CVPR 2016  
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

# Object Detection: Impact of Deep Learning

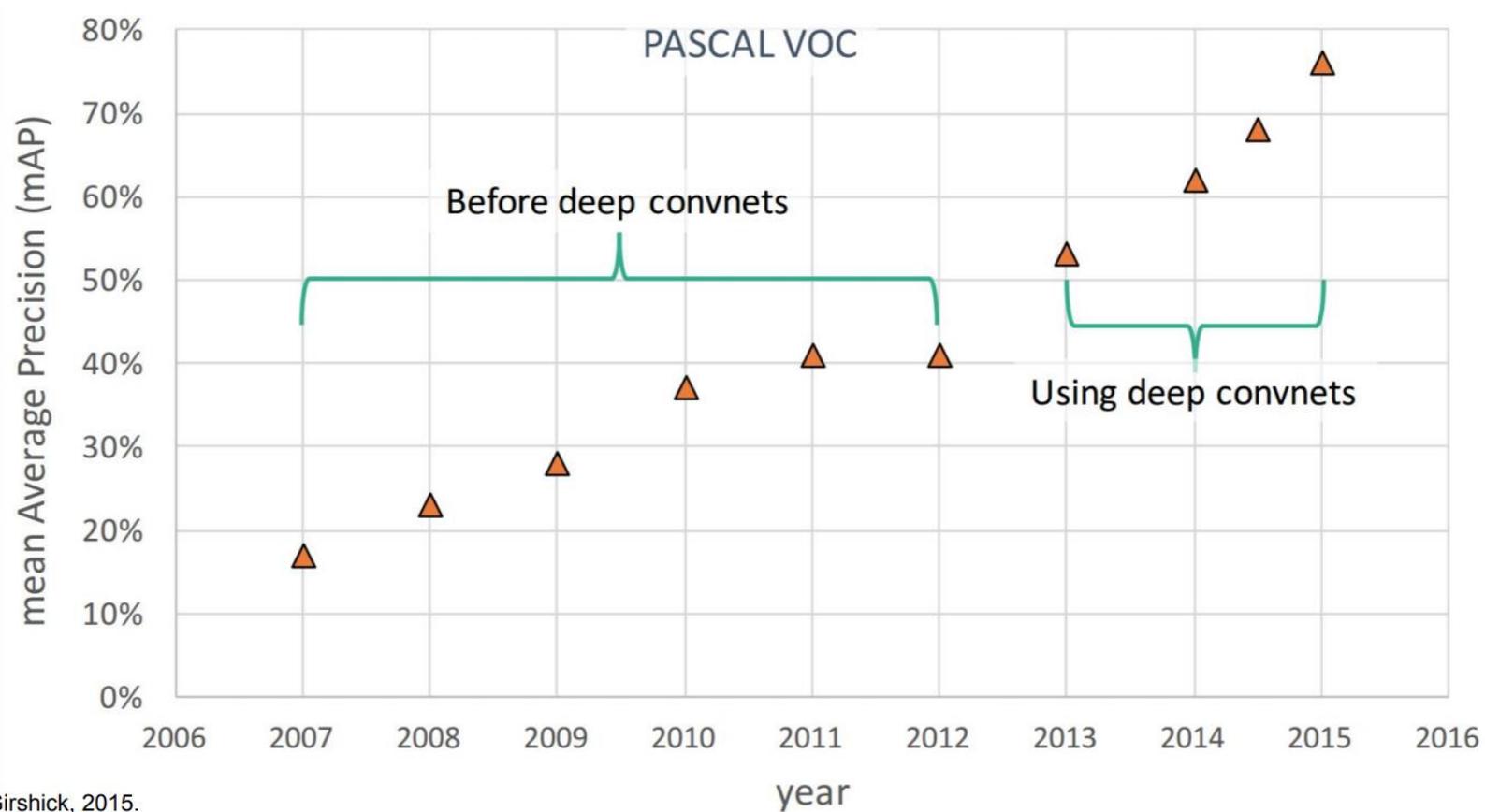


Figure copyright Ross Girshick, 2015.  
Reproduced with permission.

# Other Computer Vision Tasks

**Semantic Segmentation**



GRASS, CAT, TREE, SKY

No objects, just pixels

**Balanced equation**

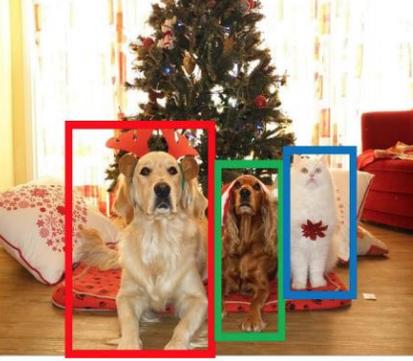
**Classification + Localization**



CAT

Single Object

**Object Detection**



DOG, DOG, CAT

Multiple Object

**Instance Segmentation**



DOG, DOG, CAT

Multiple Object

**Super-resolution**

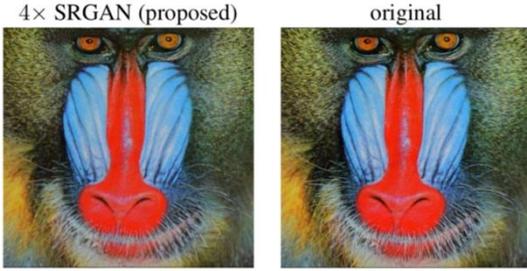


Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4x upscaling]

This image is CC0 public domain

**Over-determined**

**Over-determined**

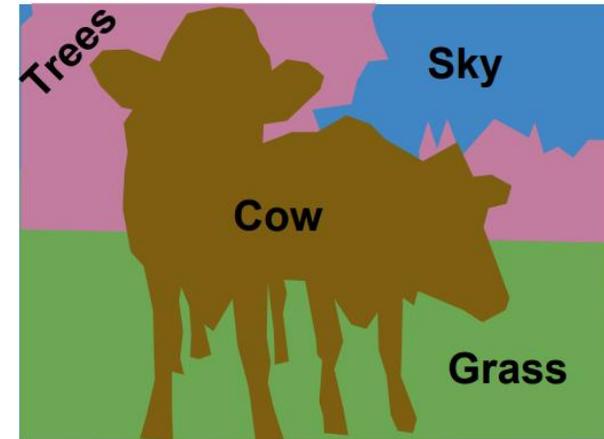
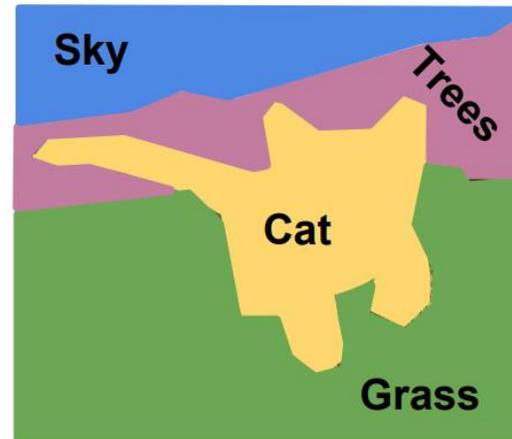
**Over-determined**

