



# Machine Learning

## Lecture 14: Computer Vision

Fall 2023

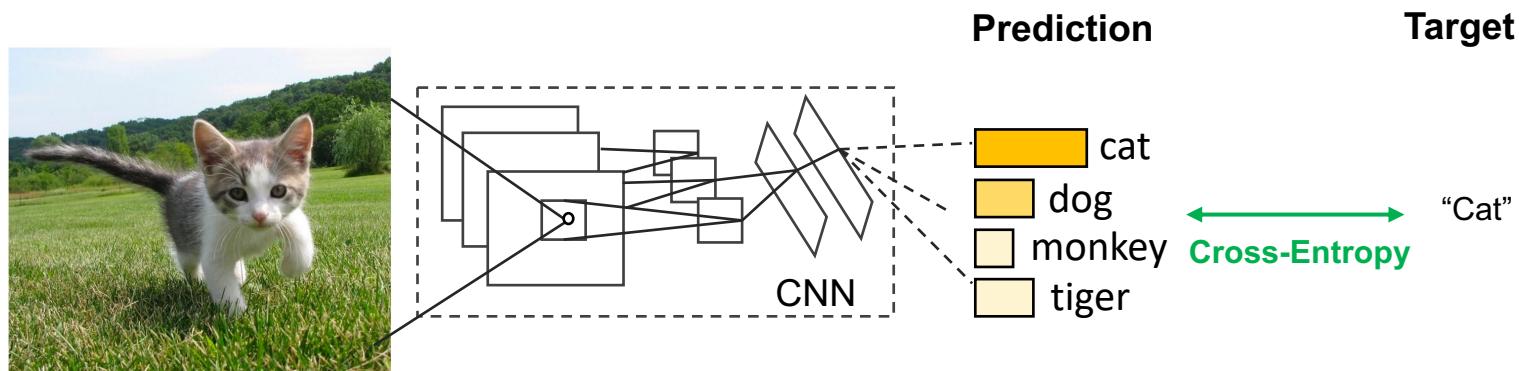
Instructor: Xiaodong Gu



# Recall: Image Classification



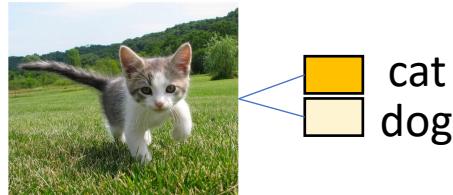
- Convolution + MLP + Softmax



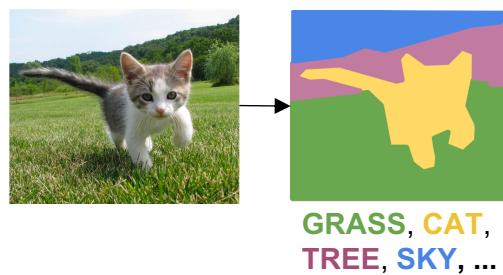
# Computer Vision Tasks



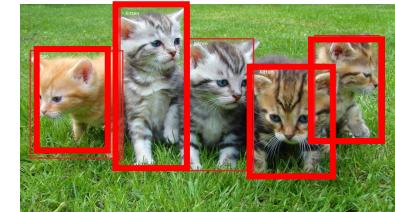
## Classification



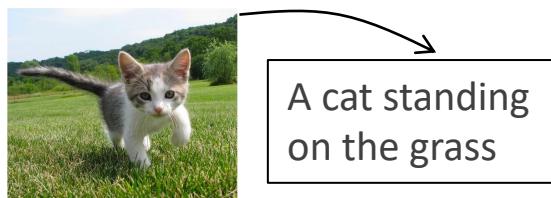
## Semantic Segmentation



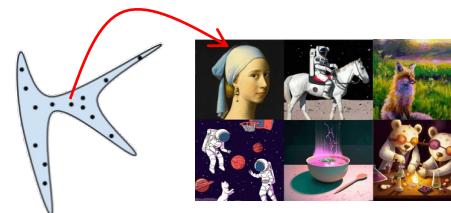
## Object Detection



## Image Captioning



## Image Generation

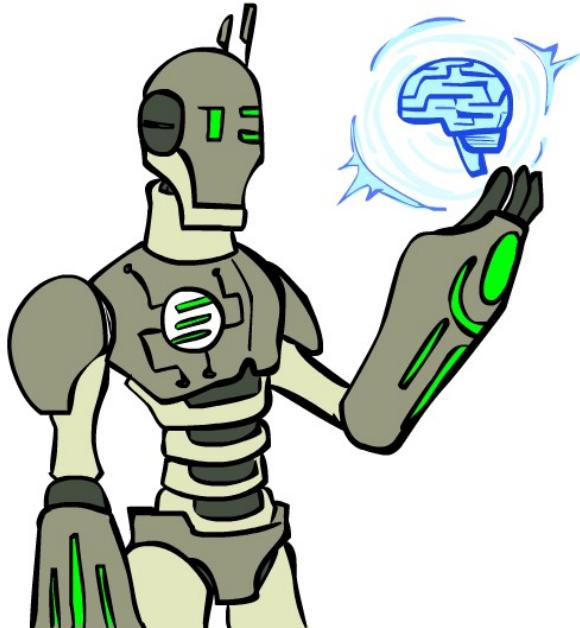


# Today

---



- Image Segmentation
- Object Detection
- Image Captioning



# Semantic Segmentation

---

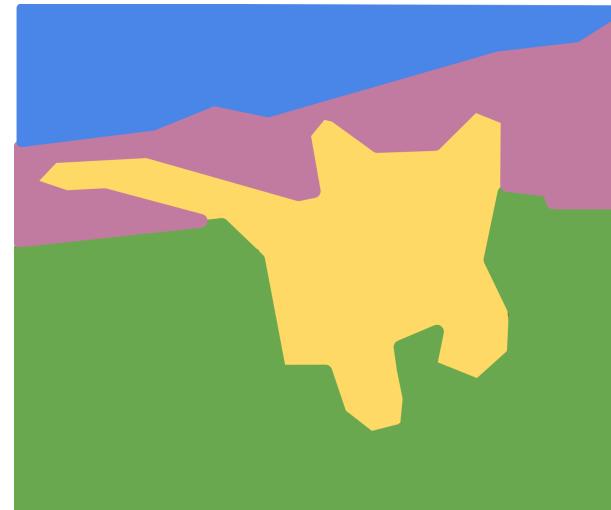


- Segment an image into disjoint regions that have independent semantics.

Label each pixel in the image with a category



**Input:** an image



**Output:** class of each pixel

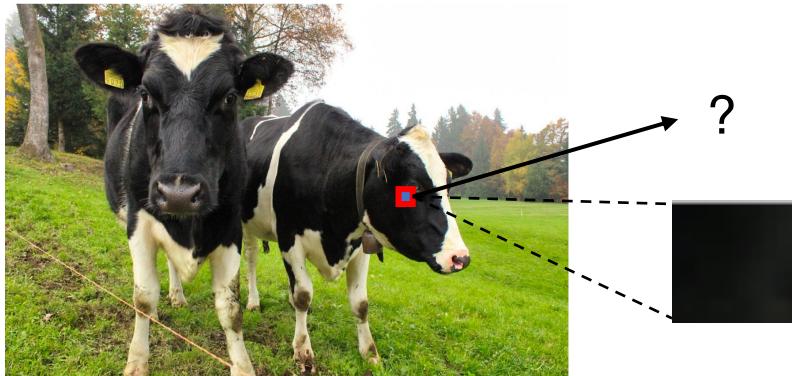
# A Naïve Idea

---



- Classifying Pixels?

Full image



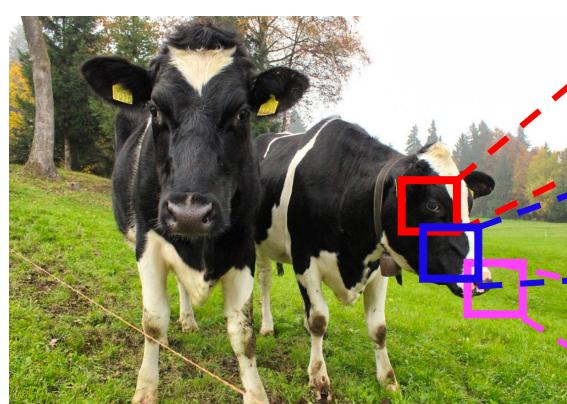
Impossible to classify without context

Q: how do we include context?



# The Idea

- Sliding Window?

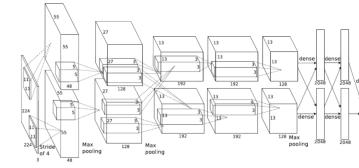


Full image

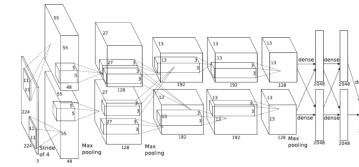
Extract patch



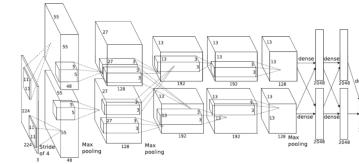
Classify center pixel  
with CNN



Cow



Cow



Grass

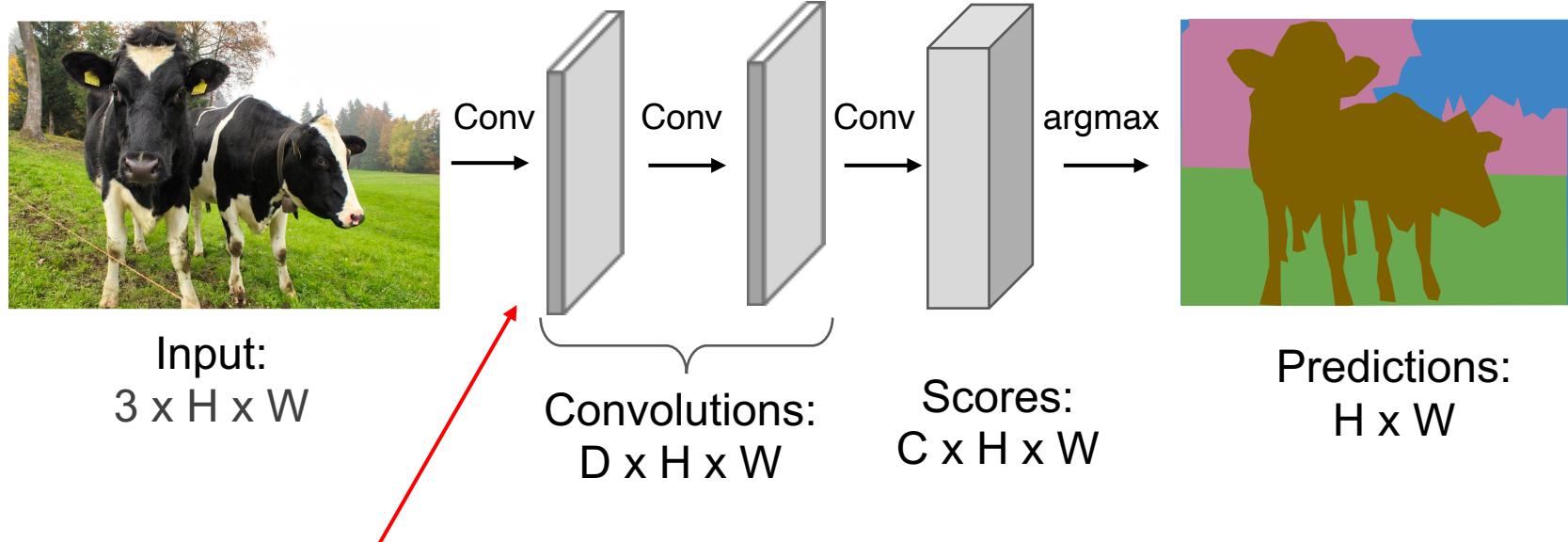
Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013  
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014



