

Deep Learning

Episode 7

Advanced computer vision



MEMBER?

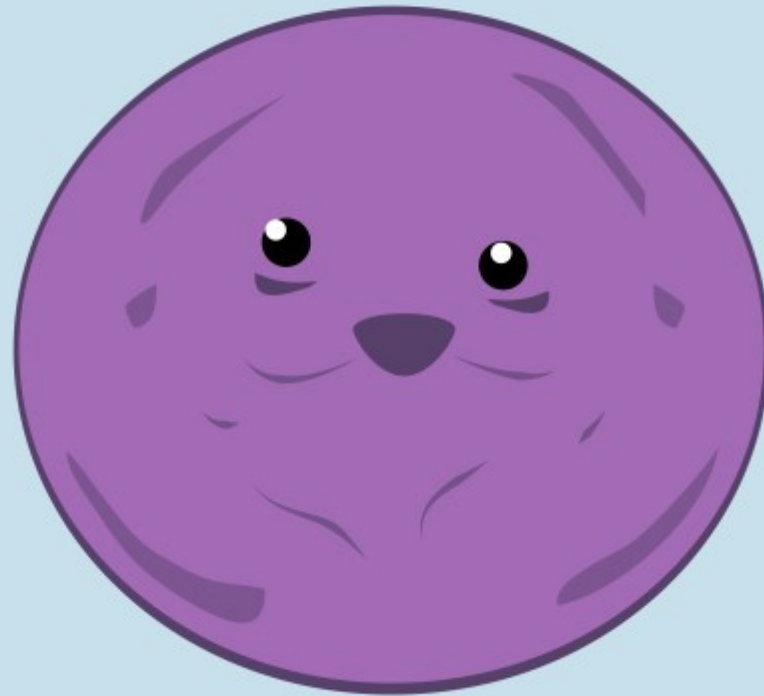


Image recognition



“Dog”

Bounding box regression

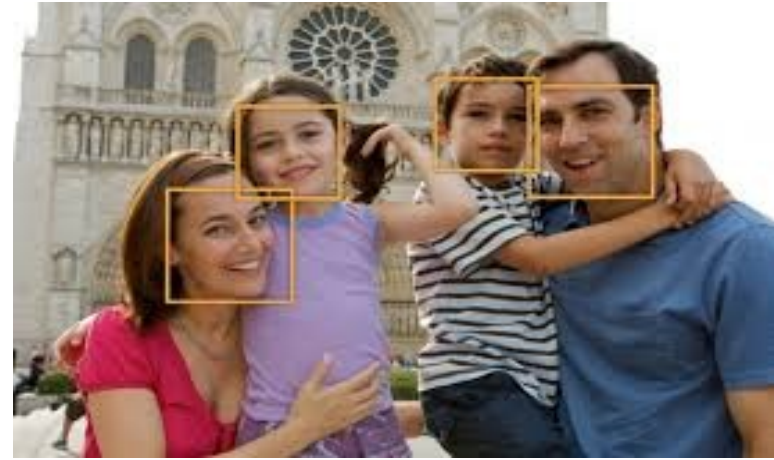
Predict object bounding box

(x_0, y_0, w, h)

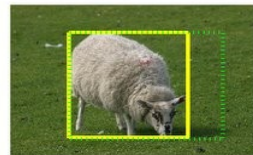
or several bounding boxes for multiple objects.

Applications examples:

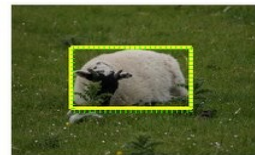
- Face detection @ cameras
- Surveillance cameras
- Self-driving cars