**Under-determined**

# Semantic Segmentation

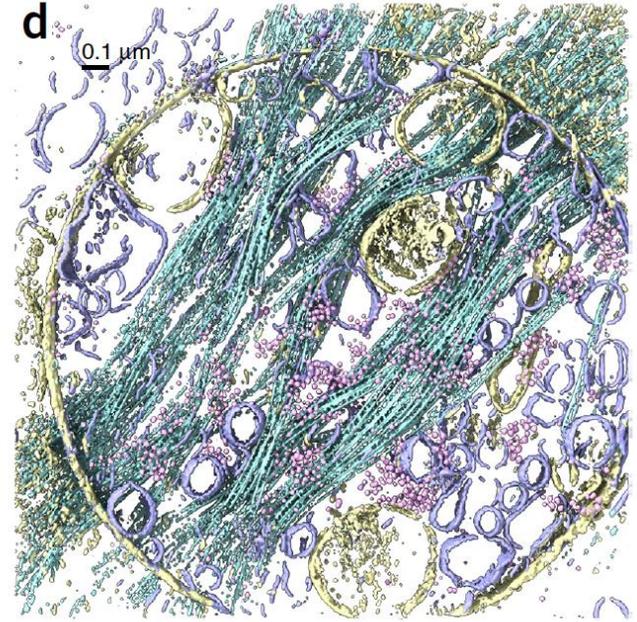
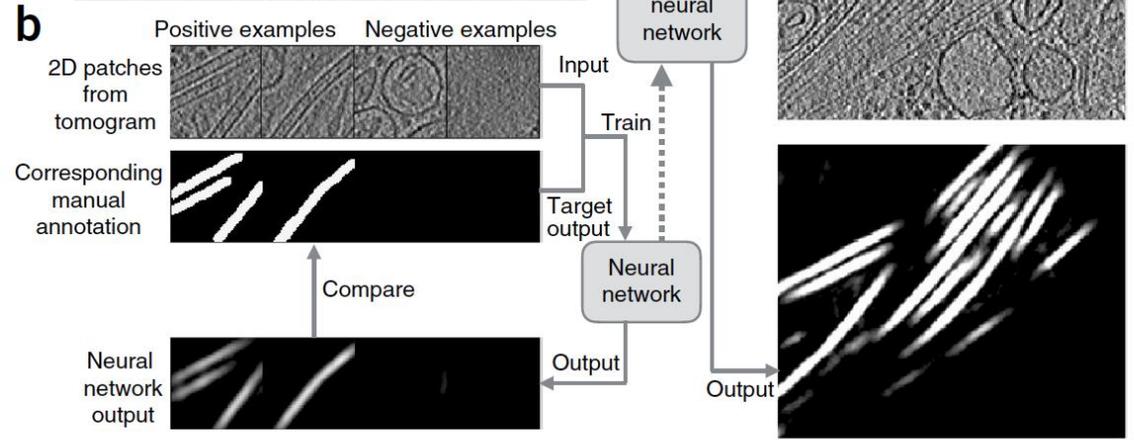
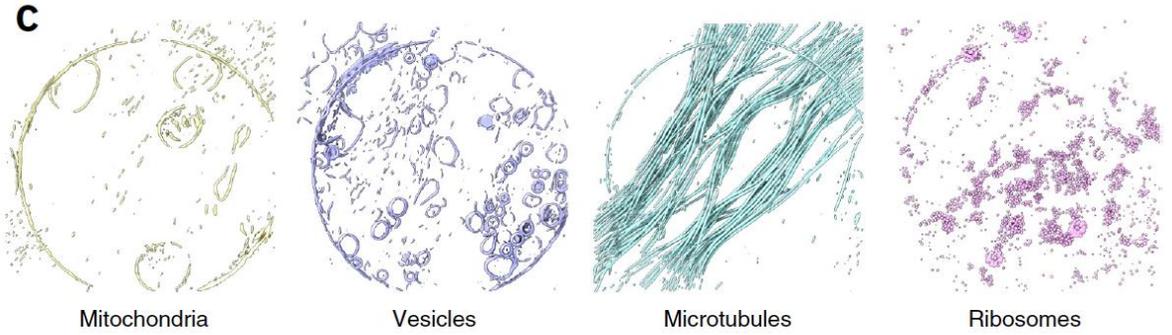
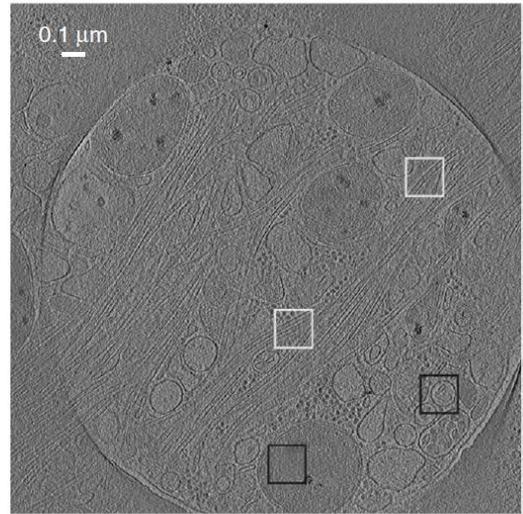
Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



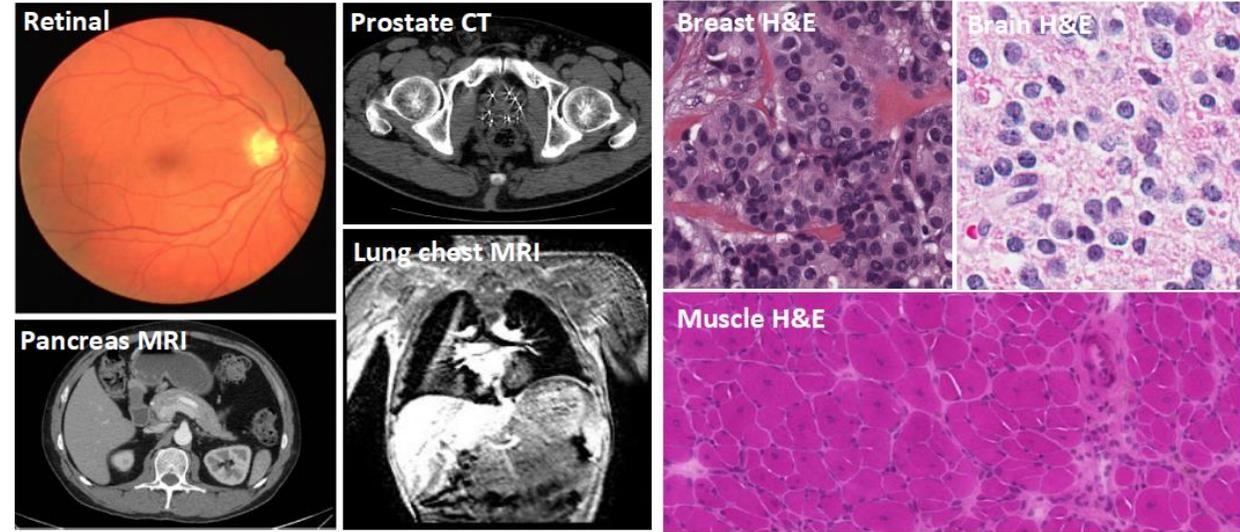
# Convolutional neural networks for automated annotation of cellular cryo-electron tomograms

Muyuan Chen<sup>1,2</sup>, Wei Dai<sup>2,4</sup>, Stella Y Sun<sup>2</sup>,  
 Darius Jonasch<sup>2</sup>, Cynthia Y He<sup>3</sup>, Michael F Schmid<sup>2</sup>,  
 Wah Chiu<sup>2</sup> & Steven J Ludtke<sup>2</sup> 

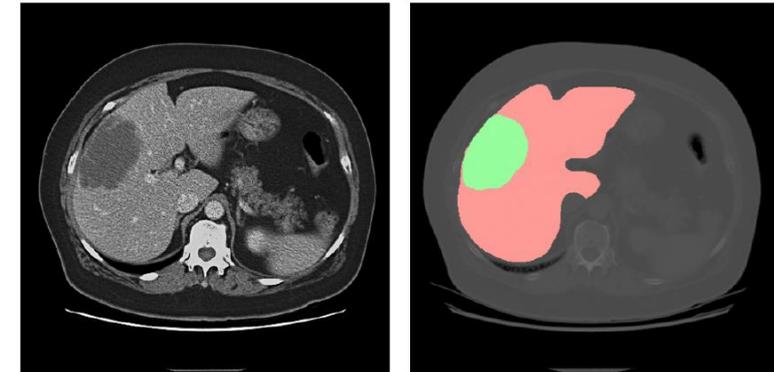
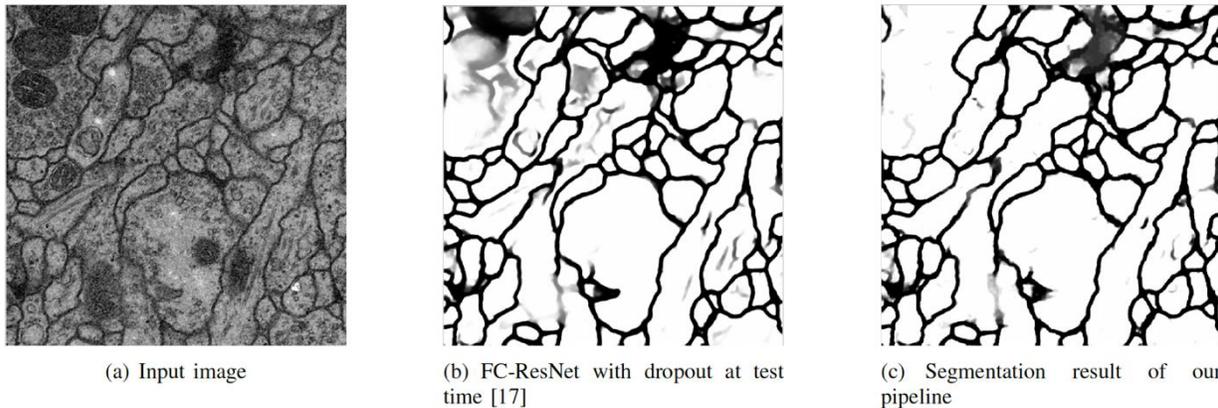