# 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 ...

# Fully Convolutional?

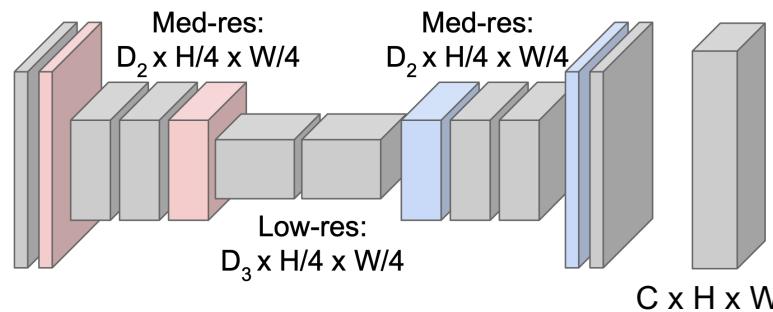


**Instead:** Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

**Downsampling:**  
Pooling, strided  
convolution



Input:  
 $3 \times H \times W$



**Upsampling:**  
???



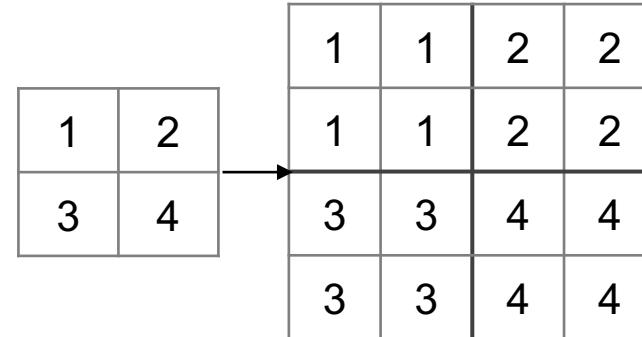
Predictions:  
 $H \times W$

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”



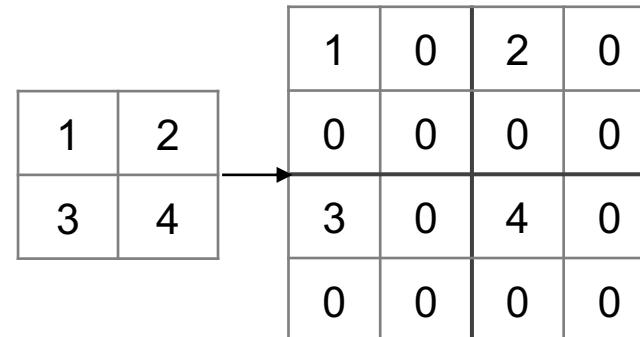
**Nearest Neighbor**



Input: 2 x 2

Output: 4 x 4

**“Bed of Nails”**



Input: 2 x 2

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

5	6
7	8

Rest of the network

## Max Unpooling

Use positions from pooling layer

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

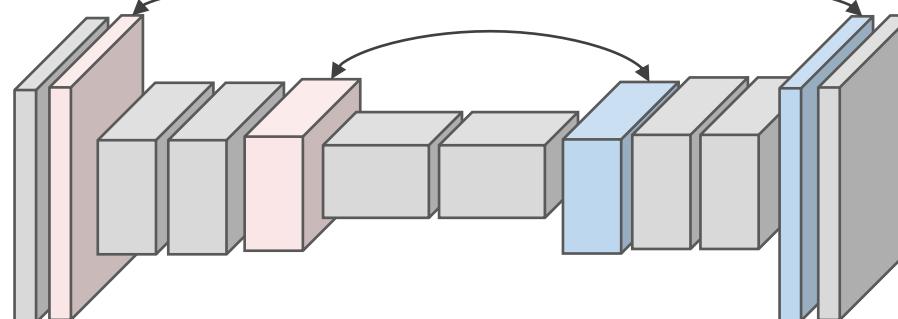
Input: 4 x 4

Output: 2 x 2

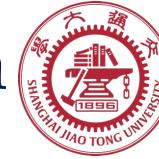
Input: 2 x 2

Output: 4 x 4

Corresponding pairs of  
downsampling and  
upsampling layers

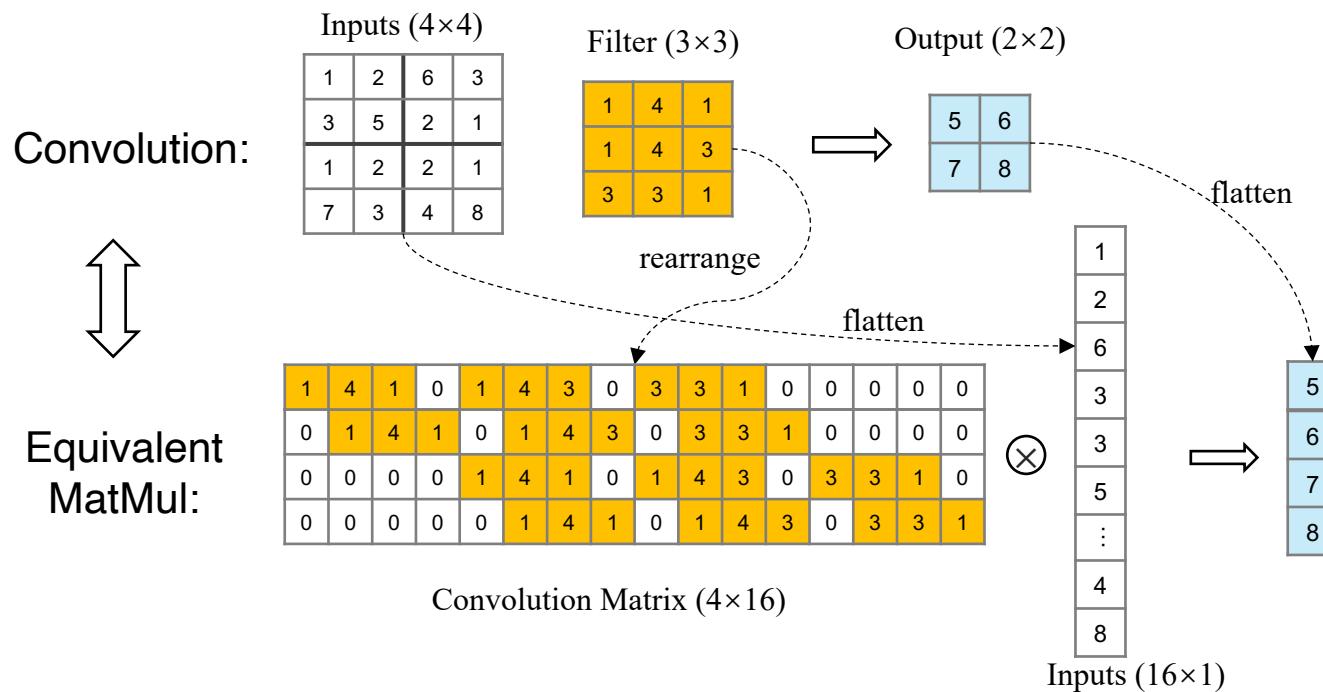


# Upsampling with Transposed Convolution



## • Convolution Matrix

We can express a convolution operation using a **matrix multiplication** by rearranging the kernel matrix into a general matrix.



# Upsampling with Transposed Convolution



- **Transposed Convolution Matrix**

Transpose the convolution matrix  $C$  ( $4 \times 16$ ) to  $C^T$  ( $16 \times 4$ ).

We can matrix-multiply  $C^T$  ( $16 \times 4$ ) with a column vector ( $4 \times 1$ ) to generate an output matrix ( $16 \times 1$ ).

**Deconvolution:**

Transposed Convolution  
Matrix ( $4 \times 16$ )

connects 1 value to 9  
values in the output.

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

Inputs ( $2 \times 2$ )

$$\otimes \begin{matrix} 5 \\ 6 \\ 7 \\ 8 \end{matrix} \Rightarrow$$

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

outputs ( $4 \times 4$ )

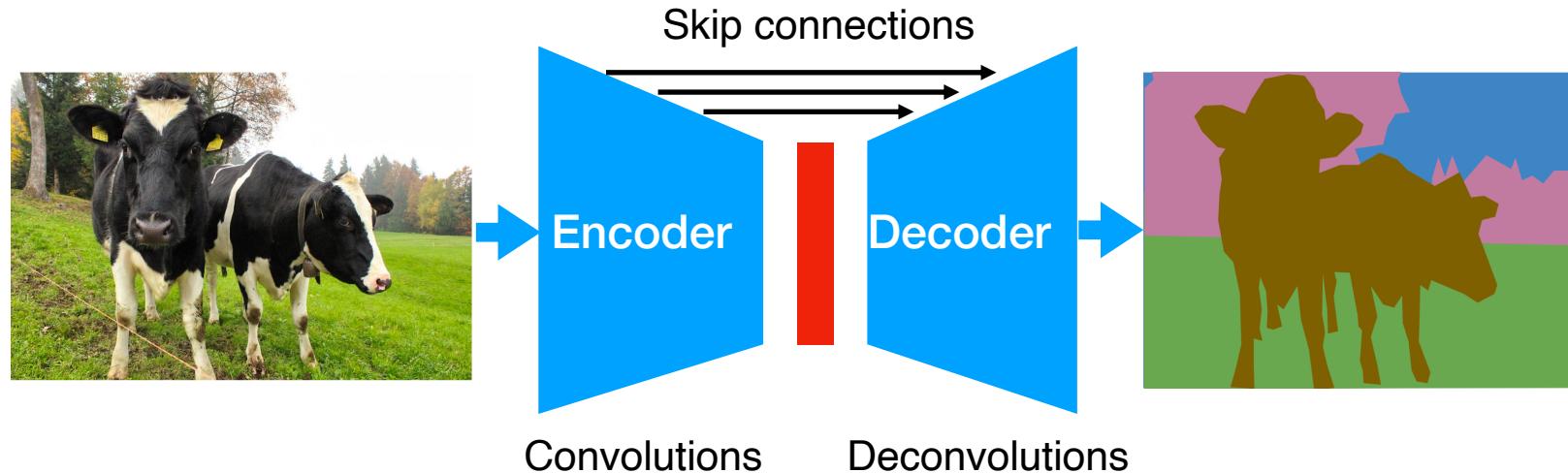
reshape →

$$\begin{matrix} 1 & 2 & 6 & 3 \\ 3 & 5 & 2 & 1 \\ 1 & 2 & 2 & 1 \\ 7 & 3 & 4 & 8 \end{matrix}$$