IM:"005194" Conf=0.835223



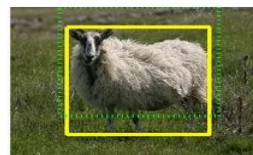
IM:"003538" Conf=0.829488



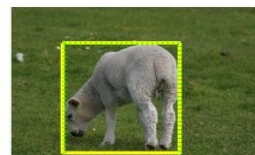
IM:"002810" Conf=0.801748



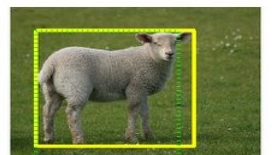
IM:"004522" Conf=0.799045



IM:"001064" Conf=0.797061



IM:"000819" Conf=0.794456



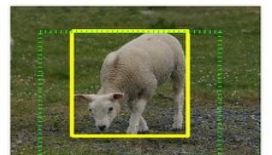
IM:"002306" Conf=0.789123



IM:"001956" Conf=0.788438



IM:"004285" Conf=0.782058

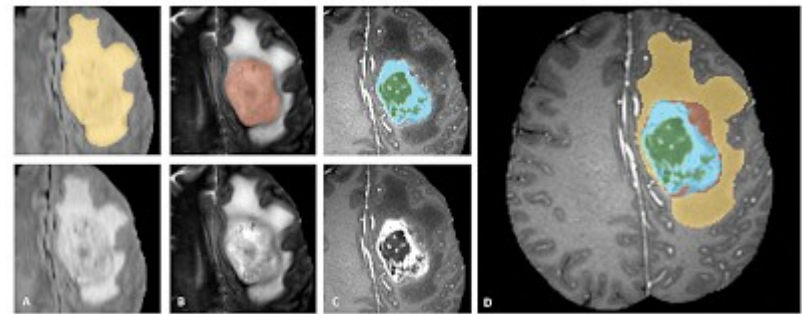
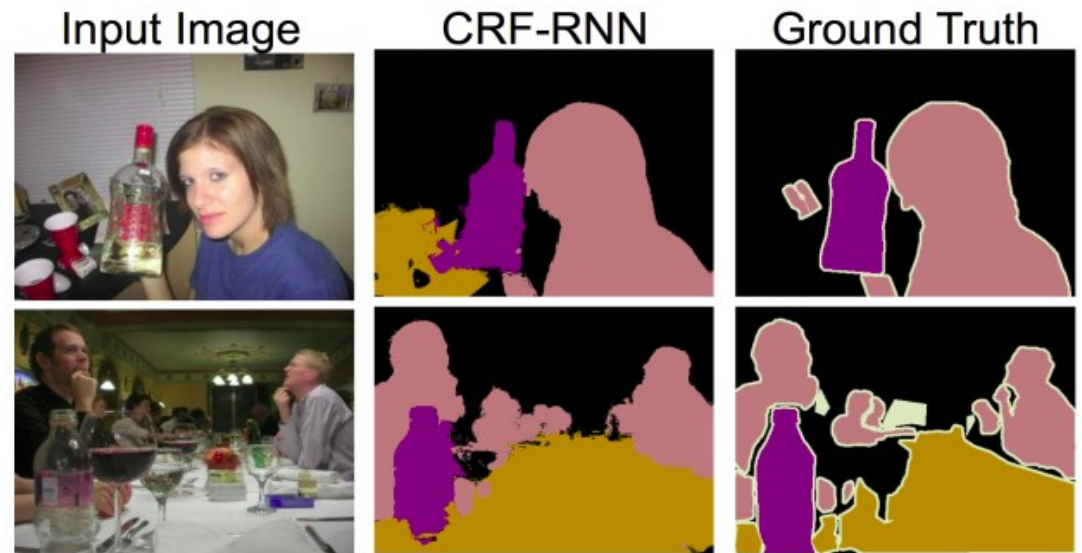


Segmentation

Predict class for each pixel
(fully-convolutional networks)

Applications examples:

- Moar surveillance
- Brain scan labeling
- Map labeling



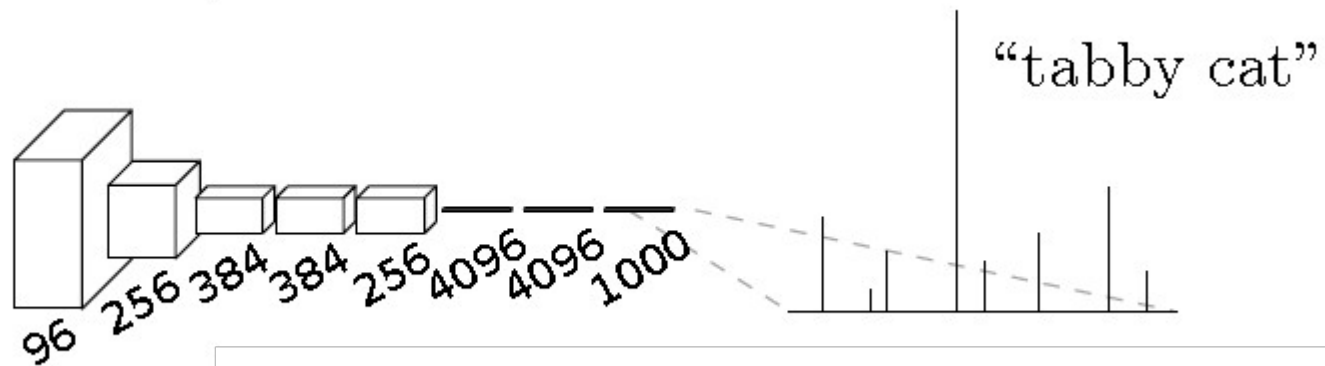
MEMBER WE HAD 10+ WEEKS UNTIL DEADLINE?