Other possible examples:



Oren Z. Kraus et al., "Classifying and Segmenting Microscopy Images Using Convolutional Multiple Instance Learning," arXiv 2015

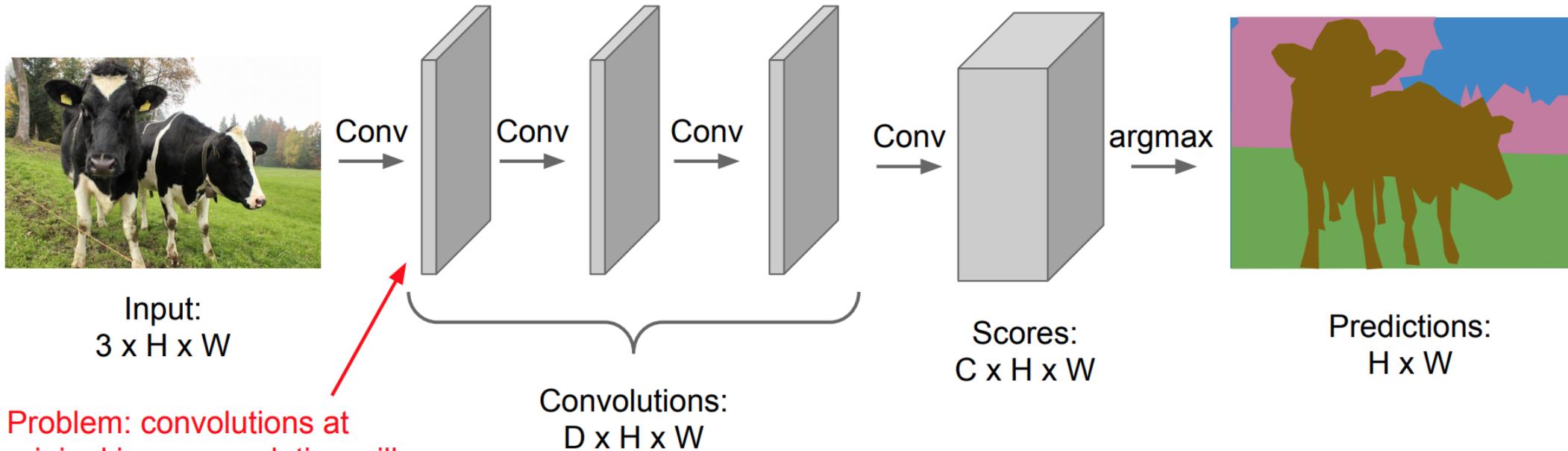


M. Drozdal et al., Learning Normalized Inputs for Iterative Estimation in Medical Image Segmentation (2017)

Z. Zhang et al., Recent Advances in the Applications of Convolutional Neural Networks to Medical Image Contour Detection (2017)<sup>35</sup>

# Semantic Segmentation Idea: Fully Convolutional ?

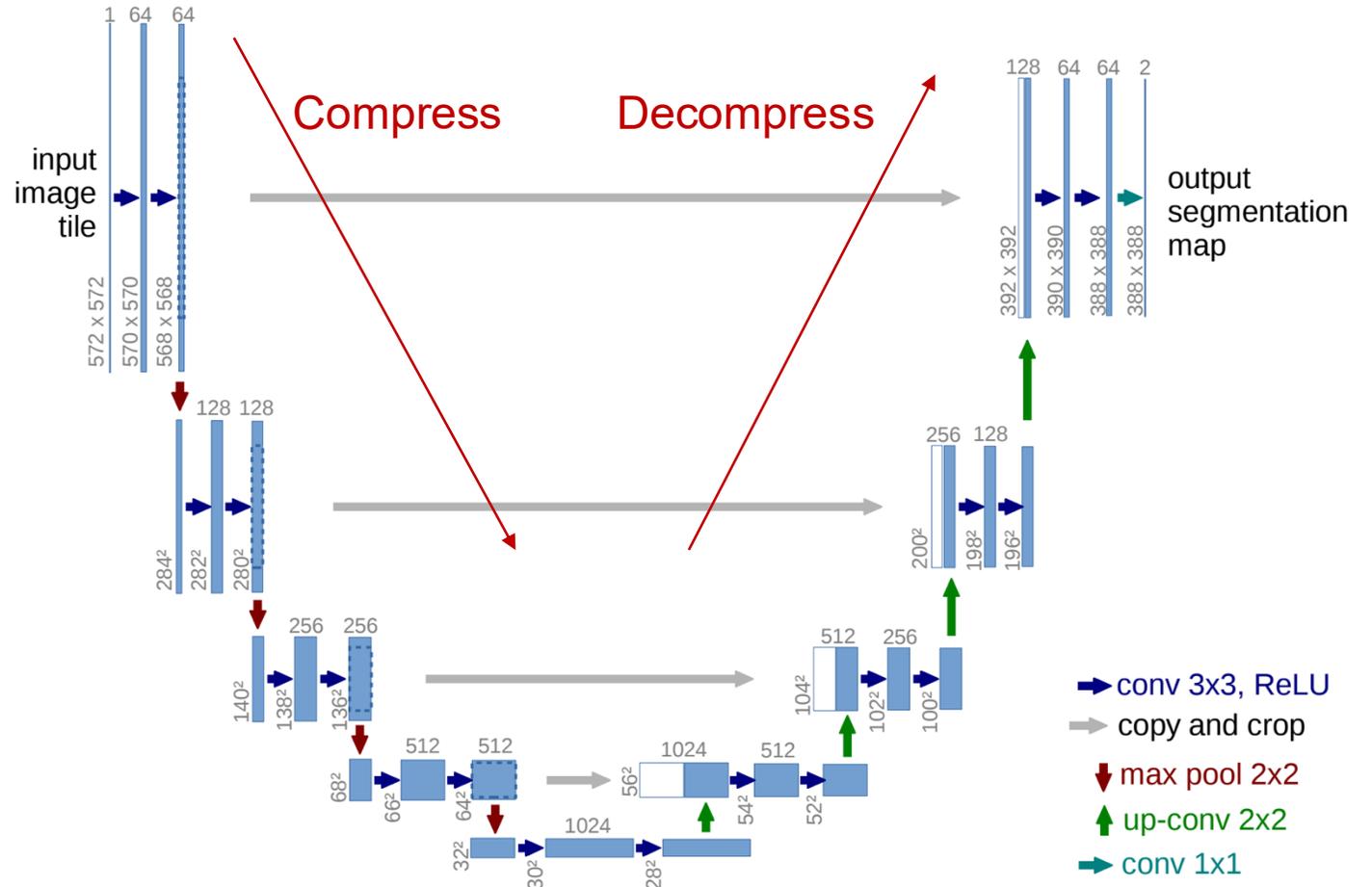
Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



# Instead, *compress x-y dimensions of input image*

- Compress spatial features into learned filters
- Then, decompress learned filters back into same spatial dimensions

## U-Net Architecture



## U-Net: Convolutional Networks for Biomedical Image Segmentation

Olaf Ronneberger, Philipp Fischer, and Thomas Brox

Computer Science Department and BIOS Centre for Biological Signalling Studies,  
University of Freiburg, Germany  
ronneber@informatik.uni-freiburg.de,  
WWW home page: <http://lmb.informatik.uni-freiburg.de/>

# Semantic Segmentation Idea: Fully Convolutional

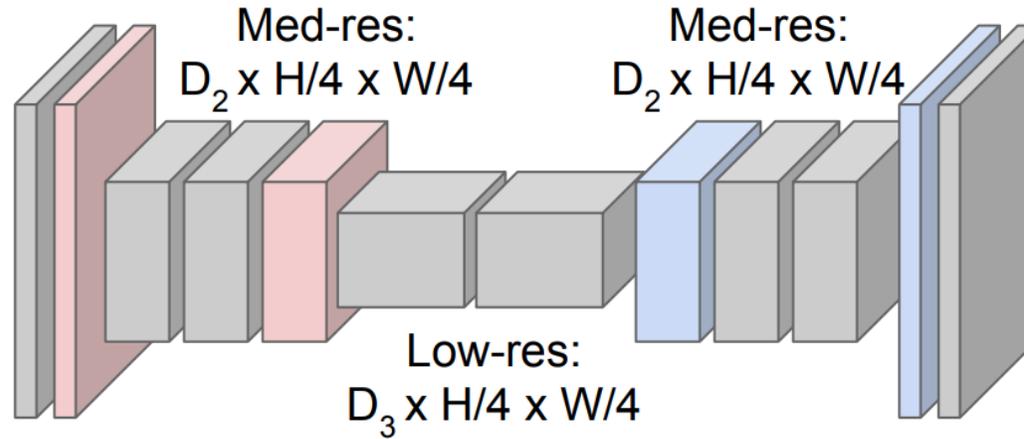
**Downsampling:**  
Pooling, strided convolution

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

**Upsampling:**  
???



Input:  
 $3 \times H \times W$



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

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



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# In-Network upsampling: “Unpooling”

**Nearest Neighbor**

1	2
3	4

Input: 2 x 2



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

Output: 4 x 4

**“Bed of Nails”**

1	2
3	4

Input: 2 x 2



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

Output: 4 x 4

# In-Network upsampling: “Max Unpooling”

## Max Pooling

Remember which element was max!