# Encoder-Decoder Architectures

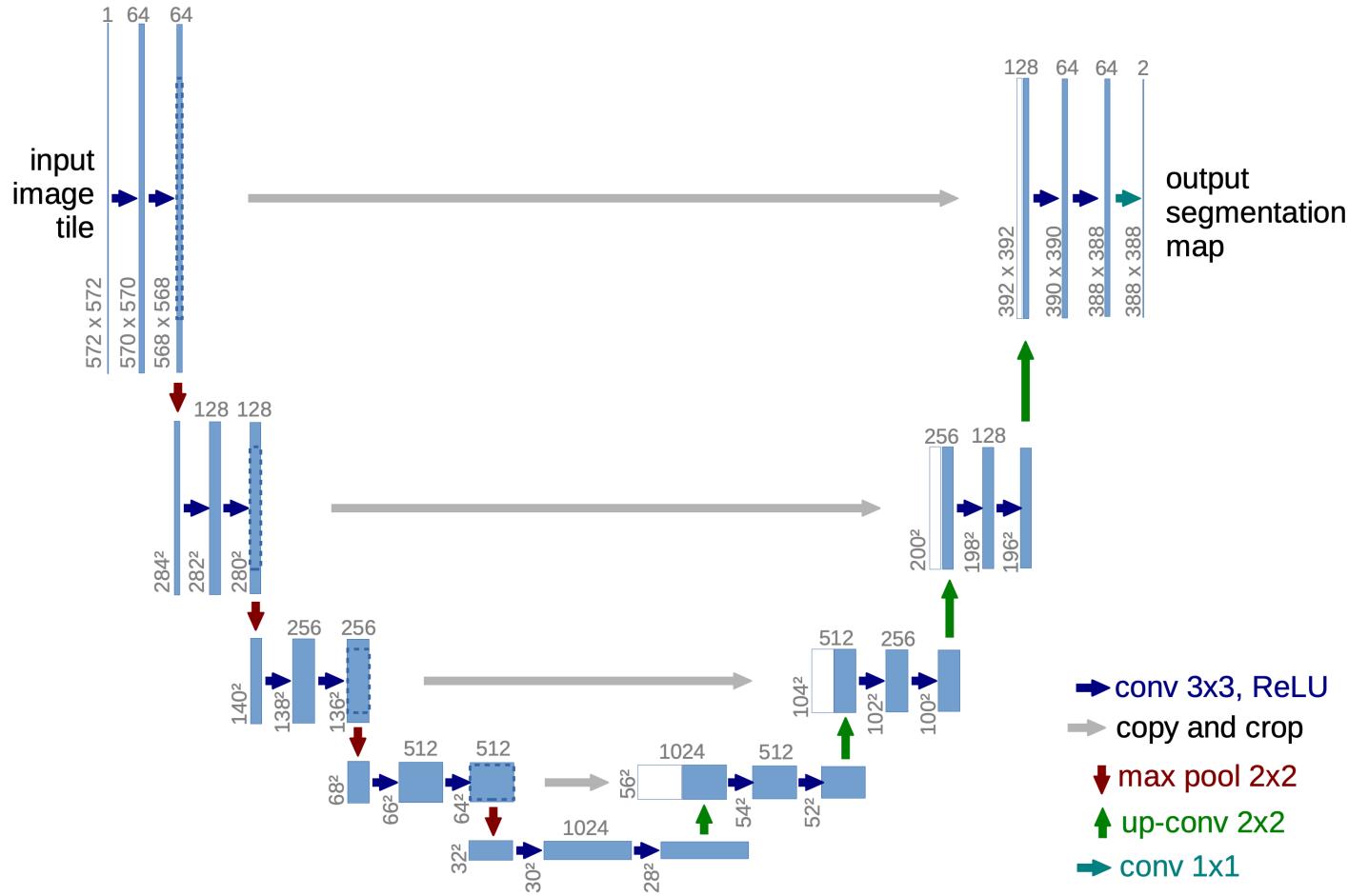


- U-Net



Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation

# U-Net



Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation

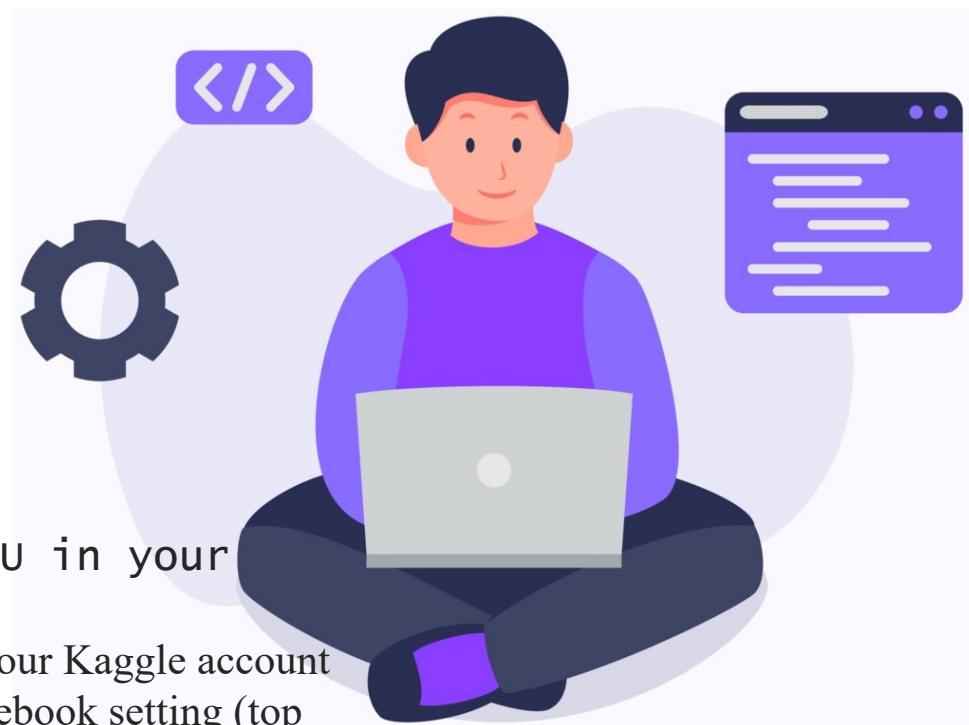
# TIME for Coding

---



- Tutorial: Image segmentation with U-Net

<https://www.kaggle.com/code/gokulkarthik/image-segmentation-with-unet-pytorch>



**NOTE:** How to enable GPU in your Kaggle kernel?

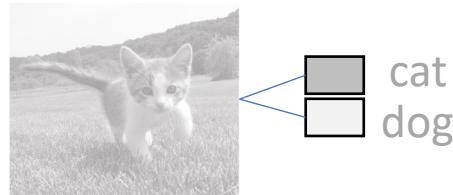
- Verify the phone number of your Kaggle account
- Click “accelerator” in the notebook setting (top right) and then select the GPU  
(<https://www.kaggle.com/general/97939>)

# Computer Vision Tasks

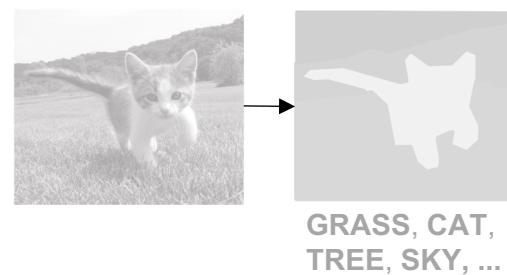
---



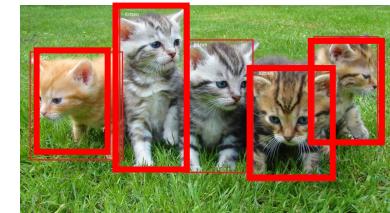
## Classification



## Semantic Segmentation



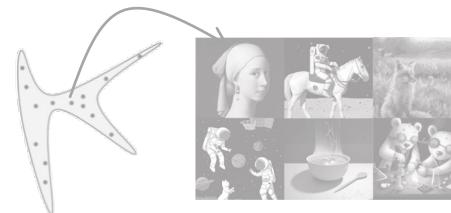
## Object Detection



## Image Captioning

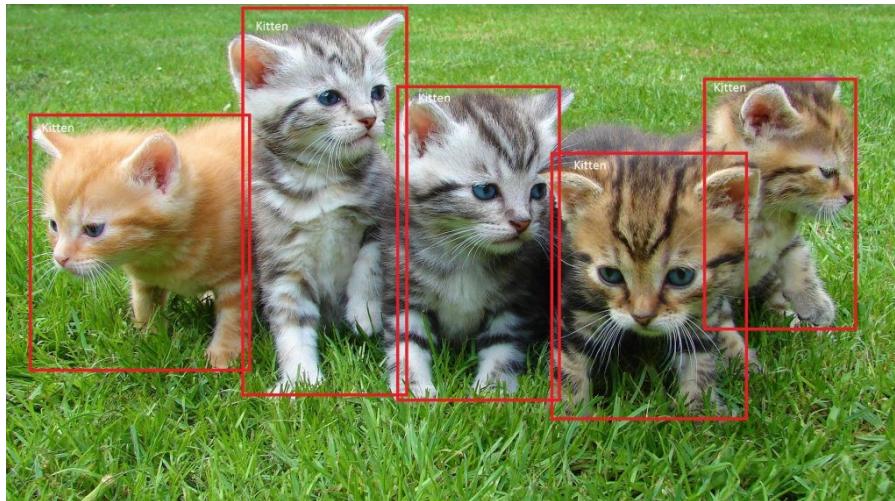


## Image Generation



# Object Detection

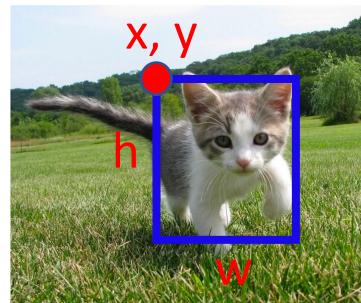
---



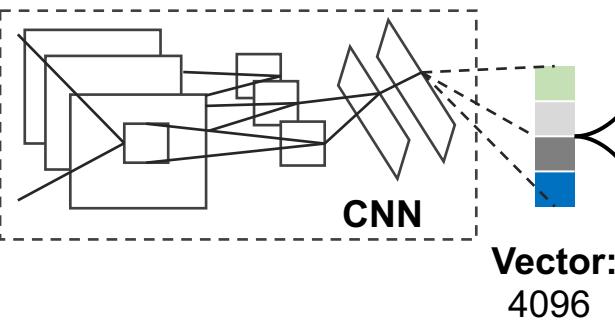
# Single Object



- Object detection = Classification + Localization



Input: an image



Vector:  
4096

Class Scores      Target

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

"Cat"

Fully  
Connected:  
4096 to 1000

Cross-Entropy

+ = Loss

Fully  
Connected:  
4096 to 4

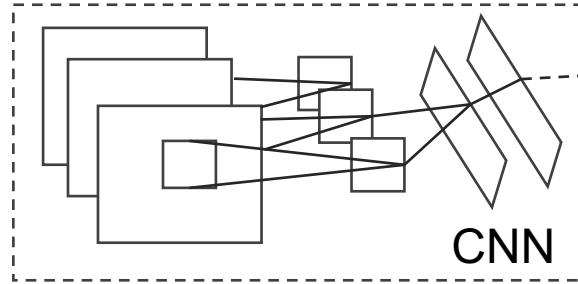
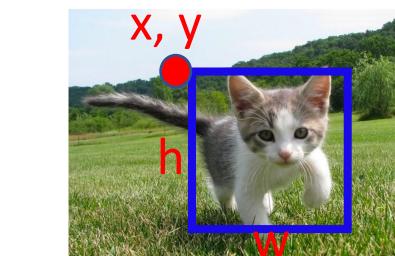
L2 Loss

( $x, y, w, h$ )

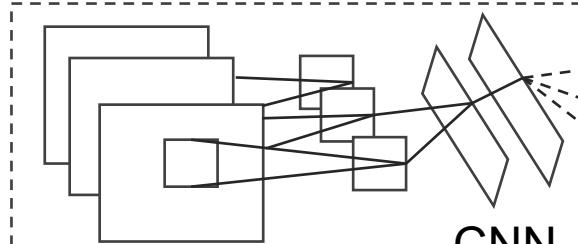
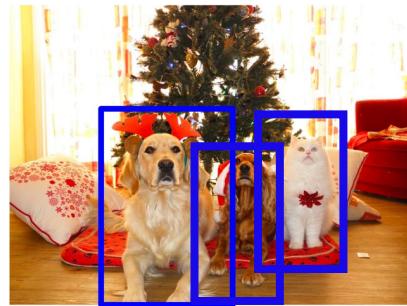
Box Coordinates      Target

$(x', y', w', h')$

# Multiple Objects?



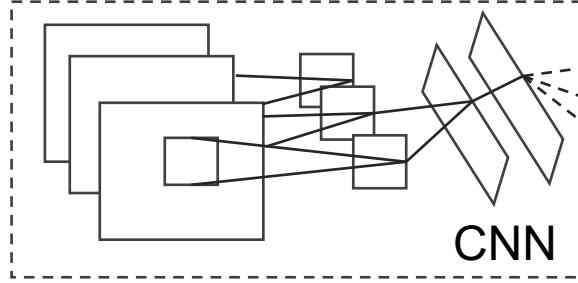
CAT:  $(x, y, w, h)$



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

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

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



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

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

....

