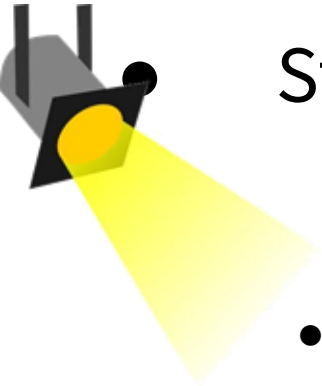


# Lecture 9: Detection, Segmentation, Visualization, and Understanding

# Administrative Announcements

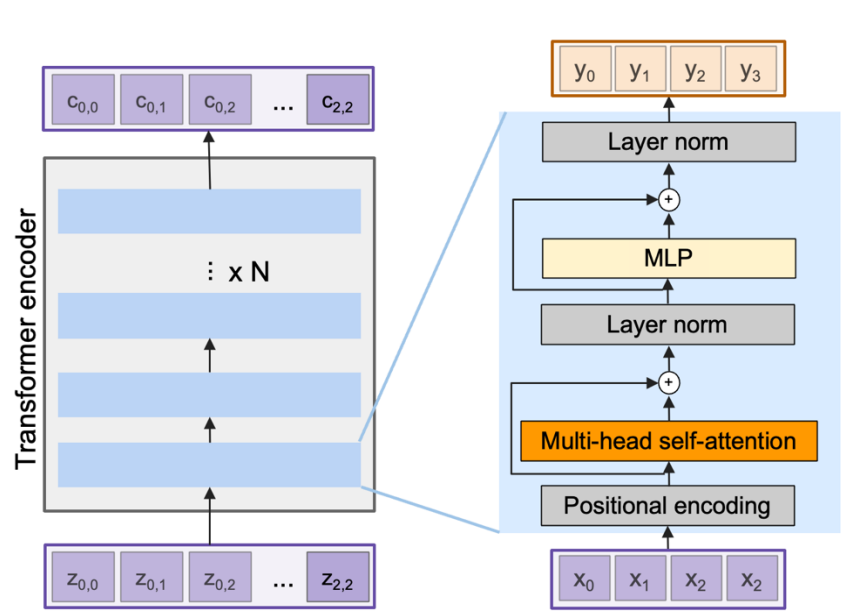
- A fix to assignment 2 notebook has been posted in course webpage. See Edstem posts for details.
- Make sure to start Assignment 2 early. It is the longest of the three assignments, and the midterm and project milestone deadlines follow closely after the Assignment 2 deadline.



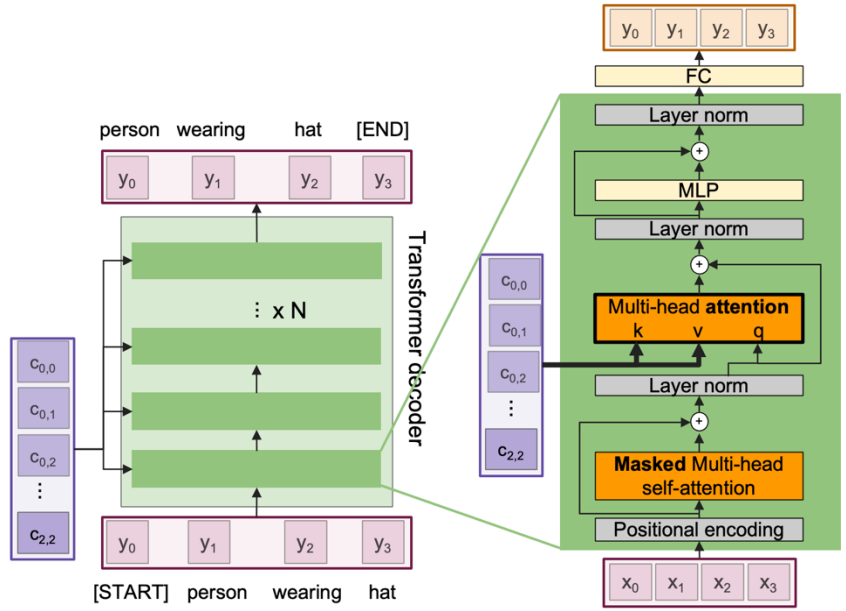
# Student Spotlights

- Thank you for your engagement and improving the course experience!
  - Ben Wengreen
  - Arunabh Sharma
  - Warren Chan
  - Eyas Taifour
  - Ali Ahmad
  - Florencio Sedano

# Last time: Transformer



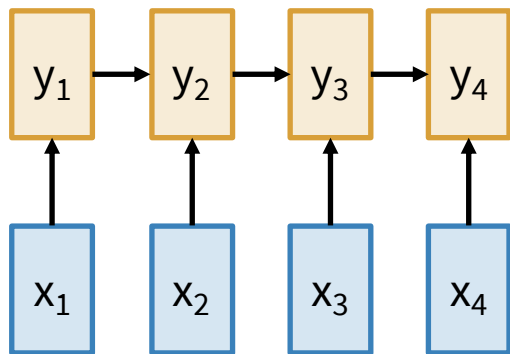
Encoder



Decoder

# Three Ways of Processing Sequences

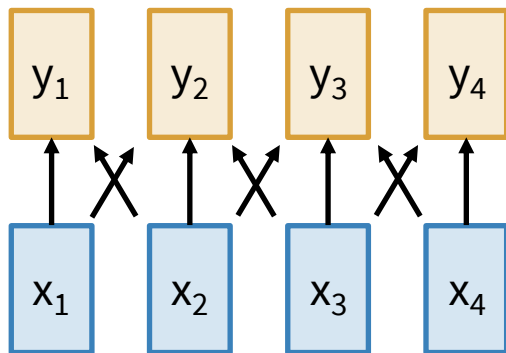
## Recurrent Neural Network



Works on 1D ordered sequences

- (+) Theoretically good at long sequences:  $O(N)$  compute and memory for a sequence of length  $N$
- (-) Not parallelizable. Need to compute hidden states sequentially

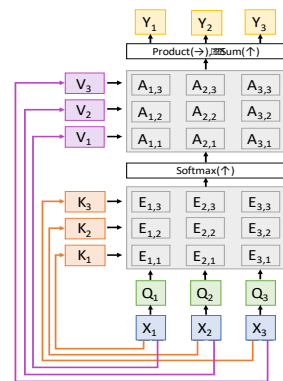
## Convolution



Works on N-dimensional grids

- (-) Bad for long sequences: need to stack many layers to build up large receptive fields
- (+) Parallelizable, outputs can be computed in parallel

## Self-Attention



Works on sets of vectors

- (+) Great for long sequences; each output depends directly on all inputs
- (+) Highly parallel, it's just 4 matmuls
- (-) Expensive:  $O(N^2)$  compute,  $O(N)$  memory for sequence of length  $N$

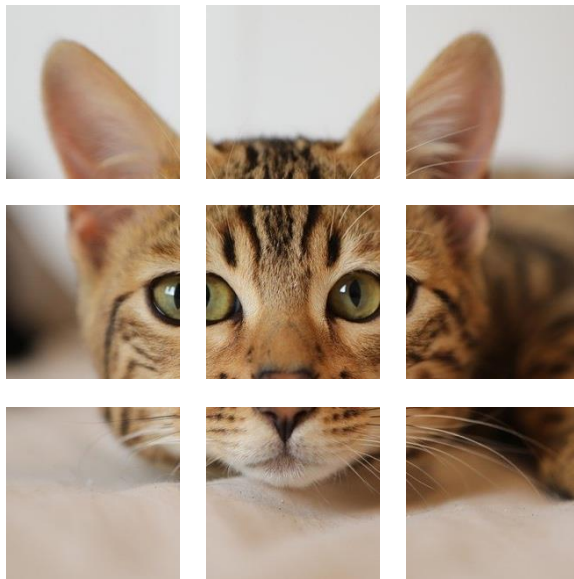
# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3

Dosovitskiy et al, "An Image is Worth  
16x16 Words: Transformers for Image  
Recognition at Scale", ICLR 2021

# Vision Transformers (ViT)

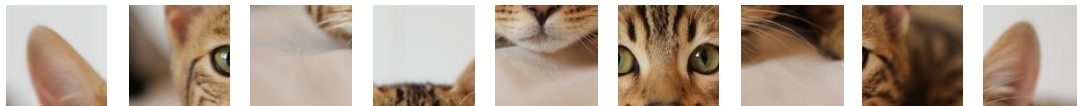


Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformers (ViT)

N input patches, each of  
shape  $3 \times 16 \times 16$



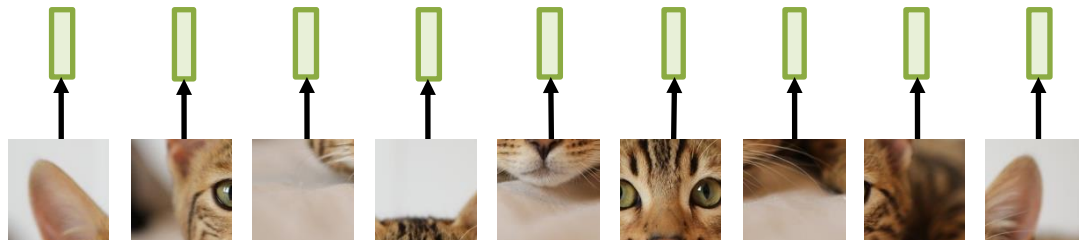
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

# Vision Transformers (ViT)

Linear projection to D-dimensional vector

N input patches, each of shape  $3 \times 16 \times 16$



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

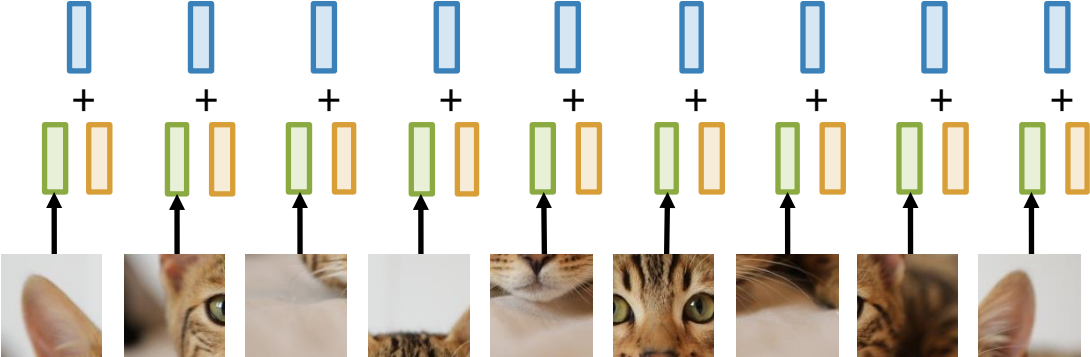
[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformers (ViT)

Add positional embedding:  
learned D-dim vector per  
position

Linear projection to D-  
dimensional vector

N input patches, each of  
shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformers (ViT)

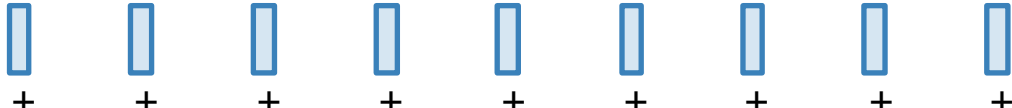
Output vectors



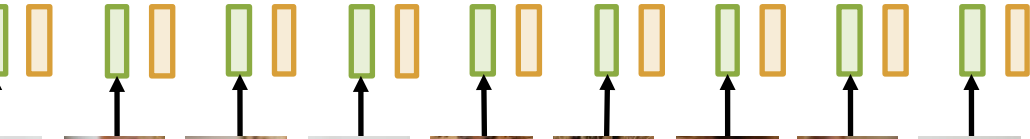
Exact same as NLP Transformer!



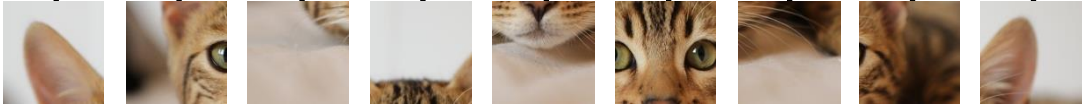
Add positional embedding:  
learned D-dim vector per  
position



Linear projection to D-  
dimensional vector



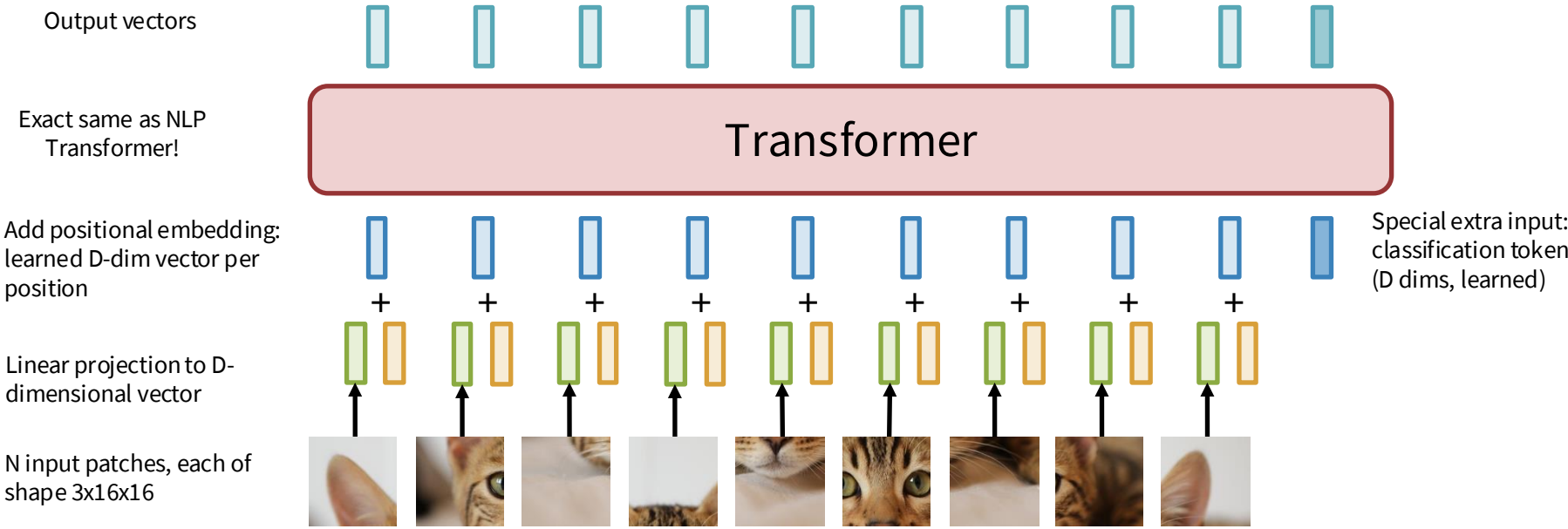
N input patches, each of  
shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

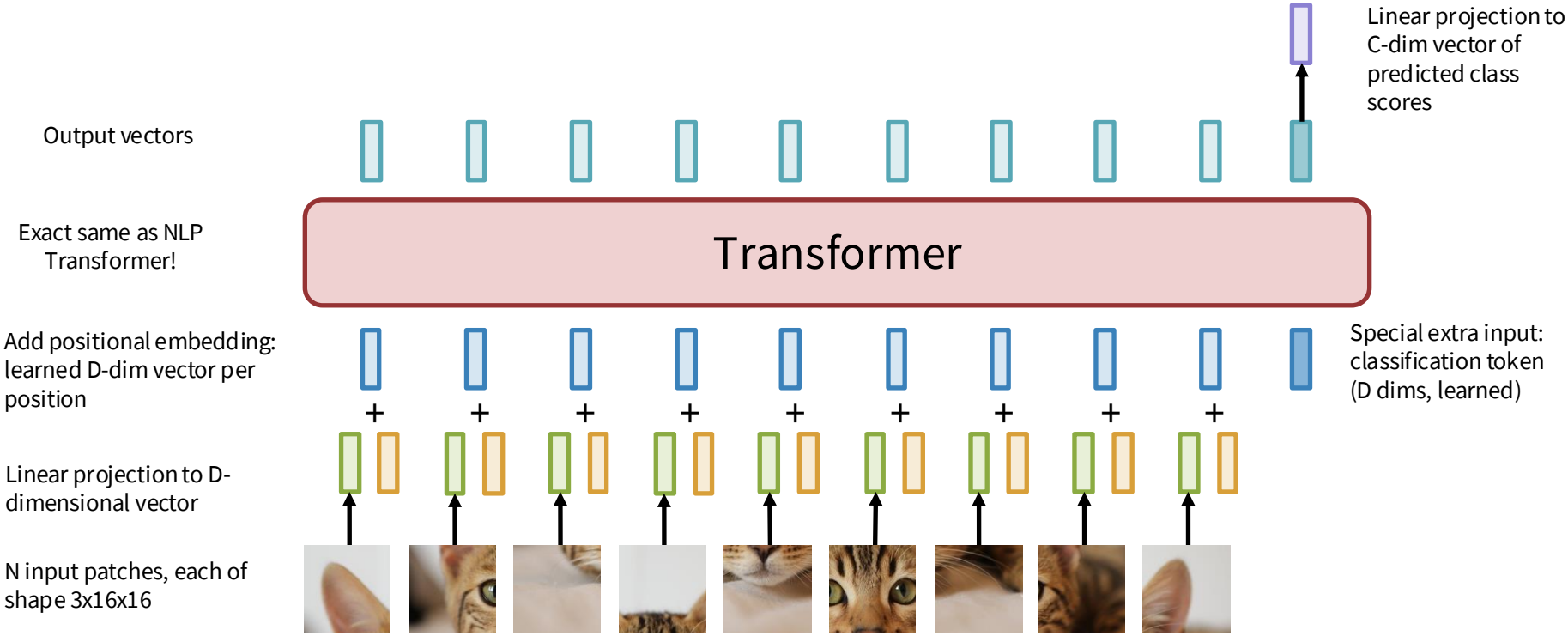
# Vision Transformers (ViT)



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

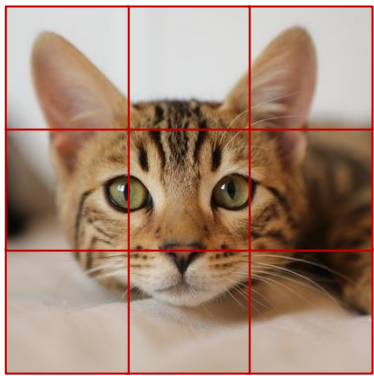
# Vision Transformers (ViT)



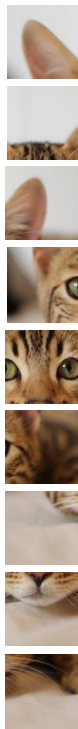
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformers (ViT) – a similar approach (different classifier)



Input image:  
e.g. 224x224x3



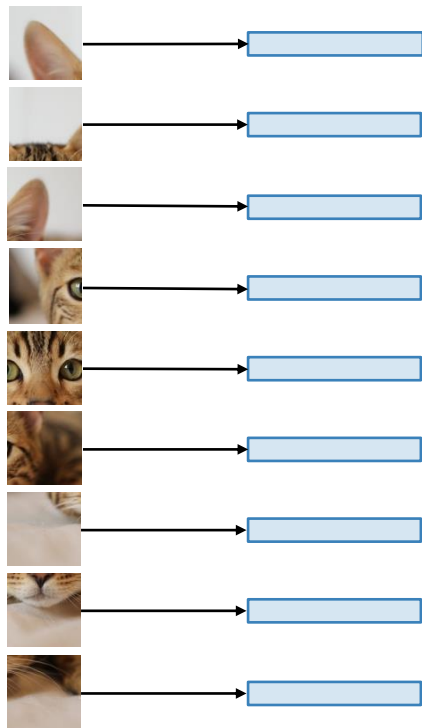
Break into patches  
e.g. 16x16x3

Dosovitskiy et al, "An Image is Worth  
16x16 Words: Transformers for Image  
Recognition at Scale", ICLR 2021

# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

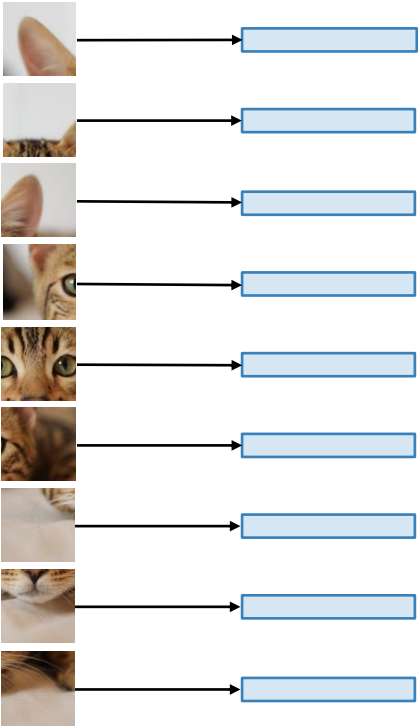
Flatten and apply a linear  
transform  $768 \Rightarrow D$

Dosovitskiy et al, "An Image is Worth  
16x16 Words: Transformers for Image  
Recognition at Scale", ICLR 2021

# Vision Transformers (ViT)

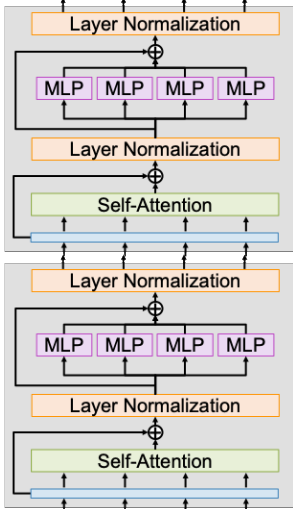


Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

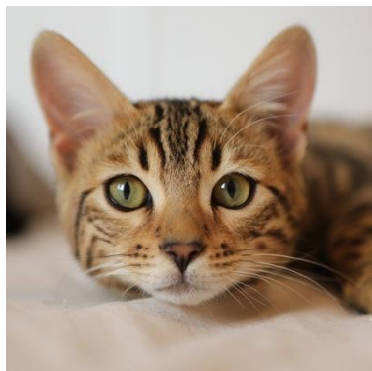
Flatten and apply a linear  
transform 768 => D



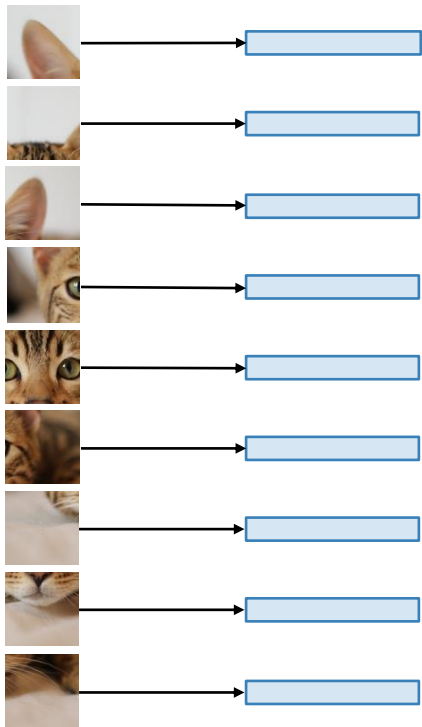
D-dim vector per patch are  
the input vectors to the  
Transformer

Dosovitskiy et al, "An Image is Worth  
16x16 Words: Transformers for Image  
Recognition at Scale", ICLR 2021

# Vision Transformers (ViT)

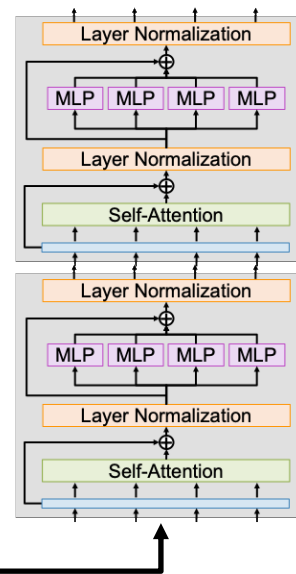


Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform 768 => D



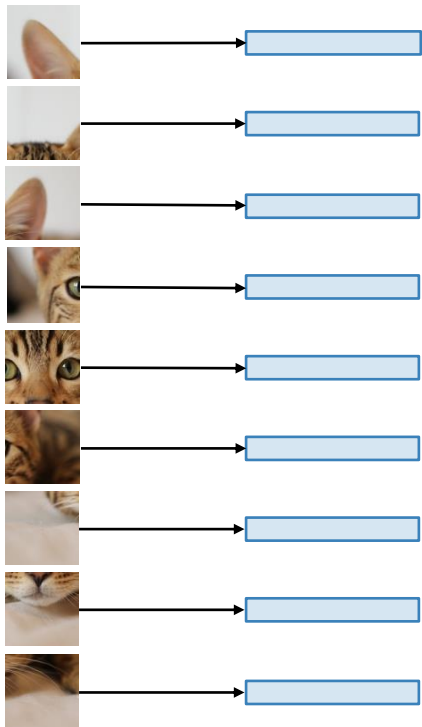
D-dim vector per patch are  
the input vectors to the  
Transformer

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

# Vision Transformers (ViT)

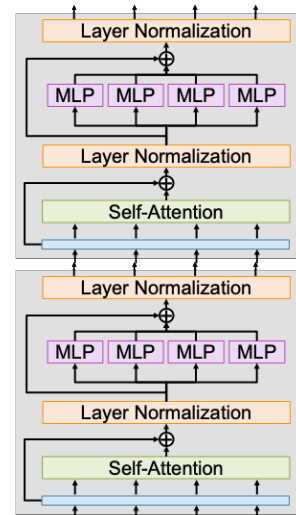


Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform 768 => D



D-dim vector per patch are  
the input vectors to the  
Transformer

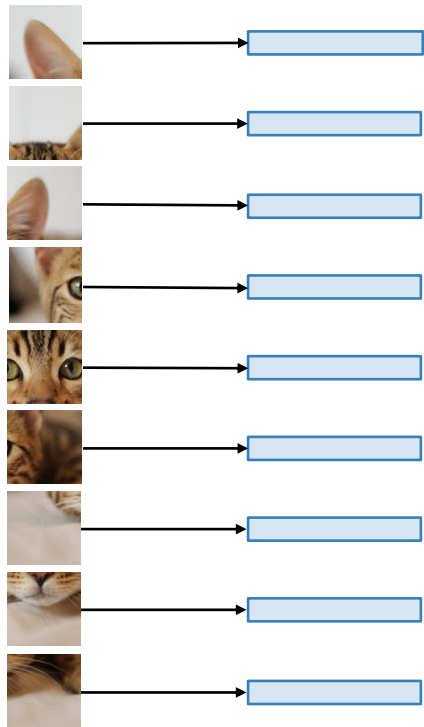
Don't use any  
masking; each  
image patch can  
look at all other  
image patches

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

# Vision Transformers (ViT)

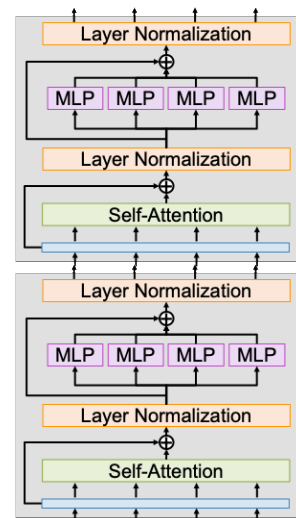


Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform 768 => D



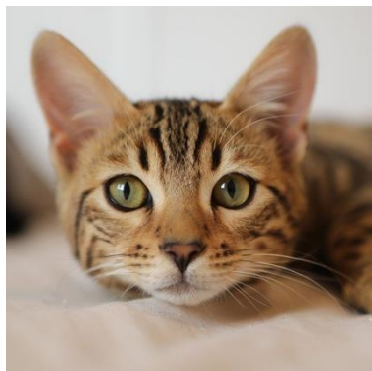
D-dim vector per patch are  
the input vectors to the  
Transformer

Transformer  
gives an output  
vector per patch

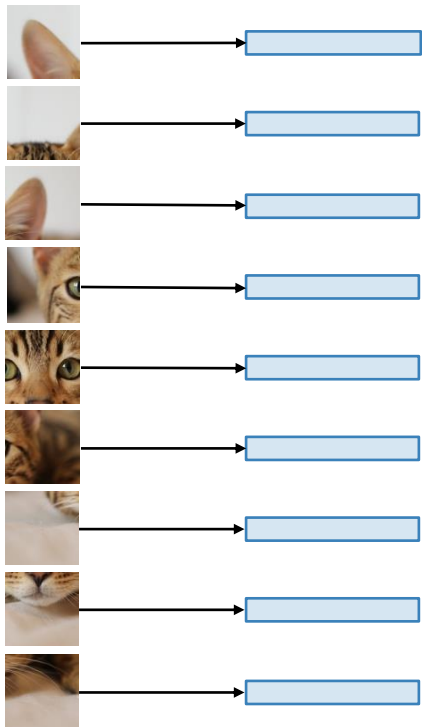
Don't use any  
masking; each  
image patch can  
look at all other  
image patches

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

# Vision Transformers (ViT)



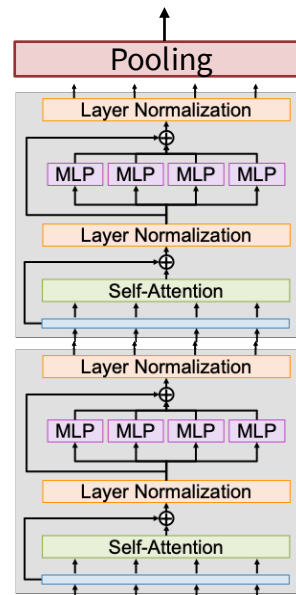
Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform 768 => D

Average pool  $N \times D$  vectors to  
 $1 \times D$ , apply a linear layer  $D \Rightarrow C$  to  
predict class scores



D-dim vector per patch are  
the input vectors to the  
Transformer

Transformer  
gives an output  
vector per patch

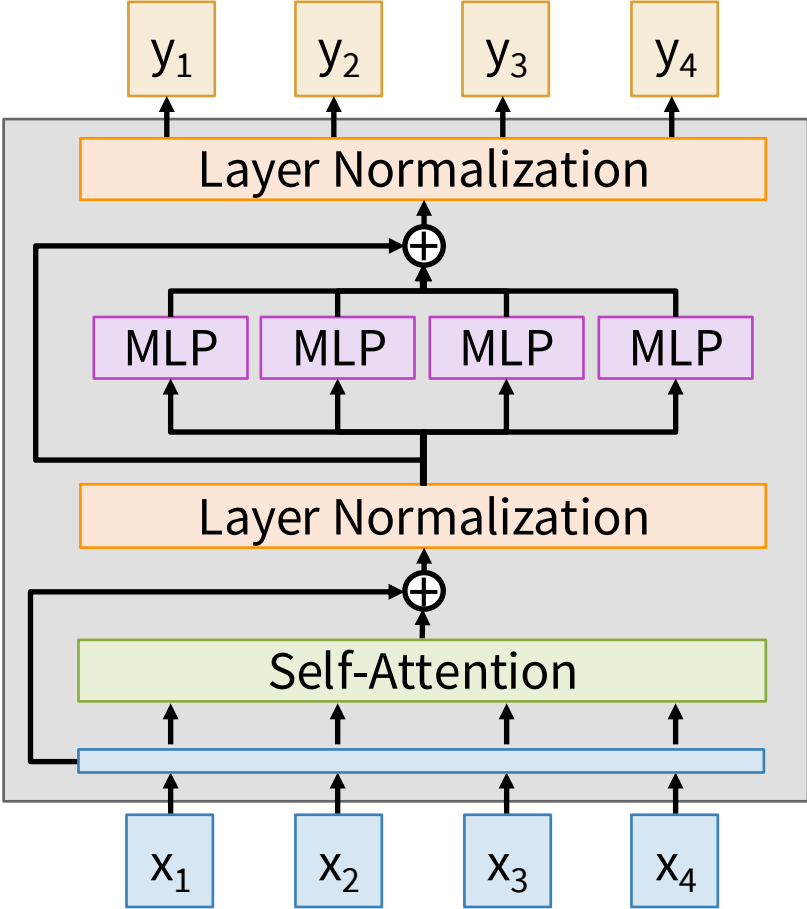
Don't use any  
masking; each  
image patch can  
look at all other  
image patches

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

# Tweaking Transformers

The Transformer architecture has not changed much since 2017.