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

Input: 4 x 4



5	6
7	8

Output: 2 x 2



Rest of the network

## Max Unpooling

Use positions from pooling layer

1	2
3	4

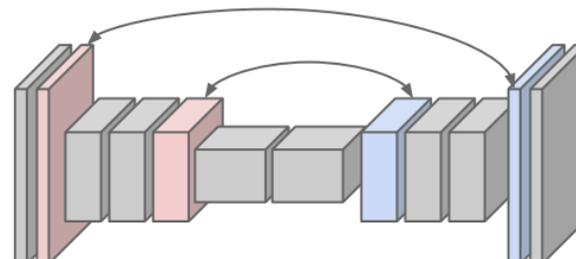
Input: 2 x 2



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

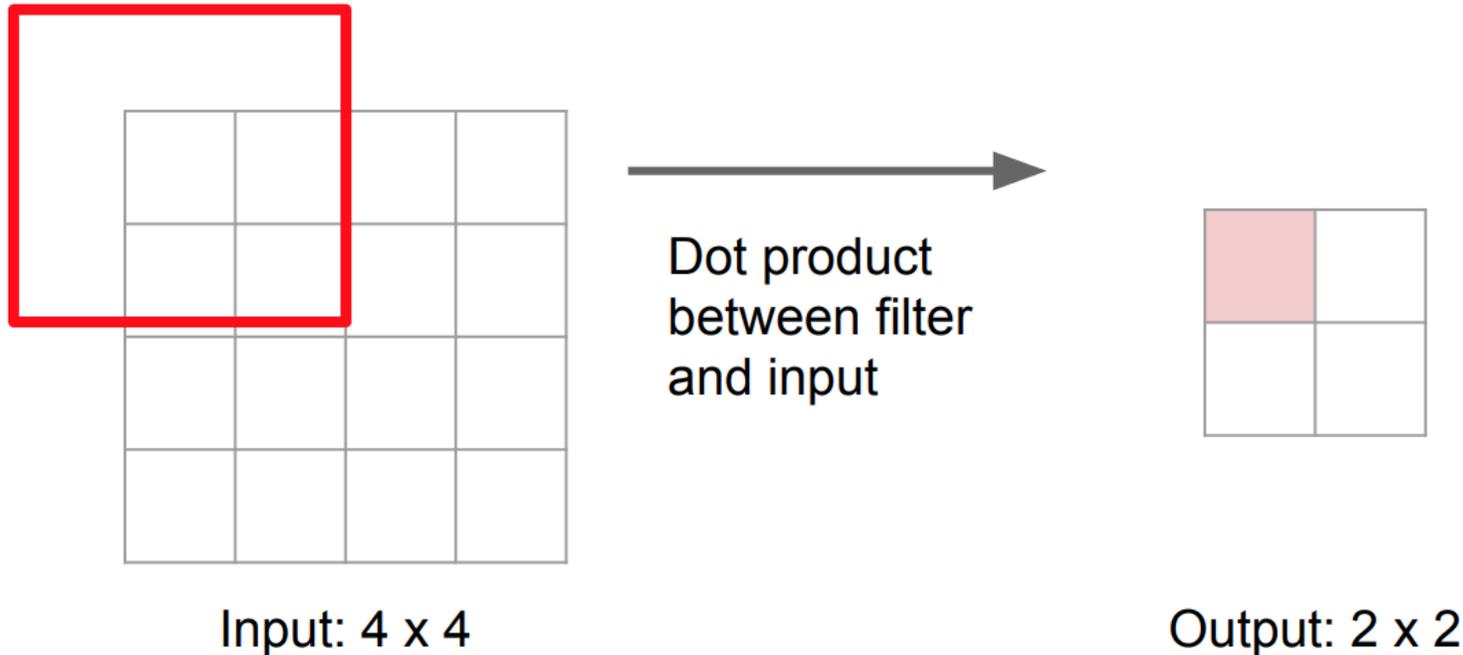
Output: 4 x 4

Corresponding pairs of downsampling and upsampling layers



# Learnable Upsampling: Transpose Convolution

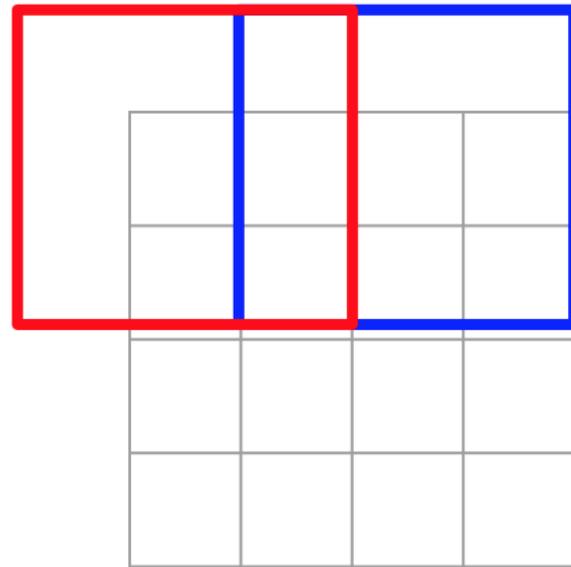
**Recall:** Normal 3 x 3 convolution, stride 2 pad 1





# Learnable Upsampling: Transpose Convolution

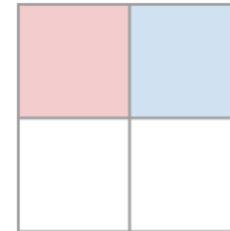
**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4