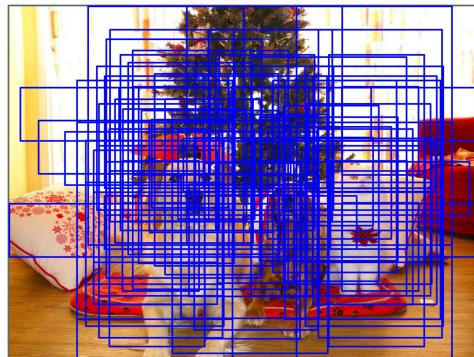
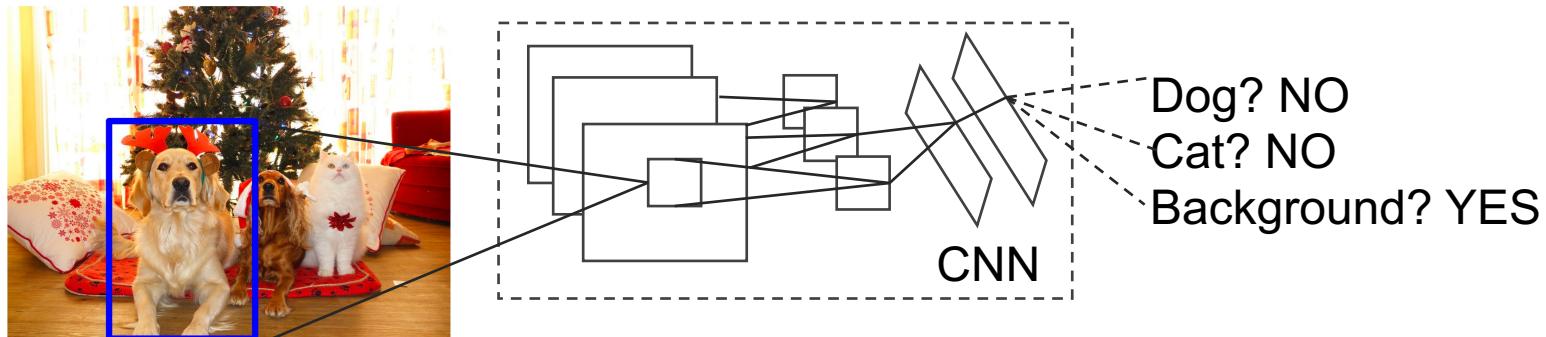
Too many predictions

- **Problem:** each image (CNN) needs a different number of outputs!



# Multiple Objects

- **Idea:** enumerate different crops of the image, apply a CNN to classify **each** crop as object or background.

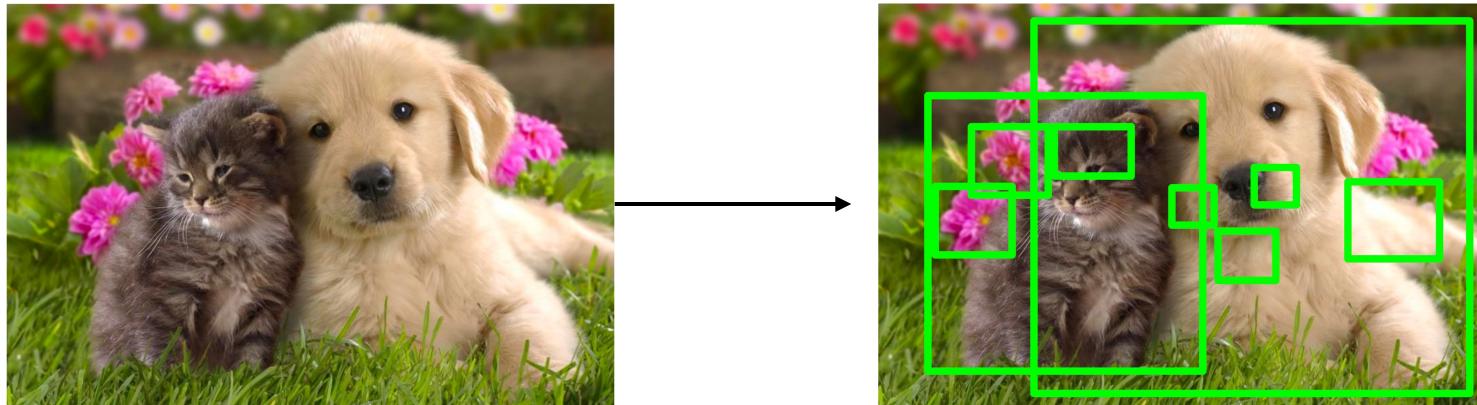


**Problem:** Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

# 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



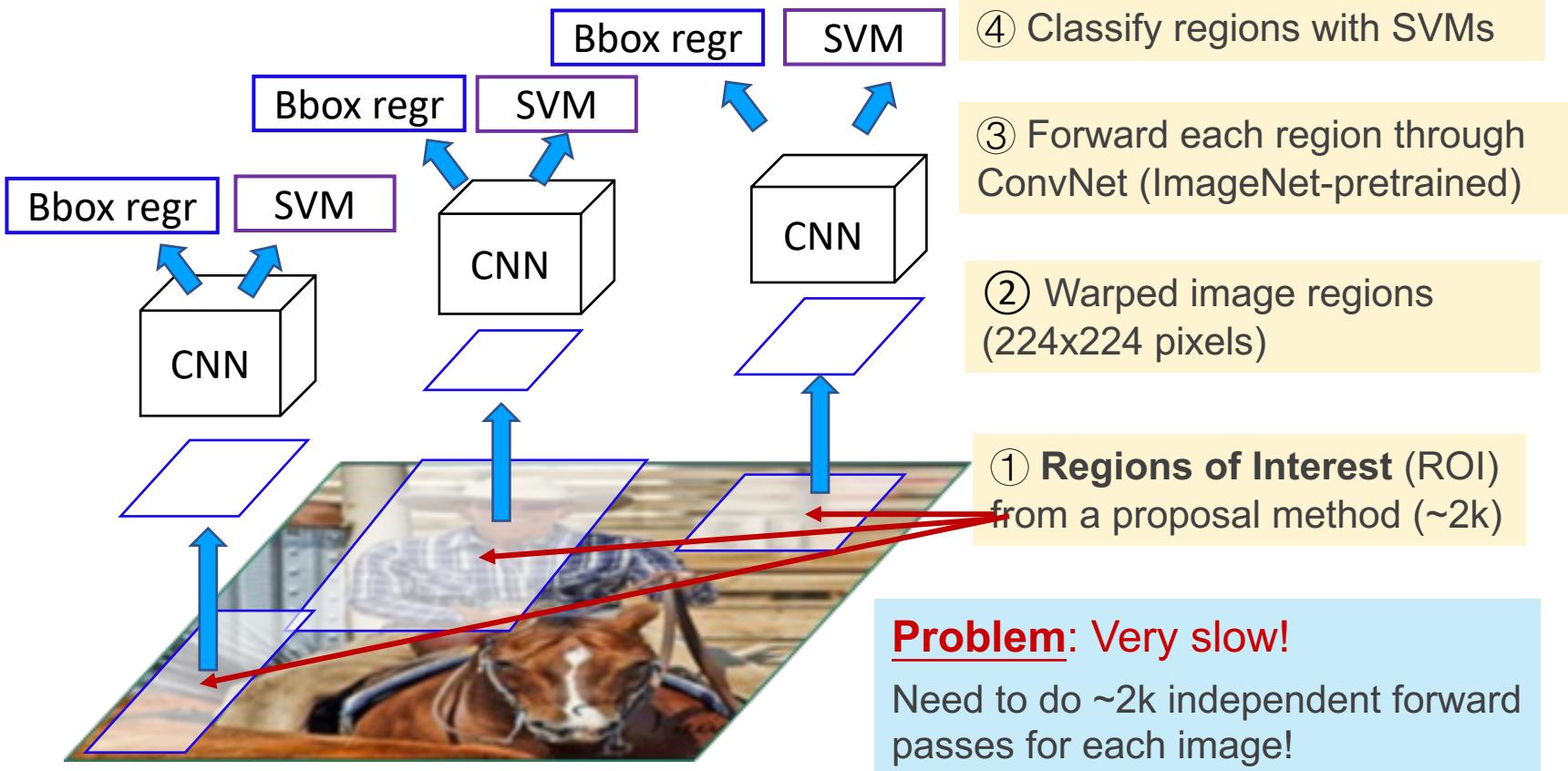
Alexe et al, “Measuring the objectness of image windows”, TPAMI 2012

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

⑤ Predict “corrections” to the ROI position: (dx, dy, dw, dh)