OH, I 'MEMBER

imgflip.com

Semantic segmentation

- Classification



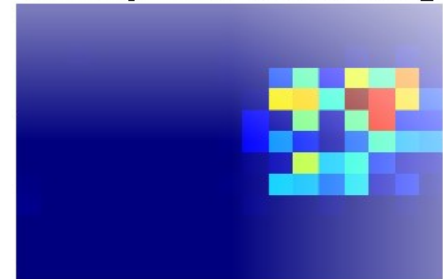
- Segmentation/Detection



???

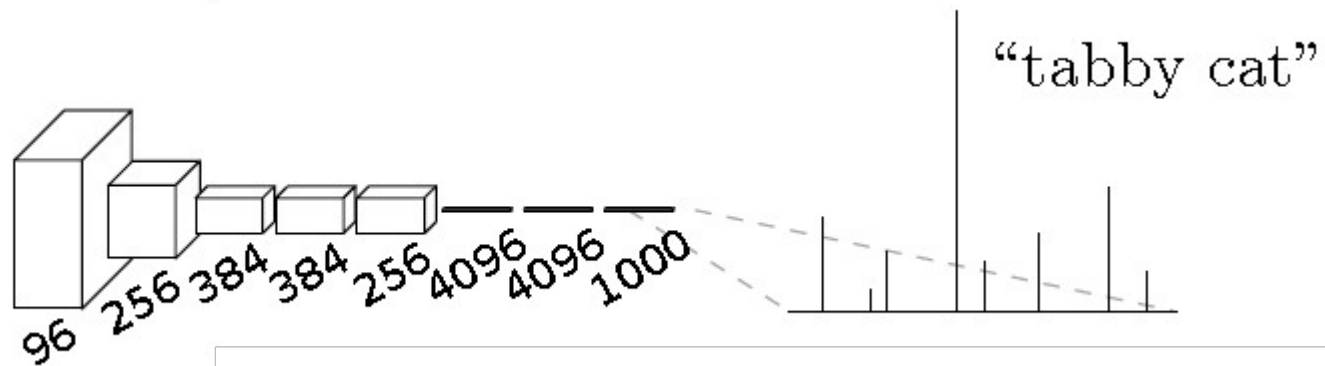
Ideas on how to
predict that?

tabby cat heatmap



Semantic segmentation

- Classification



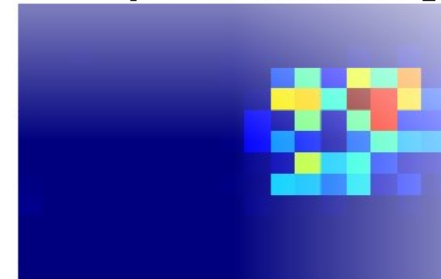
- Segmentation/Detection (naïve)