Dot product  
between filter  
and input



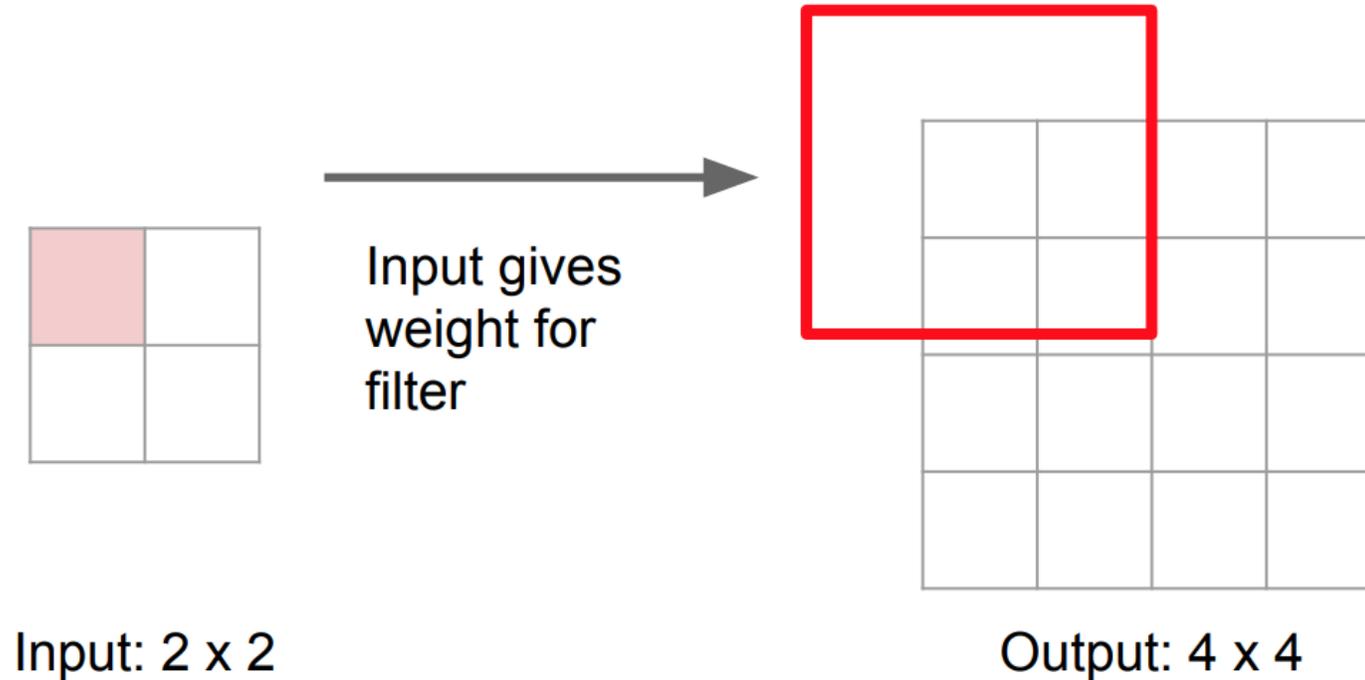
Output: 2 x 2

Filter moves 2 pixels in  
the input for every one  
pixel in the output

Stride gives ratio between  
movement in input and  
output

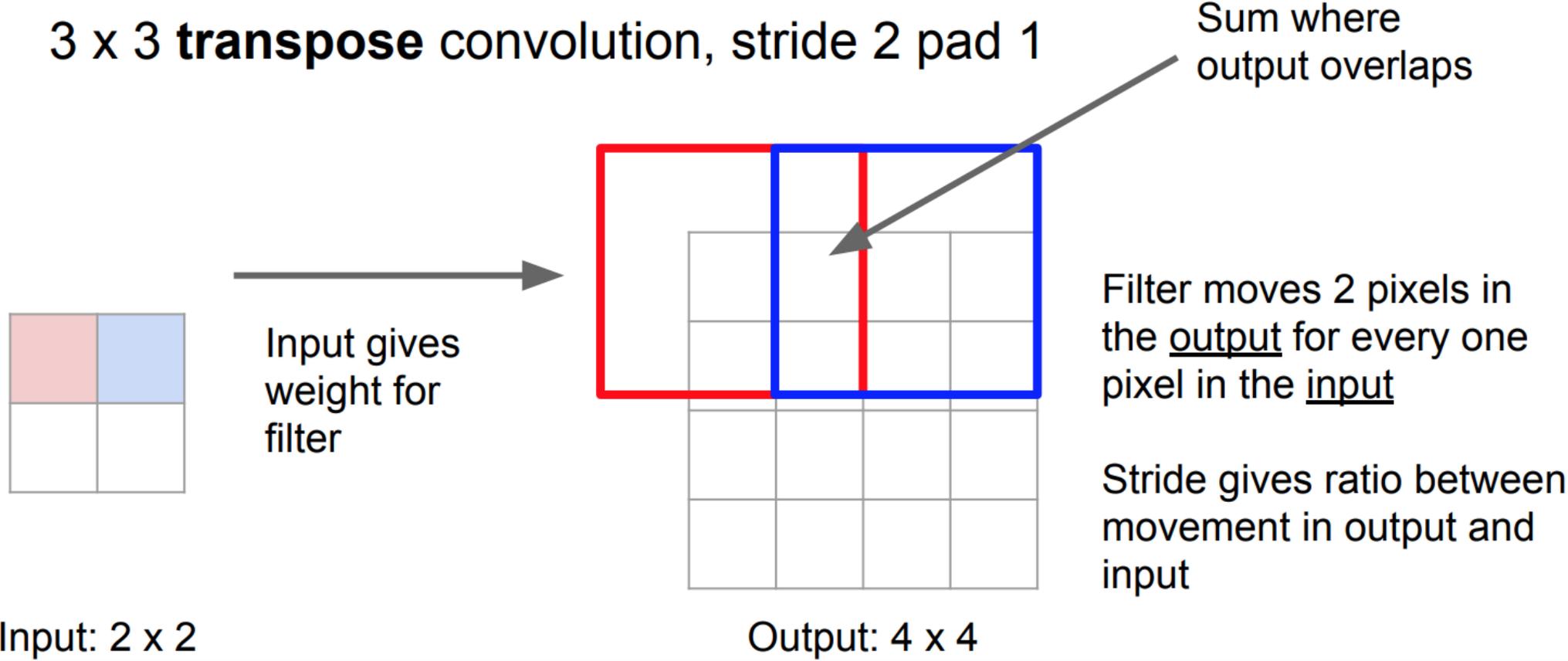
# Learnable Upsampling: Transpose Convolution

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



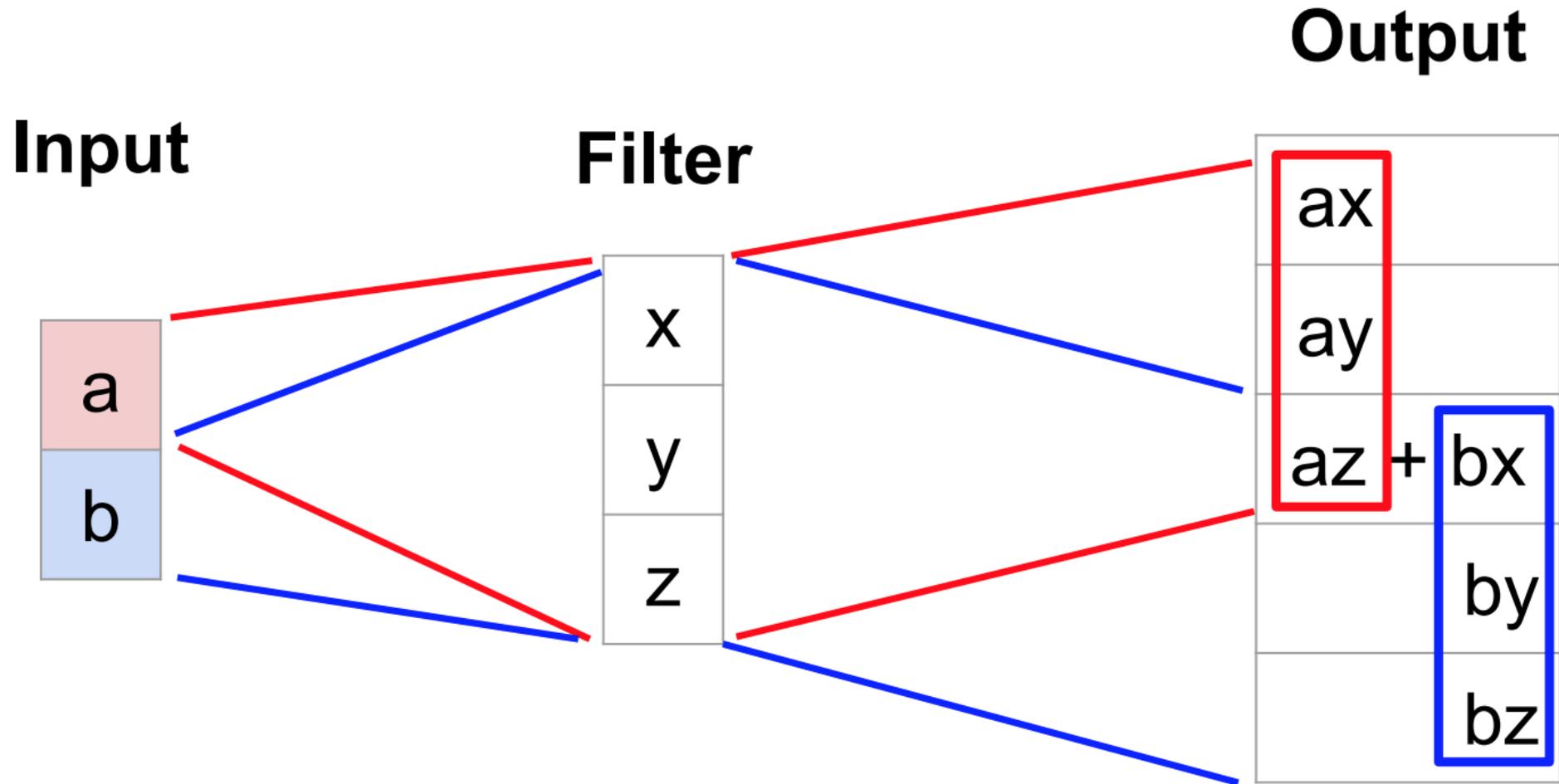
# Learnable Upsampling: Transpose Convolution