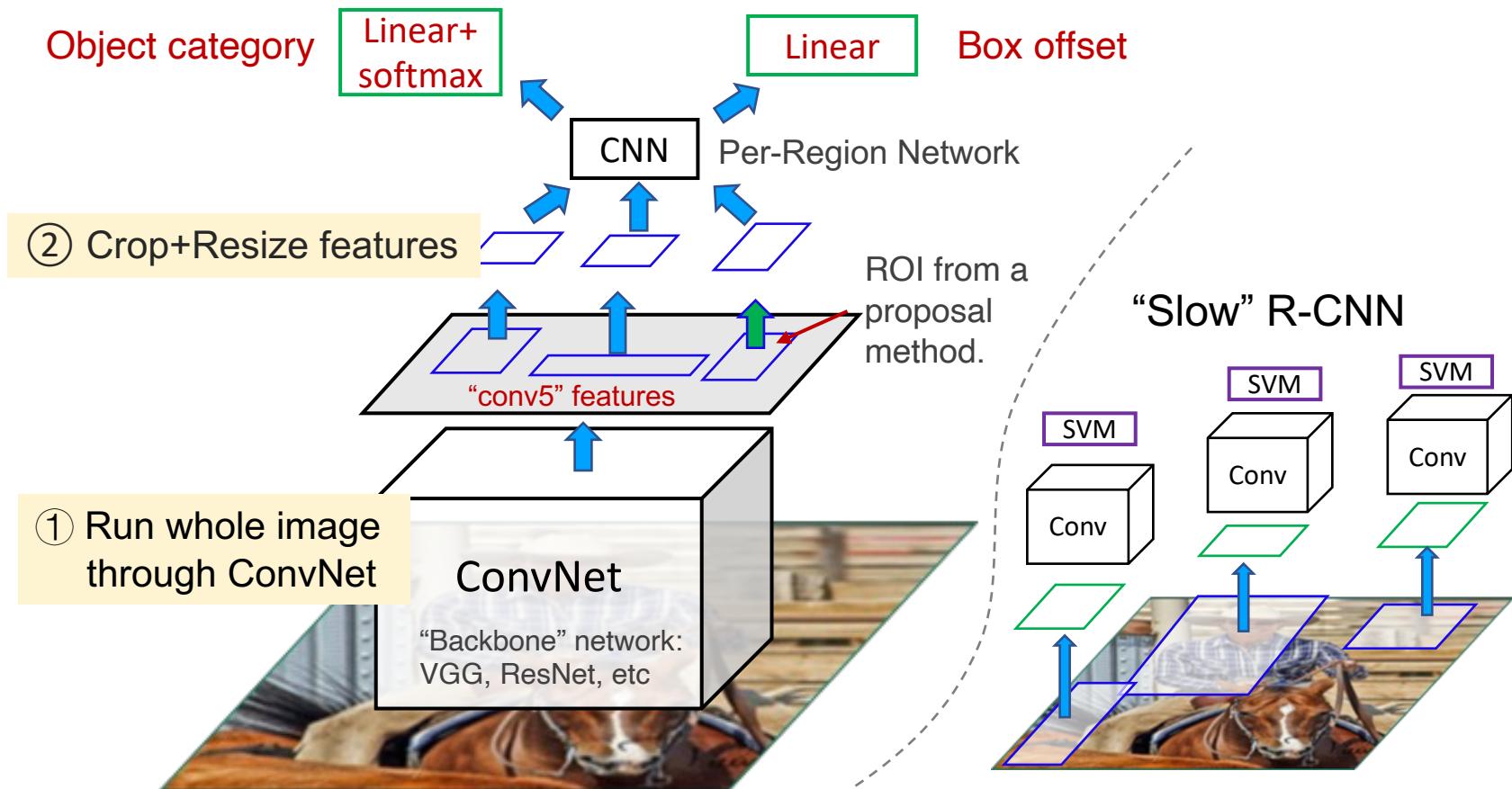


Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

# Fast R-CNN *Crop the conv feature instead!*



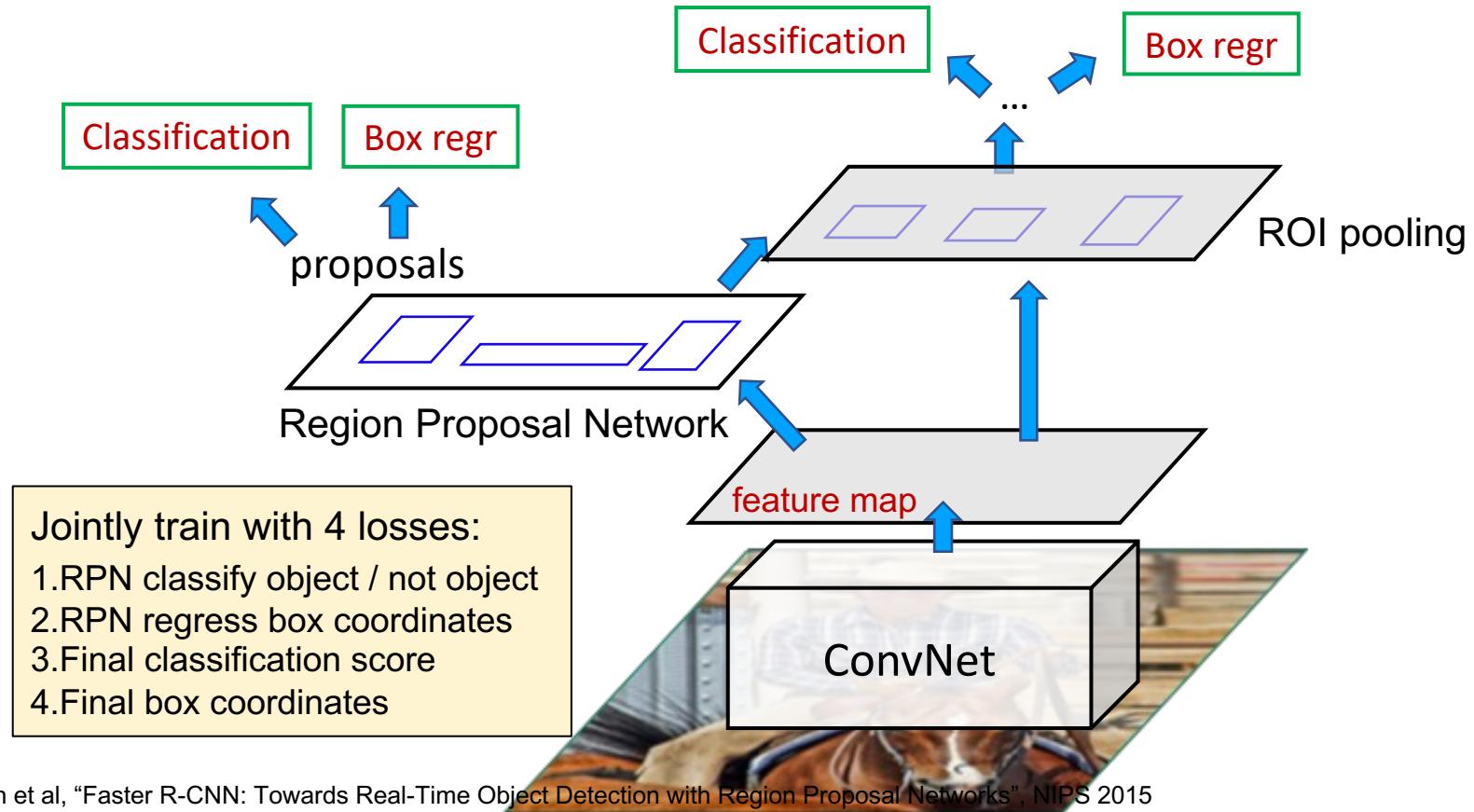
- Idea: Pass the image through convnet before cropping!



# Faster R-CNN Make CNN do proposals !



- Insert a 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

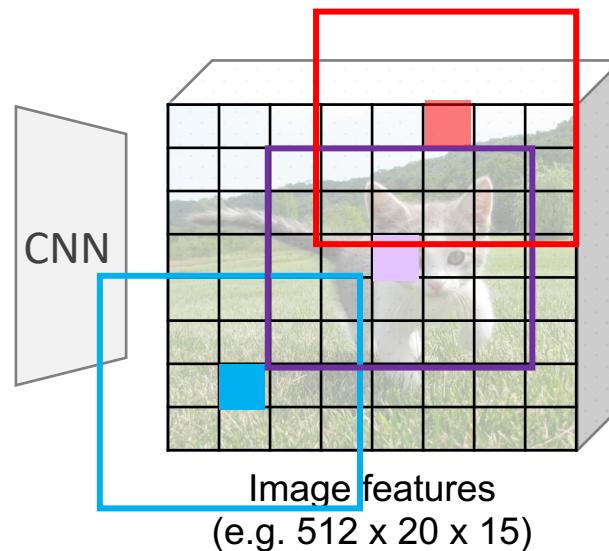
# Region Proposal Network



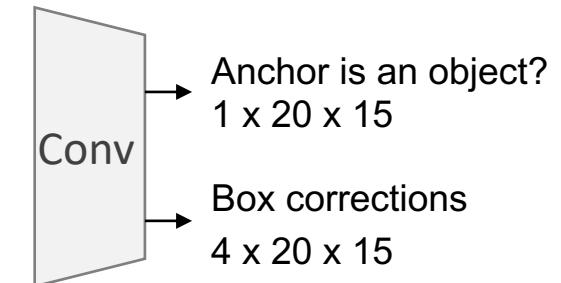
- Imagine an **anchor box** of fixed size at each point in the feature map



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



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

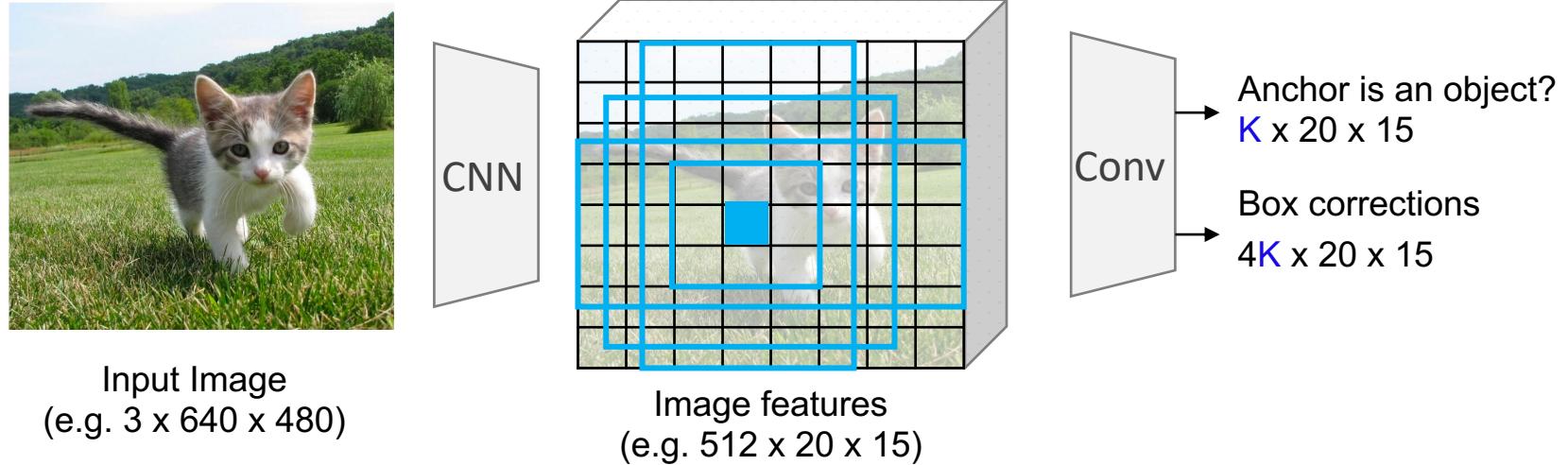


For positive boxes, also predict a corrections from the anchor to the ground-truth box (regress 4 numbers per pixel)

# Region Proposal Network



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



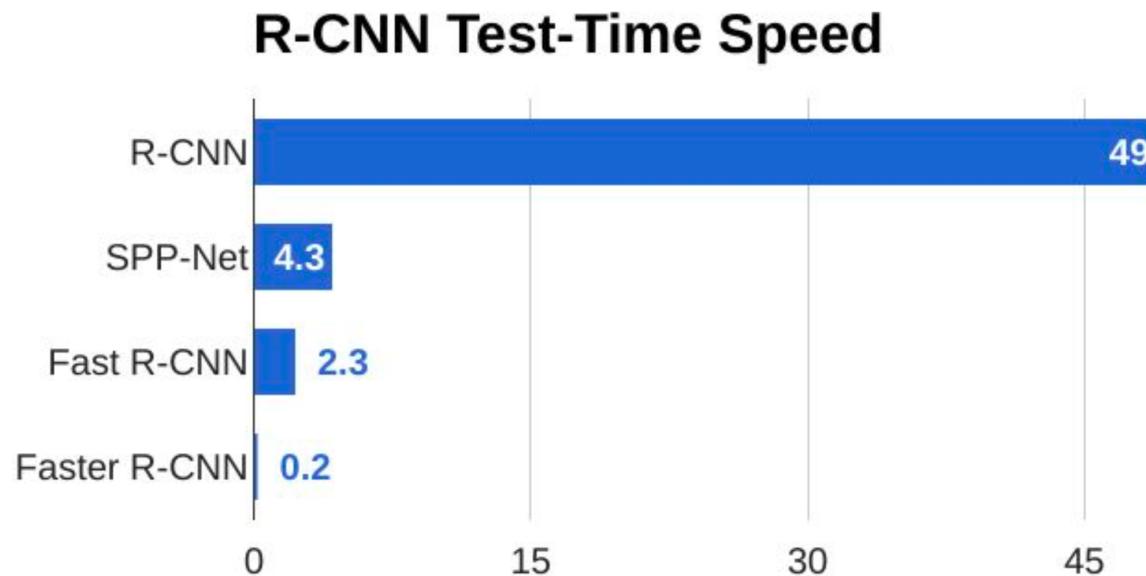
- Sort the  $K \times 20 \times 15$  boxes by their “objectness” score, take top ~300 as our proposals

# Faster R-CNN

---



- How fast?