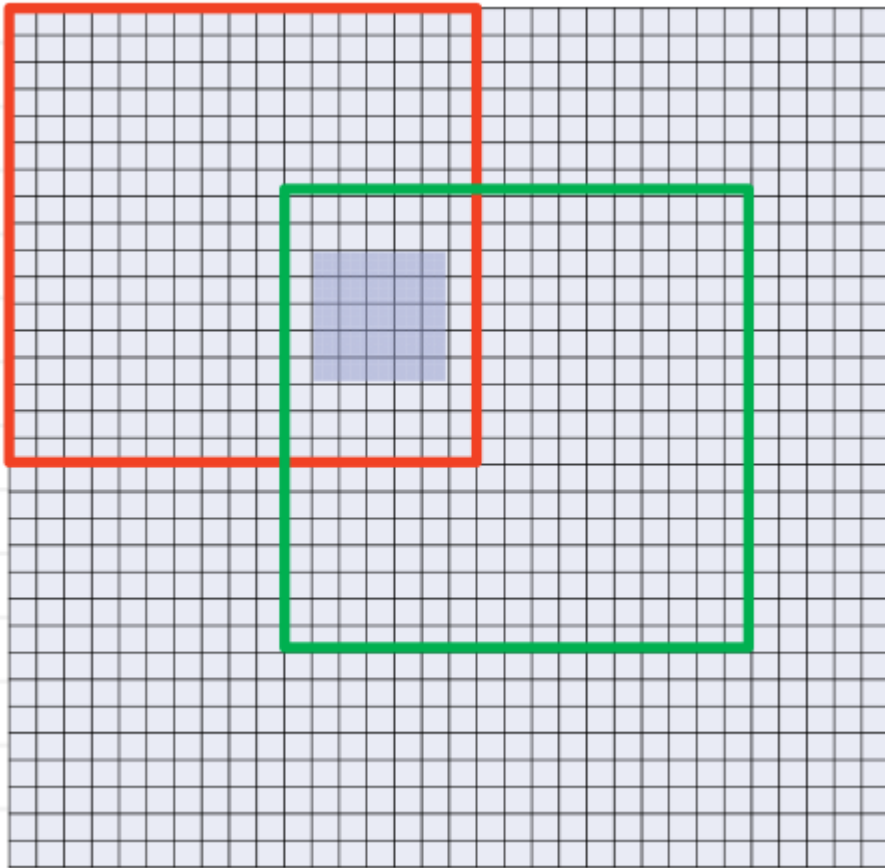
Apply to each spot

“convolution with whole
NN as a filter”