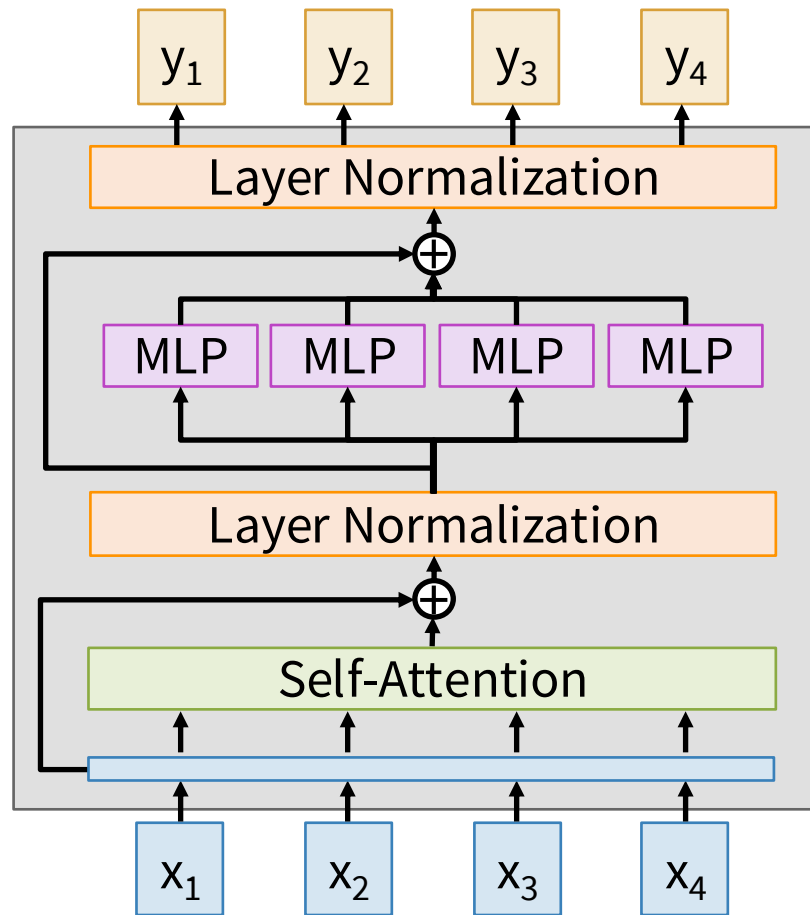
But a few changes have become common:



# Pre-Norm Transformer

Layer normalization is outside the residual connections

Kind of weird, the model can't actually learn the identity function

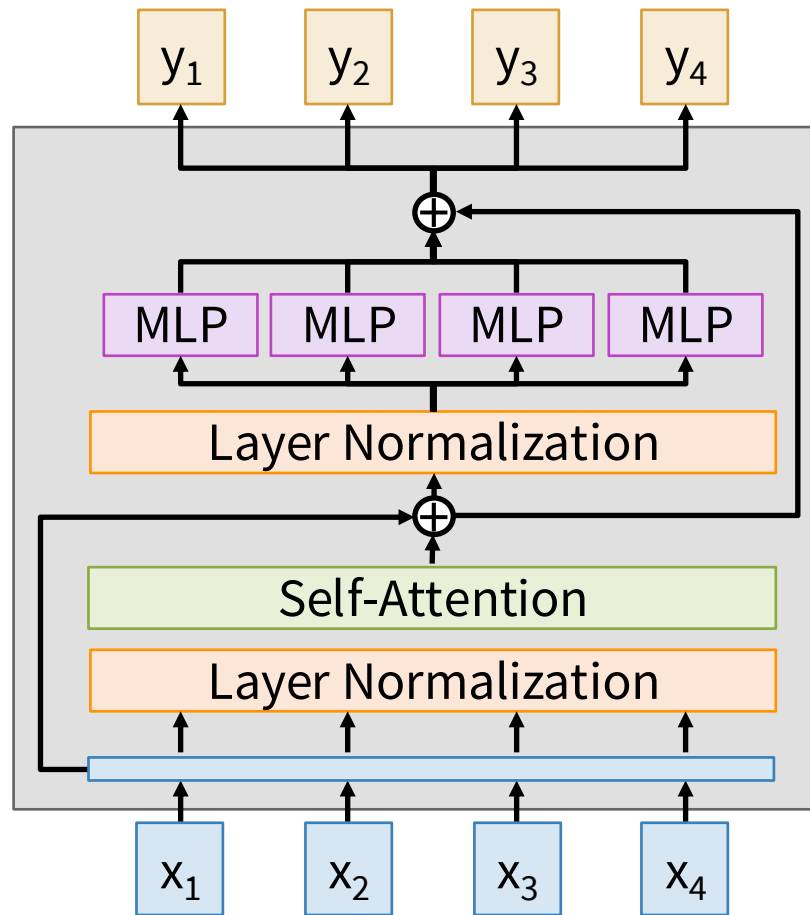


# Pre-Norm Transformer

Layer normalization is outside the residual connections

Kind of weird, the model can't actually learn the identity function

Solution: Move layer normalization before the Self-Attention and MLP, inside the residual connections. Training is more stable.



# RMSNorm

Replace Layer Normalization with Root-Mean-Square Normalization (RMSNorm)

Input:  $x$  [shape  $D$ ]

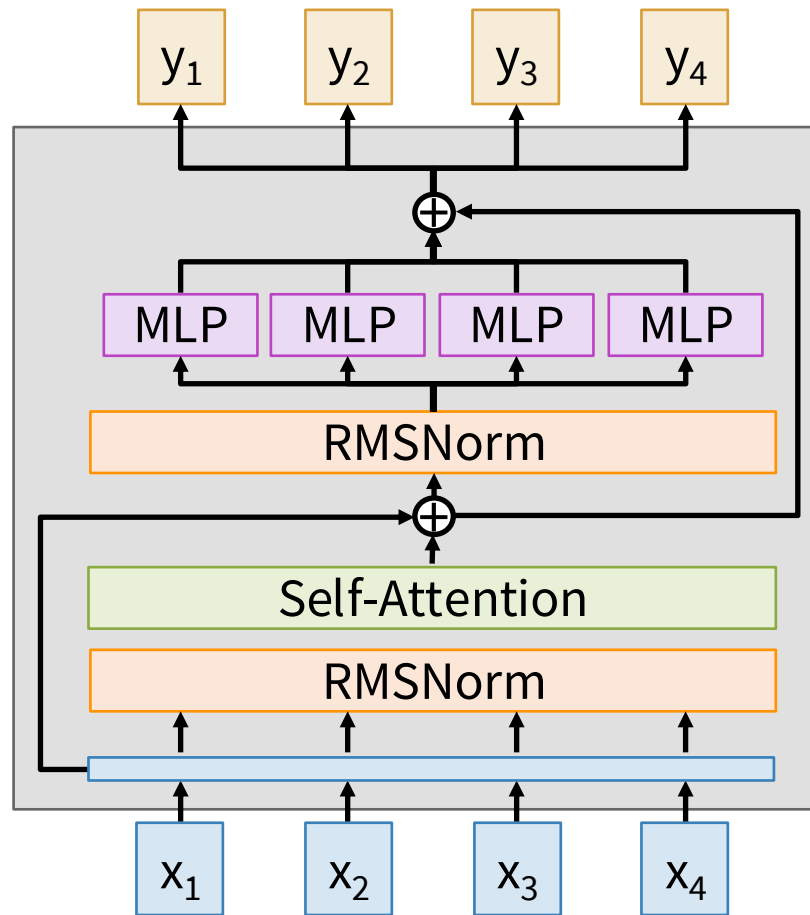
Output:  $y$  [shape  $D$ ]

Weight:  $\gamma$  [shape  $D$ ]

$$y_i = \frac{x_i}{RMS(x)} * \gamma_i$$

$$RMS(x) = \sqrt{\epsilon + \frac{1}{N} \sum_{i=1}^N x_i^2}$$

Training is a bit more stable



# SwiGLU MLP

Classic MLP:

Input:  $X [N \times D]$

Weights:  $W_1 [D \times 4D]$

$W_2 [4D \times D]$

Output:  $Y = \sigma(XW_1)W_2 [N \times D]$

SwiGLU MLP:

Input:  $X [N \times D]$

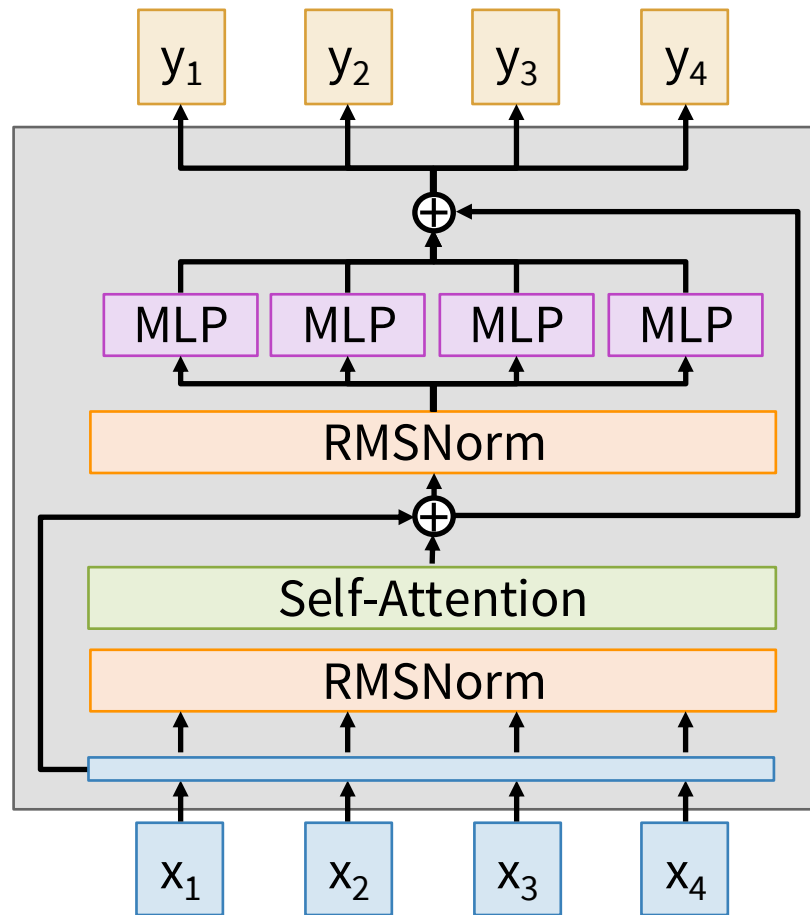
Weights:  $W_1, W_2 [D \times H]$

$W_3 [H \times D]$

Output:

$$Y = (\sigma(XW_1) \odot XW_2)W_3$$

Setting  $H = 8D/3$  keeps  
same total params

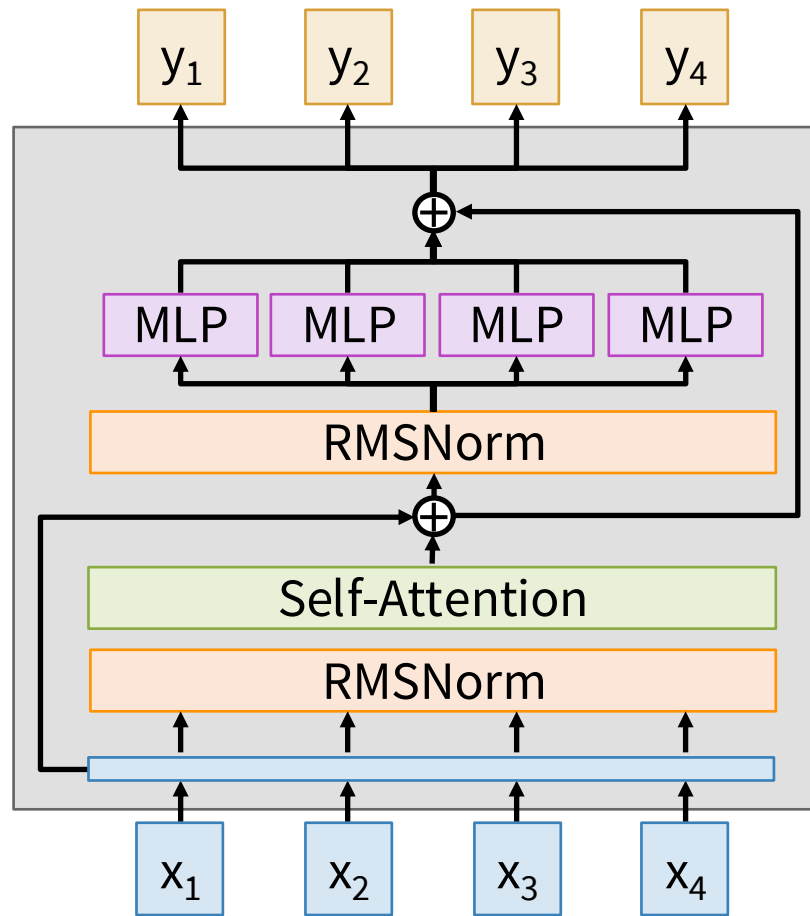


# Mixture of Experts (MoE)

Learn E separate sets of MLP weights in each block; each MLP is an expert

$W_1: [D \times 4D] \Rightarrow [E \times D \times 4D]$

$W_2: [4D \times D] \Rightarrow [E \times 4D \times D]$



Shazeer et al, "Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer", 2017

# Mixture of Experts (MoE)

Learn  $E$  separate sets of MLP weights in each block; each MLP is an expert

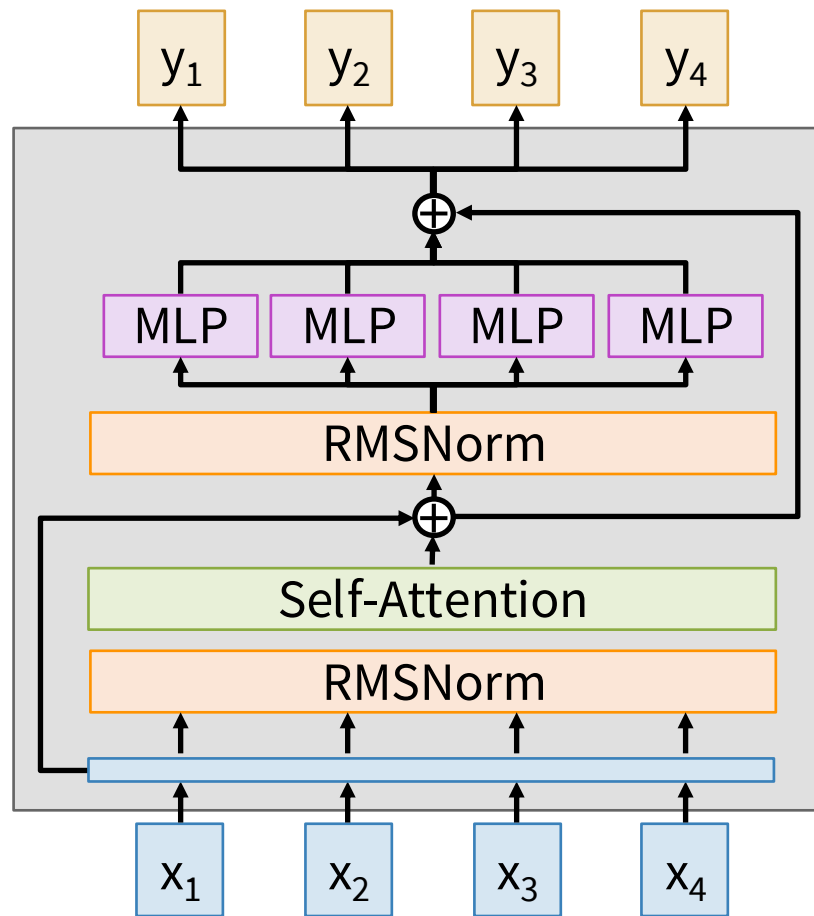
$W_1: [D \times 4D] \Rightarrow [E \times D \times 4D]$

$W_2: [4D \times D] \Rightarrow [E \times 4D \times D]$

Each token gets routed to  $A < E$  of the experts. These are the active experts.

Increases params by  $E$ ,  
But only increases compute by  $A$

All of the biggest LLMs today (e.g. GPT4o, GPT4.5, Claude 3.7, Gemini 2.5 Pro, etc) almost certainly use MoE and have  $> 1T$  params; but they don't publish details anymore

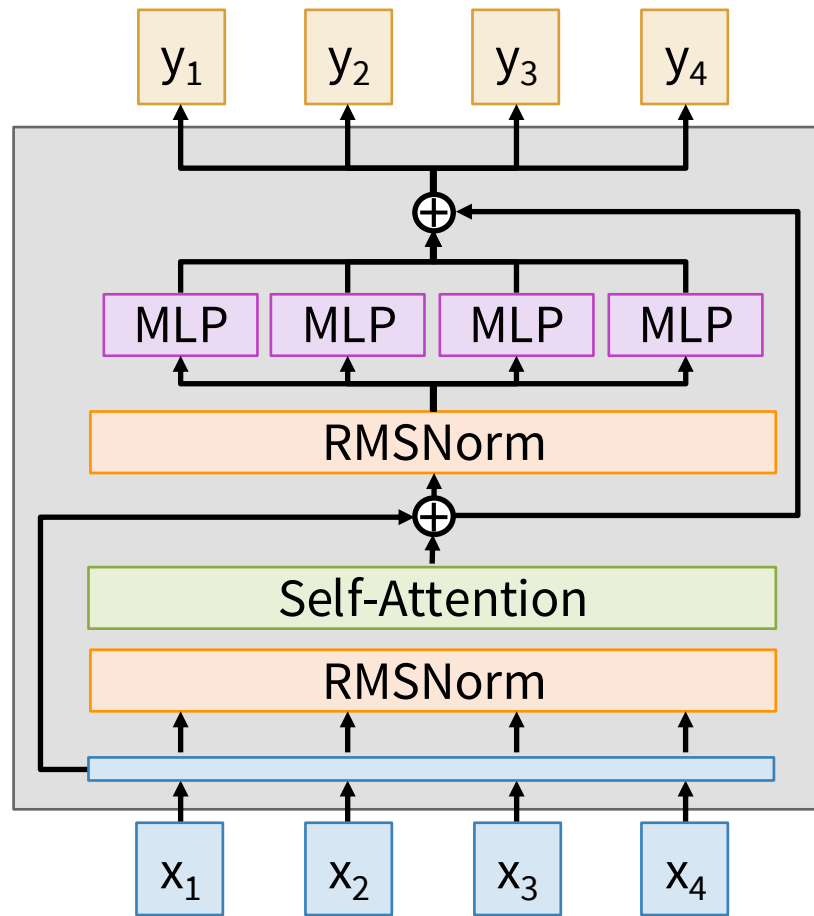


# Tweaking Transformers

The Transformer architecture has not changed much since 2017.

But a few changes have become common:

- Pre-Norm: Move normalization inside residual
- RMSNorm: Different normalization layer
- SwiGLU: Different MLP architecture
- Mixture of Experts (MoE): Learn  $E$  different MLPs, use  $A < E$  of them per token. Massively increase params, modest increase to compute cost.



# Today

- Transformers Recap
- **Computer Vision Tasks**
  - Semantic Segmentation
  - Object Detection
  - Instance Segmentation
- Visualization & Understanding
  - Model Layers Visualization
  - Saliency Maps
  - CAM & Grad-CAM

# Computer Vision Tasks

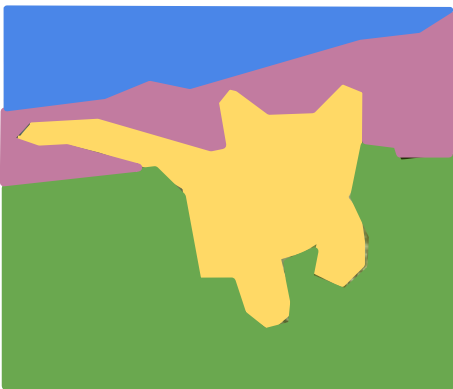
## Classification



CAT

No spatial extent

## Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

## Object Detection



DOG, DOG, CAT

Multiple Object

## Instance Segmentation



DOG, DOG, CAT

[This image is CC0 public domain](#)

# Recall: Image Classification: A core task in Computer Vision



This image by Nikita is licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)  
{dog, cat, truck, plane, ...}



cat

# Semantic Segmentation

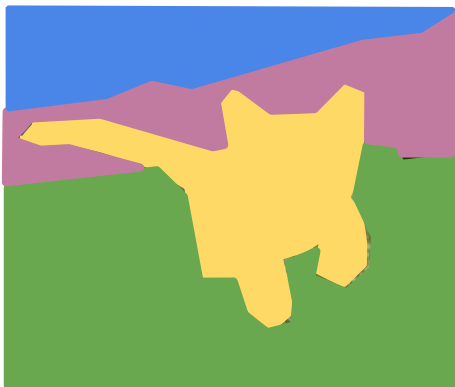
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



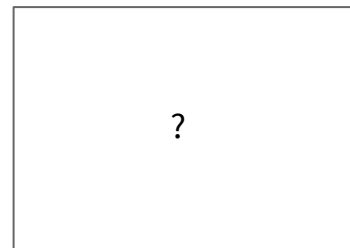
DOG, DOG, CAT

# Semantic Segmentation: The Problem



GRASS, CAT, TREE,  
SKY, ...

Paired training data: for each training image,  
each pixel is labeled with a semantic category.



At test time, classify each pixel of a new image.

# Semantic Segmentation Idea: Sliding Window

Full image



# Semantic Segmentation Idea: Sliding Window

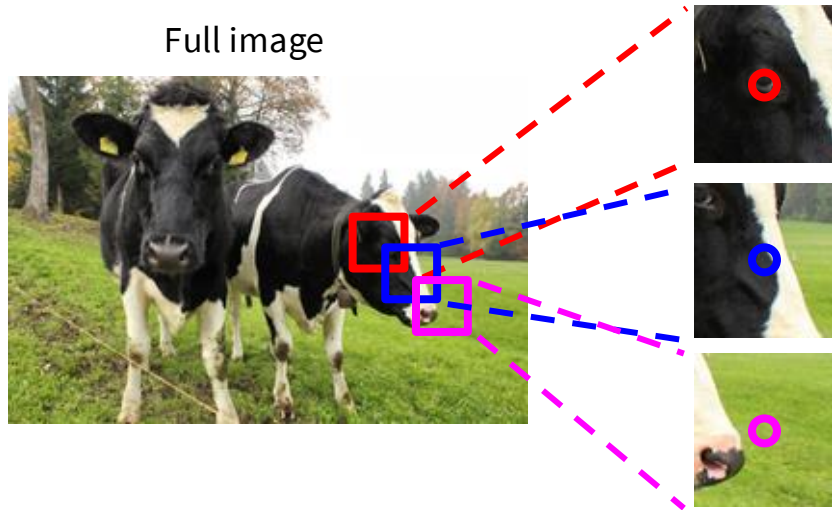
Full image



Impossible to classify without context

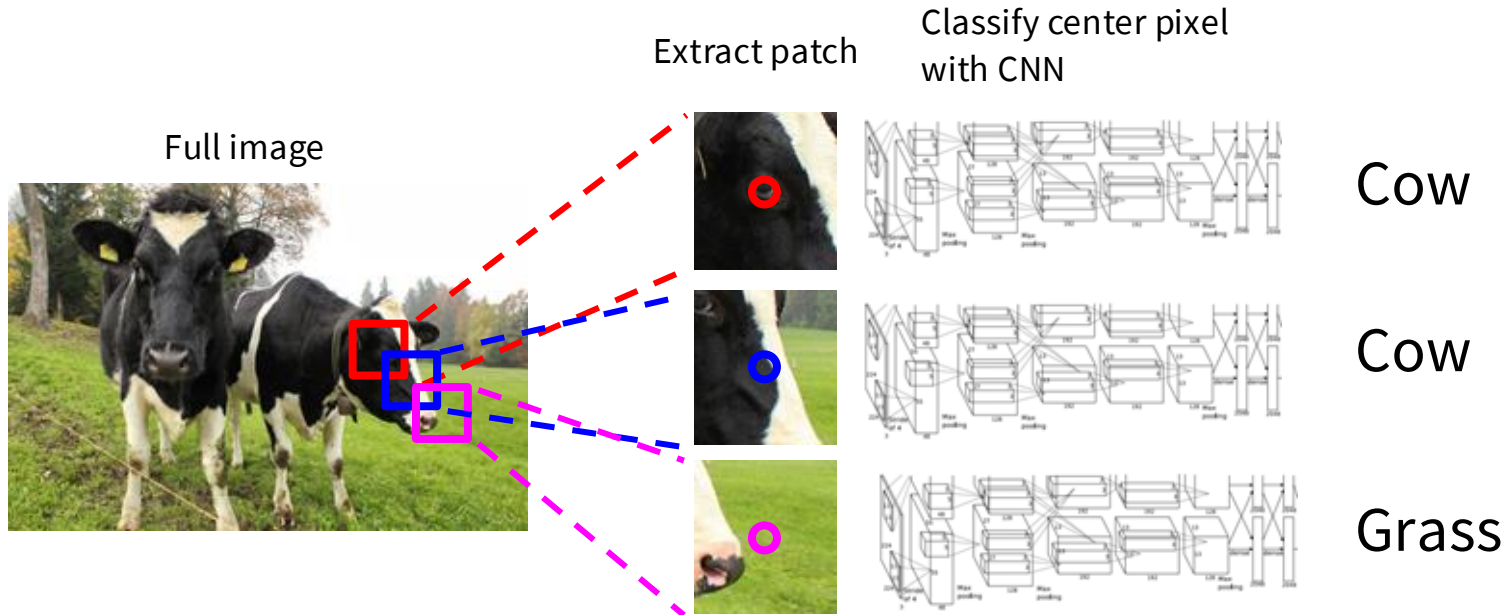
Q: how do we include context?

# Semantic Segmentation Idea: Sliding Window



Q: how do we model this?