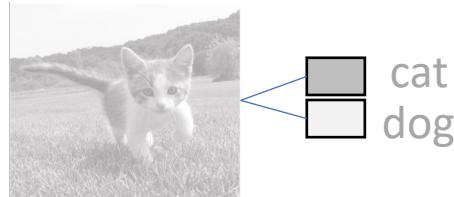


# Computer Vision Tasks

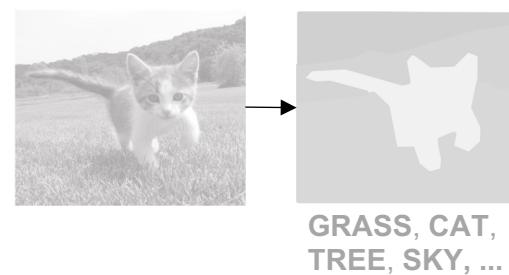
---



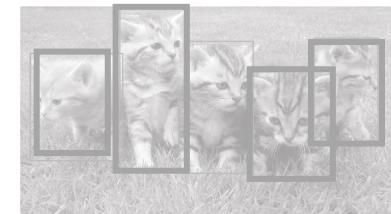
## Classification



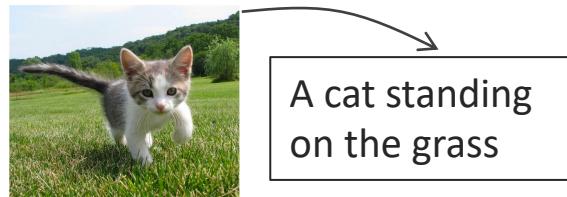
## Semantic Segmentation



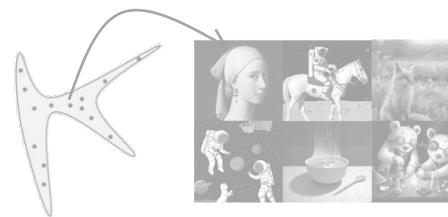
## Object Detection



## Image Captioning



## Image Generation



# Image Captioning



- Generate a natural language description for an image.



**Input:** Image I



“A person wearing a hat  
is sitting by the river”

**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$

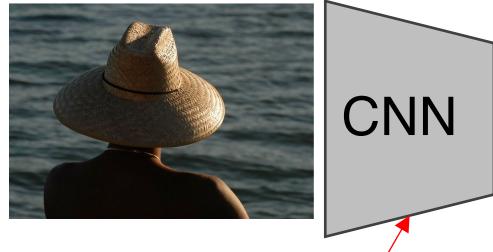


# Idea: Encoder-Decoder?

**Encoder:**  $c = f_w(z)$

where  $z$  is spatial CNN features  
 $f_w(\cdot)$  is an MLP

**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c)$   
where context vector  $c$  is often  $h_0 = c$



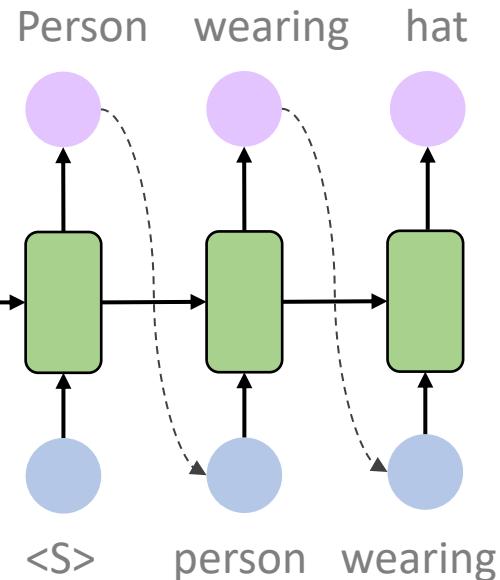
Extract spatial features from a pretrained CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

MLP

$c$



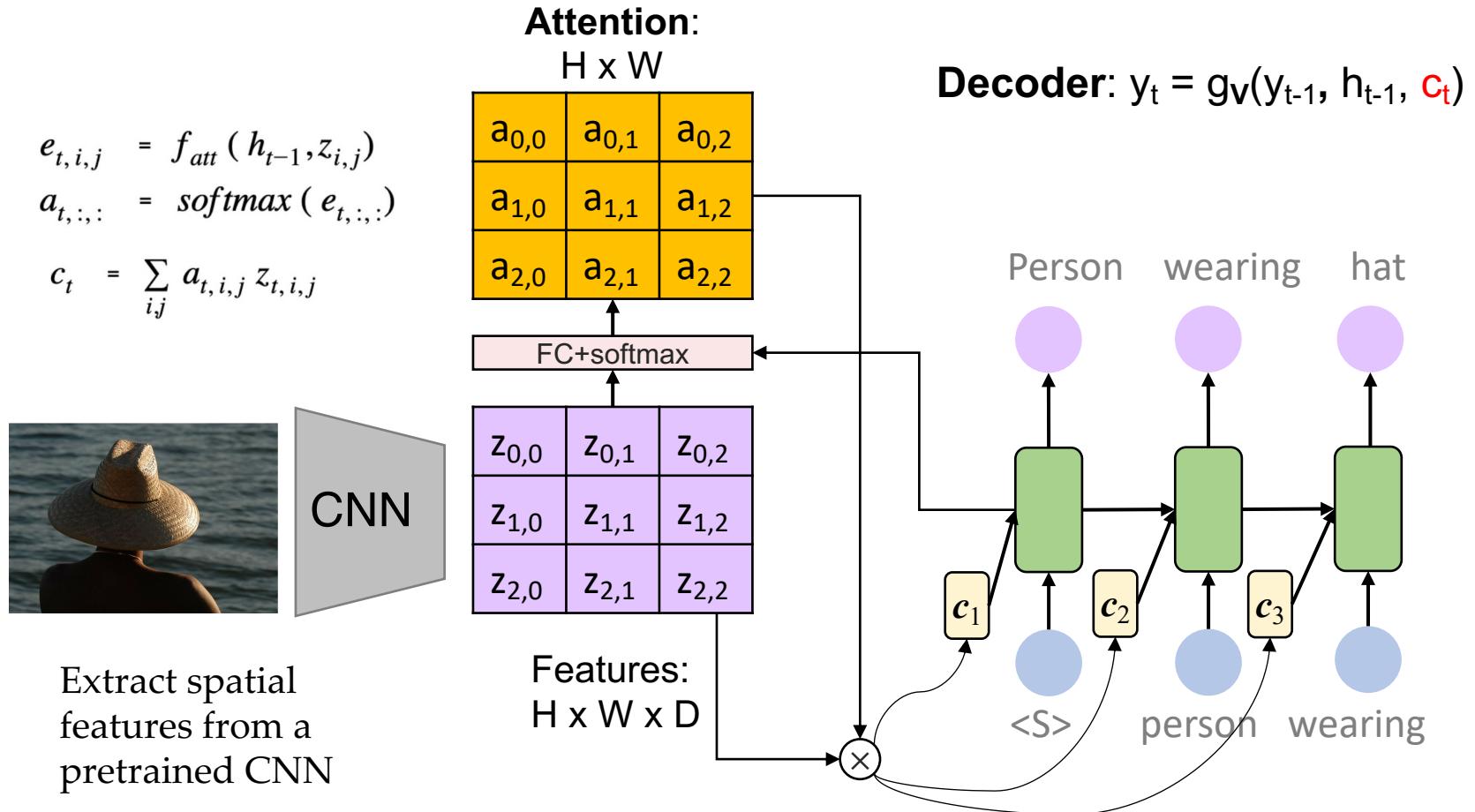
**Problem: Input is "bottlenecked" through  $c$**

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

# Encoder-Decoder with Attention?



- Each context vector will attend to different image regions



# Image Captioning with Attention



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



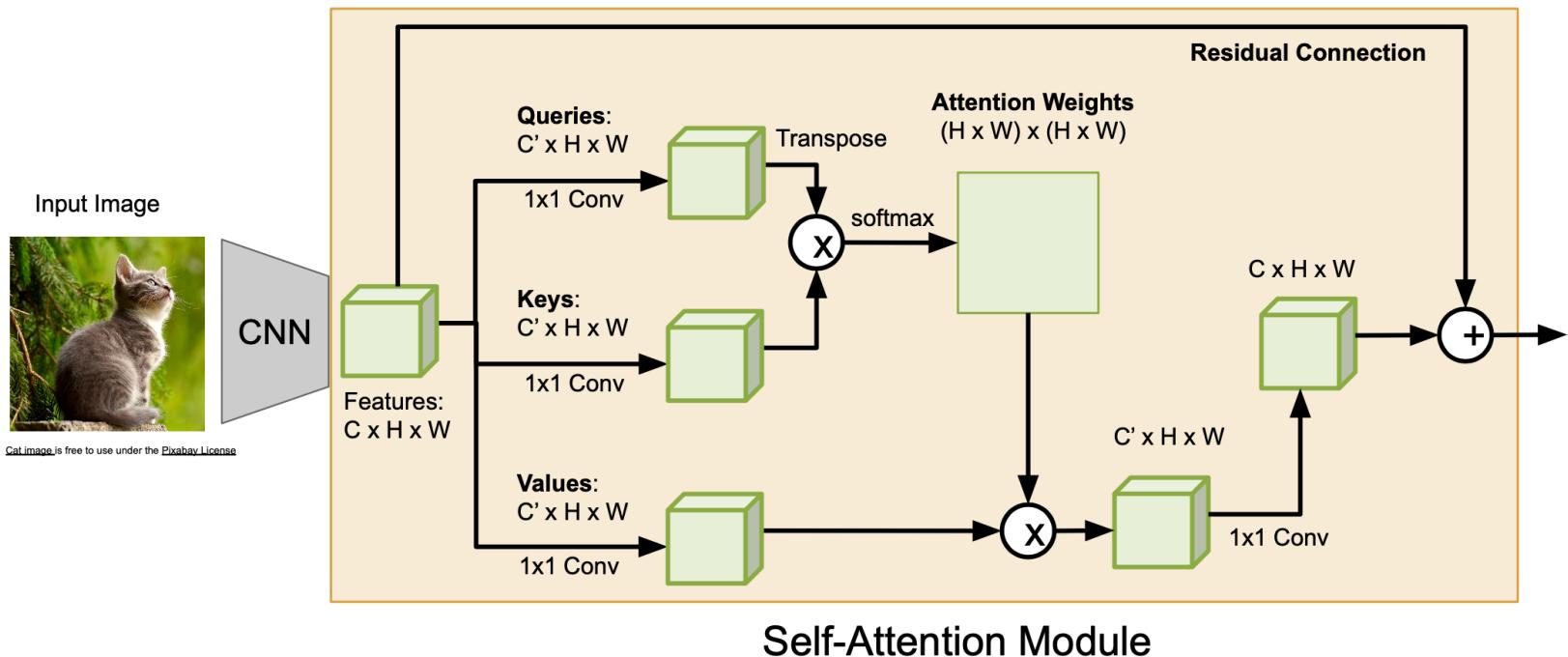
A giraffe standing in a forest with trees in the background.

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

# CNN with Self-Attention



- We can also apply **self-attention** to CNN features.