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





# Learnable Upsampling: 1D Example



Output contains copies of the filter weighted by the input, summing at where it overlaps in the output

Need to crop one pixel from output to make output exactly 2x input



# Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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



# Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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

Convolution transpose multiplies by the transpose of the same matrix:

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

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

When stride=1, convolution transpose is just a regular convolution (with different padding rules)



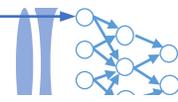
# Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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



# Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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

Convolution transpose multiplies by the transpose of the same matrix:

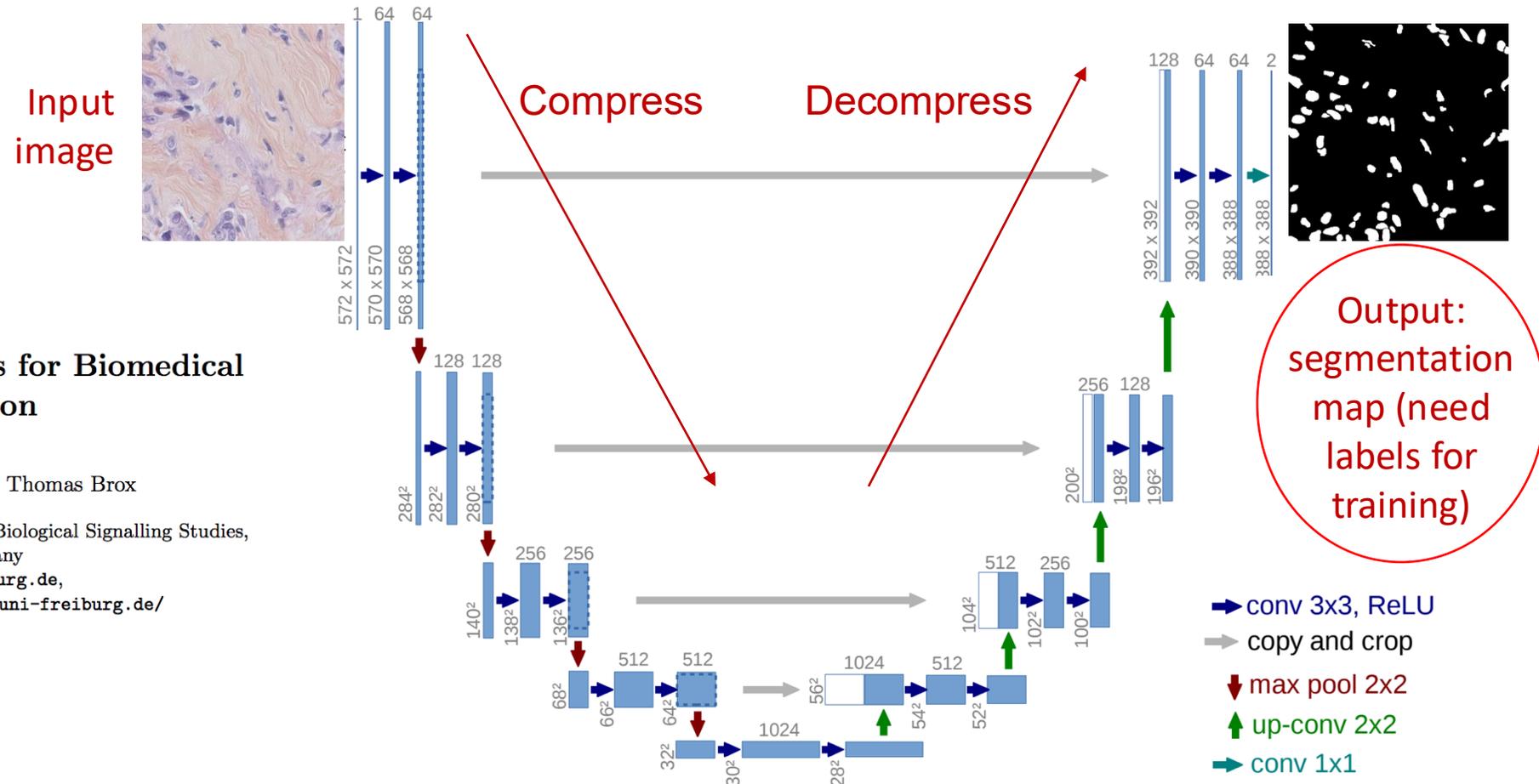
$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

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

When stride>1, convolution transpose is no longer a normal convolution!