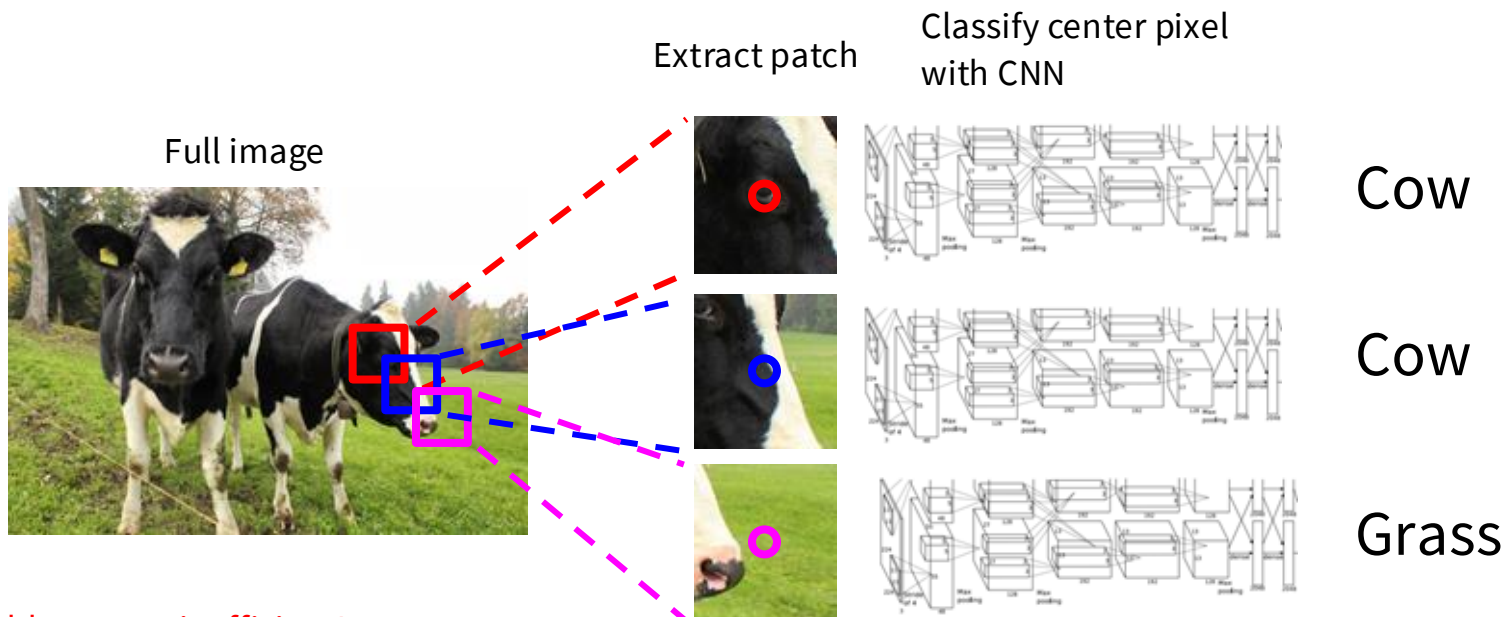
# Semantic Segmentation Idea: Sliding Window



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

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

# Semantic Segmentation Idea: Sliding Window



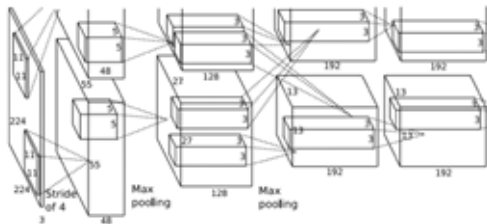
Problem: Very inefficient! Not reusing shared features between overlapping patches

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

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

# Semantic Segmentation Idea: Convolution

Full image

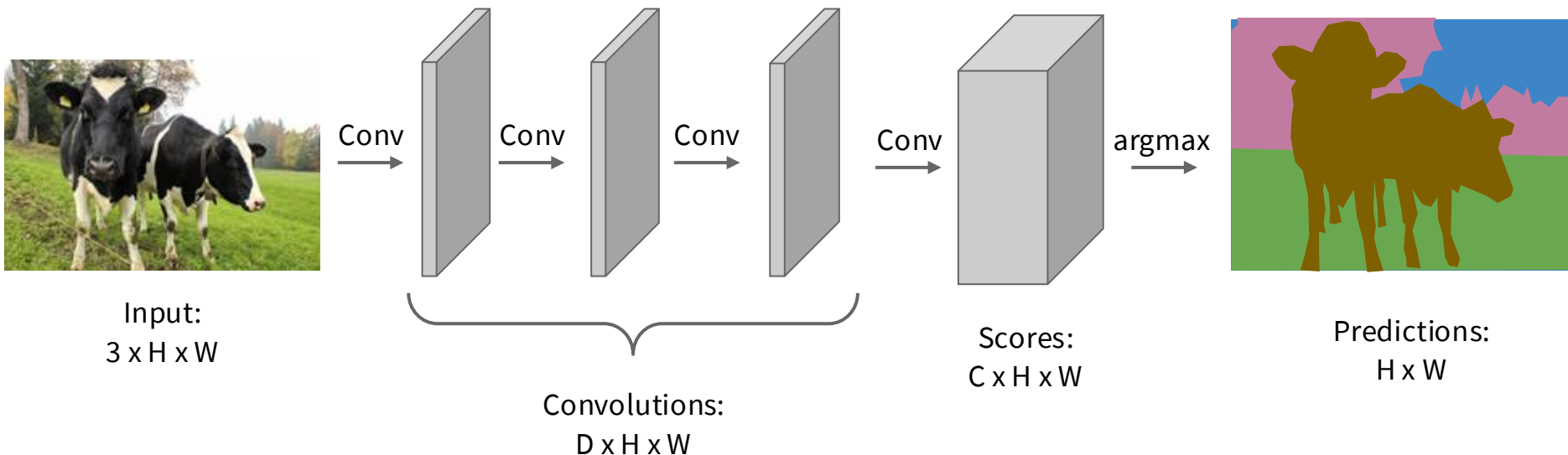


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

**Problem:** classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

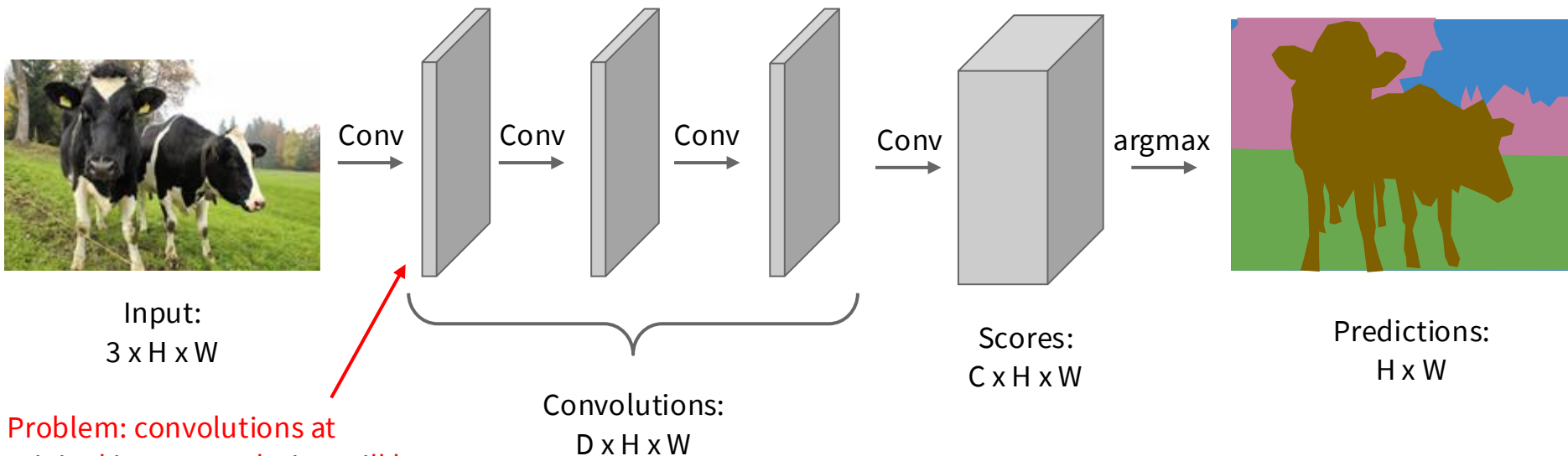
# Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



# Semantic Segmentation Idea: Fully Convolutional

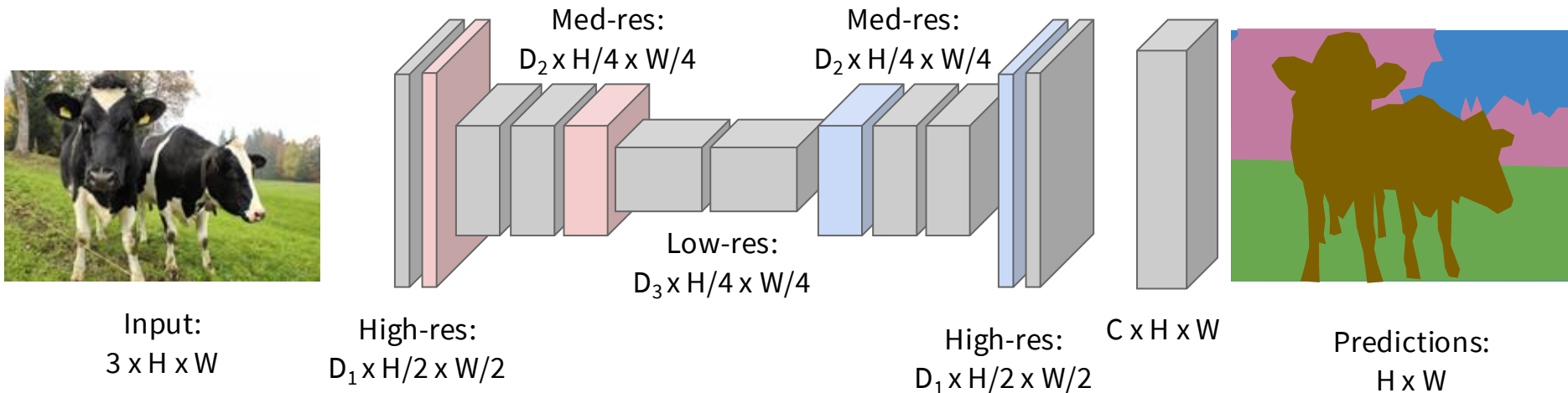
Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



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

# Semantic Segmentation Idea: Fully Convolutional

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



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

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

# Semantic Segmentation Idea: Fully Convolutional

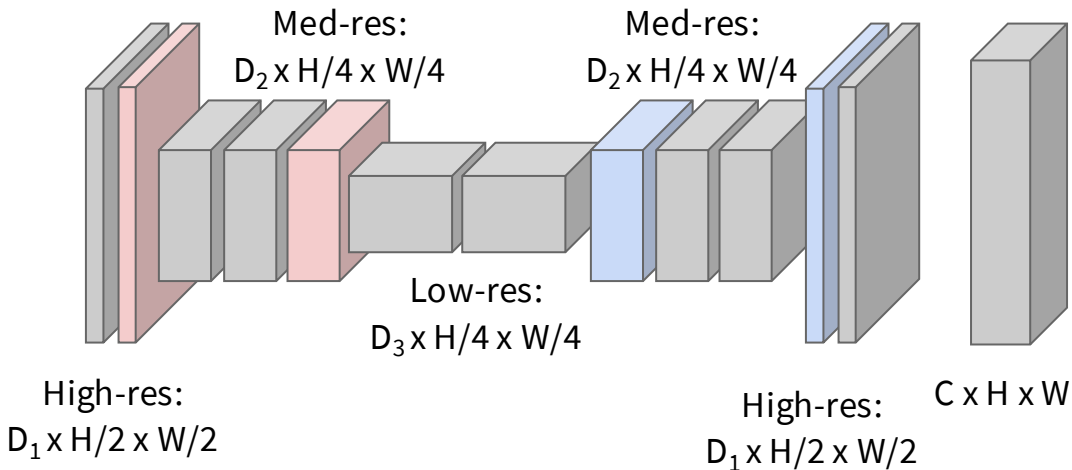
Downsampling:  
Pooling, strided  
convolution

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

Upsampling:  
???



Input:  
 $3 \times H \times W$



Predictions:  
 $H \times W$

# 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

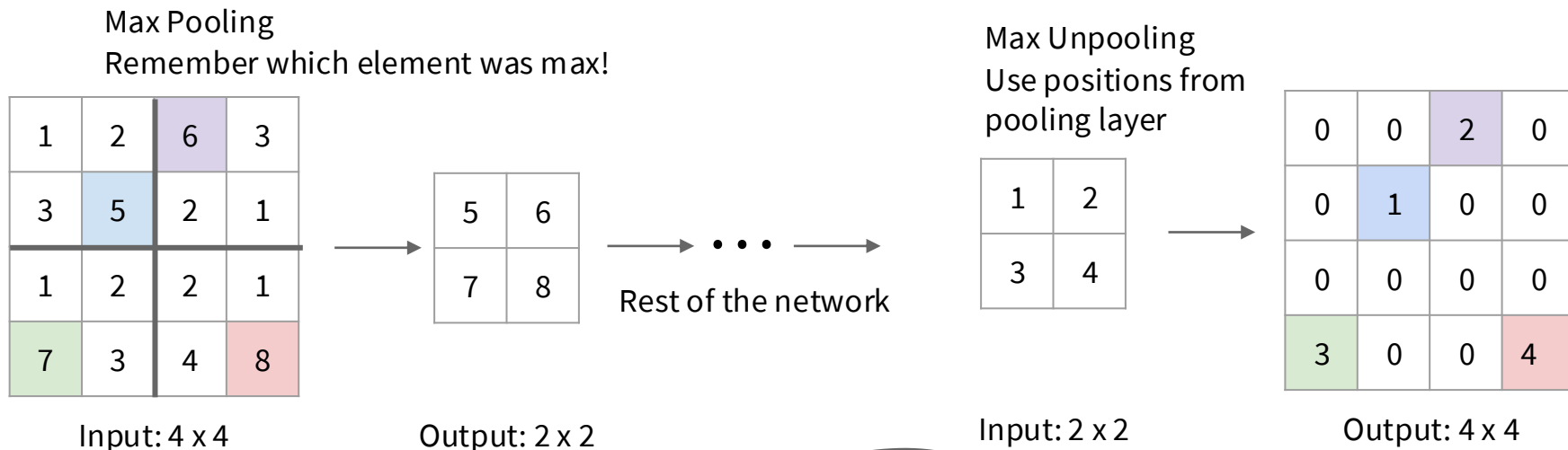
Input: 2 x 2



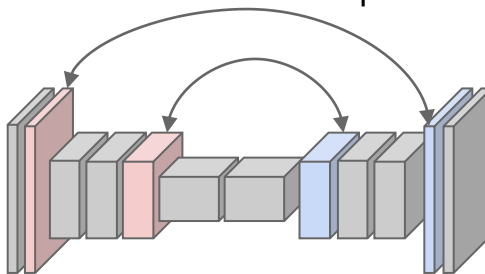
1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output: 4 x 4

# In-Network upsampling: “Max Unpooling”

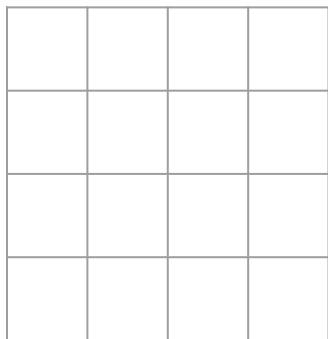


Corresponding pairs of downsampling and upsampling layers

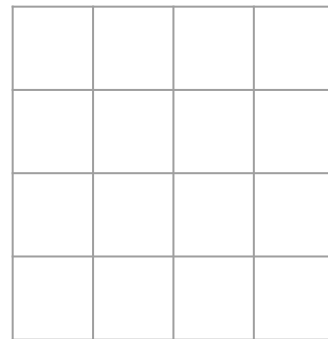


# Learnable Upsampling

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



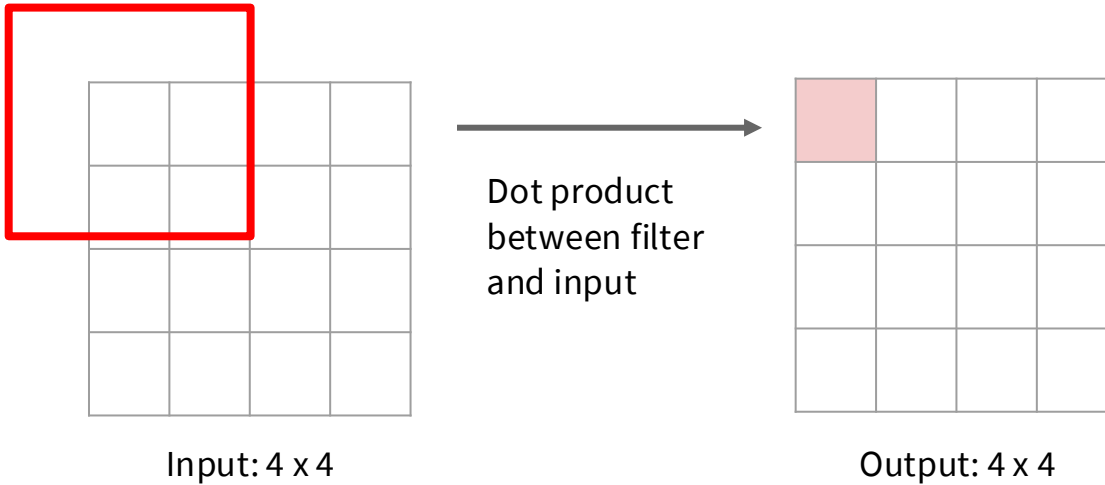
Input: 4 x 4



Output: 4 x 4

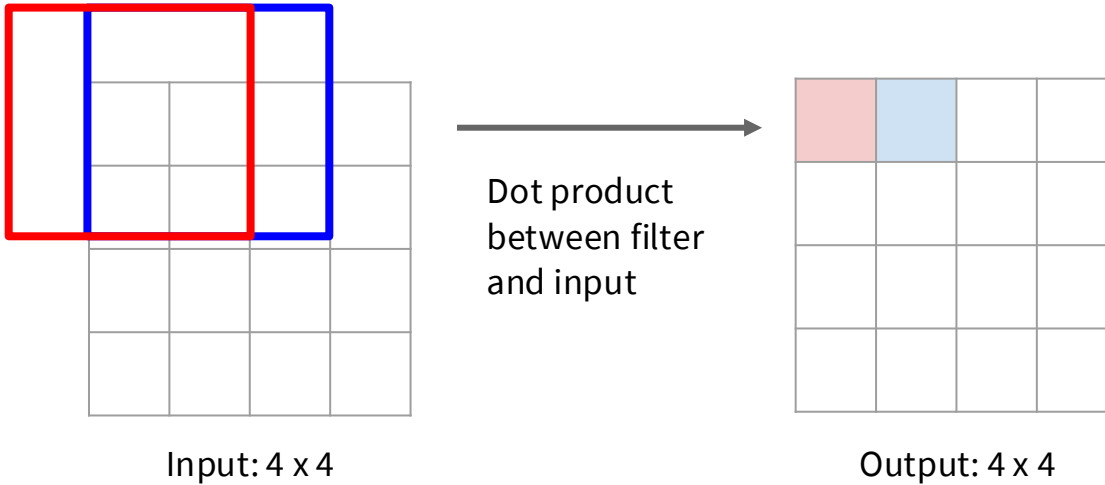
# Learnable Upsampling

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



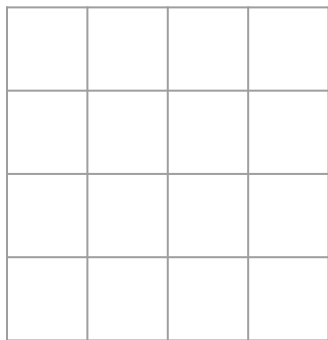
# Learnable Upsampling

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

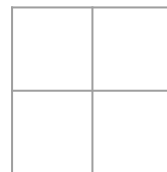


# Learnable Upsampling

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



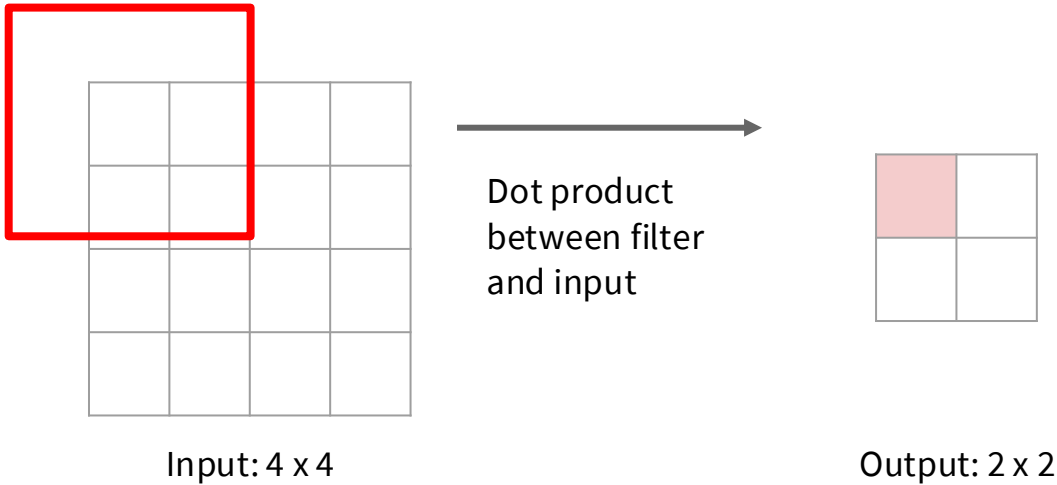
Input: 4 x 4



Output: 2 x 2

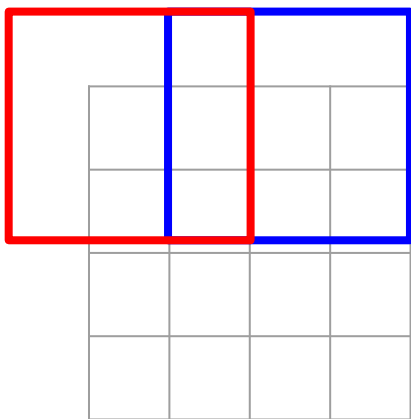
# Learnable Upsampling

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



# Learnable Upsampling

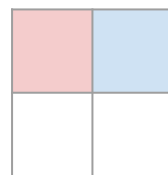
Recall: Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4



Dot product  
between filter  
and input



Output: 2 x 2

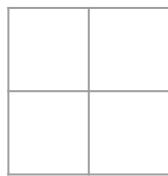
Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

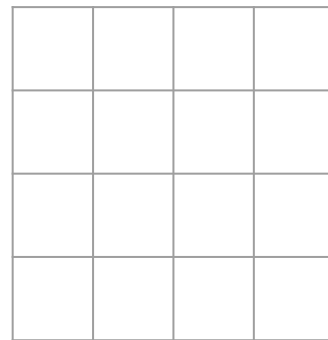
We can interpret strided convolution as “learnable downsampling”.

# Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1



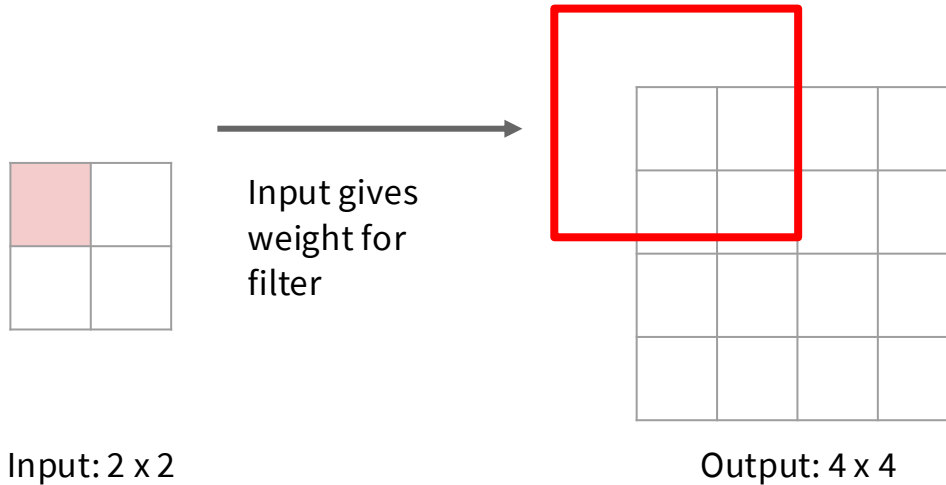
Input: 2 x 2



Output: 4 x 4

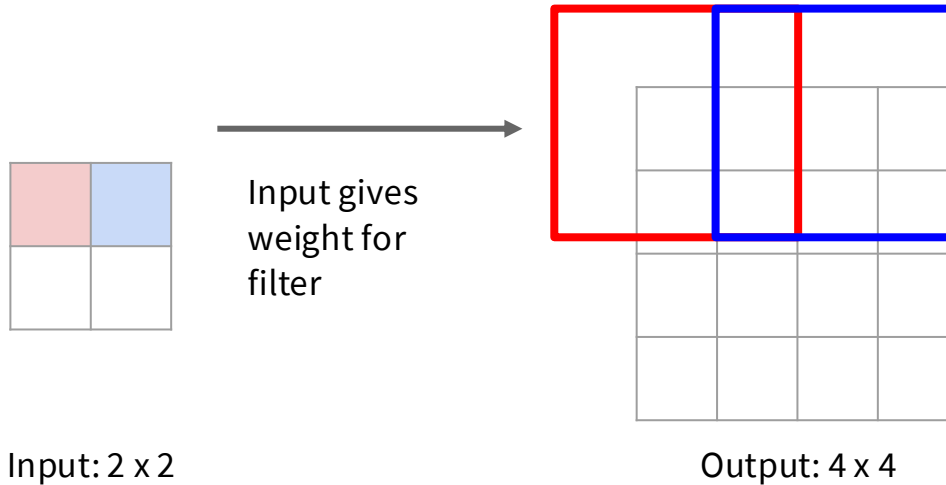
# Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1



# Learnable Upsampling: Transposed Convolution

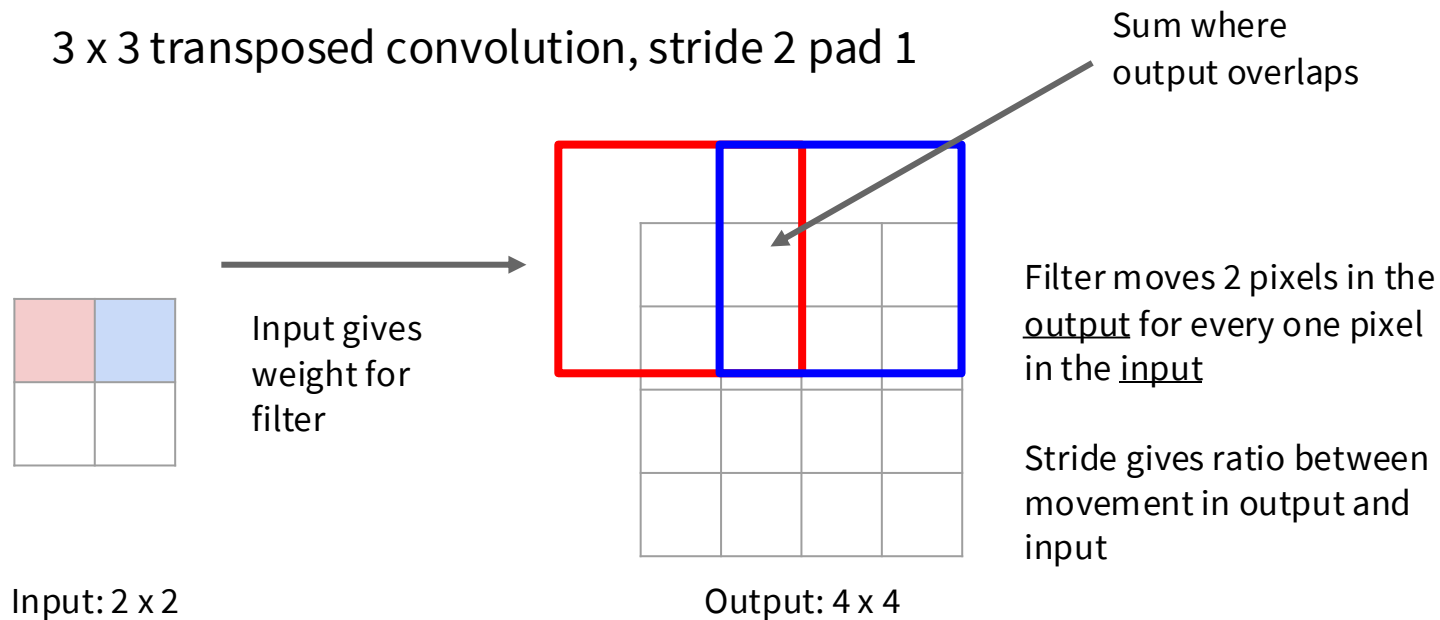
3 x 3 transposed convolution, stride 2 pad 1



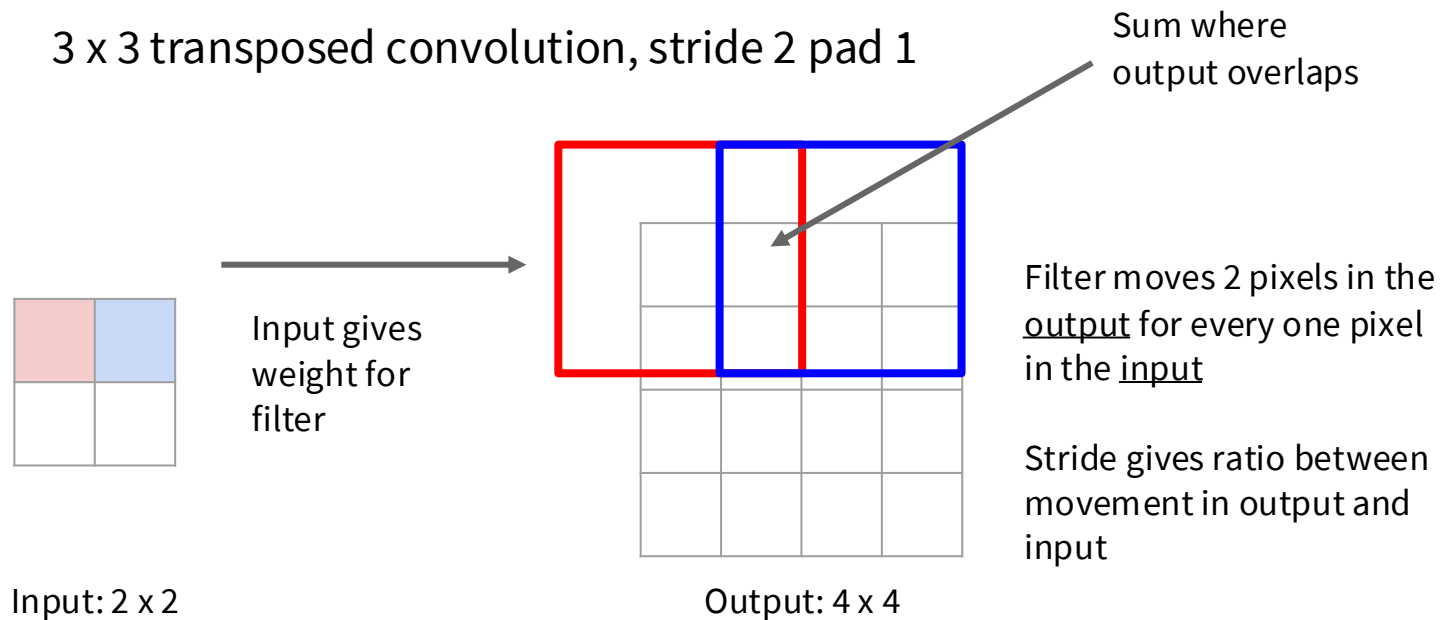
Filter moves 2 pixels in the output for every one pixel in the input

Stride gives ratio between movement in output and input

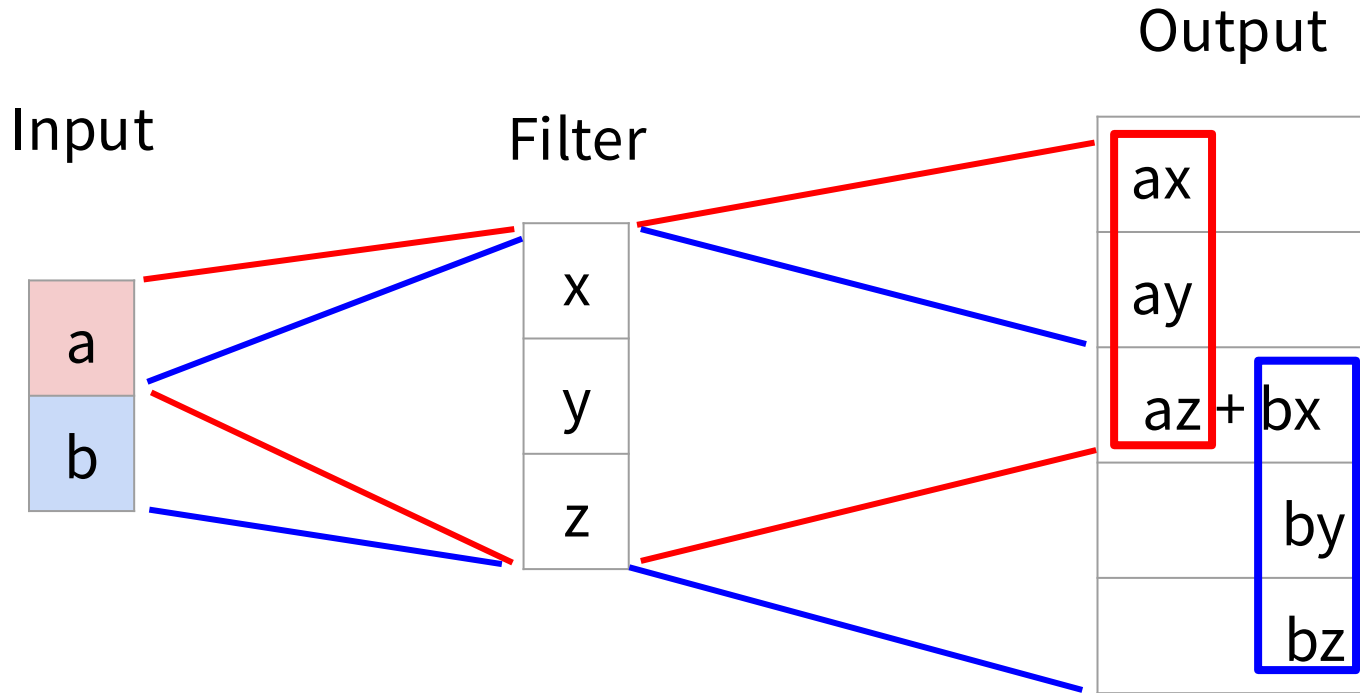
# Learnable Upsampling: Transposed Convolution