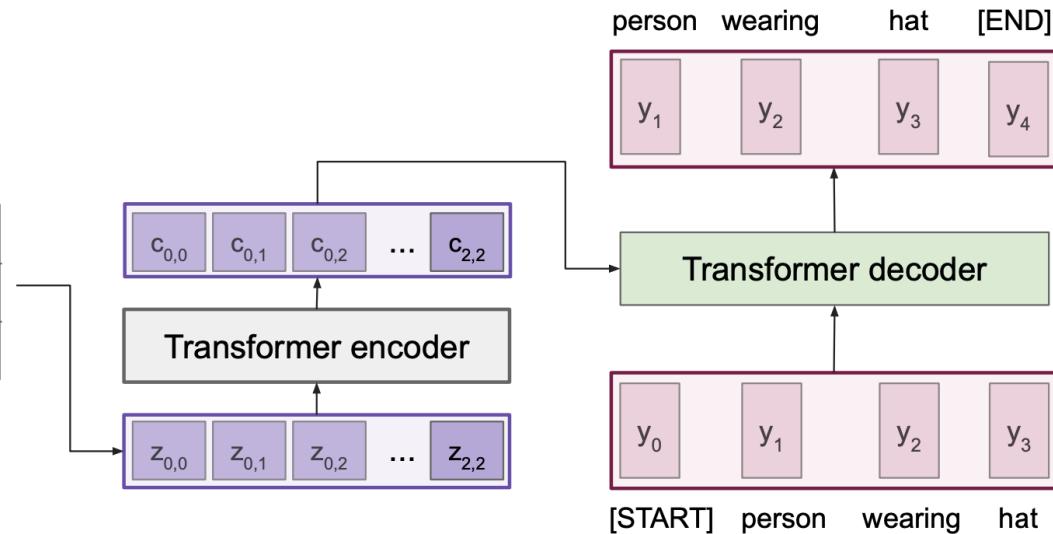
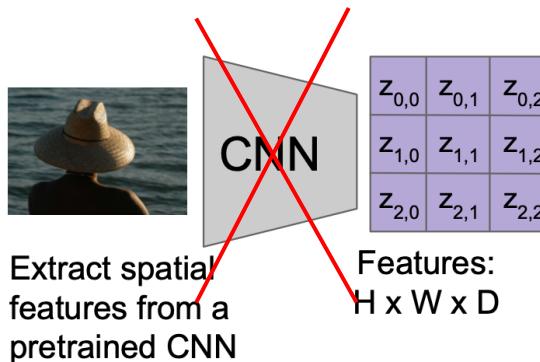


# Image Captioning ~~with~~ using Transformers



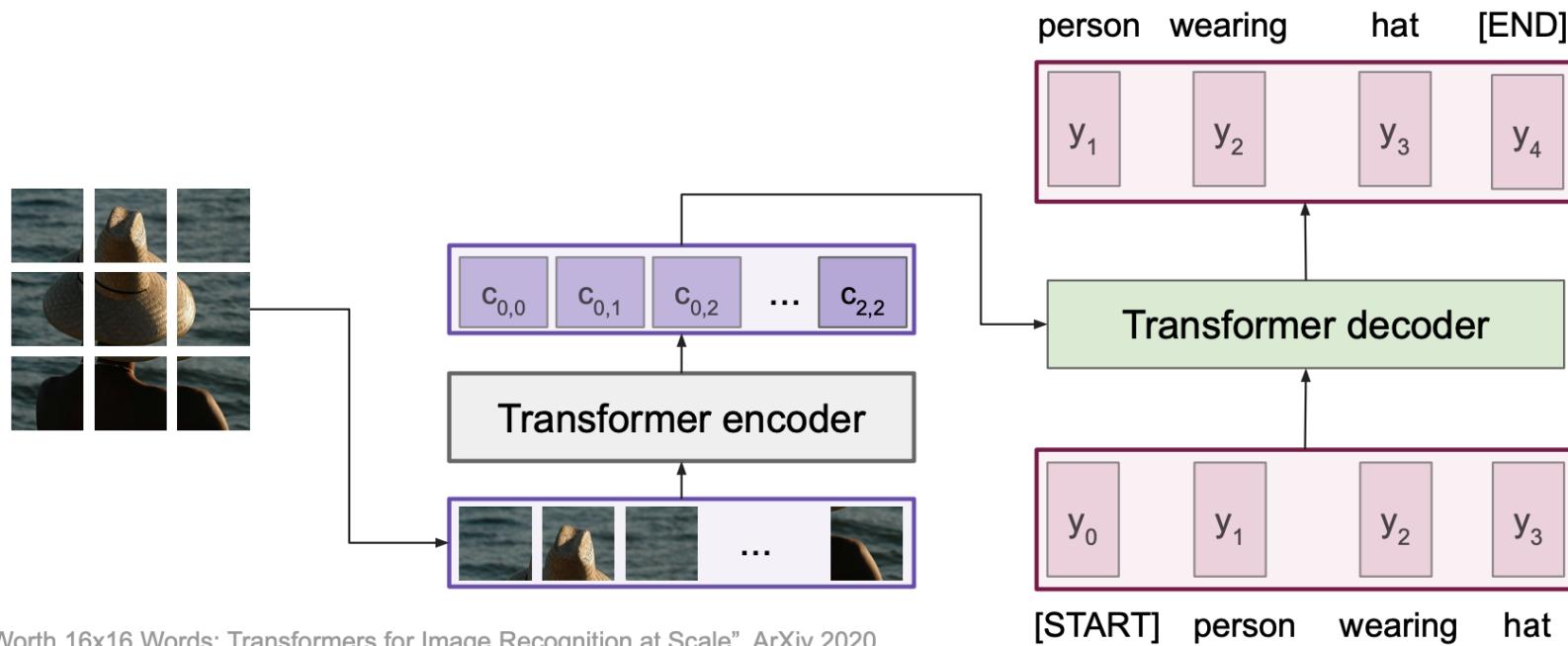
Encoder:  $c = T_w(z)$

where  $z$  is spatial CNN features  
 $T_w(\cdot)$  is the transformer encoder



Perhaps we don't need convolutions at all?

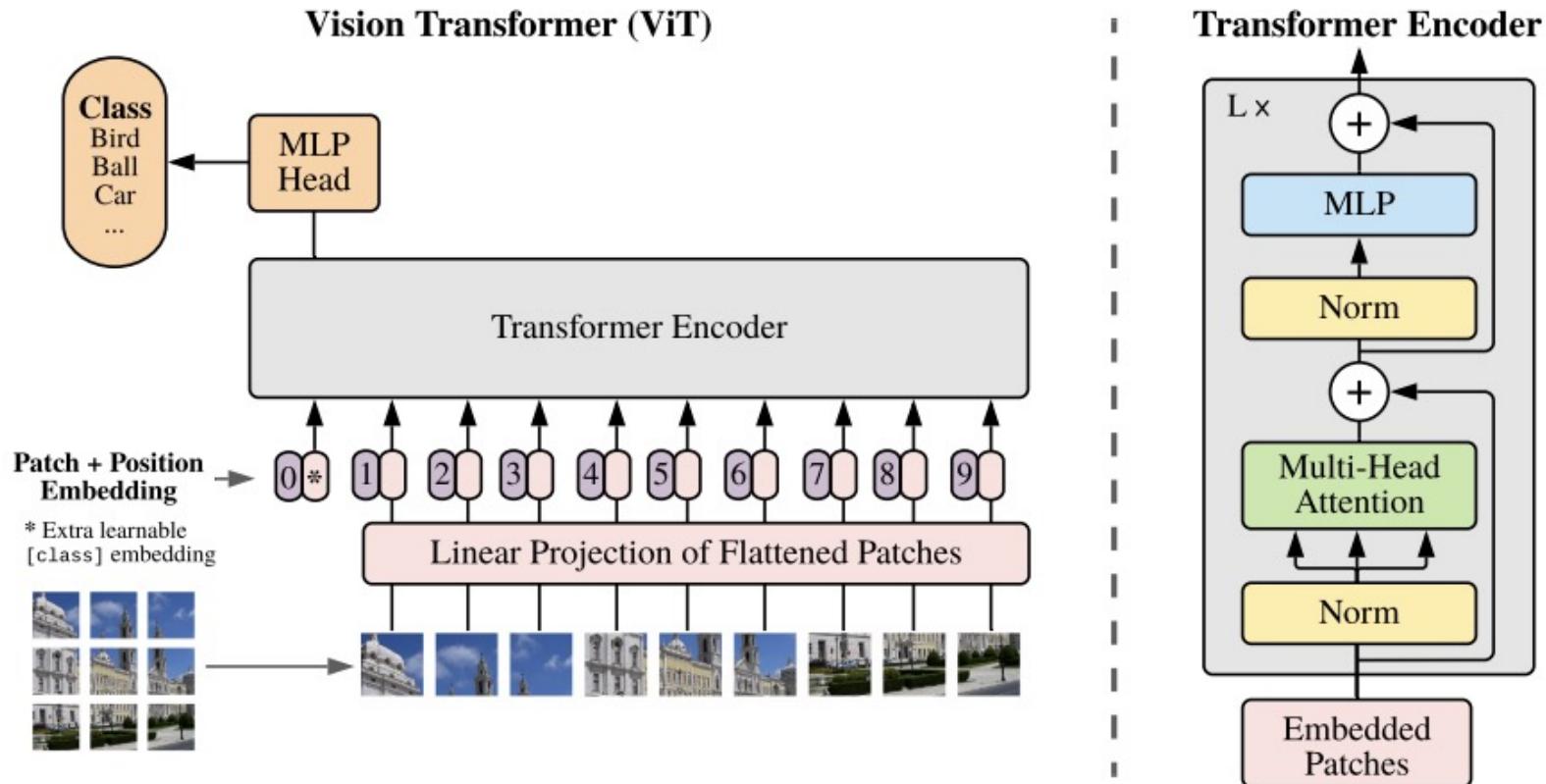
# Image Captioning using ONLY Transformers