# Segmentation: need a map of classes for label

## U-Net Architecture

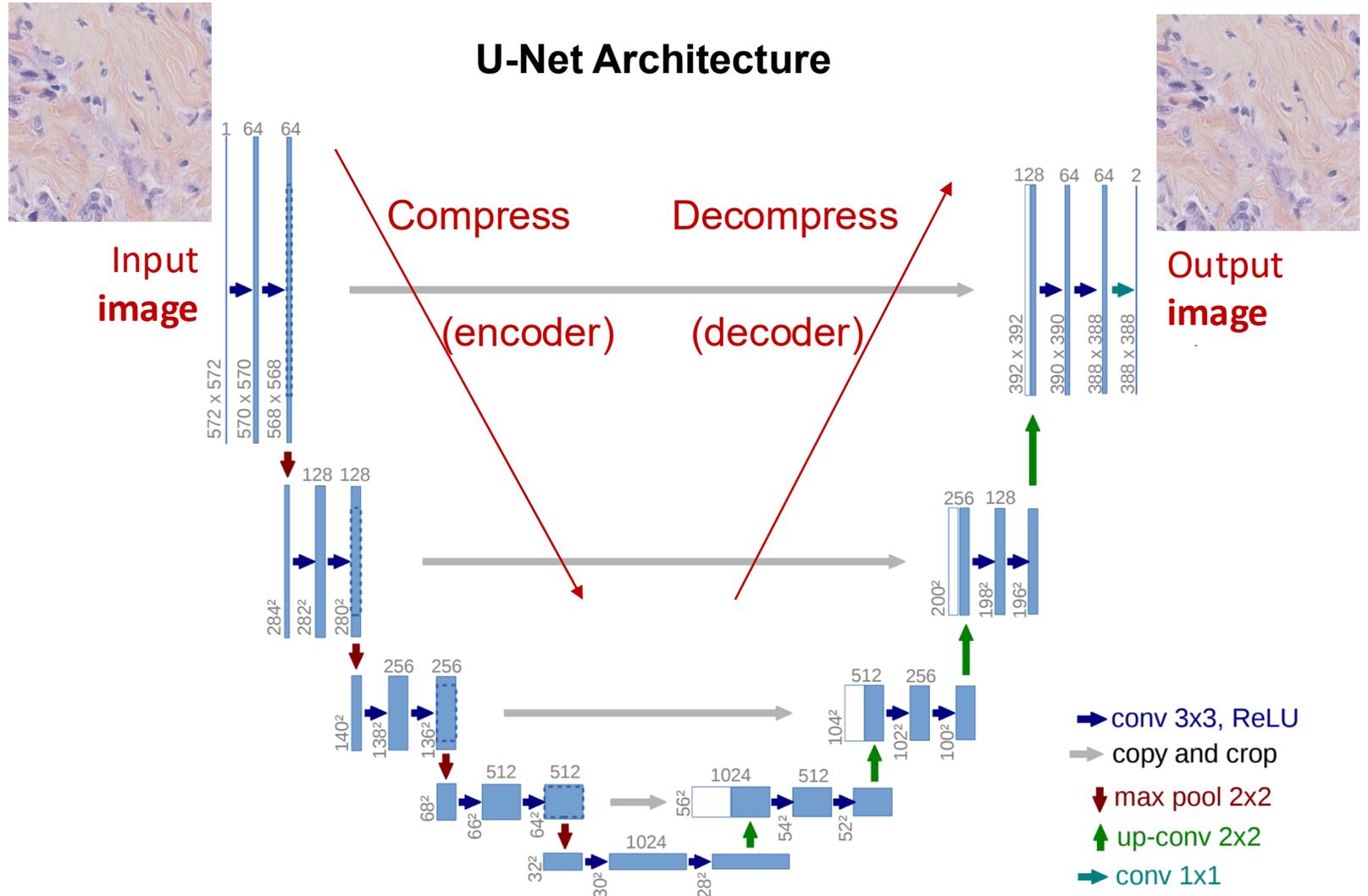


### U-Net: Convolutional Networks for Biomedical Image Segmentation

Olaf Ronneberger, Philipp Fischer, and Thomas Brox

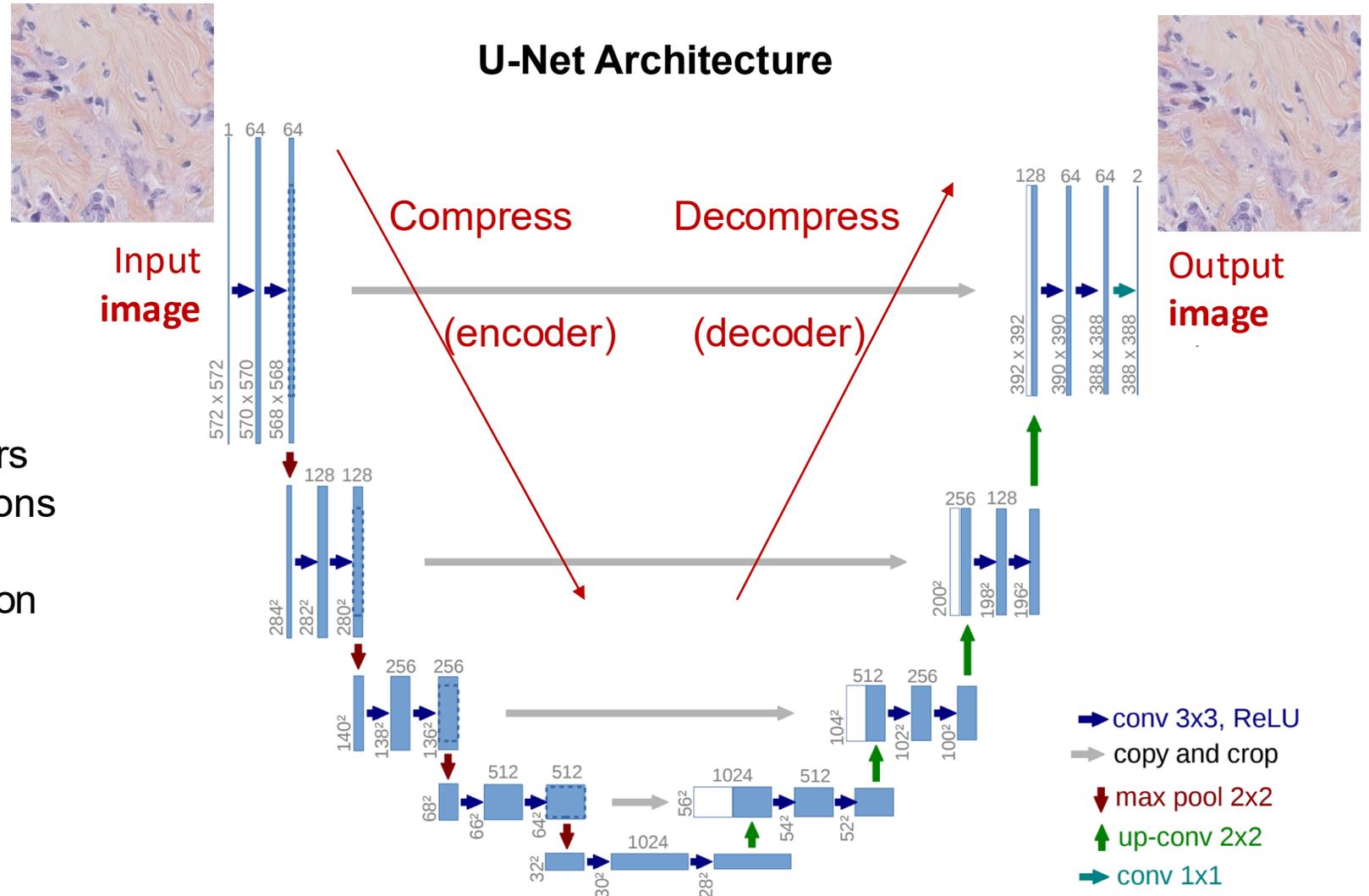
Computer Science Department and BIOS Centre for Biological Signalling Studies,  
 University of Freiburg, Germany  
 ronneber@informatik.uni-freiburg.de,  
 WWW home page: <http://lmb.informatik.uni-freiburg.de/>

# Instead, *compress x-y dimensions of input image*



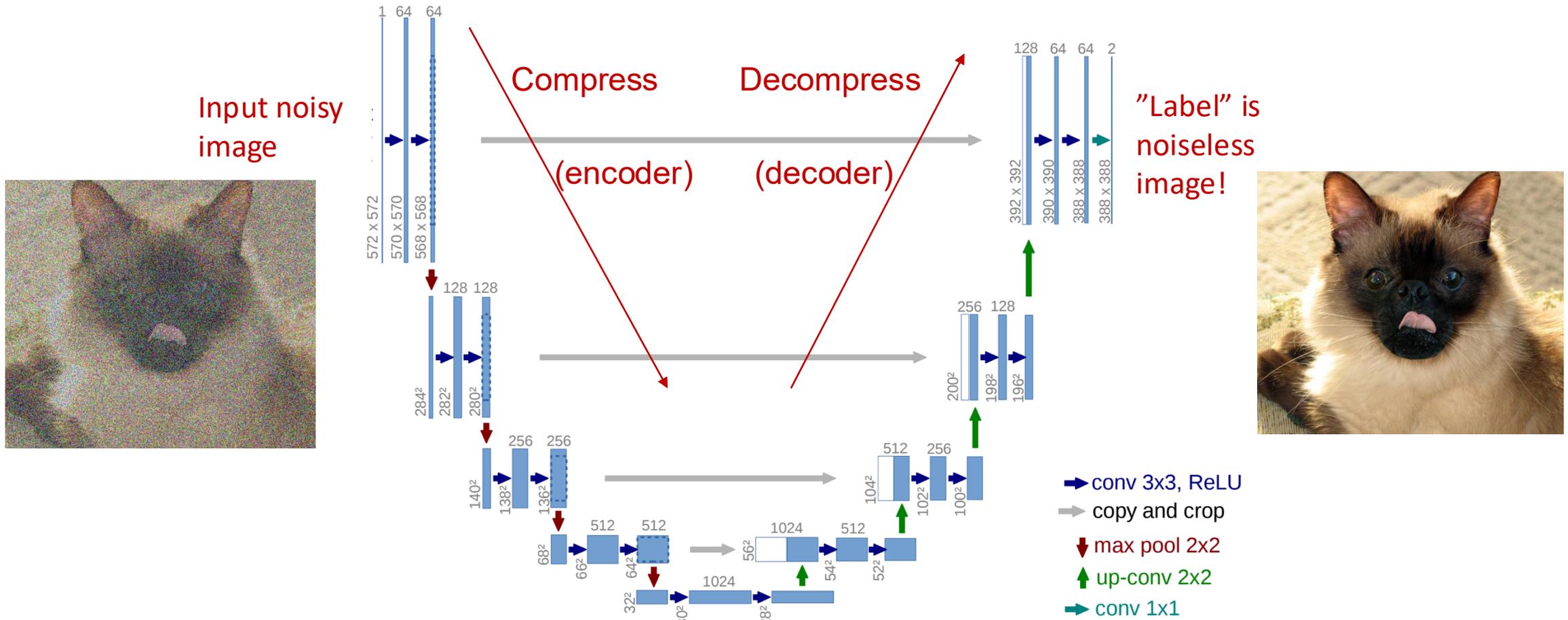
# Instead, *compress x-y dimensions of input image*

- Compress spatial features into learned filters
- Then, decompress learned filters back into same spatial dimensions
- **Termed an autoencoder**
- Analogous to image compression
- A pretty powerful idea...