# Learnable Upsampling: Transposed Convolution



# Learnable Upsampling: 1D Example



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

# Semantic Segmentation Idea: Fully Convolutional

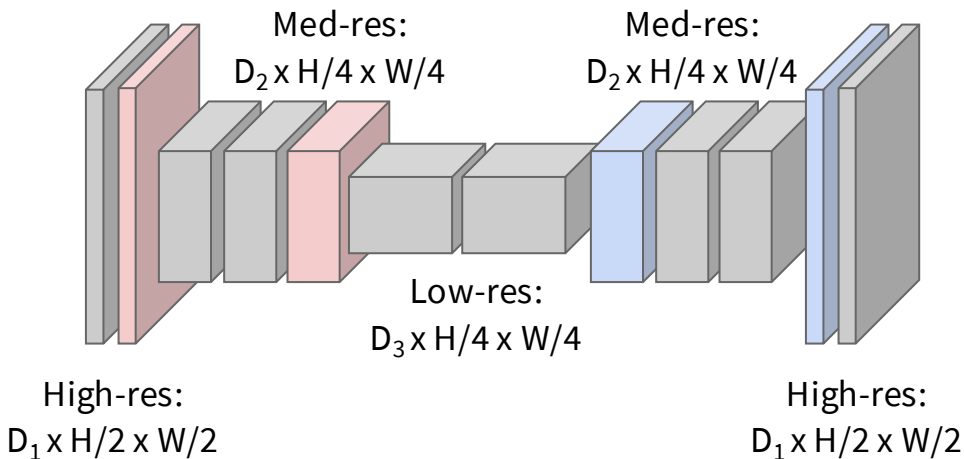
Downsampling:  
Pooling, strided  
convolution

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

Upsampling:  
Unpooling or strided  
transposed convolution



Input:  
 $3 \times H \times W$

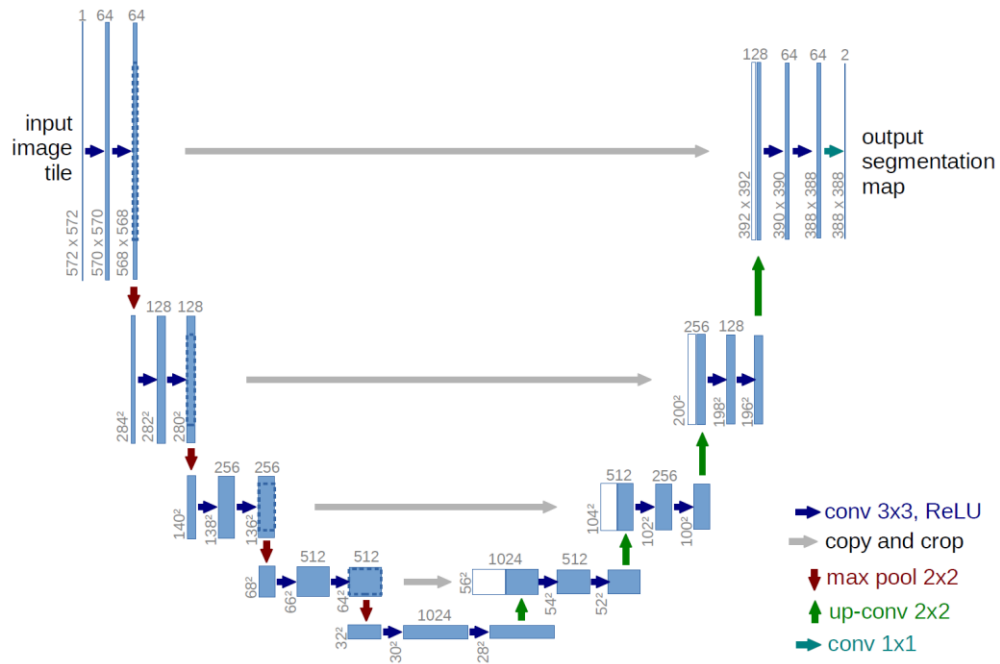


Predictions:  
 $H \times W$

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

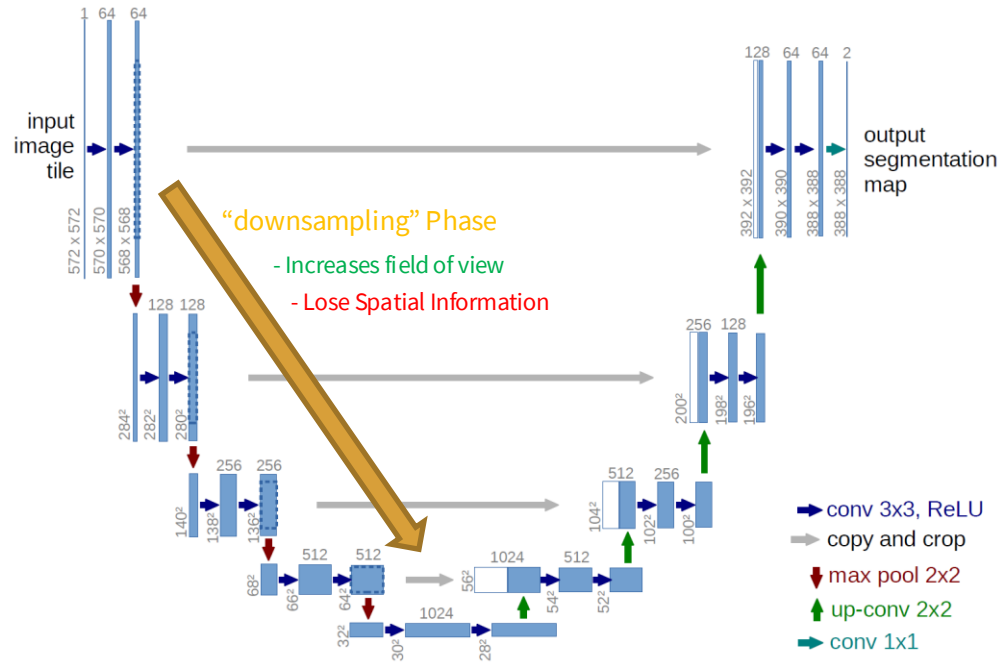
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# U-Net



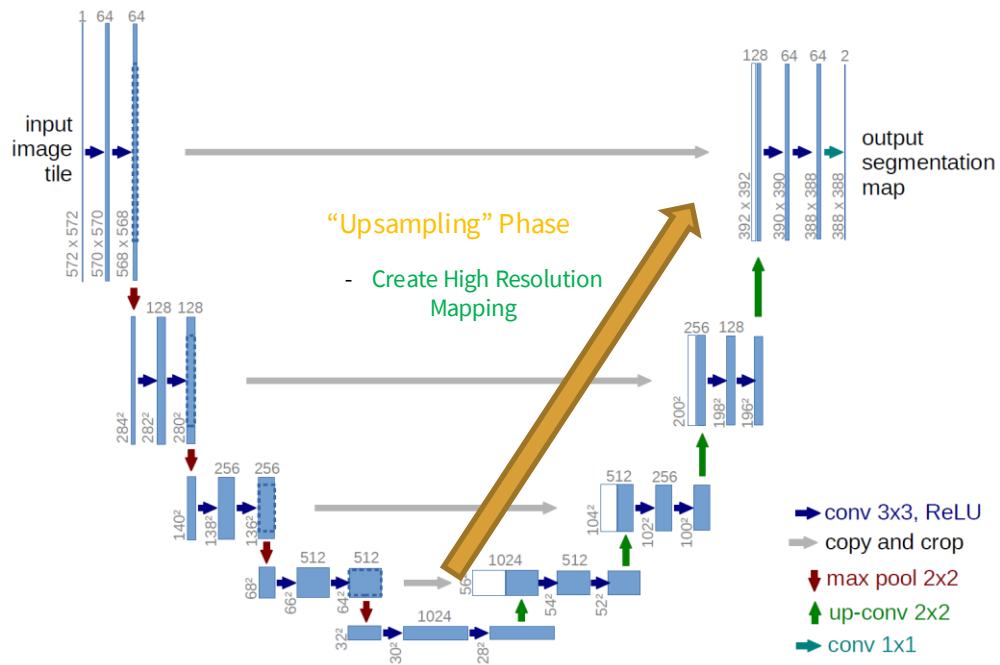
Ronneberger et al. (2015) U-net Architecture

# U-Net



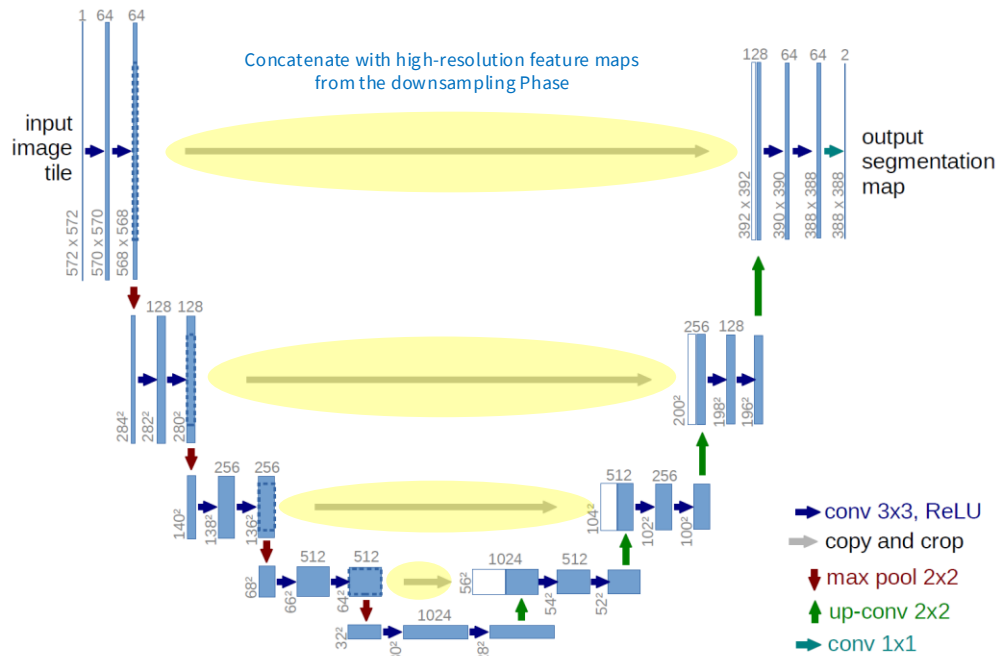
Ronneberger et al. (2015) U-net Architecture

# U-Net



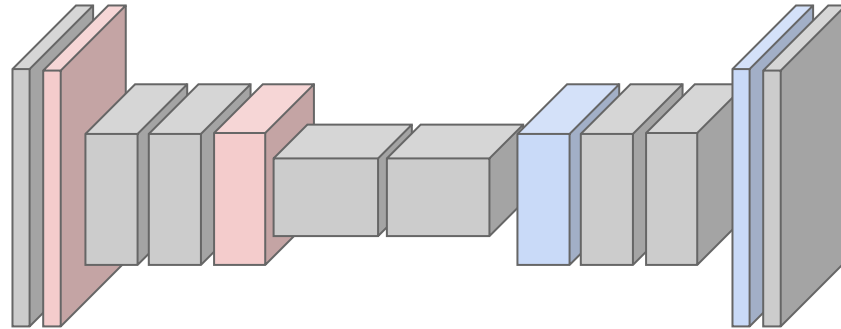
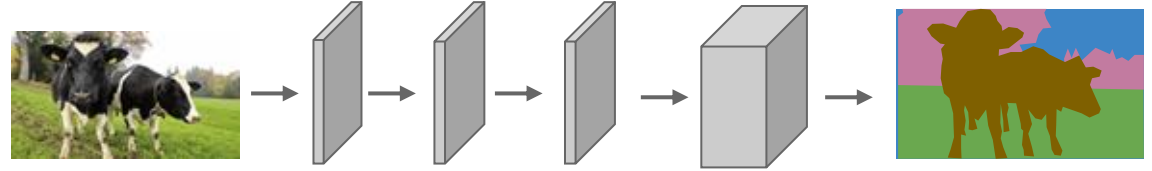
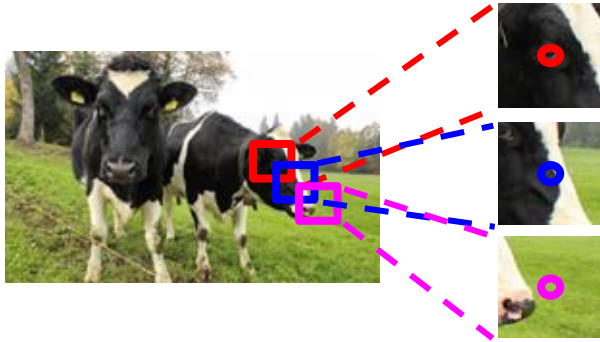
Ronneberger et al. (2015) U-net Architecture

# U-Net



Ronneberger et al. (2015) U-net Architecture

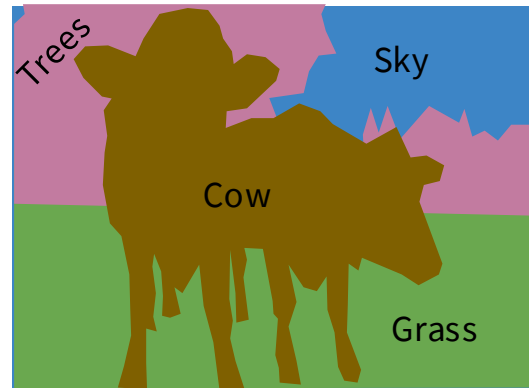
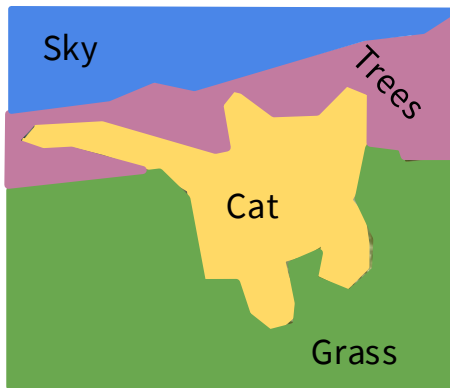
# Semantic Segmentation: Summary



# Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



# Object Detection

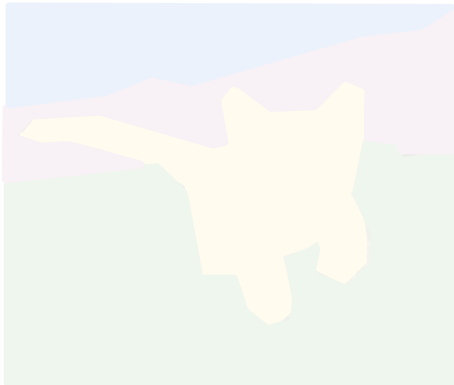
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

# Object Detection

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

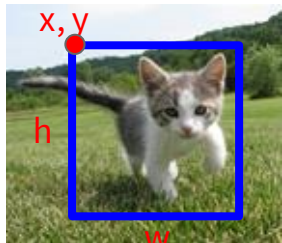
Instance Segmentation



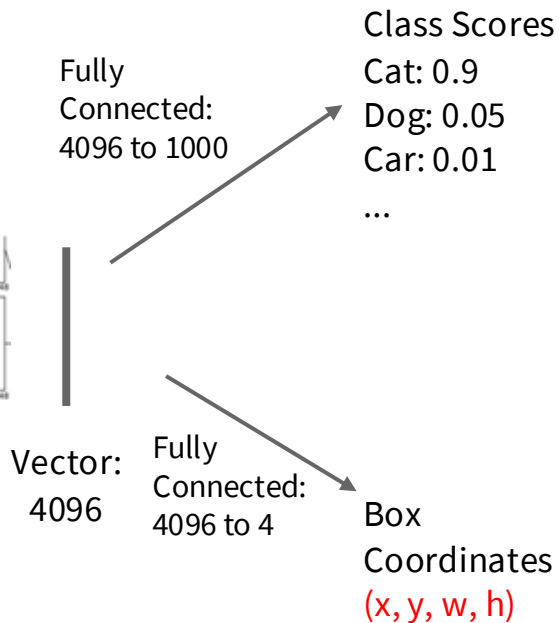
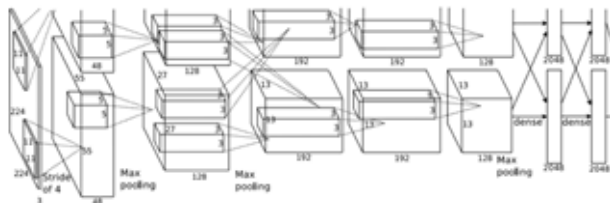
DOG, DOG, CAT

# Object Detection: Single Object

(Classification + Localization)

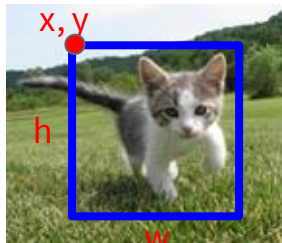


[This image is CC0 public domain](#)

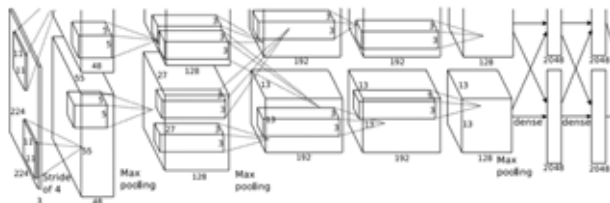


# Object Detection: Single Object

(Classification + Localization)



[This image is CC0 public domain](#)



Vector:  
4096

Fully  
Connected:  
4096 to 1000

Class Scores

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

Correct label:  
Cat

Softmax  
Loss

Fully  
Connected:  
4096 to 4

Box  
Coordinates  
(x, y, w, h)

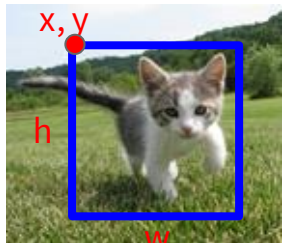
L2 Loss

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

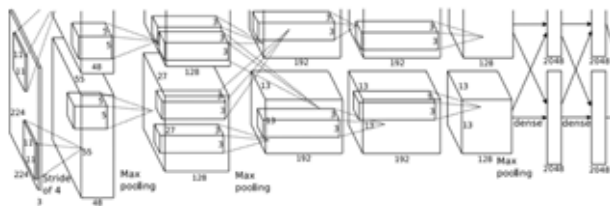
Treat localization as a  
regression problem!

# Object Detection: Single Object

(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

**Multitask Loss**

+ → Loss

Vector:  
4096

Fully Connected:  
4096 to 4

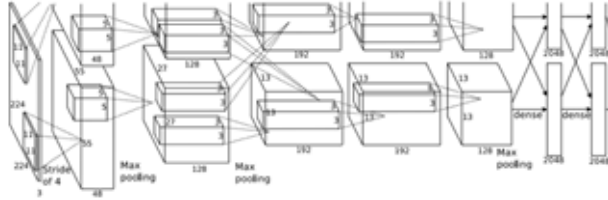
Box  
Coordinates  
(x, y, w, h)

L2 Loss

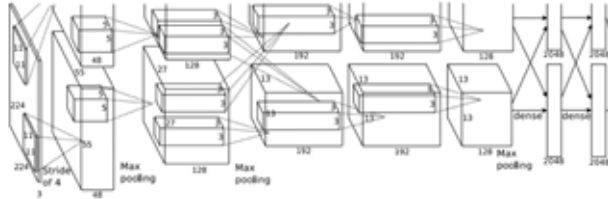
Treat localization as a  
regression problem!

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

# Object Detection: Multiple Objects



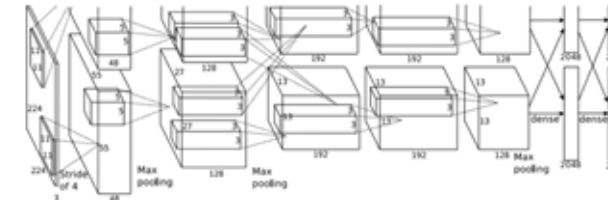
CAT:  $(x, y, w, h)$



DOG:  $(x, y, w, h)$

DOG:  $(x, y, w, h)$

CAT:  $(x, y, w, h)$



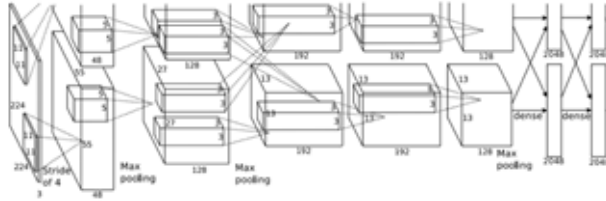
DUCK:  $(x, y, w, h)$

DUCK:  $(x, y, w, h)$

....

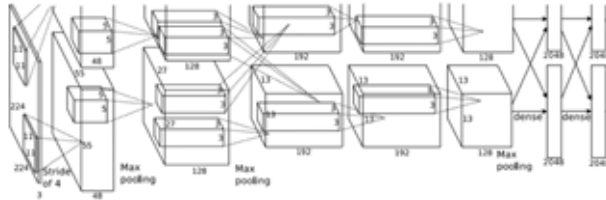
# Object Detection: Multiple Objects

Each image needs a different number of outputs!



CAT:  $(x, y, w, h)$

4 numbers

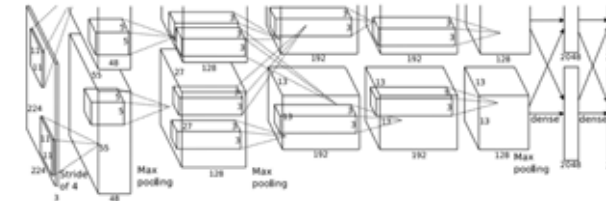


DOG:  $(x, y, w, h)$

DOG:  $(x, y, w, h)$

CAT:  $(x, y, w, h)$

12 numbers



DUCK:  $(x, y, w, h)$

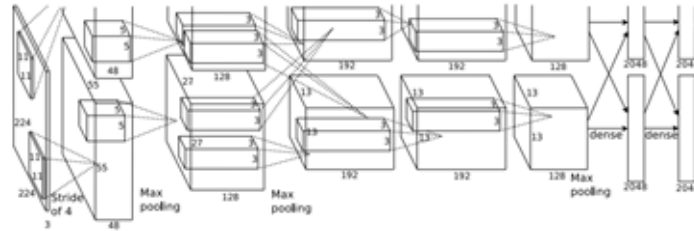
DUCK:  $(x, y, w, h)$

....

Many numbers!

# Object Detection: Multiple Objects

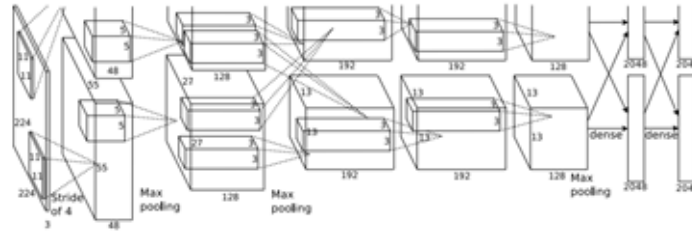
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? NO  
Background? YES

# Object Detection: Multiple Objects

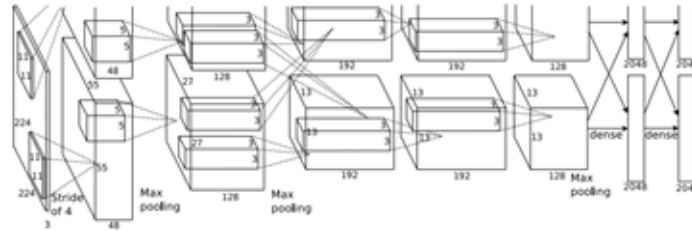
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Multiple Objects

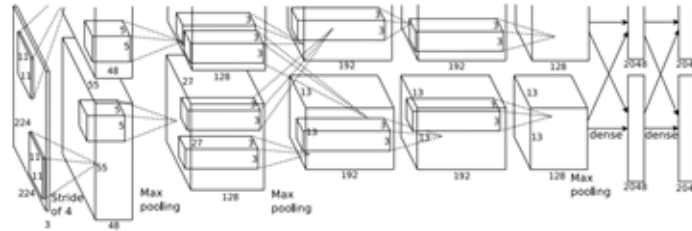
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Multiple Objects

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

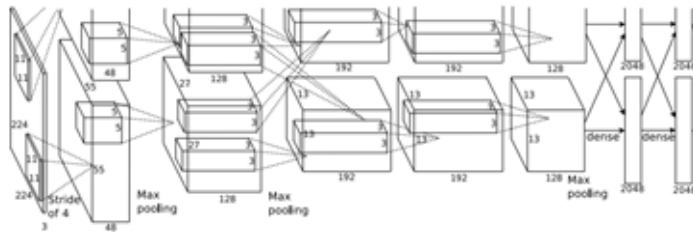
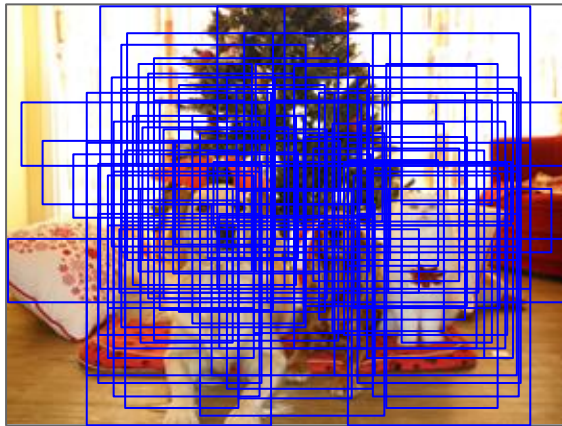


Dog? NO  
Cat? YES  
Background? NO

Q: What's the problem with this approach?

# Object Detection: Multiple Objects

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

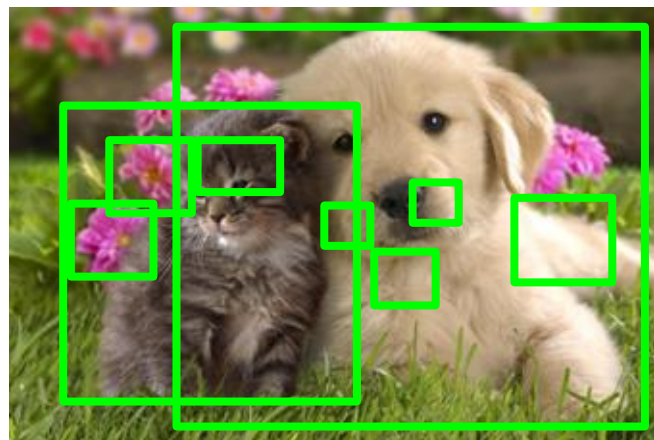


Dog? NO  
Cat? YES  
Background? NO

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

# 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



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



Input image

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

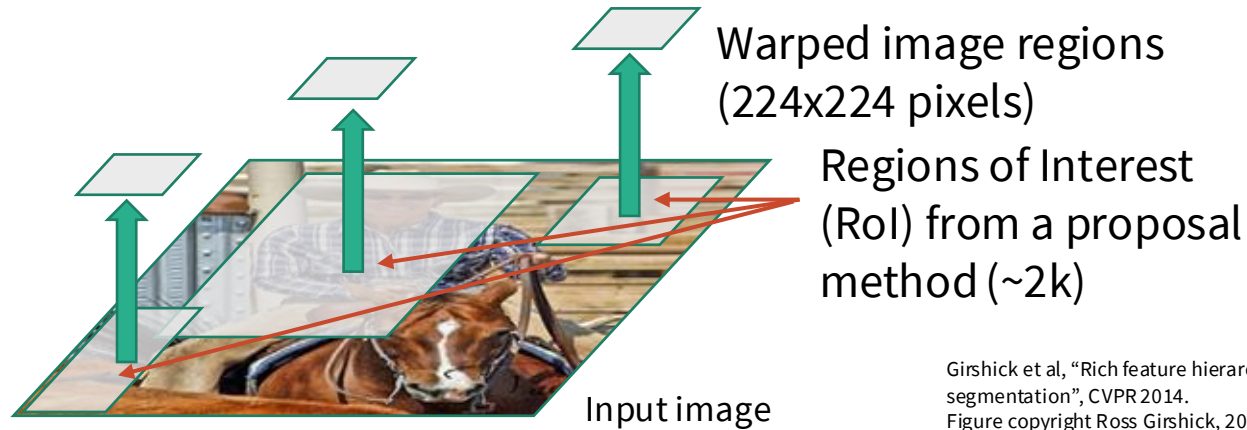
# R-CNN



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

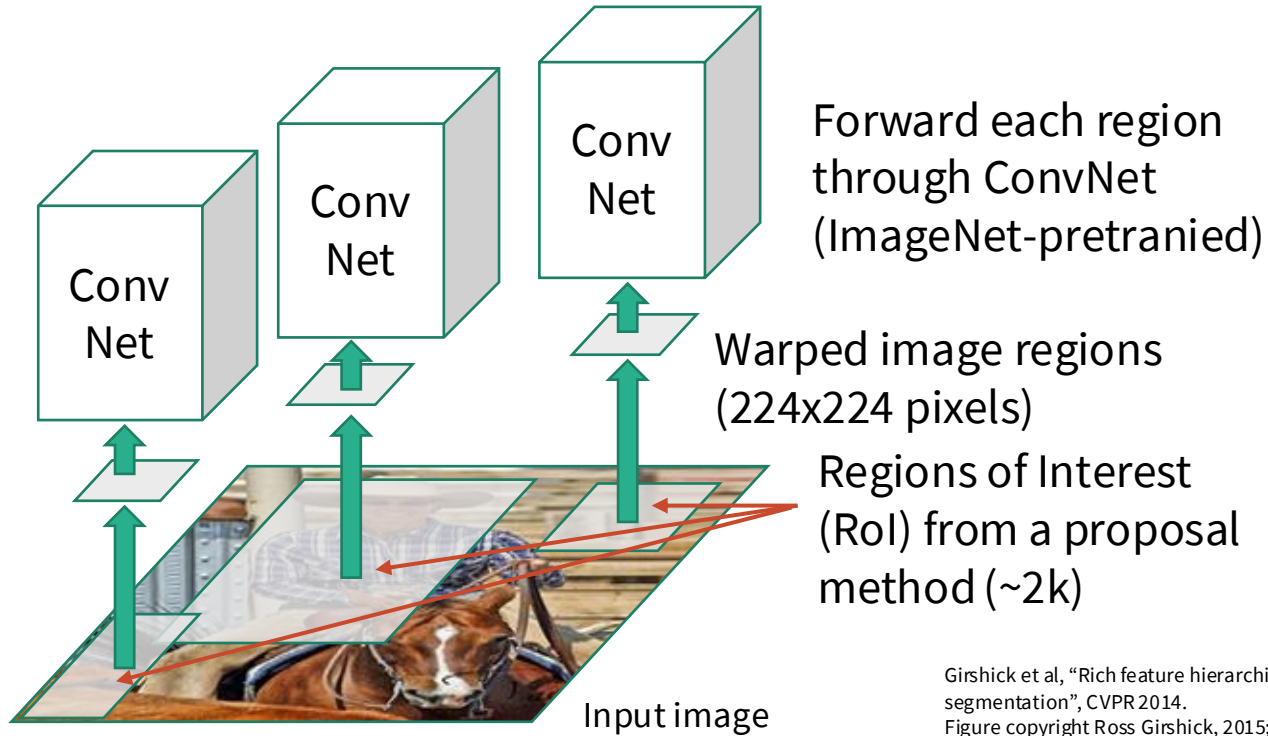
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN