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ArXiv 2020

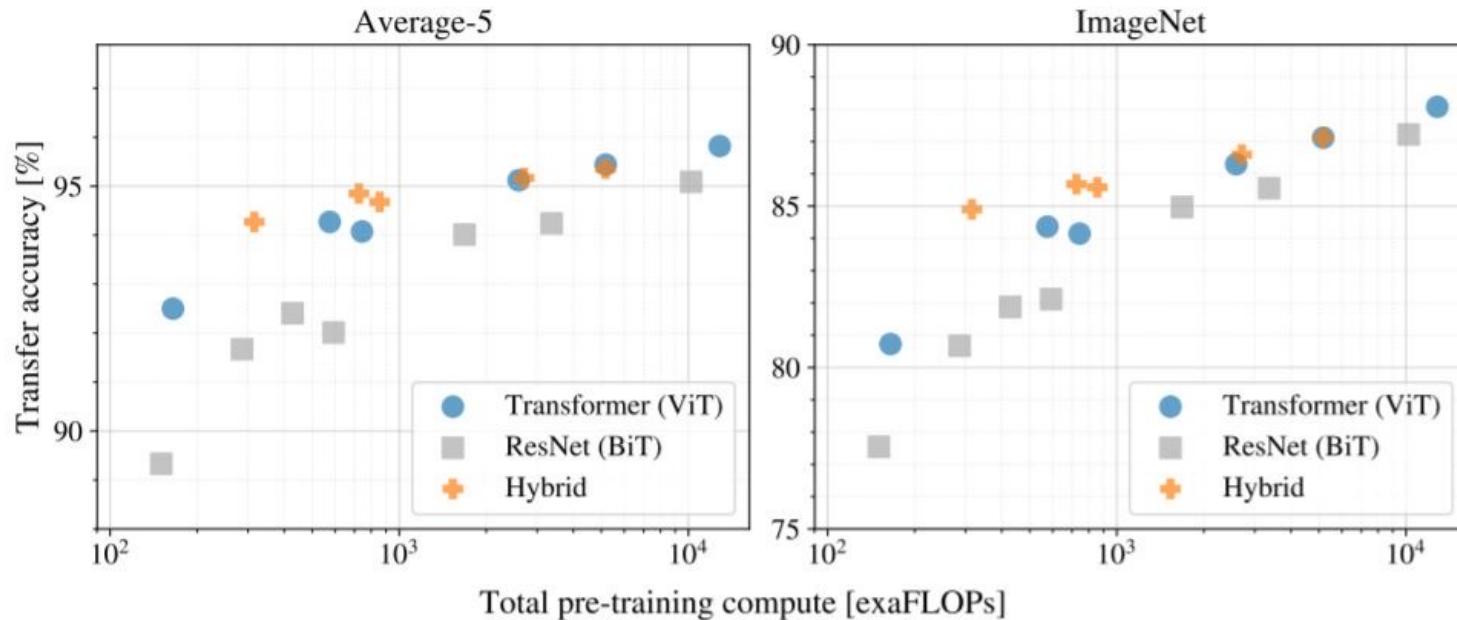


# Vision Transformer (ViT)



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ArXiv 2020

# Vision Transformers vs. ResNets



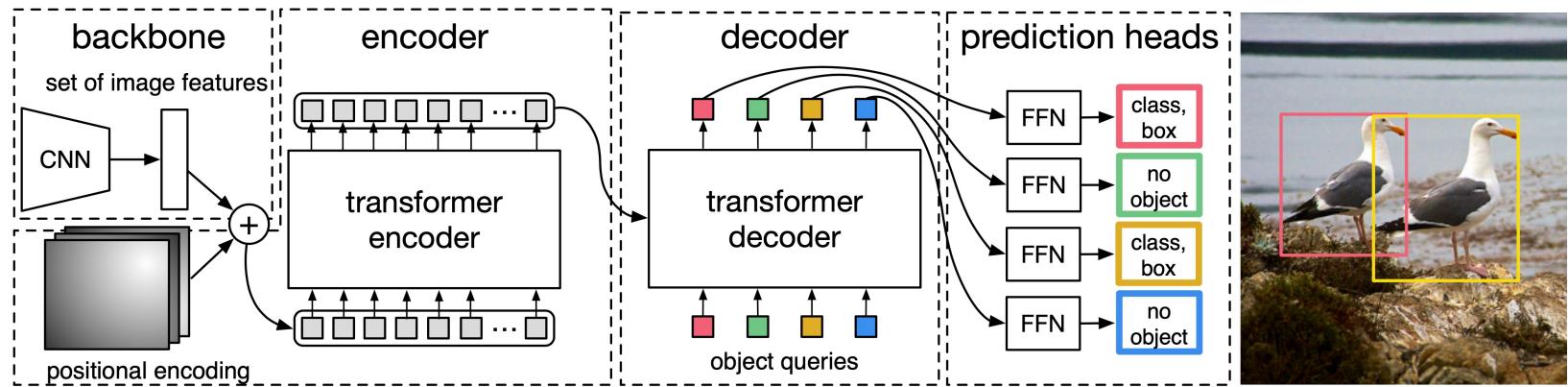
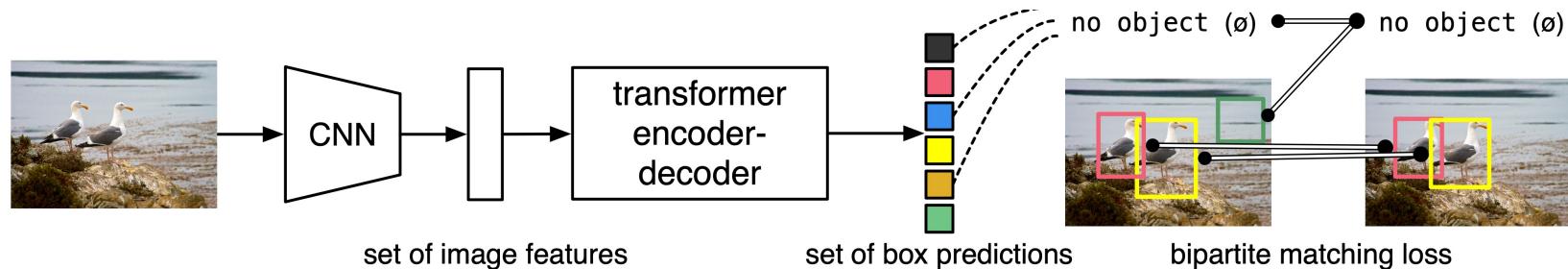
Farewell CNN ?



# Vision Transformers for Object Detection



- DE:TR

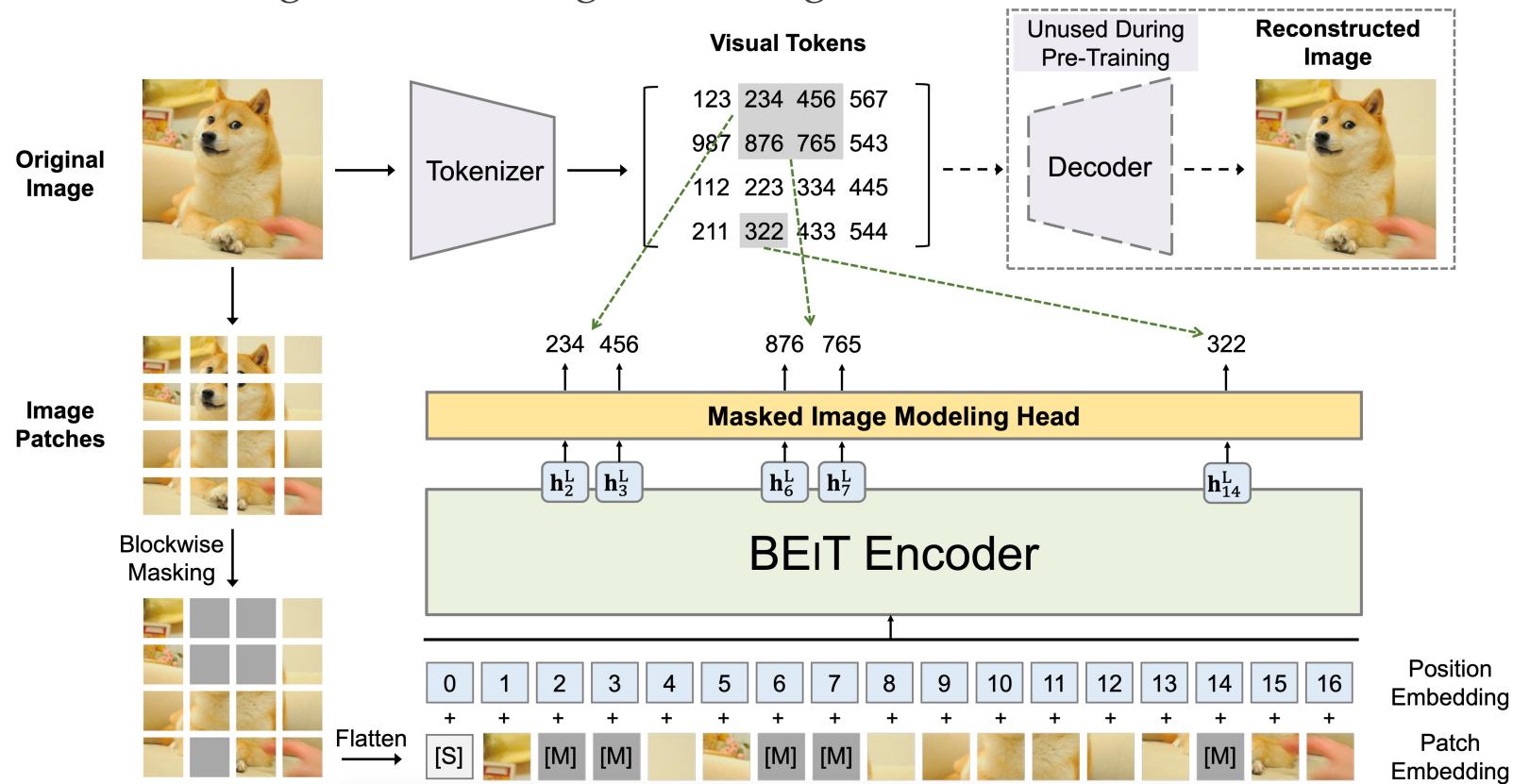


Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020

# BEiT: Pre-Trained Models for Images



- Backbone Network: Standard Transformer
- Pre-Training: Masked Image Modeling



BEiT: BERT Pre-Training of Image Transformers



- State-of-the-art performance on CV tasks

Models	Model Size	Resolution	ImageNet
<i>Training from scratch (i.e., random initialization)</i>			
ViT <sub>384</sub> -B [DBK <sup>+</sup> 20]	86M	384 <sup>2</sup>	77.9
ViT <sub>384</sub> -L [DBK <sup>+</sup> 20]	307M	384 <sup>2</sup>	76.5
DeiT-B [TCD <sup>+</sup> 20]	86M	224 <sup>2</sup>	81.8
DeiT <sub>384</sub> -B [TCD <sup>+</sup> 20]	86M	384 <sup>2</sup>	83.1
<i>Supervised Pre-Training on ImageNet-22K (using labeled data)</i>			
ViT <sub>384</sub> -B [DBK <sup>+</sup> 20]	86M	384 <sup>2</sup>	84.0
ViT <sub>384</sub> -L [DBK <sup>+</sup> 20]	307M	384 <sup>2</sup>	85.2
<i>Self-Supervised Pre-Training on ImageNet-1K (without labeled data)</i>			
iGPT-1.36B <sup>†</sup> [CRC <sup>+</sup> 20]	1.36B	224 <sup>2</sup>	66.5
ViT <sub>384</sub> -B-JFT300M <sup>†</sup> [DBK <sup>+</sup> 20]	86M	384 <sup>2</sup>	79.9
MoCo v3-B [CXH21]	86M	224 <sup>2</sup>	83.2
MoCo v3-L [CXH21]	307M	224 <sup>2</sup>	84.1
DINO-B [CTM <sup>+</sup> 21]	86M	224 <sup>2</sup>	82.8
BEiT-B (ours)	86M	224 <sup>2</sup>	83.2
BEiT <sub>384</sub> -B (ours)	86M	384 <sup>2</sup>	84.6
BEiT-T (ours)	307M	224 <sup>2</sup>	85.2
BEiT <sub>384</sub> -L (ours)	307M	384 <sup>2</sup>	<b>86.3</b>

ADE20K semantic segmentation

Top-1 accuracy on ImageNet-1K

Models	mIoU (%)	Multi-Scale mIoU (%)
<i>Supervised Pre-Training on ImageNet-22K (using labeled data)</i>		
Swin-B [LLC <sup>+</sup> 21]	50.0	51.7
Swin-L [LLC <sup>+</sup> 21]	52.1	53.5
<i>Self-Supervised Pre-Training, and Intermediate Fine-Tuning on ImageNet-22K</i>		
BEiT-B <sup>+</sup> (ours)	53.6	54.2
BEiT-L <sup>+</sup> (ours)	56.7	57.0
<i>Self-Supervised Pre-Training, and Intermediate Fine-Tuning on In-House-70M</i>		
BEiT-L <sup>+</sup> (ours)	<b>57.9</b>	<b>58.4</b>

# What's Next?



## • Unsupervised Learning

物以类聚

### Clustering

Learn hidden structures

- K-means
- Gaussian Mixture

無中生有

### Generative

density estimation

- GAN
- VAE
- Diffusion

化繁為簡

### Dimensional Reduction

Learning data compression

- Principal component analysis
- Auto-encoder