# Another example: Denoising Autoencoder

## U-Net Architecture



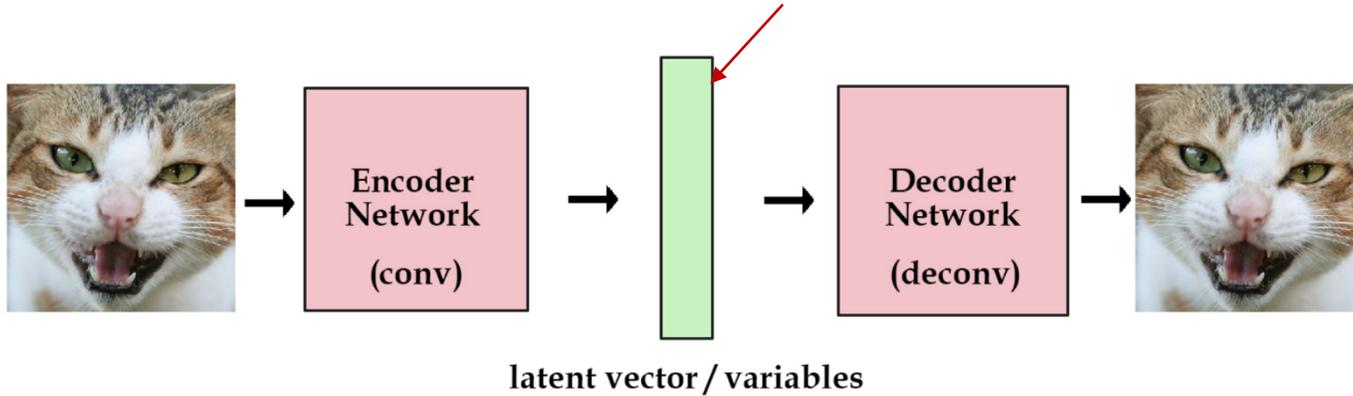
Code review: See the following:

[Jupyter Notebook: A simple Autoencoder in Tensorflow/Keras](#)

<https://deepimaging.github.io/lectures/>

# Example: Variational Autoencoder (VAE)

Force this vector to follow a Gaussian PDF



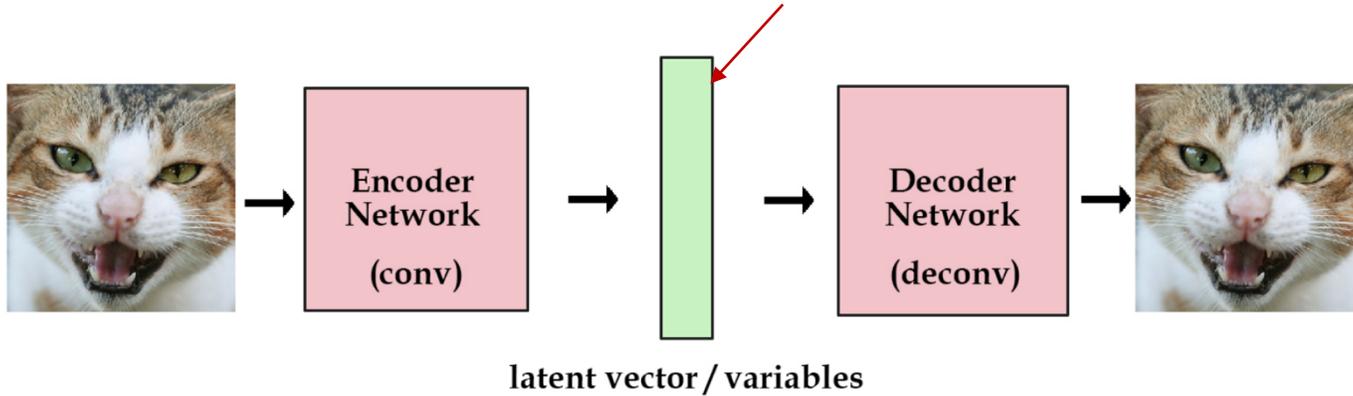
- Good *generative model*
- Have a clean probability distribution to select from to generate new examples

Minimize (KL) distance between latent vector and Gaussian normal



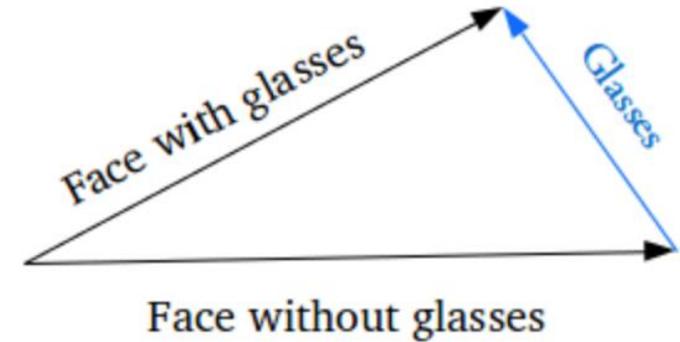
# Example: Variational Autoencoder (VAE)

Force this vector to follow a Gaussian PDF



Minimize (KL) distance between latent vector and Gaussian normal

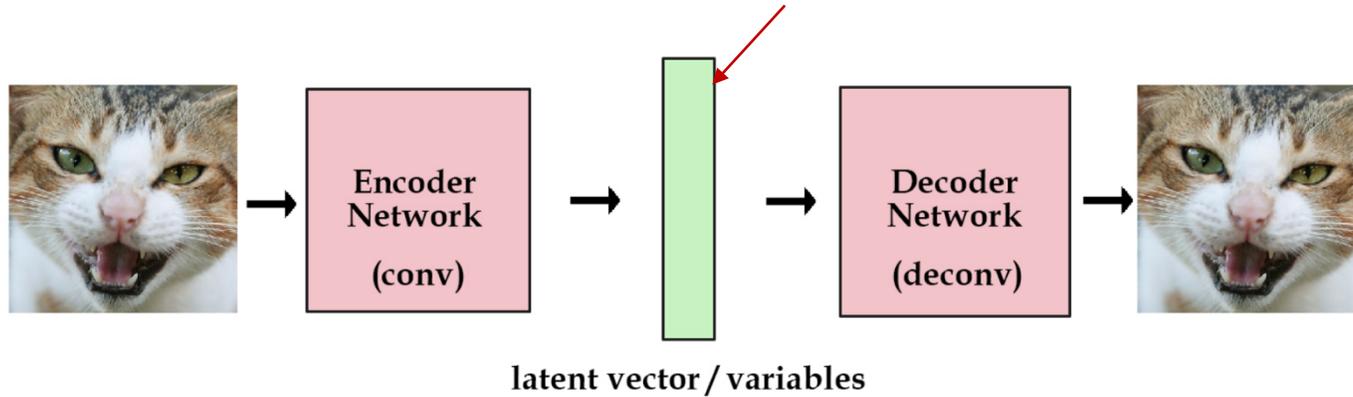
- With Gaussian PDF, can start to add/subtract latent vector in a normalized vector space



Adding new features to samples

# Example: Variational Autoencoder (VAE)

Force this vector to follow a Gaussian PDF

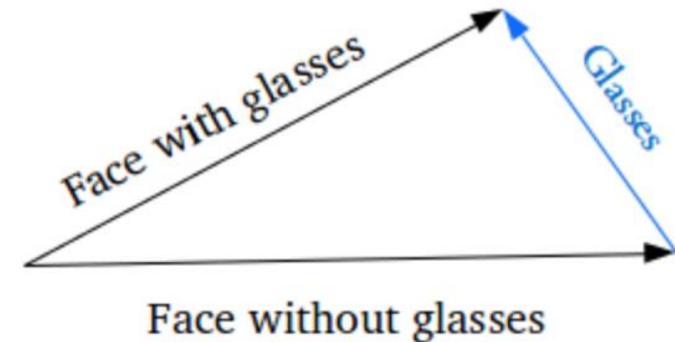


Minimize (KL) distance between latent vector and Gaussian normal

Generative Example (once trained):

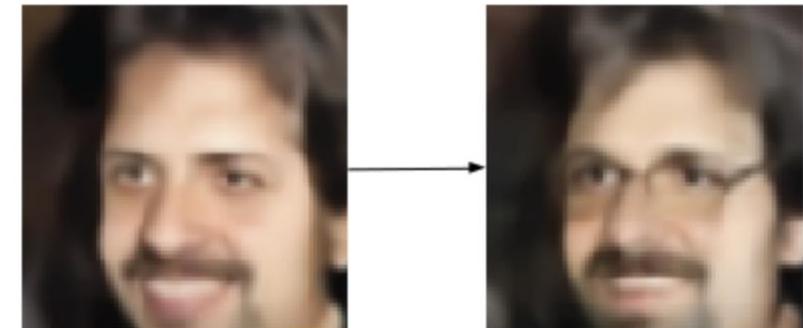
- Encode image with glasses, obtain latent vector PDF  $P_g$
- Encode image without glasses, obtain PDF  $P_{ng}$
- Compute  $\mathbf{diff} = P_g - P_{ng}$
- Encode new image to obtain  $P_{new}$ , add in  $\mathbf{diff}$
- Decode  $P_{new} + \mathbf{diff}$  to get guy with glasses!

- With Gaussian PDF, can start to add/subtract latent vector in a normalized vector space



Adding new features to samples

Glasses



Exploring a specific variation of input data[1]