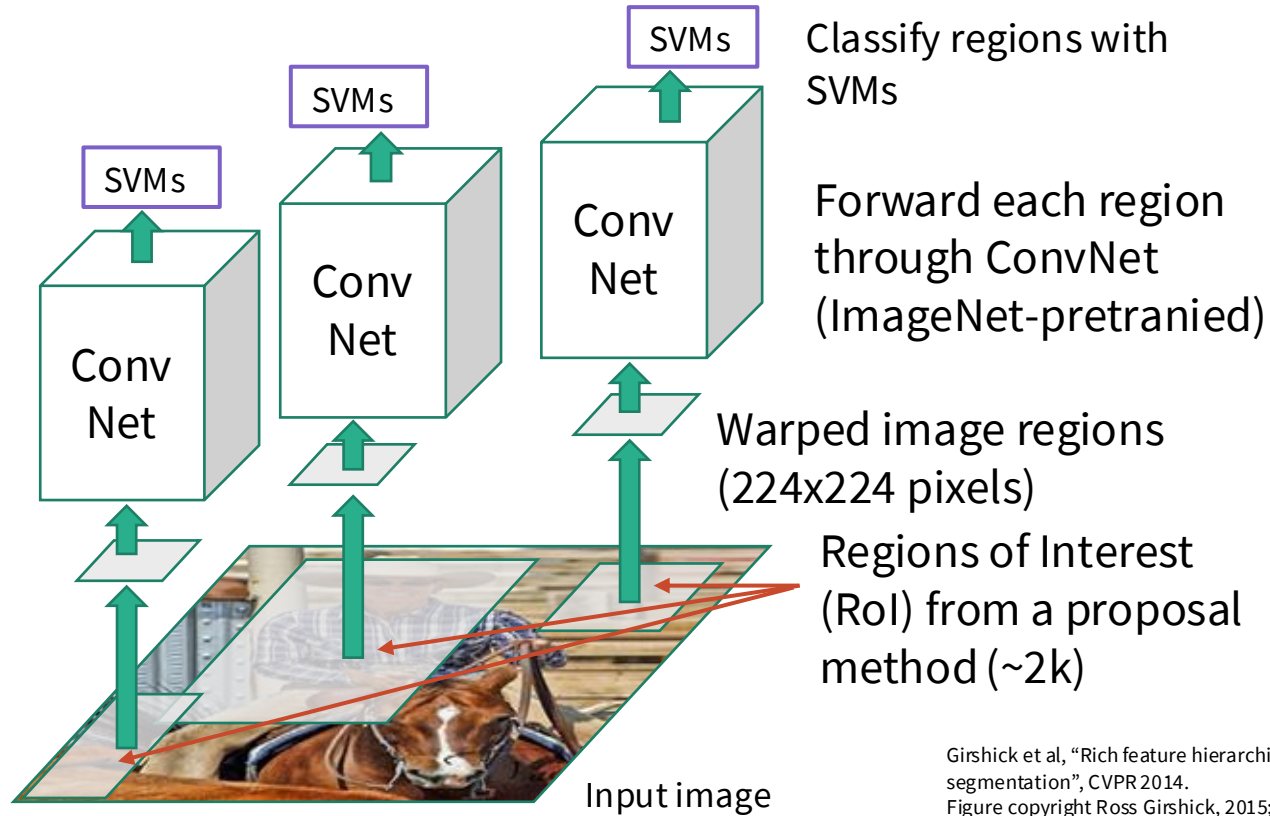
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN



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

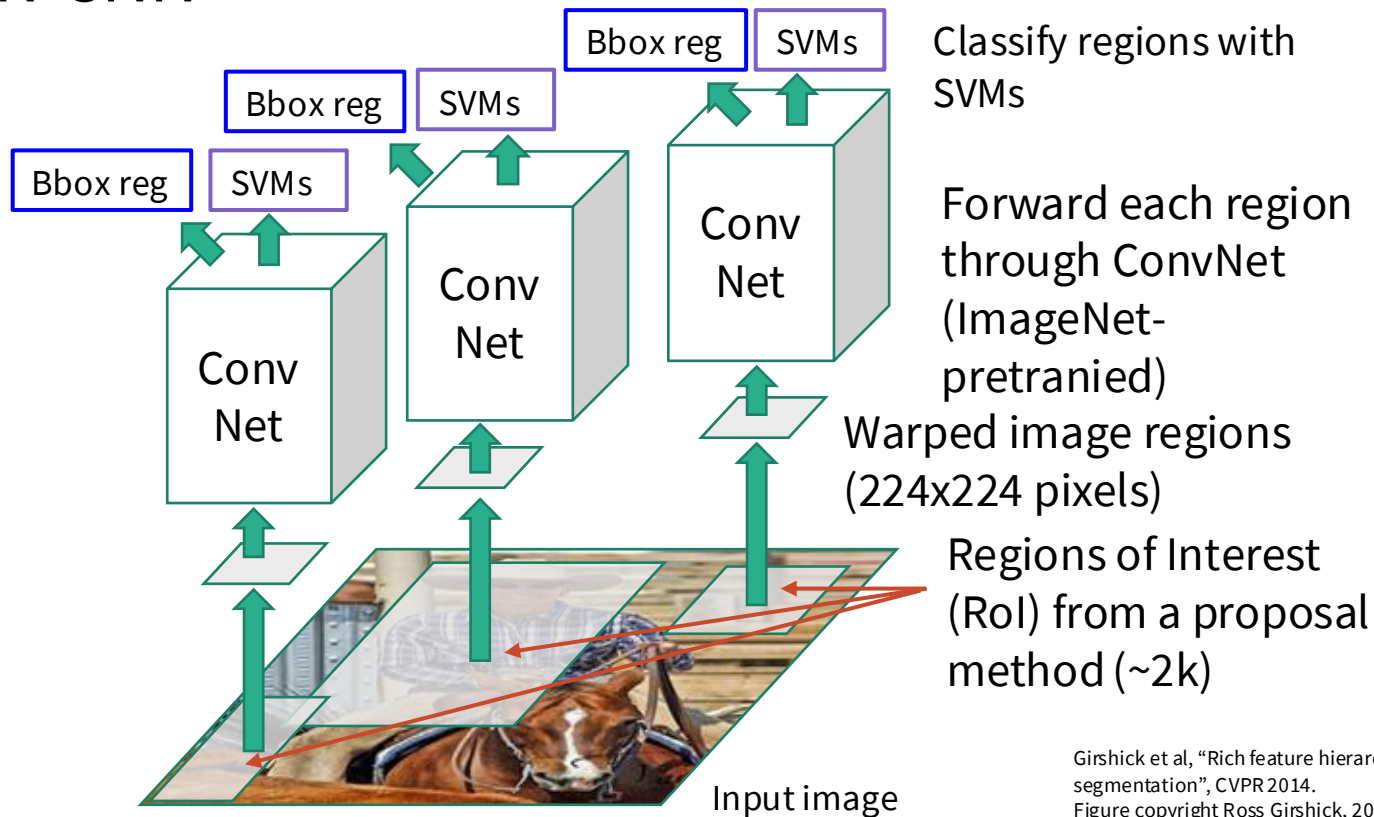
# R-CNN



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

# R-CNN

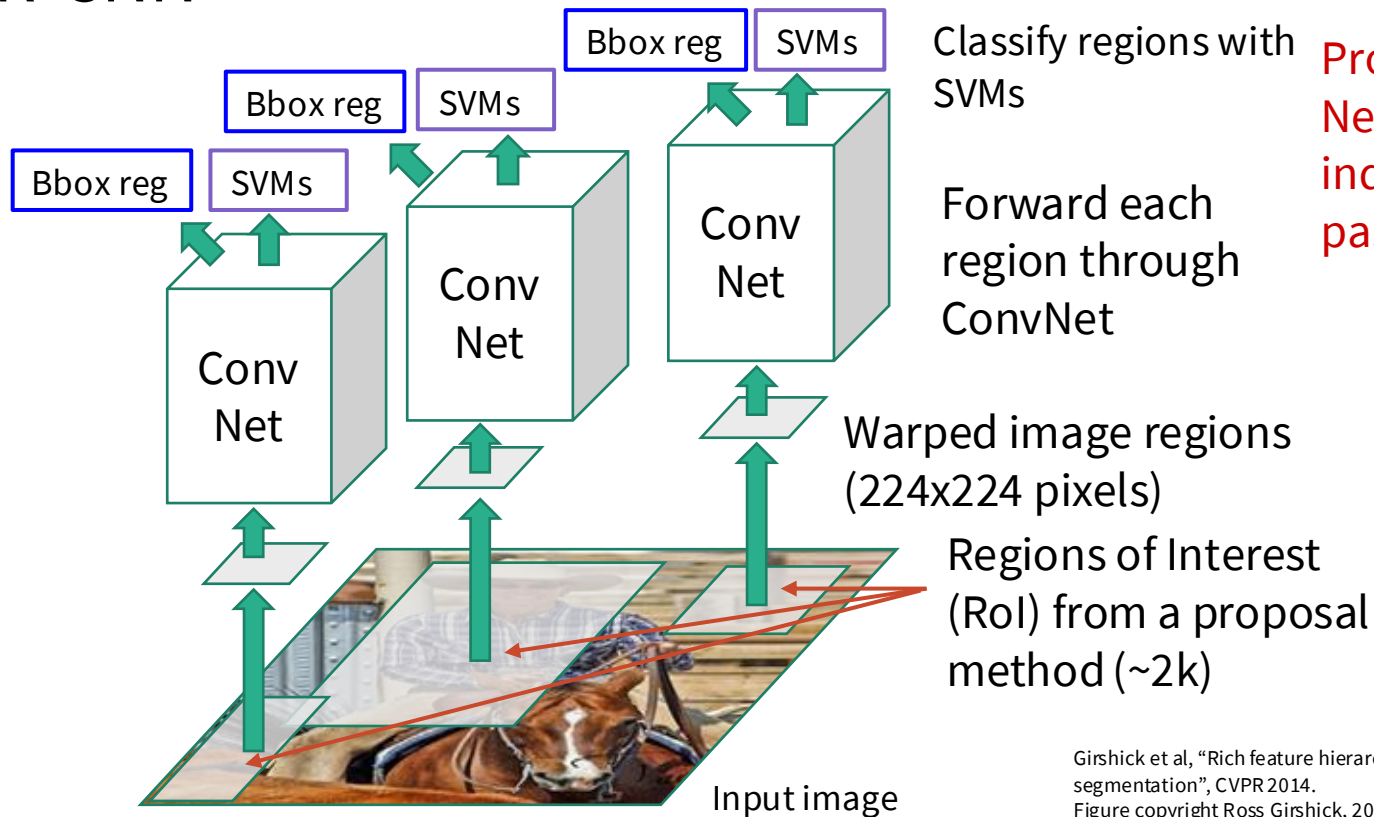
Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



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

# R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



Classify regions with SVMs

Forward each region through ConvNet

**Problem: Very slow!**  
**Need to do ~2k independent forward passes for each image!**

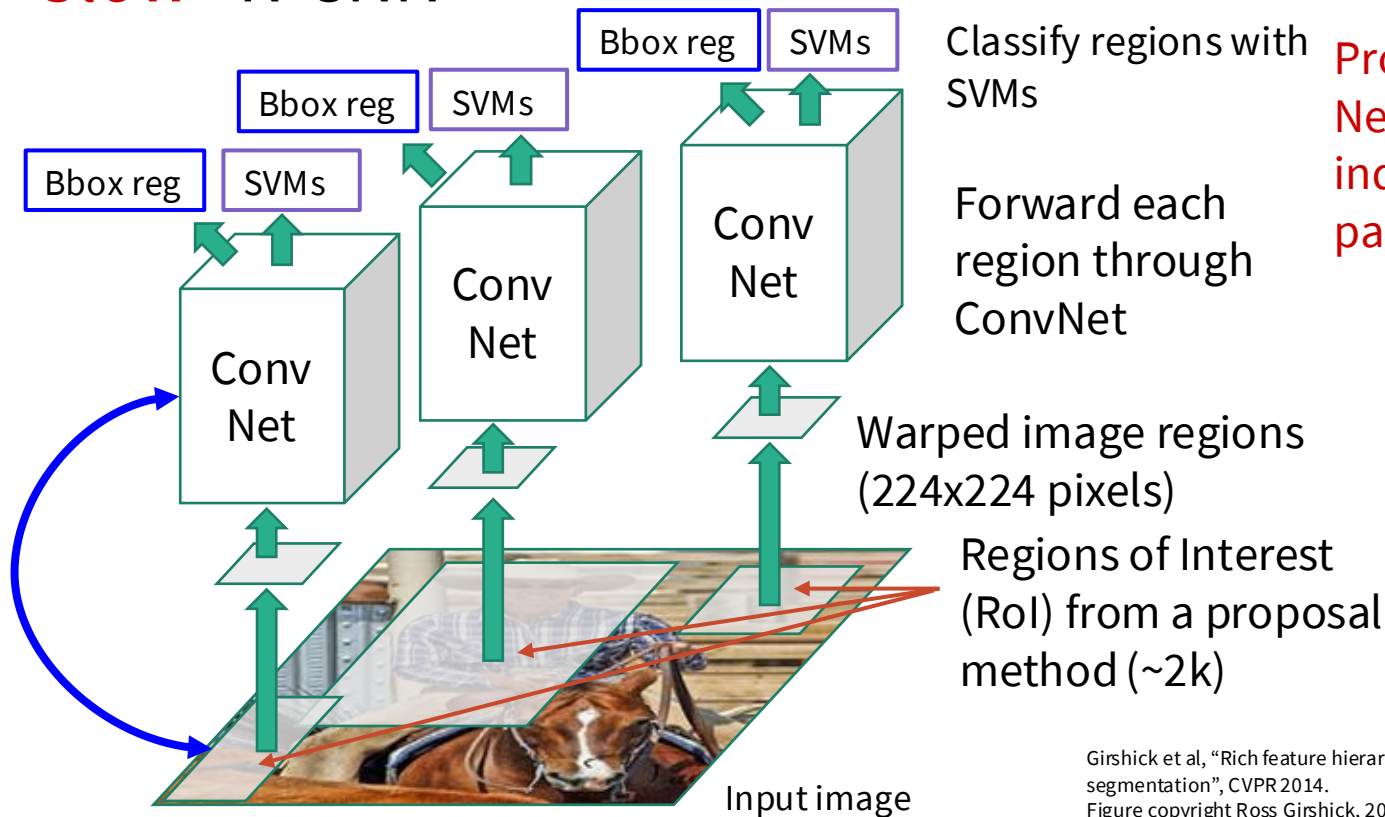
Warped image regions (224x224 pixels)

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.

# “Slow” R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



Classify regions with SVMs

Forward each region through ConvNet

Problem: Very slow!  
Need to do ~2k independent forward passes for each image!

Idea: Pass the image through convnet before cropping! Crop the conv feature instead!

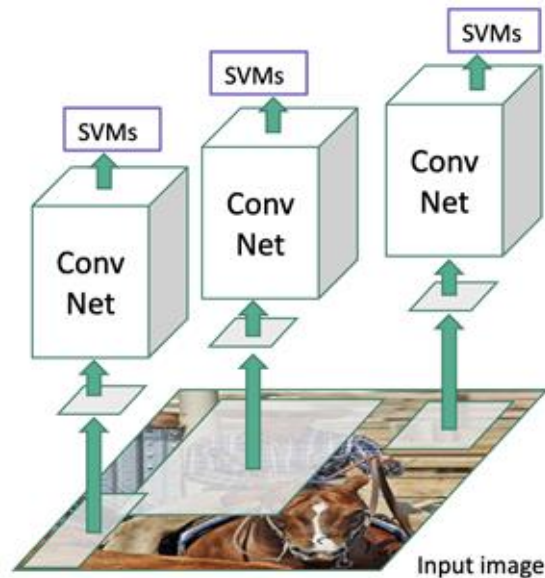
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



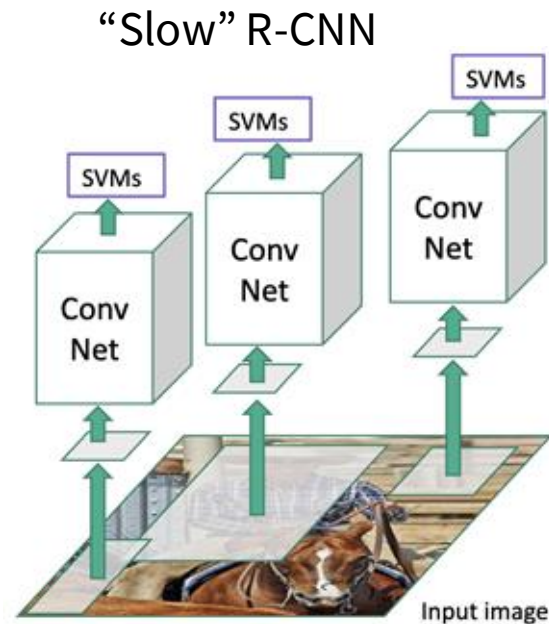
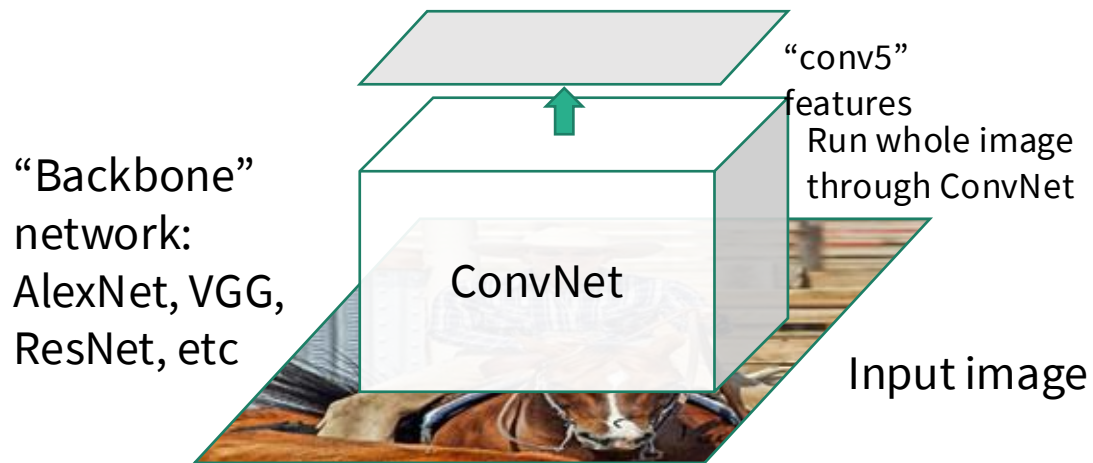
Input image

## “Slow” R-CNN



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

# Fast R-CNN

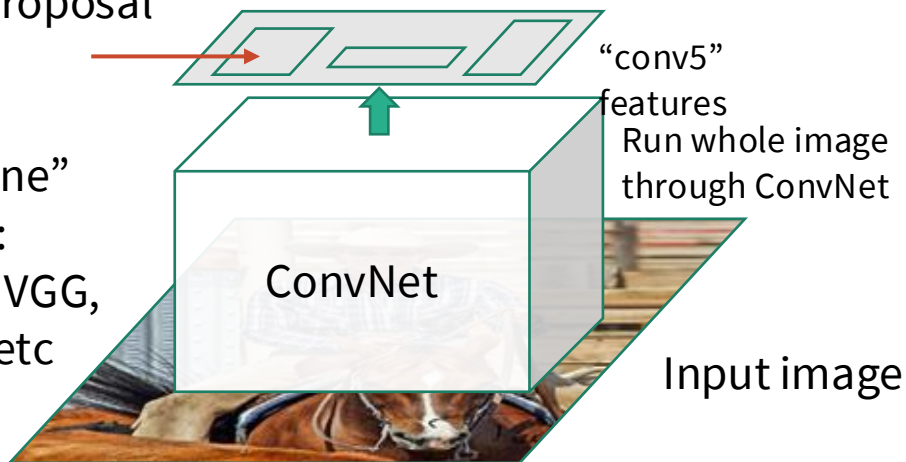


Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

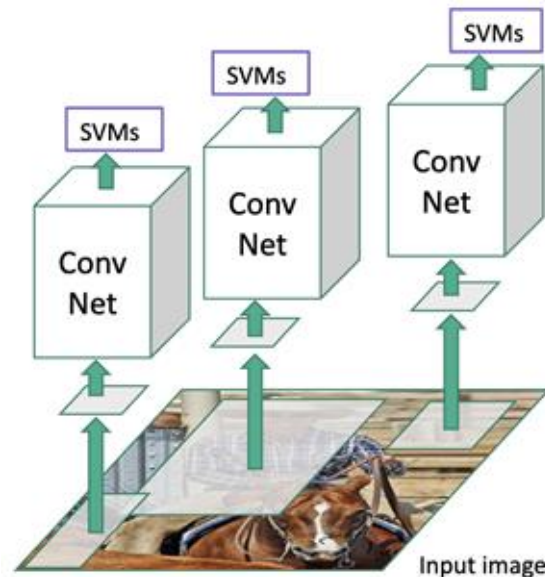
# Fast R-CNN

Regions of Interest (Rois) from a proposal method

“Backbone” network: AlexNet, VGG, ResNet, etc



“Slow” R-CNN

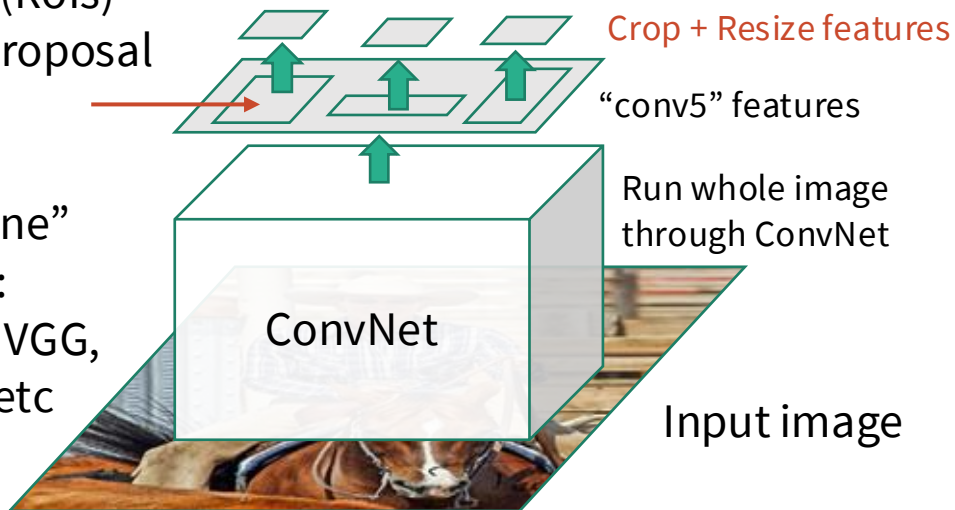


Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

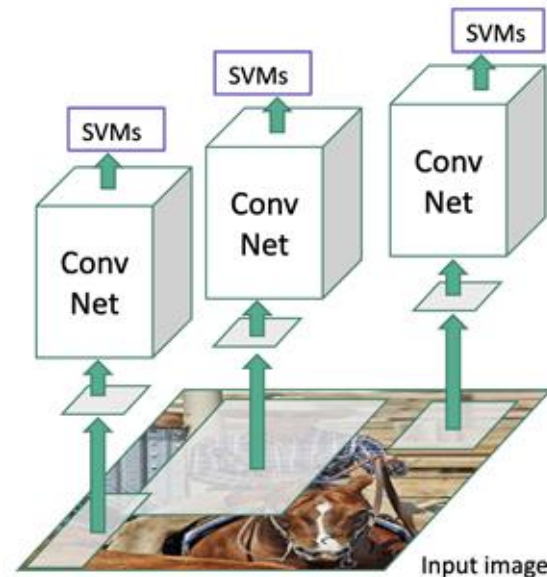
# Fast R-CNN

Regions of Interest (Rois) from a proposal method

“Backbone” network:  
AlexNet, VGG,  
ResNet, etc

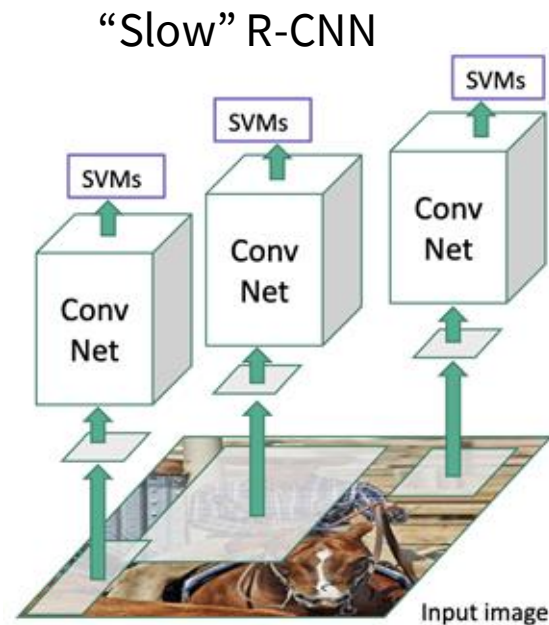
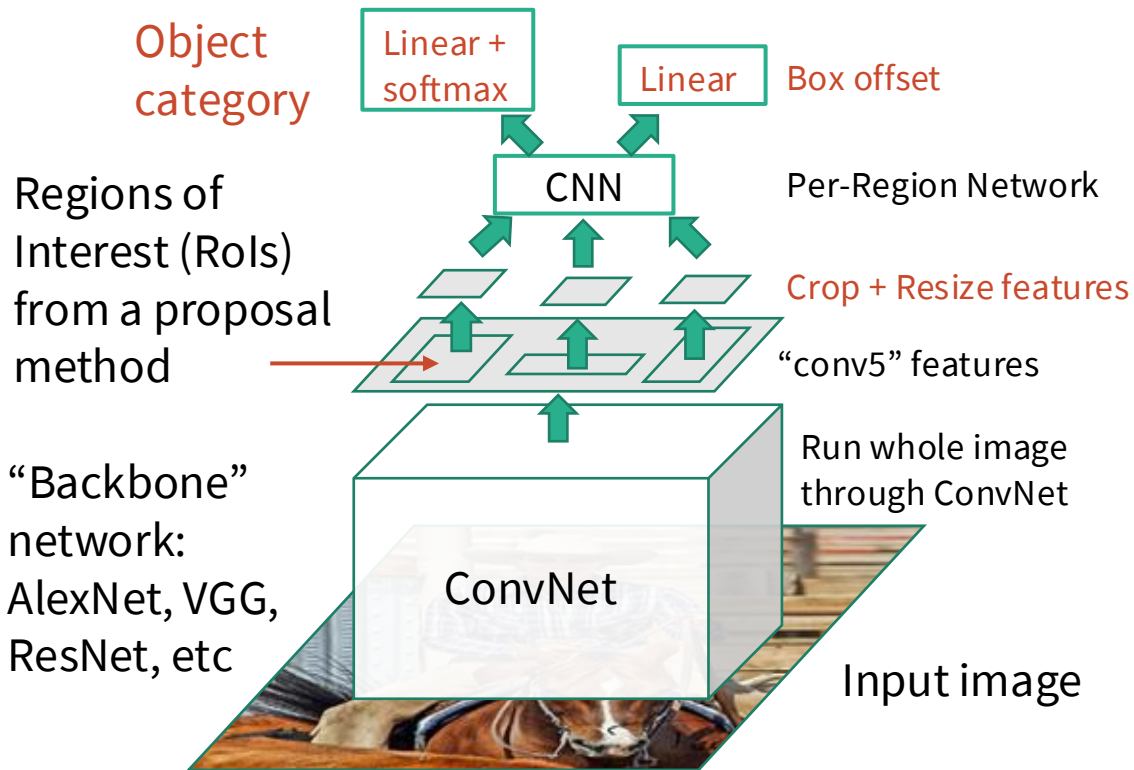


“Slow” R-CNN



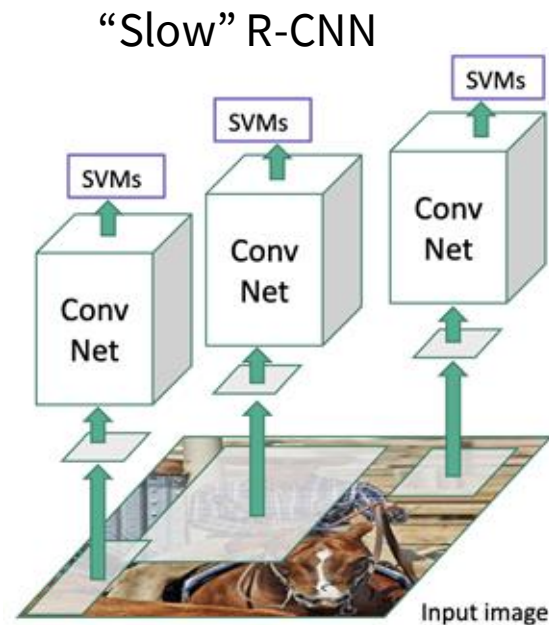
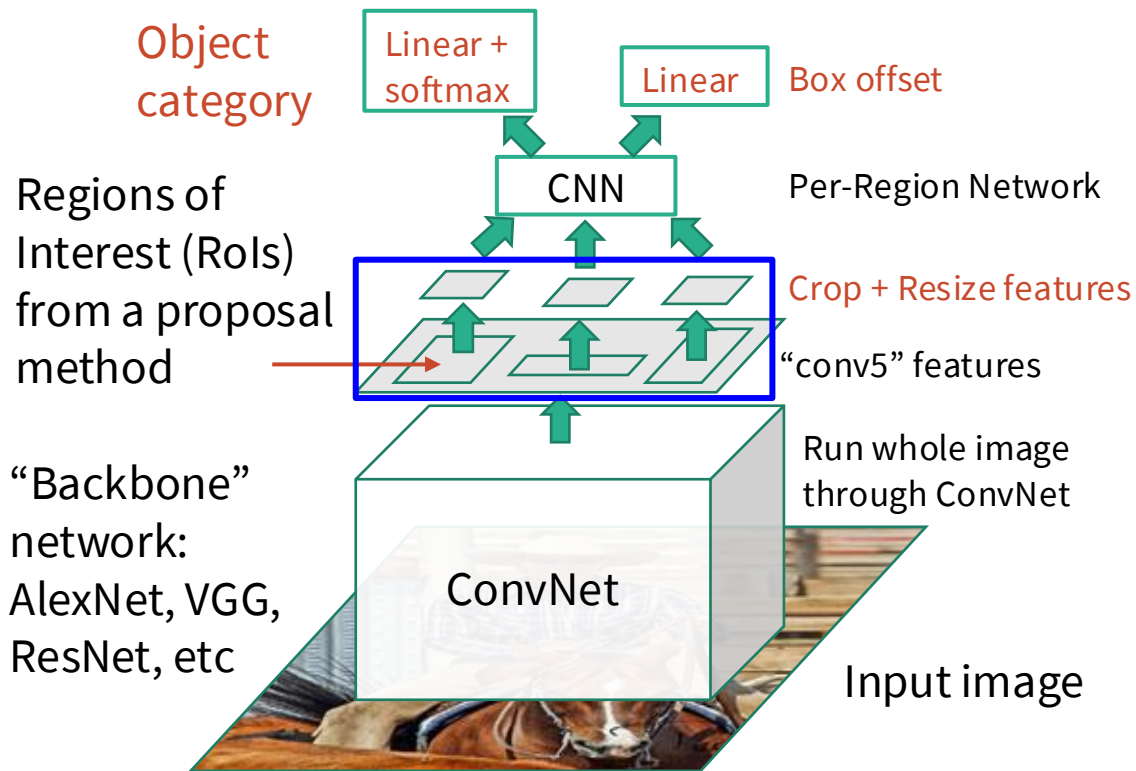
Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



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

# Fast R-CNN



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

# Region Proposal Network



Input Image  
(e.g. 3 x 640 x 480)

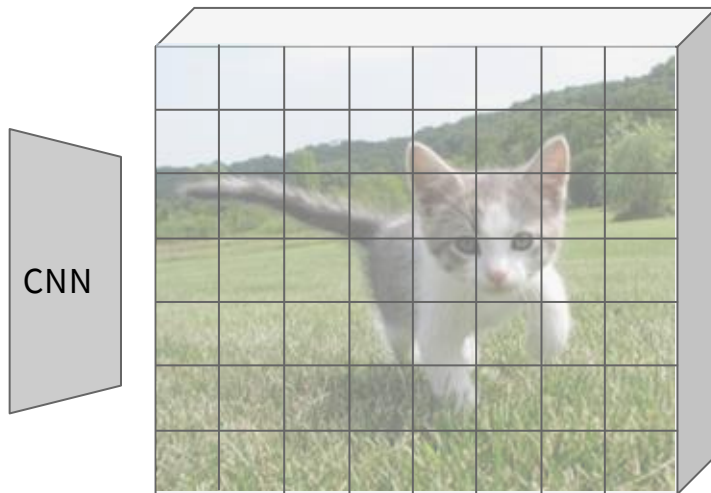


Image features  
(e.g. 512 x 20 x 15)

# Region Proposal Network

Imagine an anchor box of fixed size at each point in the feature map



Input Image  
(e.g. 3 x 640 x 480)

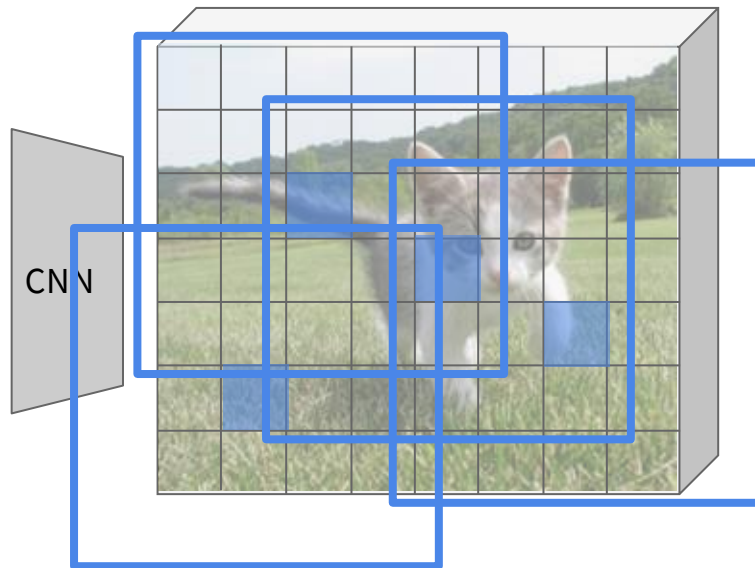


Image features  
(e.g. 512 x 20 x 15)

# Region Proposal Network



Input Image  
(e.g.  $3 \times 640 \times 480$ )

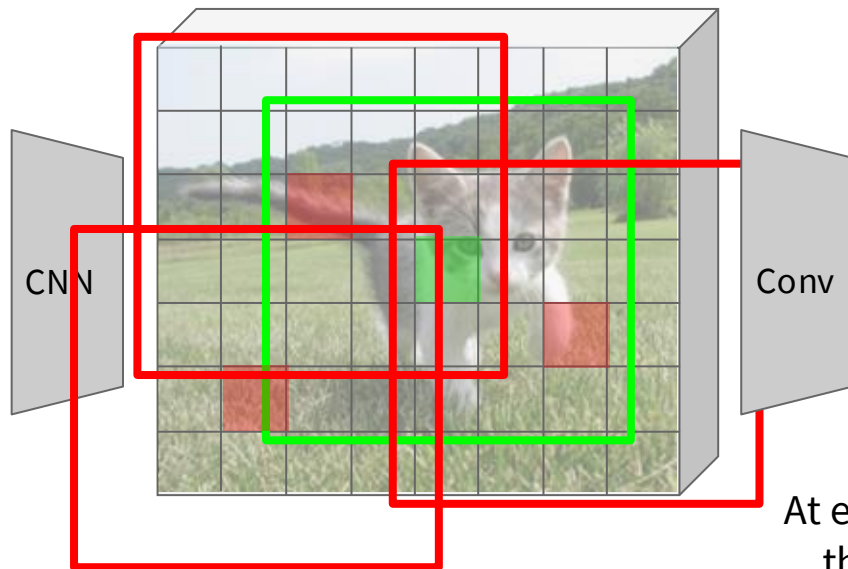


Image features  
(e.g.  $512 \times 20 \times 15$ )

Imagine an anchor box of fixed size at each point in the feature map

Anchor is an object?  
 $1 \times 20 \times 15$

At each point, predict whether the corresponding anchor contains an object (binary classification)

# Region Proposal Network



Input Image  
(e.g.  $3 \times 640 \times 480$ )

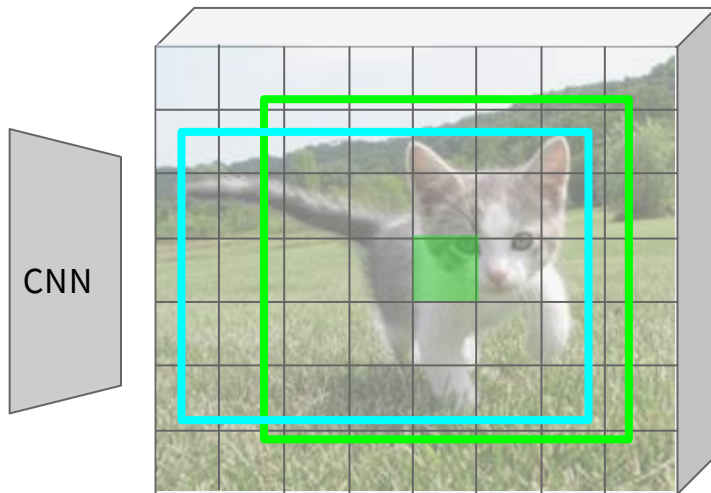
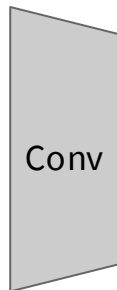


Image features  
(e.g.  $512 \times 20 \times 15$ )

Imagine an anchor box of fixed size at each point in the feature map



Anchor is an object?  
 $1 \times 20 \times 15$

Box corrections  
 $4 \times 20 \times 15$

For positive boxes, also predict a corrections from the anchor to the ground-truth box (regress 4 numbers per pixel)

# Region Proposal Network

In practice use  $K$  different anchor boxes of different size / scale at each point



Input Image  
(e.g.  $3 \times 640 \times 480$ )

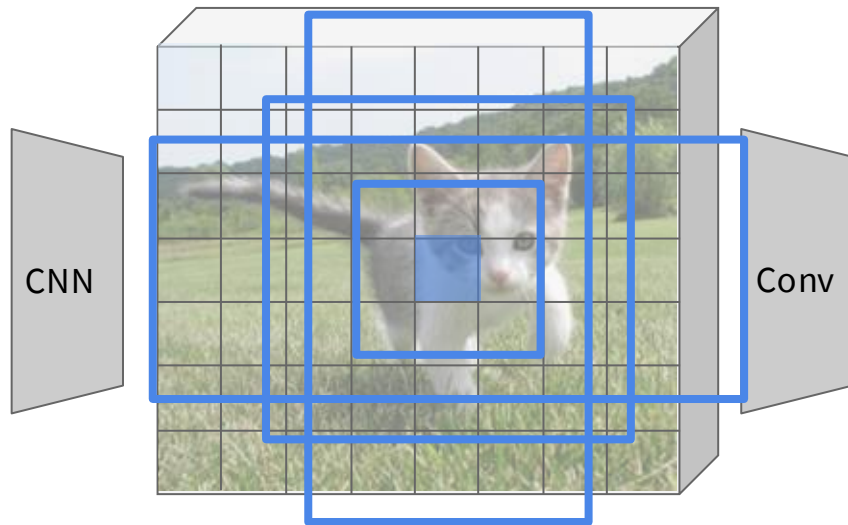


Image features  
(e.g.  $512 \times 20 \times 15$ )

Anchor is an object?  
 $K \times 20 \times 15$

Box transforms  
 $4K \times 20 \times 15$

# Region Proposal Network

In practice use  $K$  different anchor boxes of different size / scale at each point



Input Image  
(e.g.  $3 \times 640 \times 480$ )

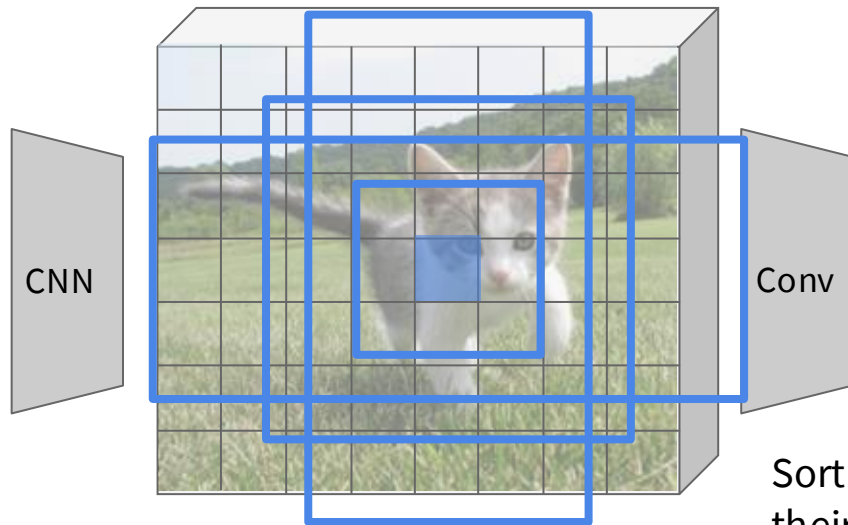


Image features  
(e.g.  $512 \times 20 \times 15$ )

Anchor is an object?  
 $K \times 20 \times 15$

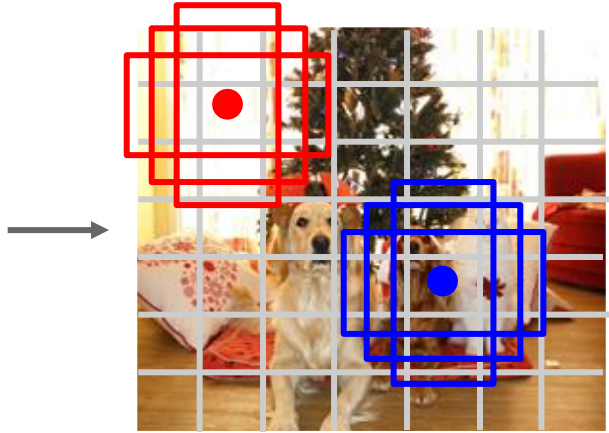
Box transforms  
 $4K \times 20 \times 15$

Sort the  $K \times 20 \times 15$  boxes by their “objectness” score, take top  $\sim 300$  as our proposals

# Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image  
 $3 \times H \times W$



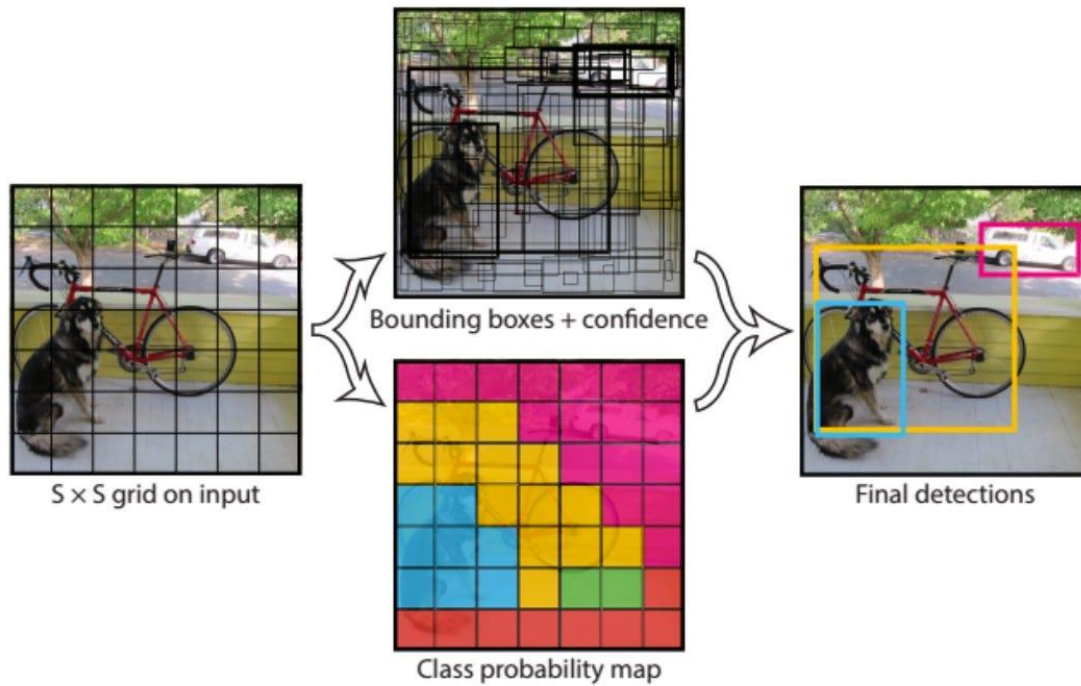
Divide image into grid  
 $7 \times 7$   
Image a set of base boxes  
centered at each grid cell  
Here  $B = 3$

- 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)
  - Looks a lot like RPN, but category-specific!

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  
Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

# YOLO (You Only Look Once) real-time object detection



Redmon et al. "You only look once: unified, real-time object detection (2015)."

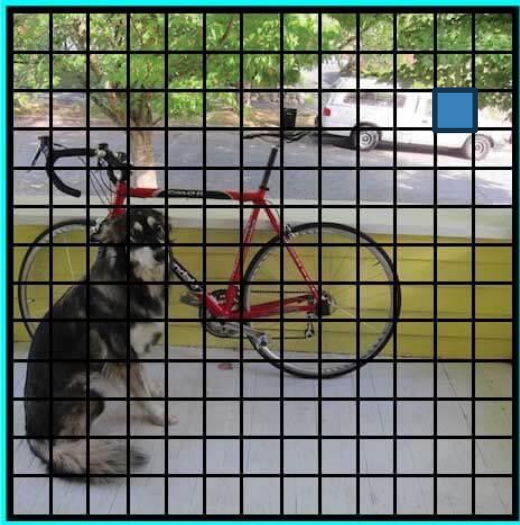
# YOLO



SxS Grid

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



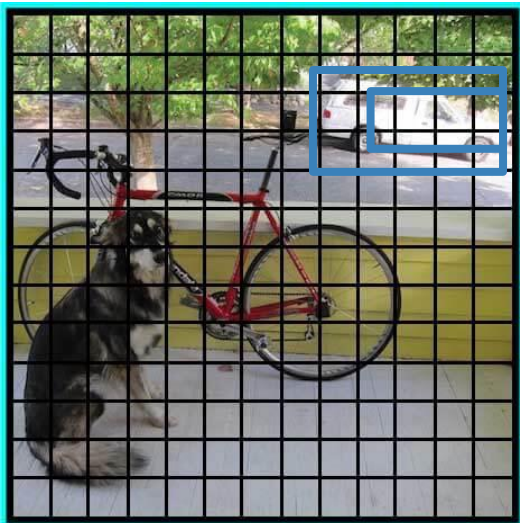
SxS Grid

For each box output:

- $P(\text{object})$ : probability that the box contains an object
- $B$  bounding boxes  $(x, y, h, w)$
- $P(\text{class})$ : probability of belonging to a class

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



SxS Grid

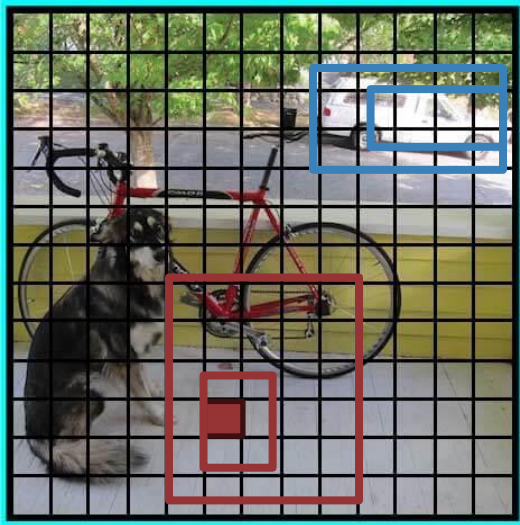
For each box output:

- $P(\text{object})$ : probability that the box contains an object
- $B$  bounding boxes  $(x, y, h, w)$
- $P(\text{class})$ : probability of belonging to a class

$B=2$

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



SxS Grid

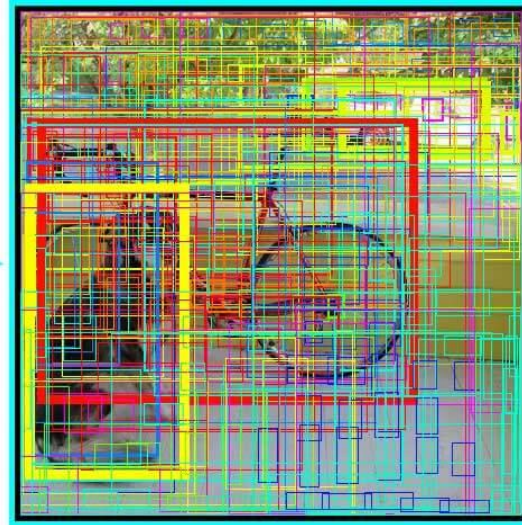
For each box output:

- $P(\text{object})$ : probability that the box contains an object
- $B$  bounding boxes  $(x, y, h, w)$
- $P(\text{class})$ : probability of belonging to a class

$B=2$

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO

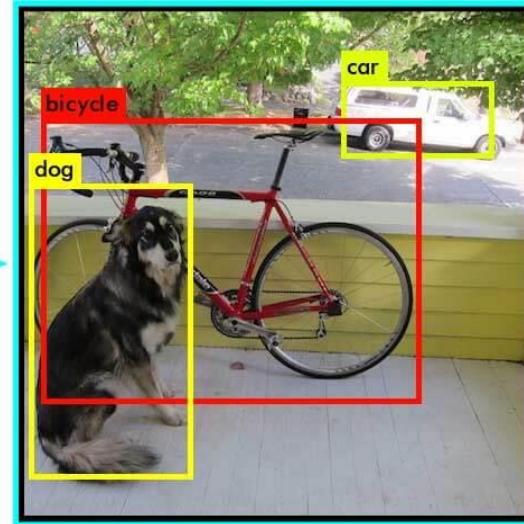
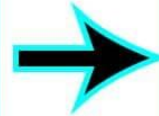
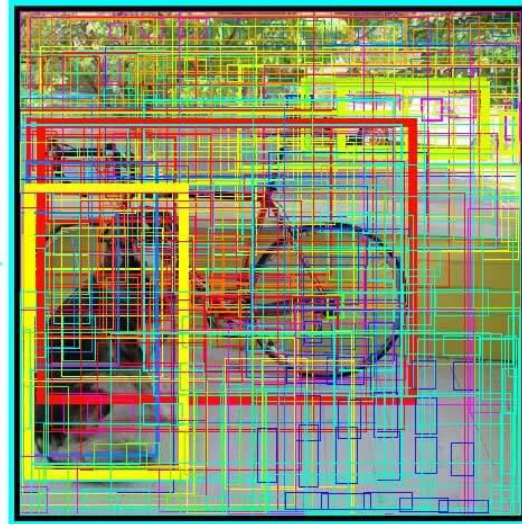
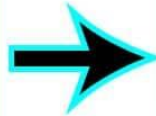
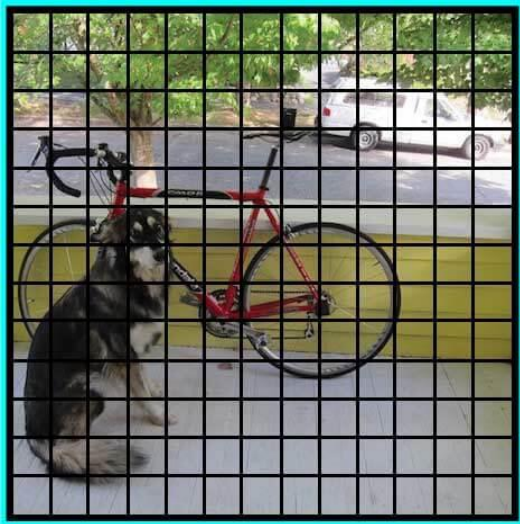


Many bounding boxes with different object probabilities

SxS Grid

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



SxS Grid

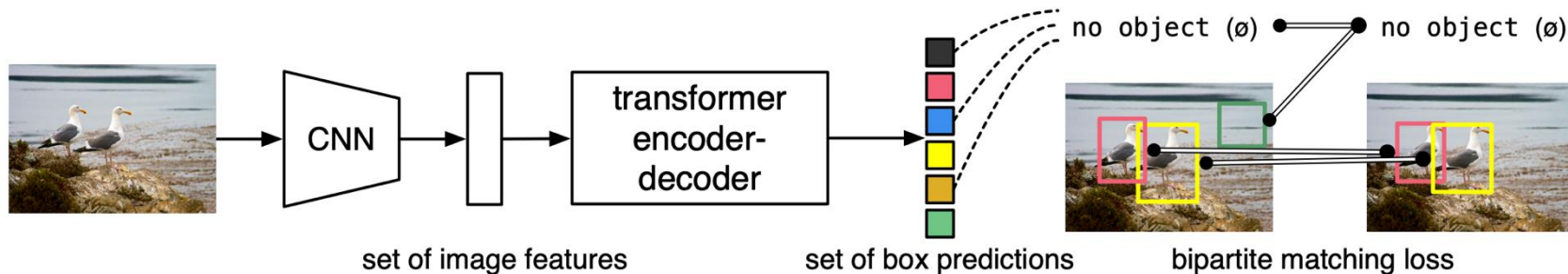
Redmon et al. "You only look once: unified, real-time object detection (2015)."

# Object **D**etection with **T**ransformers: DETR

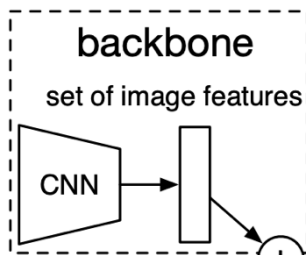
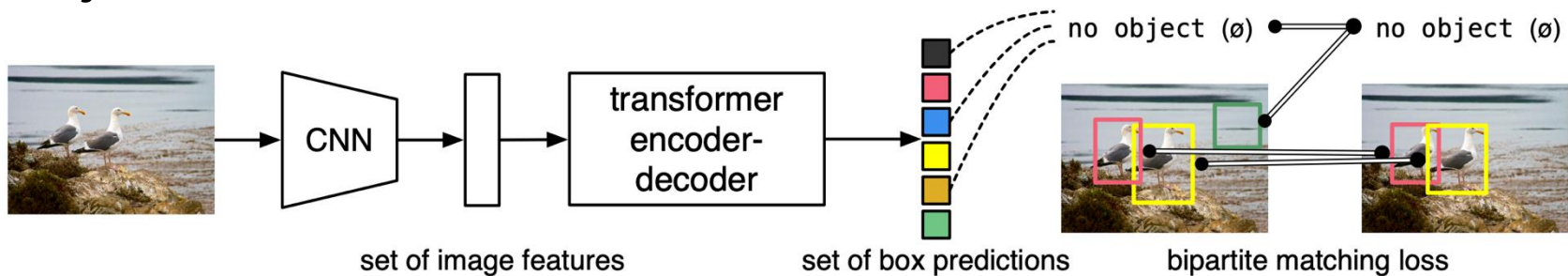
Simple object detection pipeline: directly output a set of boxes from a Transformer

No anchors, no regression of box transforms

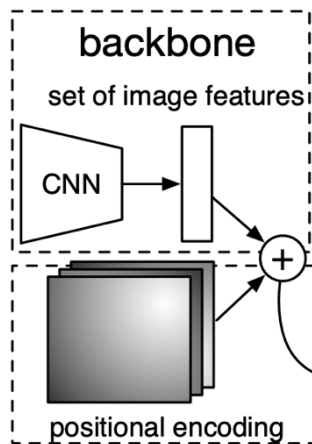
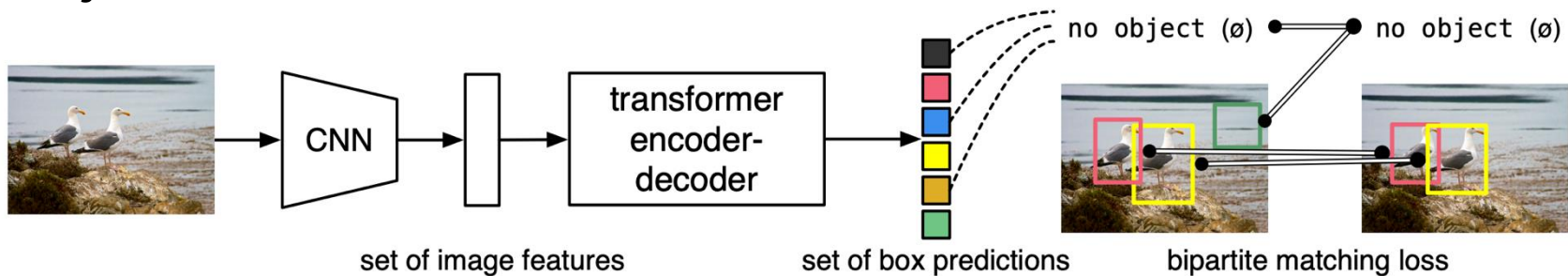
Match predicted boxes to GT boxes with bipartite matching; train to regress box coordinates



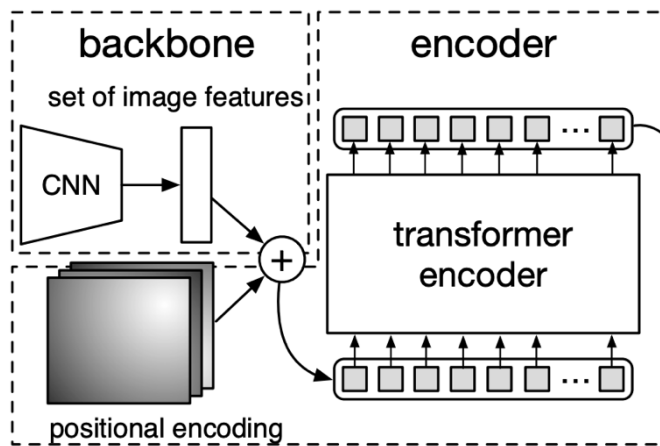
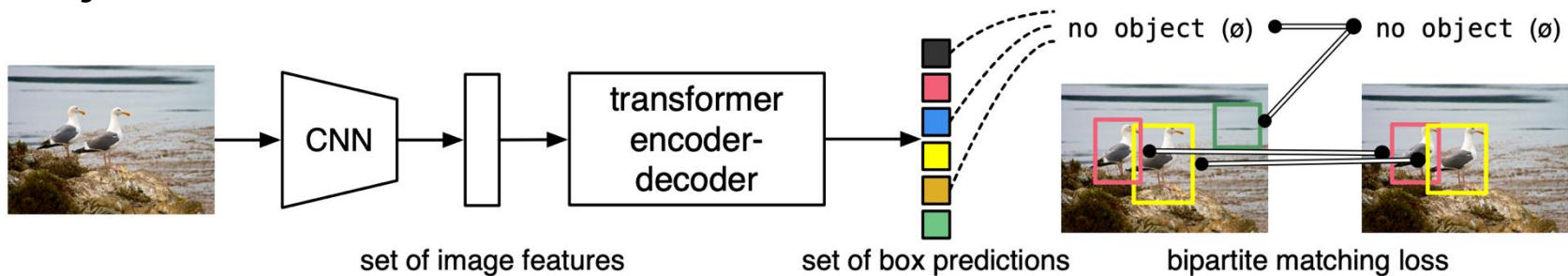
# Object Detection with Transformers: DETR



# Object Detection with Transformers: DETR

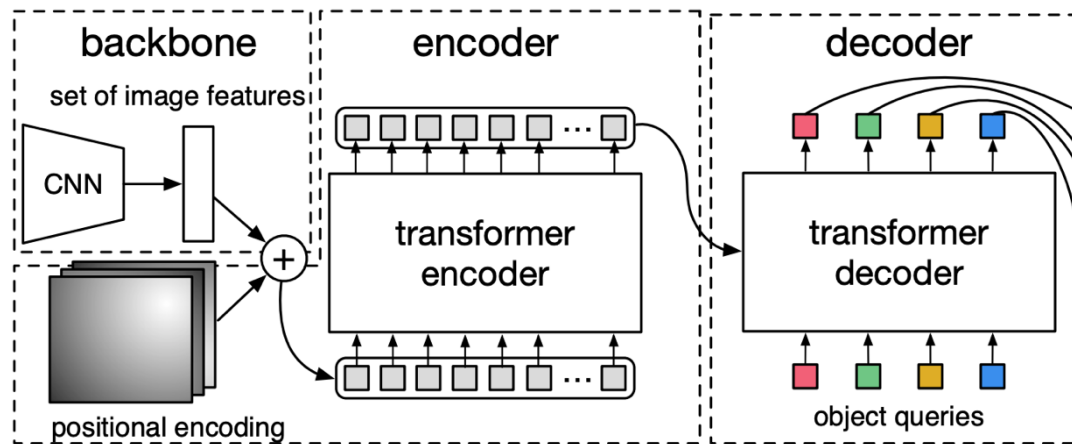
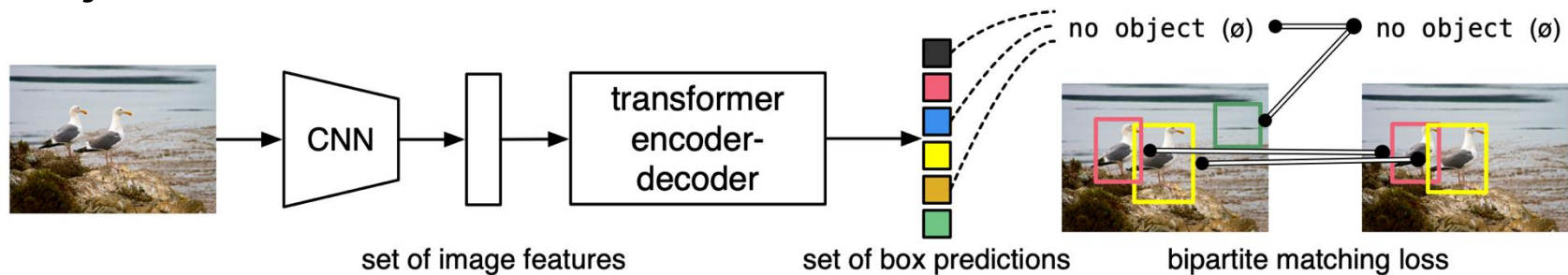


# Object Detection with Transformers: DETR



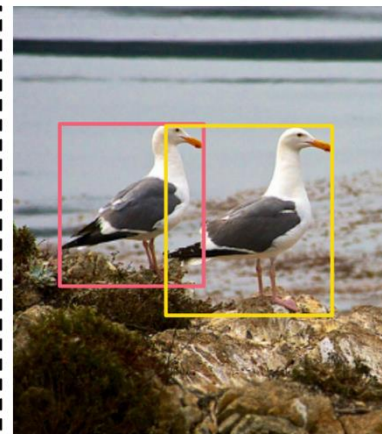
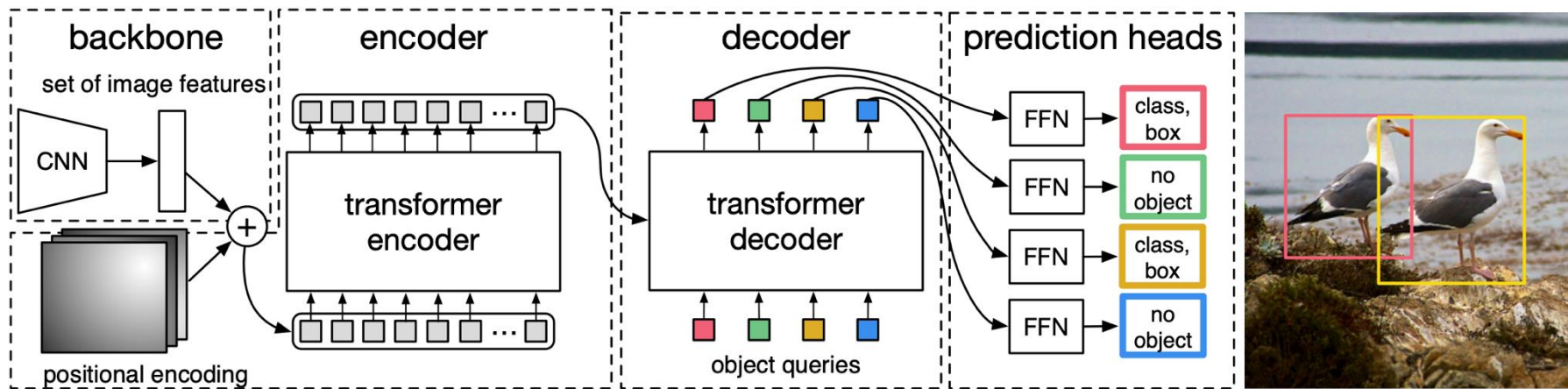
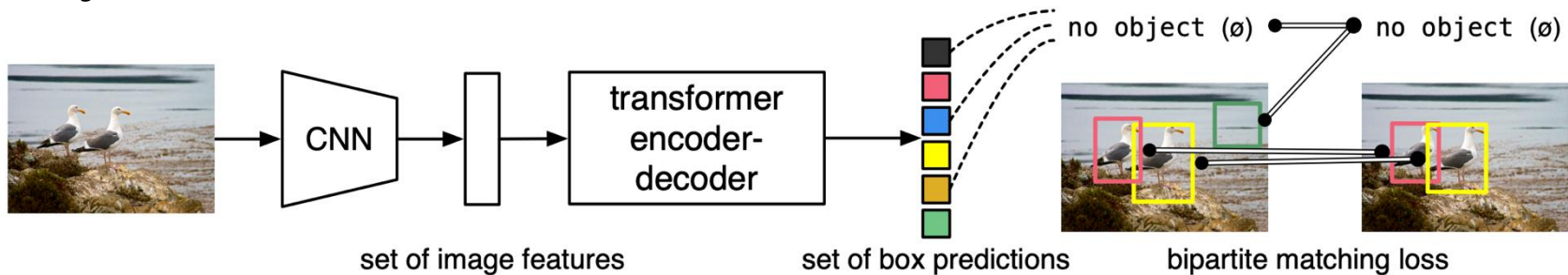
Carion et al, "End-to-End Object Detection with Transformers", ECCV2020

# Object Detection with Transformers: DETR



Carion et al, "End-to-End Object Detection with Transformers", ECCV2020

# Object Detection with Transformers: DETR



# Instance Segmentation

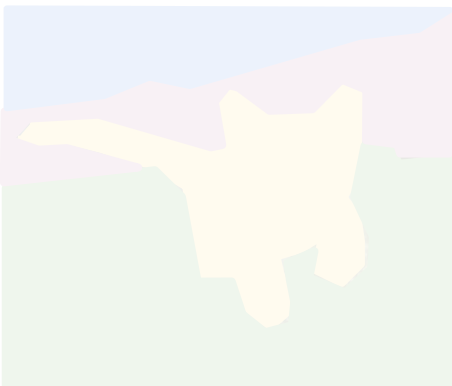
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

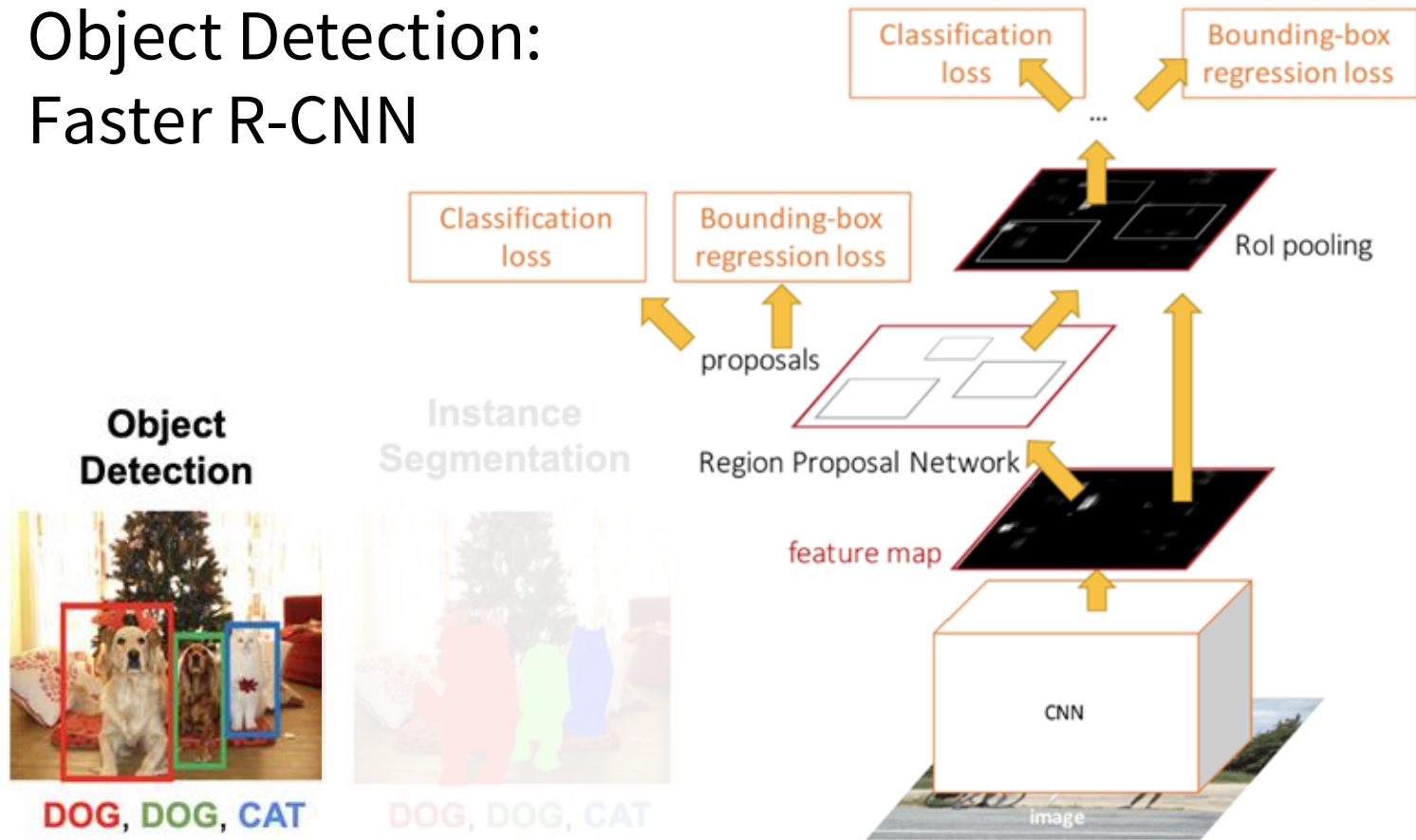
Multiple Object

Instance Segmentation

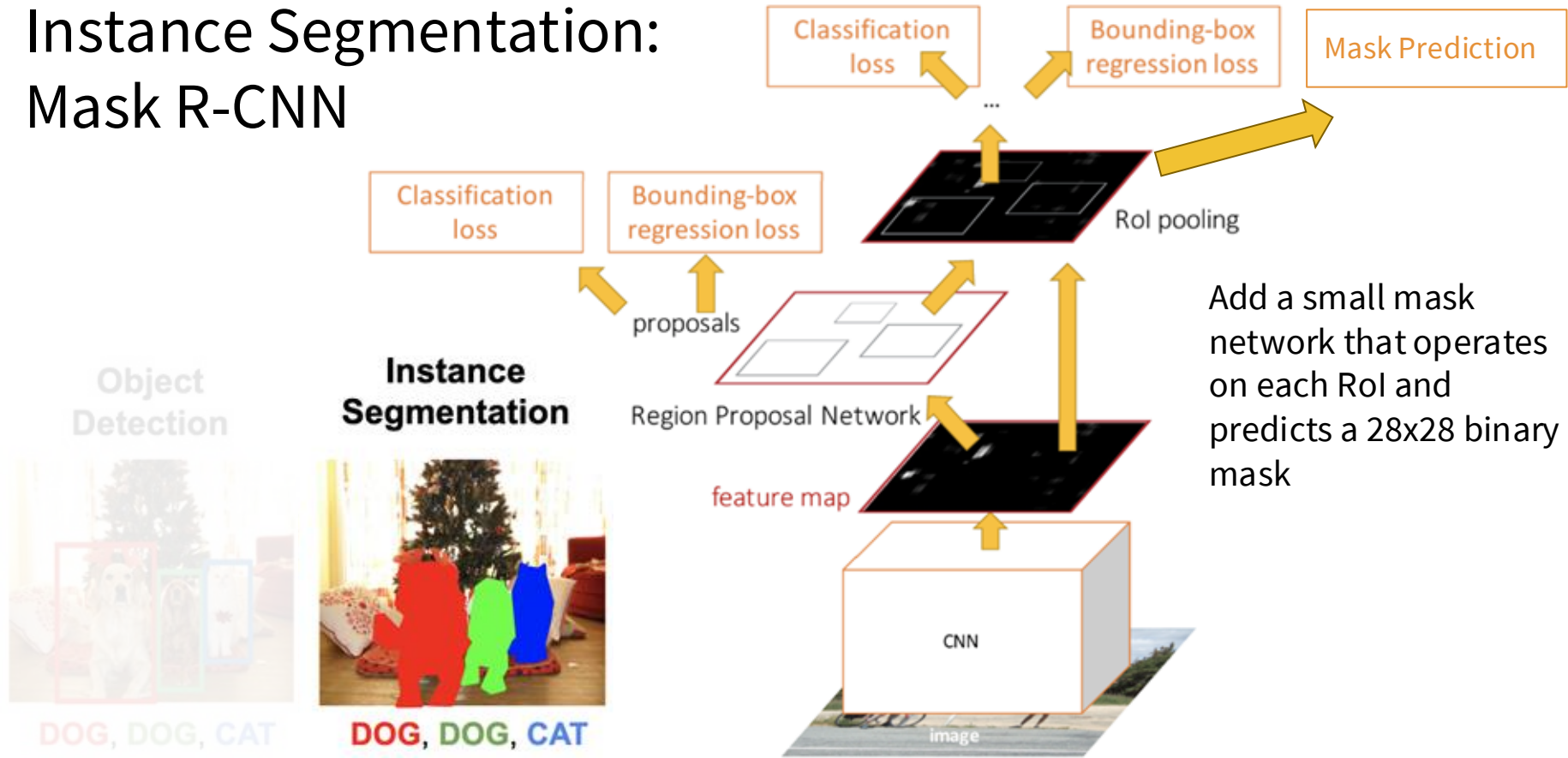


DOG, DOG, CAT

# Object Detection: Faster R-CNN

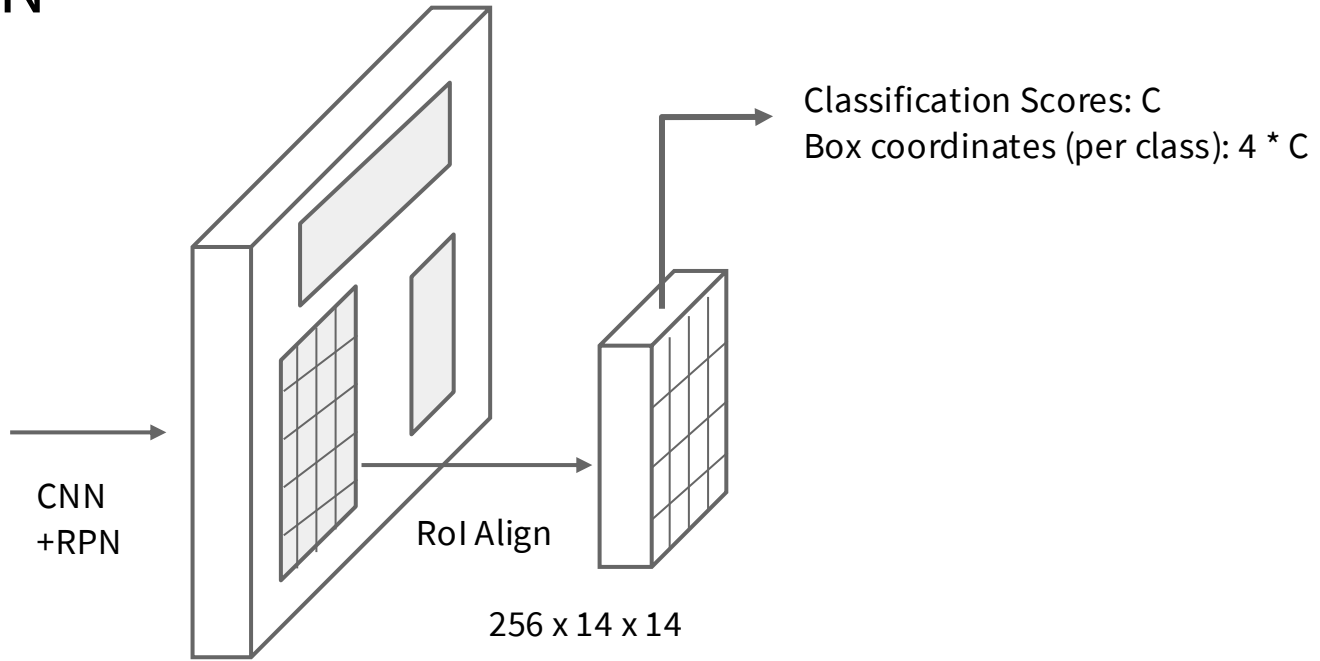
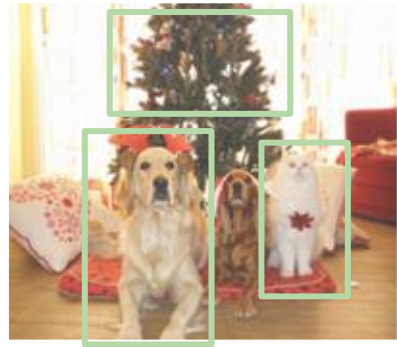


# Instance Segmentation: Mask R-CNN



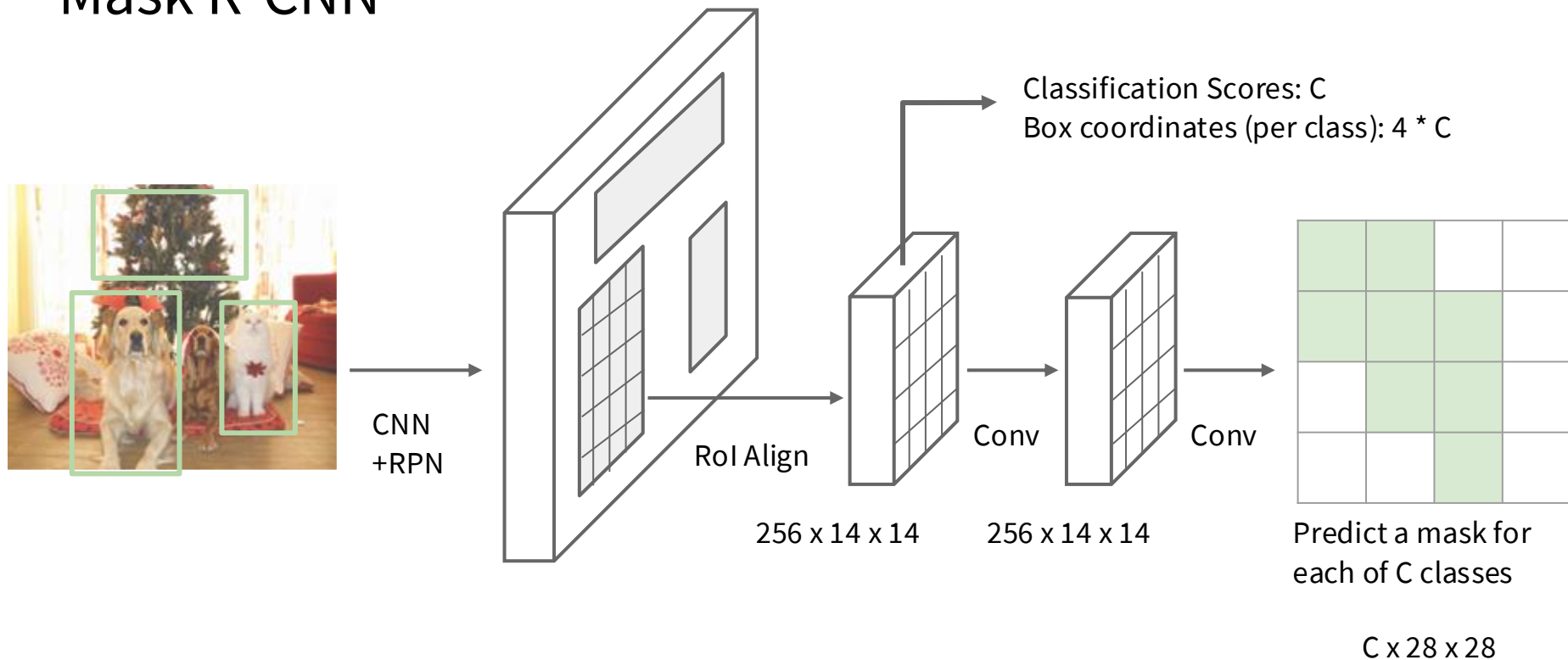
He et al, "Mask R-CNN", ICCV 2017

# Mask R-CNN



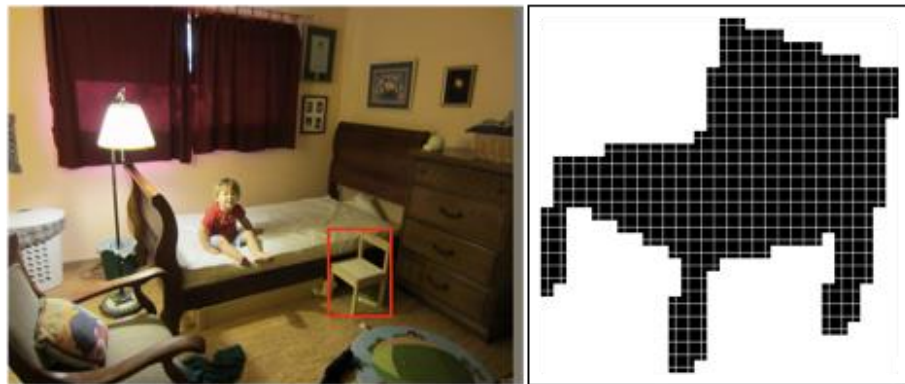
He et al, "Mask R-CNN", arXiv 2017

# Mask R-CNN

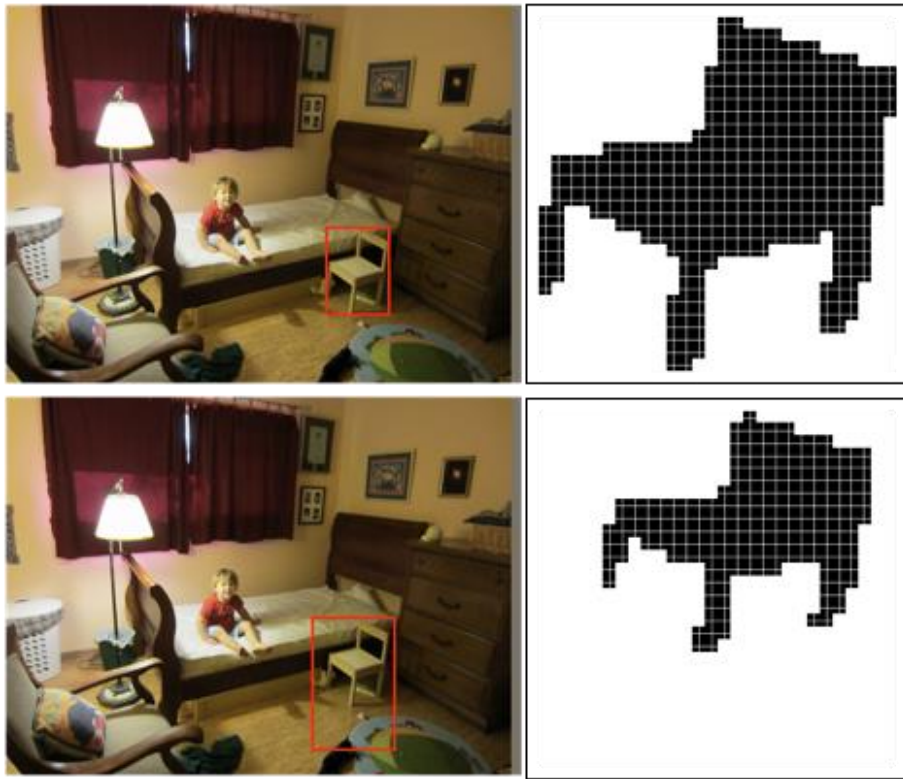


He et al, "Mask R-CNN", arXiv 2017

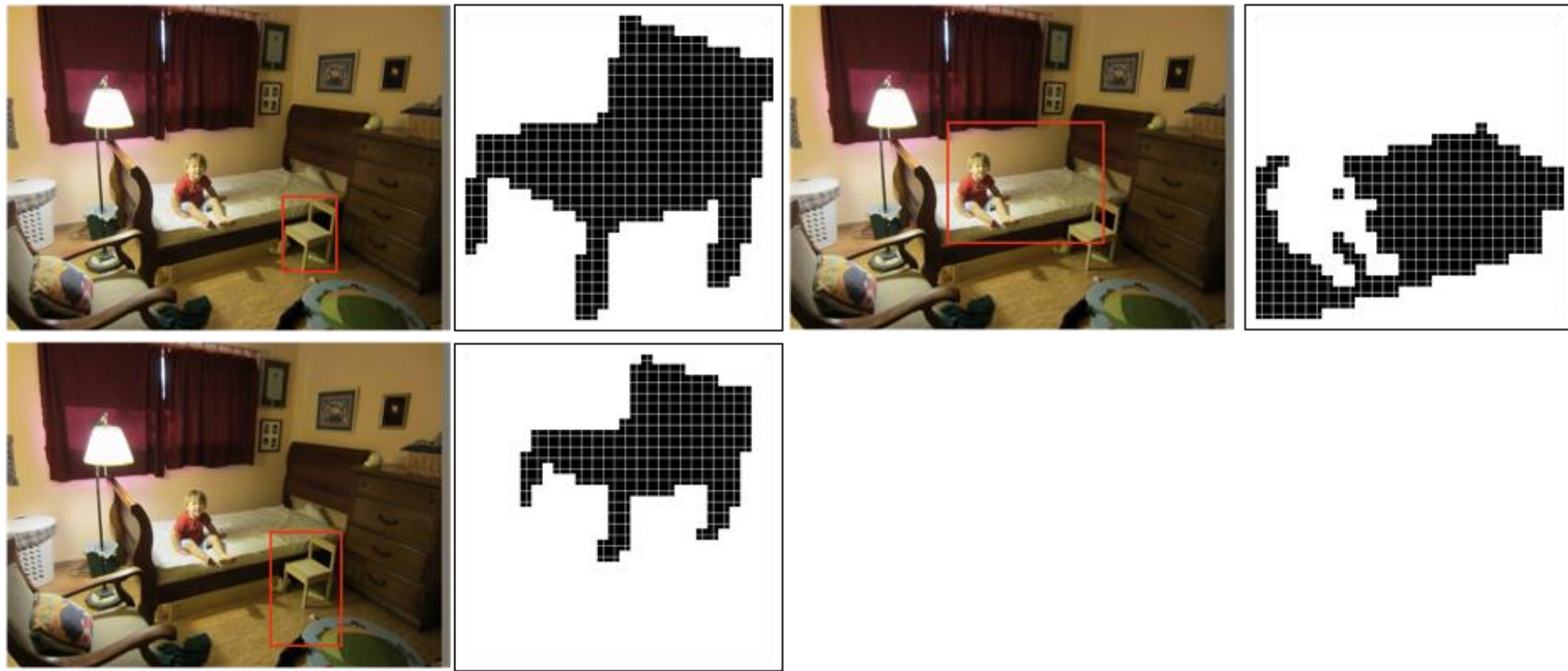
# Mask R-CNN: Example Mask Training Targets



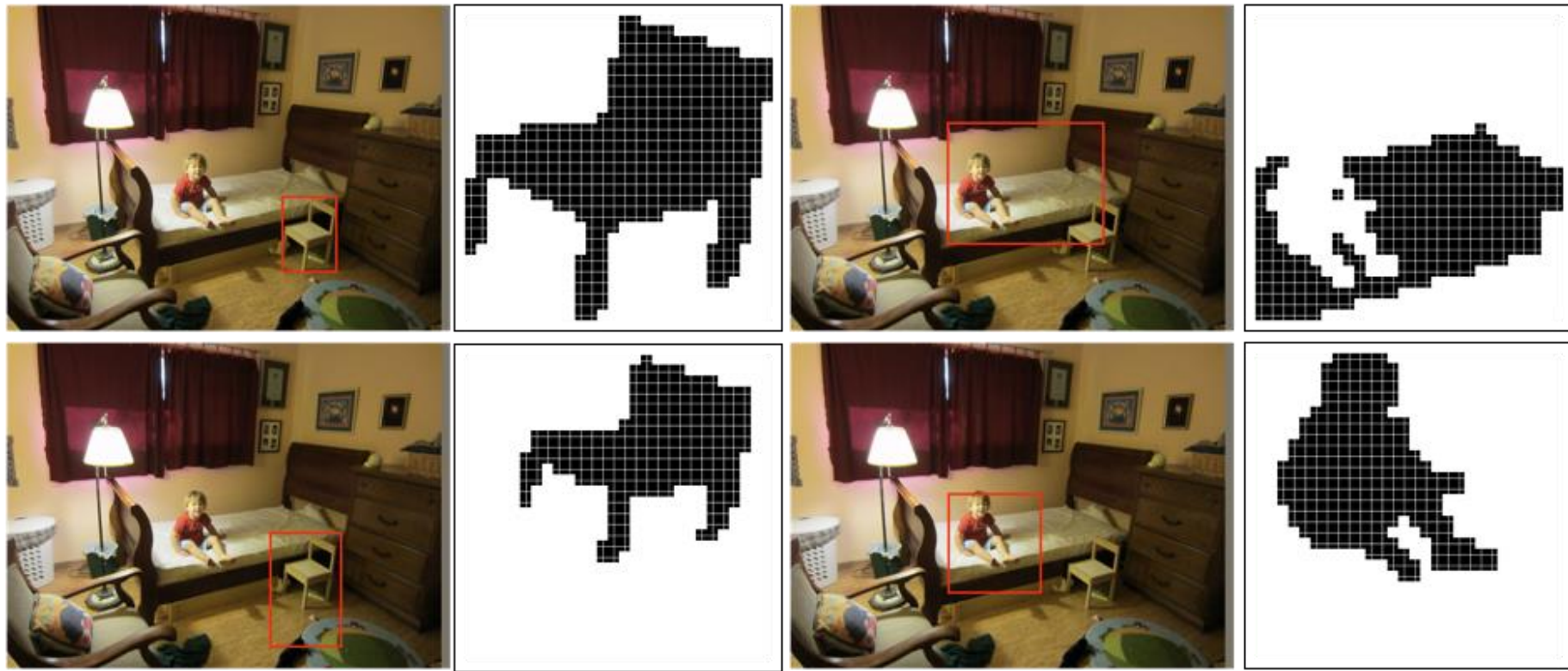
# Mask R-CNN: Example Mask Training Targets



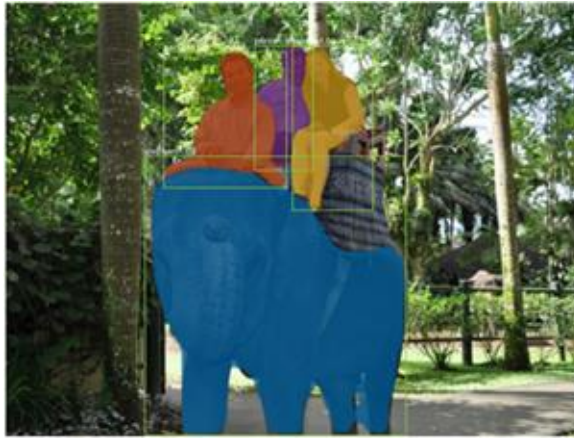
# Mask R-CNN: Example Mask Training Targets



# Mask R-CNN: Example Mask Training Targets



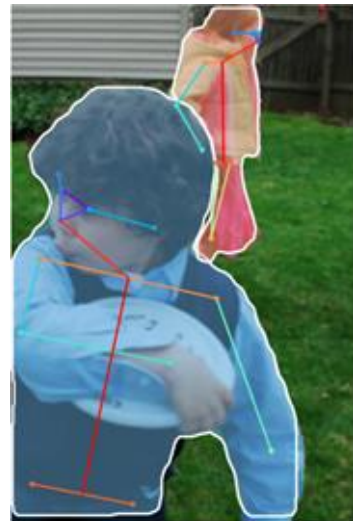
# Mask R-CNN: Very Good Results!



He et al, "Mask R-CNN", ICCV 2017

# Mask R-CNN

## Also does pose



He et al, "Mask R-CNN", ICCV 2017

# Recap: Lots of computer vision tasks!

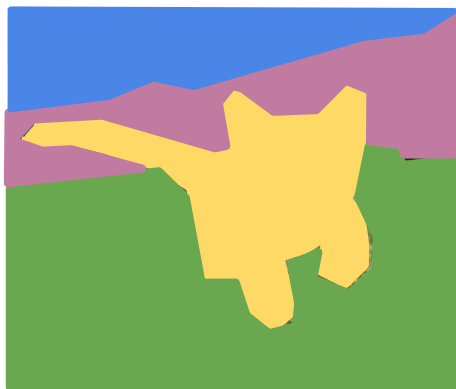
Classification



CAT

No spatial extent

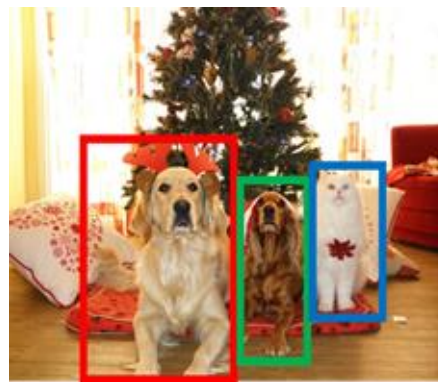
Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



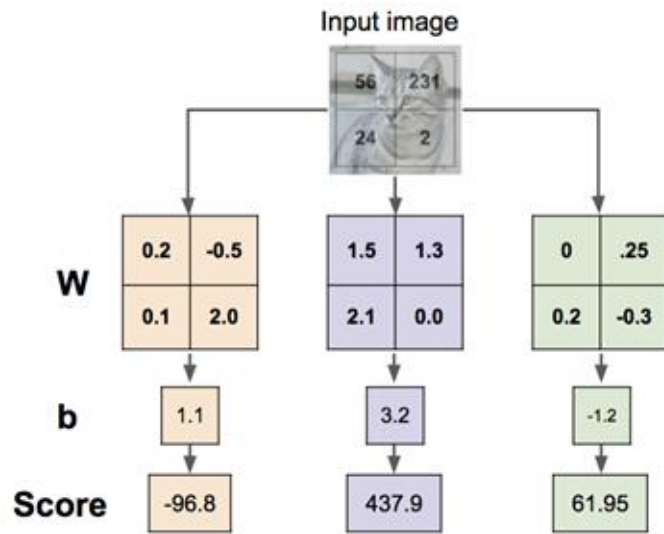
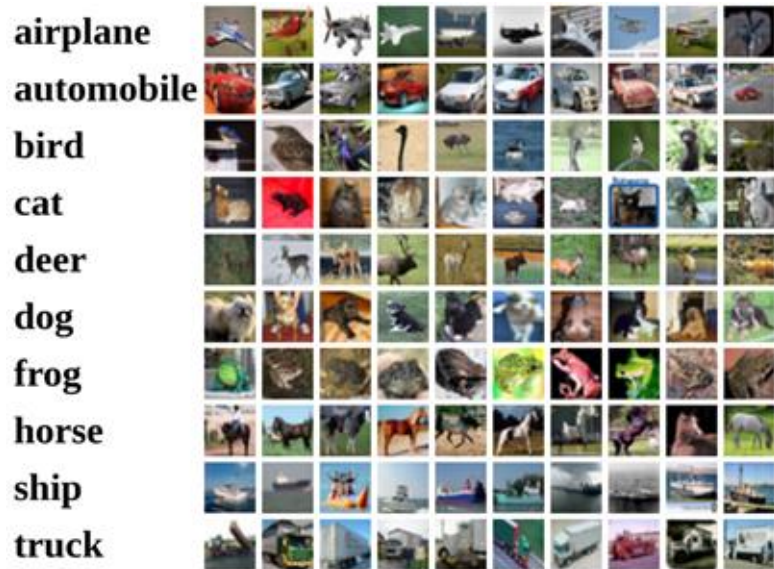
DOG, DOG, CAT

[This image is CC0 public domain](#)

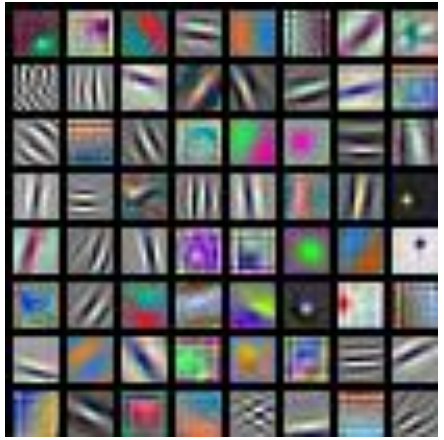
# Today

- Transformers Recap
- **Computer Vision Tasks**
  - Semantic Segmentation
  - Object Detection
  - Instance Segmentation
- Visualization & Understanding
  - Model Layers Visualization
  - Saliency Maps
  - CAM & Grad-CAM

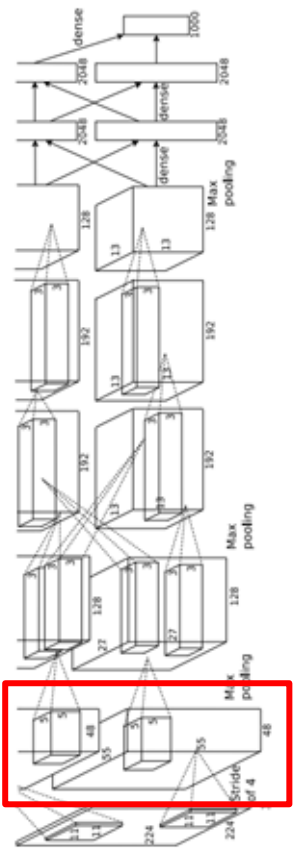
# Interpreting a Linear Classifier: Visual Viewpoint



# First Layer: Visualize Filters

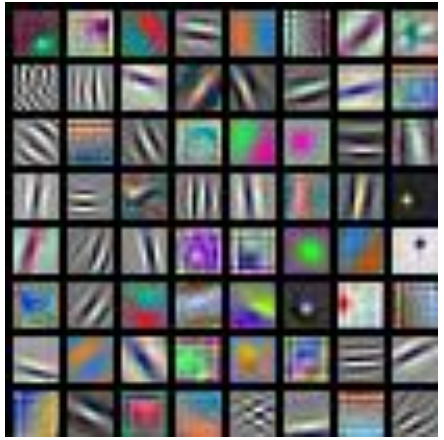


AlexNet:  
64 x 3 x 11 x 11



Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014  
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016  
Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

# First Layer: Visualize Filters



AlexNet:  
64 x 3 x 11 x 11



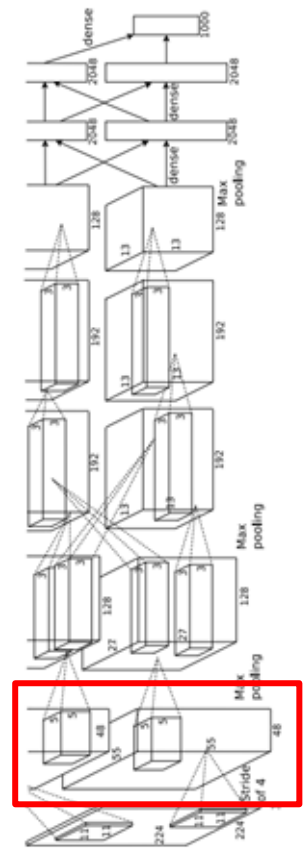
ResNet-18:  
64 x 3 x 7 x 7



ResNet-101:  
64 x 3 x 7 x 7



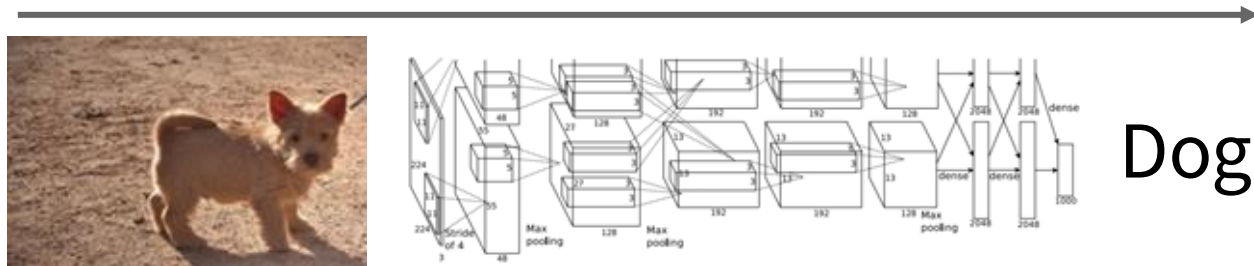
DenseNet-121:  
64 x 3 x 7 x 7



Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014  
 He et al, "Deep Residual Learning for Image Recognition", CVPR 2016  
 Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

# Which pixels matter: Saliency via Backprop

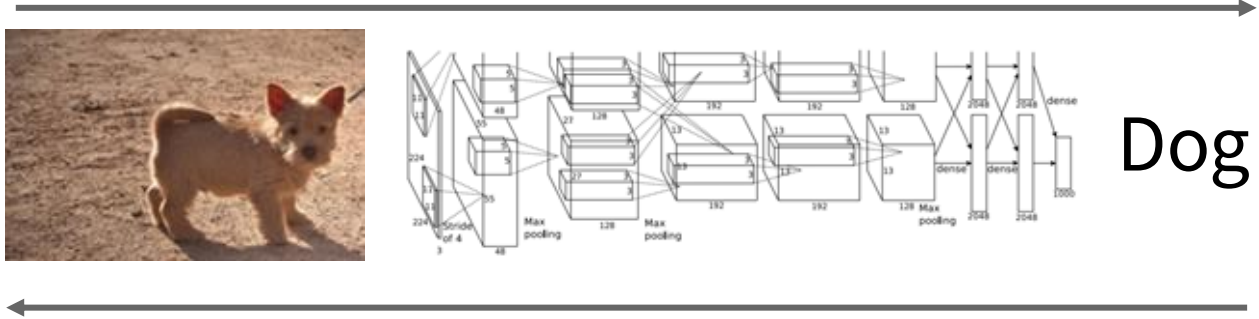
Forward pass: Compute probabilities



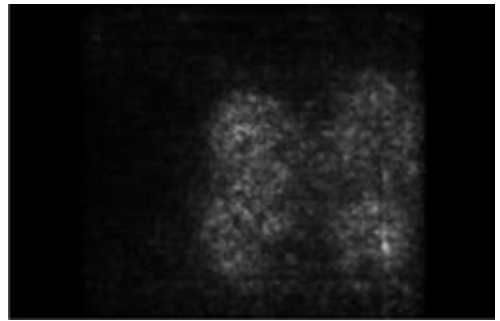
Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.  
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

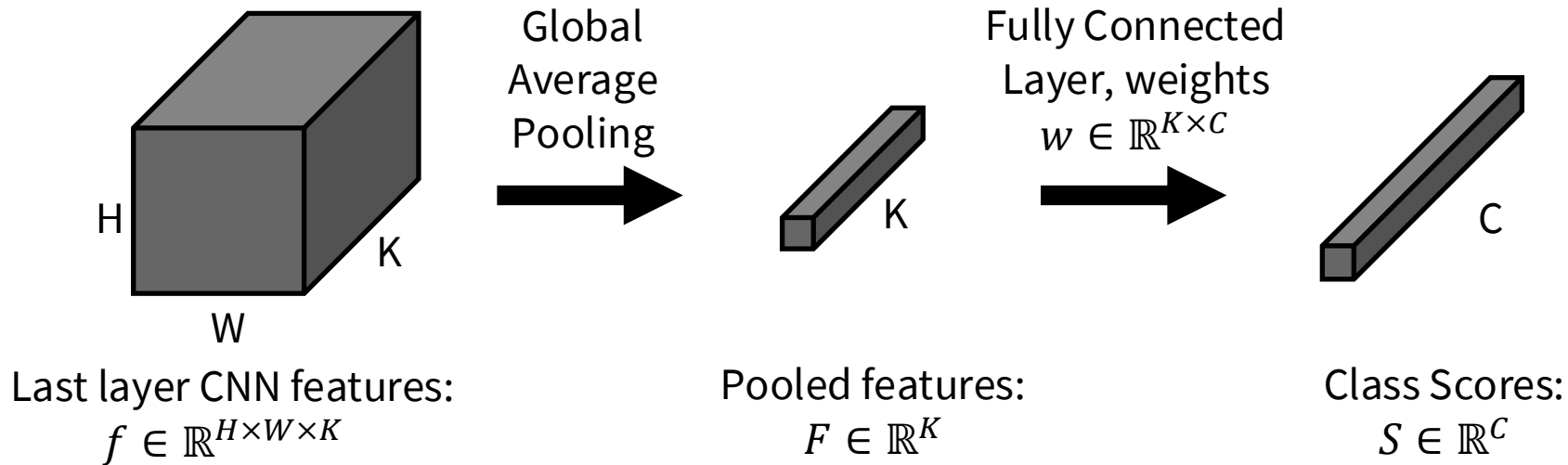
# Saliency Maps



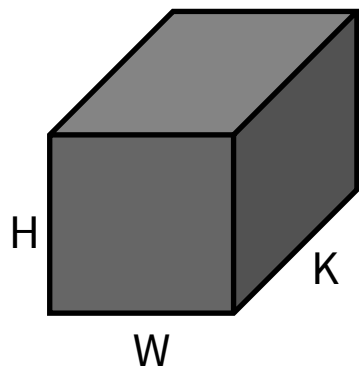
Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Class Activation Mapping (CAM)



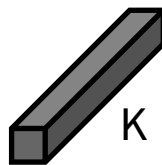
# Class Activation Mapping (CAM)



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Global  
Average  
Pooling

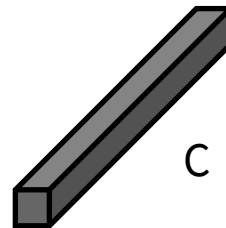


Pooled features:

$$F \in \mathbb{R}^K$$

Fully Connected  
Layer, weights

$$w \in \mathbb{R}^{K \times C}$$

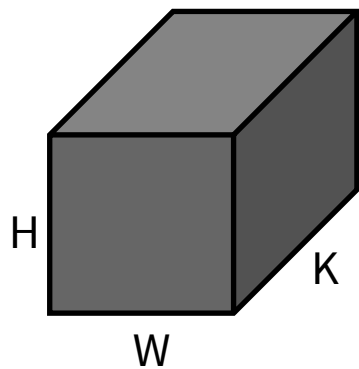


Class Scores:

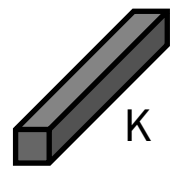
$$S \in \mathbb{R}^C$$

$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k}$$

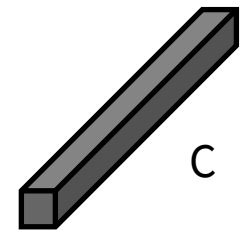
# Class Activation Mapping (CAM)



Global Average Pooling  
→



Fully Connected Layer, weights  
 $w \in \mathbb{R}^{K \times C}$   
→



Last layer CNN features:  
 $f \in \mathbb{R}^{H \times W \times K}$

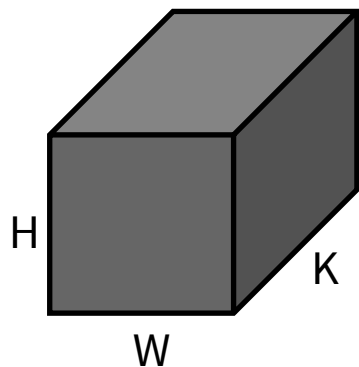
Pooled features:  
 $F \in \mathbb{R}^K$

Class Scores:  
 $S \in \mathbb{R}^C$

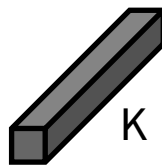
$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k$$

Zhou et al, "Learning Deep Features for Discriminative Localization", CVPR 2016

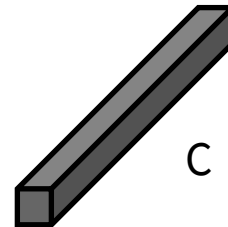
# Class Activation Mapping (CAM)



Global  
Average  
Pooling



Fully Connected  
Layer, weights  
 $w \in \mathbb{R}^{K \times C}$



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Pooled features:

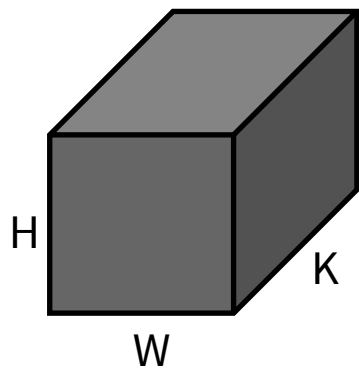
$$F \in \mathbb{R}^K$$

Class Scores:

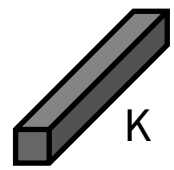
$$S \in \mathbb{R}^C$$

$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k}$$

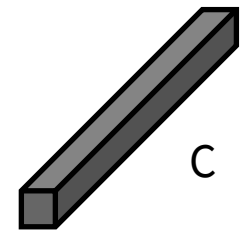
# Class Activation Mapping (CAM)



Global  
Average  
Pooling  
→



Fully Connected  
Layer, weights  
 $w \in \mathbb{R}^{K \times C}$   
→



Last layer CNN features:  
 $f \in \mathbb{R}^{H \times W \times K}$

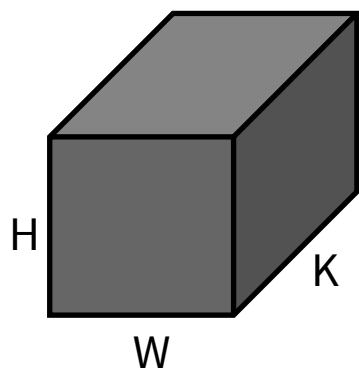
Pooled features:  
 $F \in \mathbb{R}^K$

Class Scores:  
 $S \in \mathbb{R}^C$

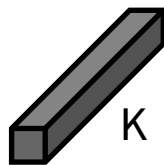
$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k} = \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k}$$

Zhou et al, "Learning Deep Features for Discriminative Localization", CVPR 2016

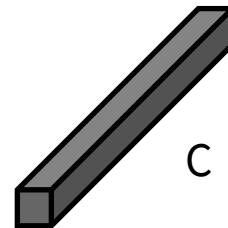
# Class Activation Mapping (CAM)



Global  
Average  
Pooling



Fully Connected  
Layer, weights  
 $w \in \mathbb{R}^{K \times C}$



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Pooled features:

$$F \in \mathbb{R}^K$$

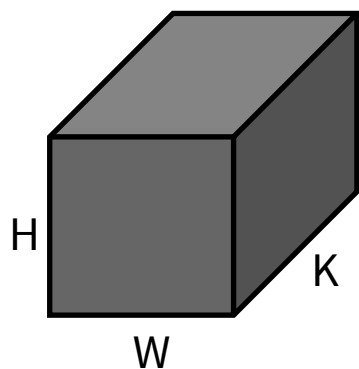
Class Scores:

$$S \in \mathbb{R}^C$$

$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k} \\ = \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k}$$

Zhou et al, "Learning Deep Features for Discriminative Localization", CVPR 2016

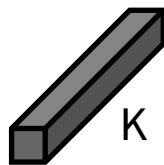
# Class Activation Mapping (CAM)



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Global  
Average  
Pooling

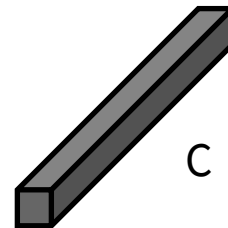


Pooled features:

$$F \in \mathbb{R}^K$$

Fully Connected  
Layer, weights

$$w \in \mathbb{R}^{K \times C}$$



Class Scores:

$$S \in \mathbb{R}^C$$

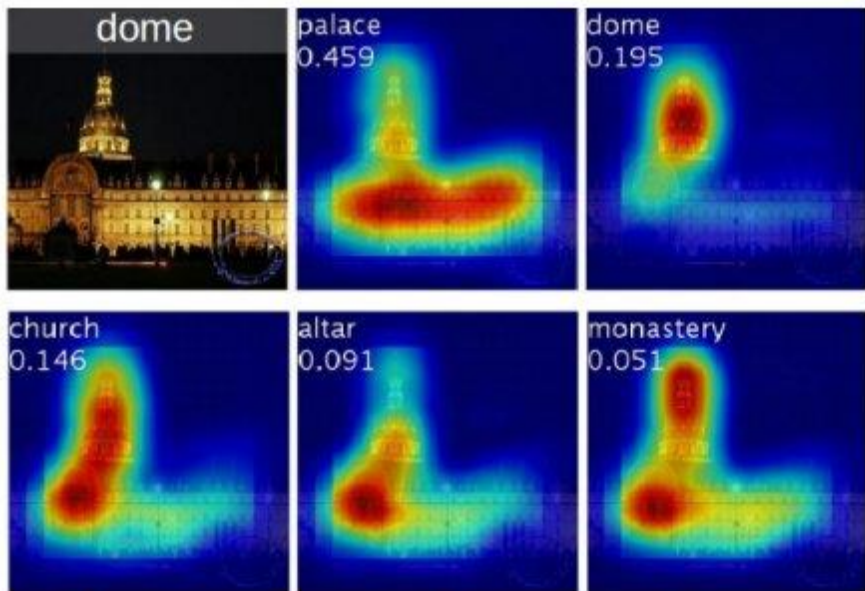
$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k} \\ = \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k}$$

Class Activation Maps:

$$M \in \mathbb{R}^{C,H,W}$$

$$M_{c,h,w} = \sum_k w_{k,c} f_{h,w,k}$$

# Class Activation Mapping (CAM)



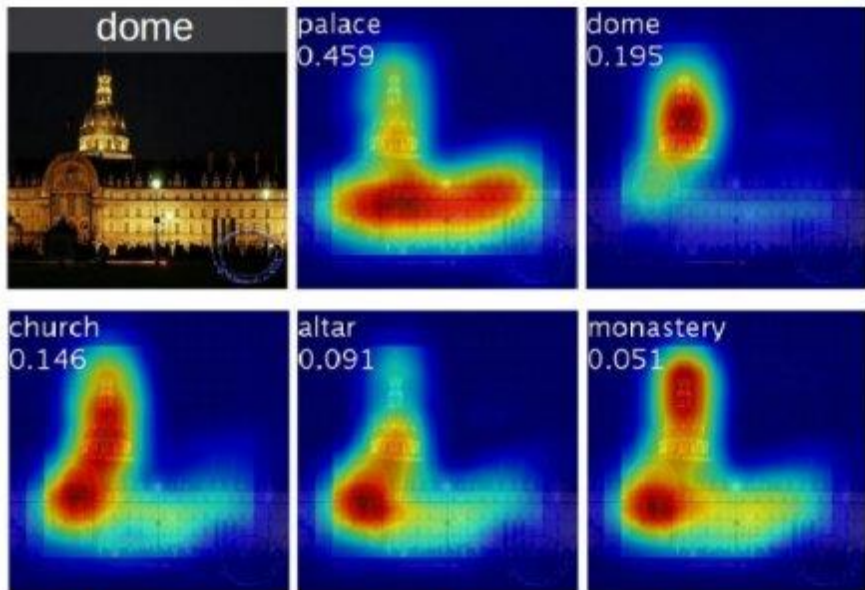
Class activation maps of top 5 predictions



Class activation maps for one object class

# Class Activation Mapping (CAM)

Problem: Can only apply to last conv layer



Class activation maps of top 5 predictions



Class activation maps for one object class

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$
2. Compute gradient of class score  $S_c$  with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$
2. Compute gradient of class score  $S_c$  with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

3. Global Average Pool the gradients to get weights  $\alpha \in \mathbb{R}^K$ :

$$\alpha_k = \frac{1}{HW} \sum_{h,w} \frac{\partial S_c}{\partial A_{h,w,k}}$$

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$
2. Compute gradient of class score  $S_c$  with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

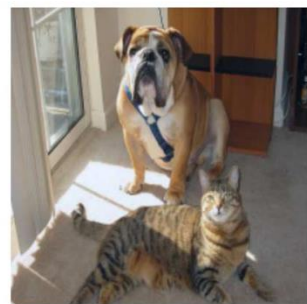
3. Global Average Pool the gradients to get weights  $\alpha \in \mathbb{R}^K$ :

$$\alpha_k = \frac{1}{HW} \sum_{h,w} \frac{\partial S_c}{\partial A_{h,w,k}}$$

4. Compute activation map  $M^c \in \mathbb{R}^{H,W}$ :

$$M_{h,w}^c = \text{ReLU} \left( \sum_k \alpha_k A_{h,w,k} \right)$$

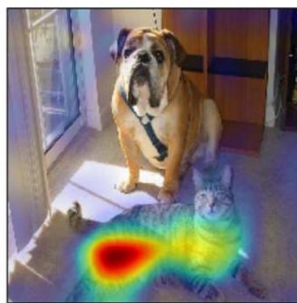
# Gradient-Weighted Class Activation Mapping (Grad-CAM)



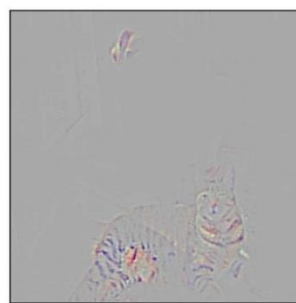
(a) Original Image



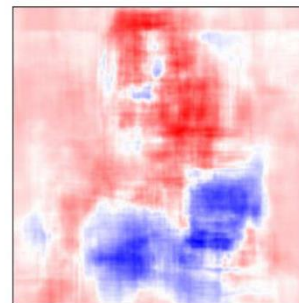
(b) Guided Backprop 'Cat'



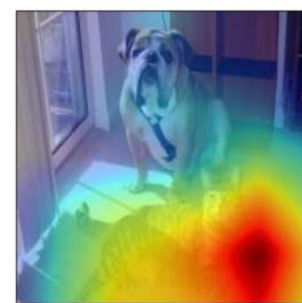
(c) Grad-CAM 'Cat'



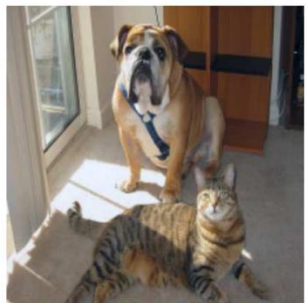
(d) Guided Grad-CAM 'Cat'



(e) Occlusion map for 'Cat'



(f) ResNet Grad-CAM 'Cat'



(g) Original Image



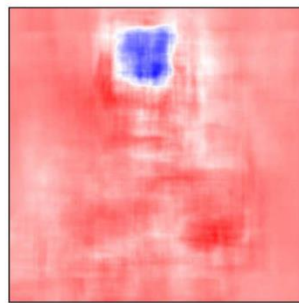
(h) Guided Backprop 'Dog'



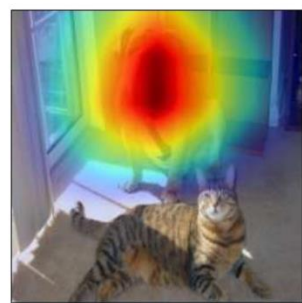
(i) Grad-CAM 'Dog'



(j) Guided Grad-CAM 'Dog'

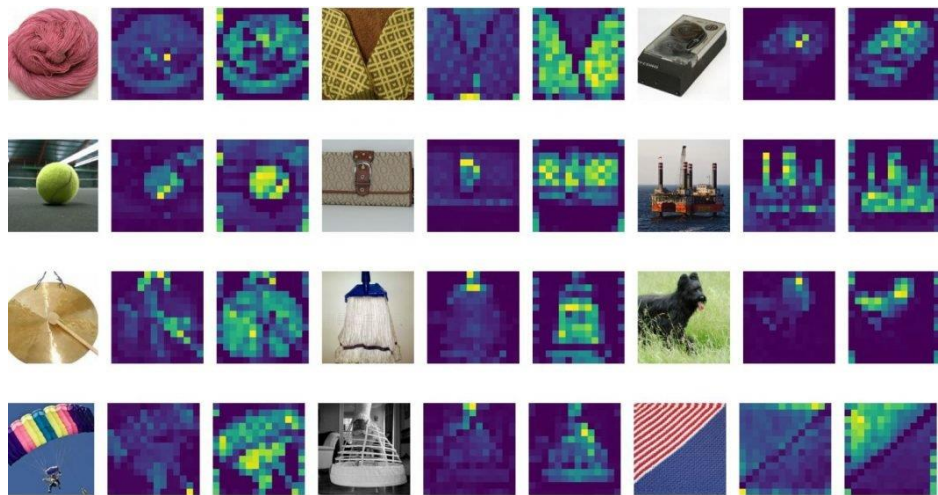


(k) Occlusion map for 'Dog'



(l) ResNet Grad-CAM 'Dog'

# Visualizing ViT features



High-confidence



Chen et al., When Vision Transformers Outperform Resnets Without Pre-training Or Strong Data Augmentations, ICLR 2022; Paul and Chen, Vision Transformers are Robust Learners, AAAI 2022. Reproduced for educational purposes.

# Today

- Transformers Recap
- **Computer Vision Tasks**
  - Semantic Segmentation
  - Object Detection
  - Instance Segmentation
- Visualization & Understanding
  - Model Layers Visualization
  - Saliency Maps
  - CAM & Grad-CAM

Next time: Video Understanding