tabby cat heatmap

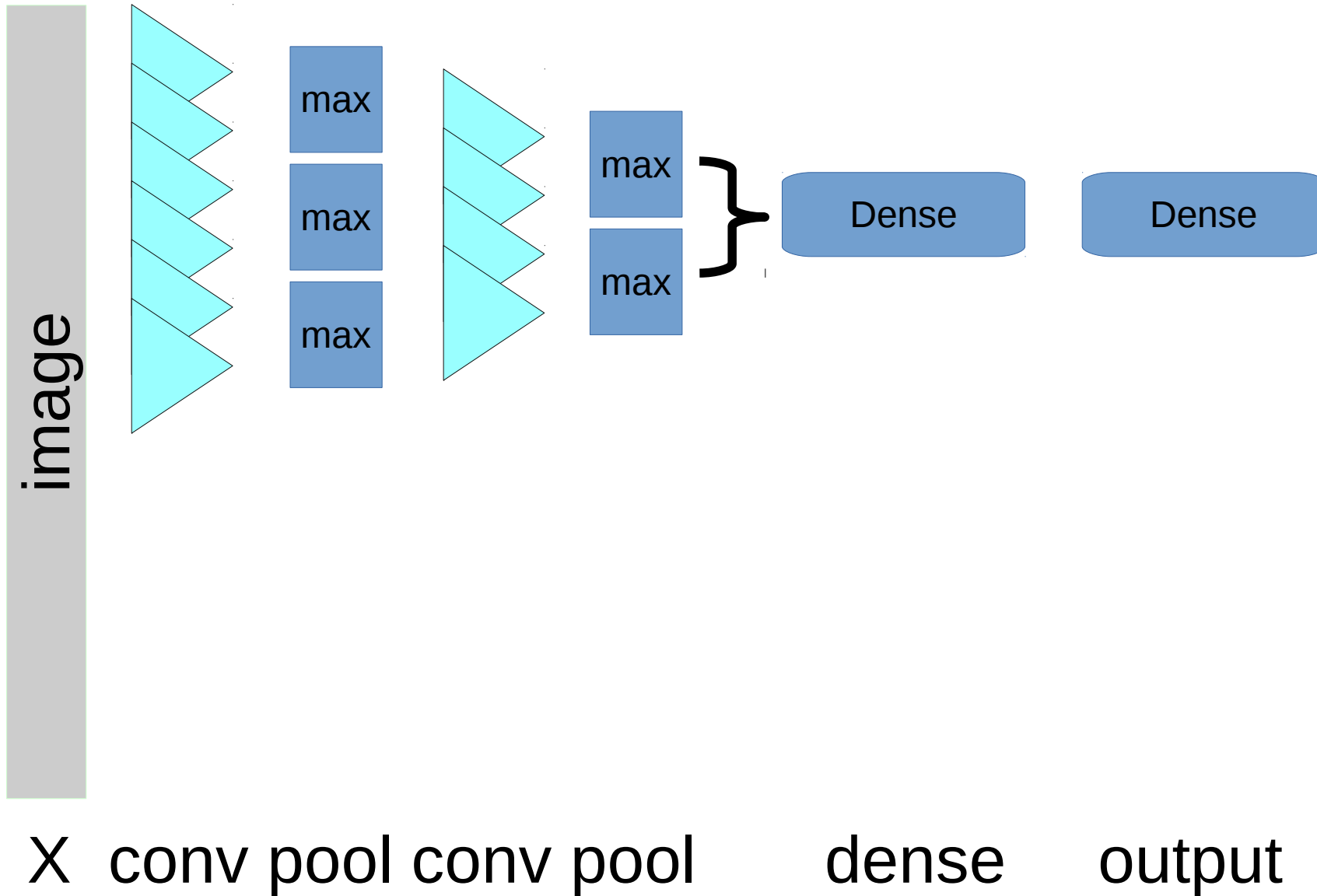


Convolutionalization

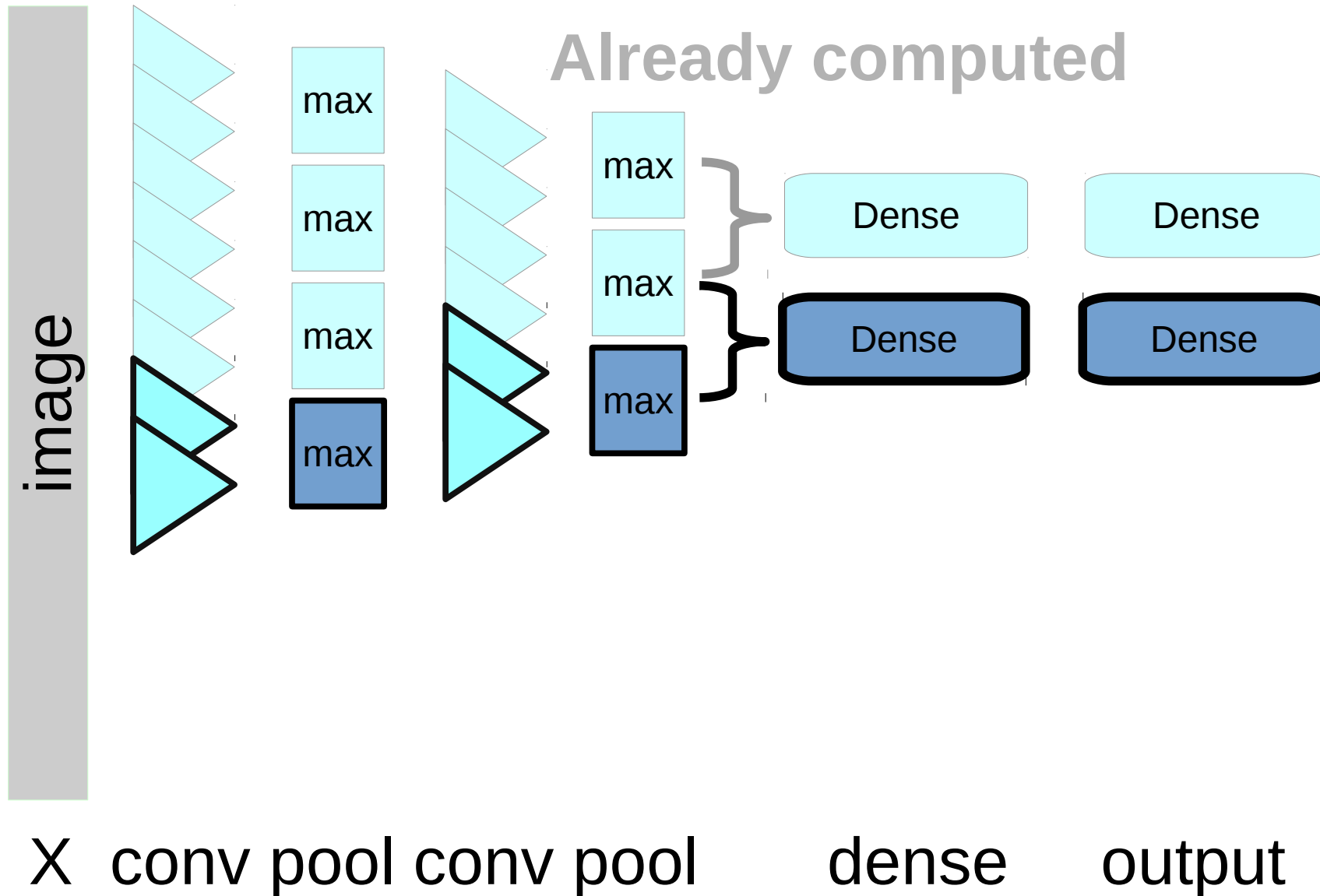


- Idea: if you apply the NN to neighboring patches, you'll have to compute same convolutions again
- Instead let us precompute convolutions for the whole image in advance.

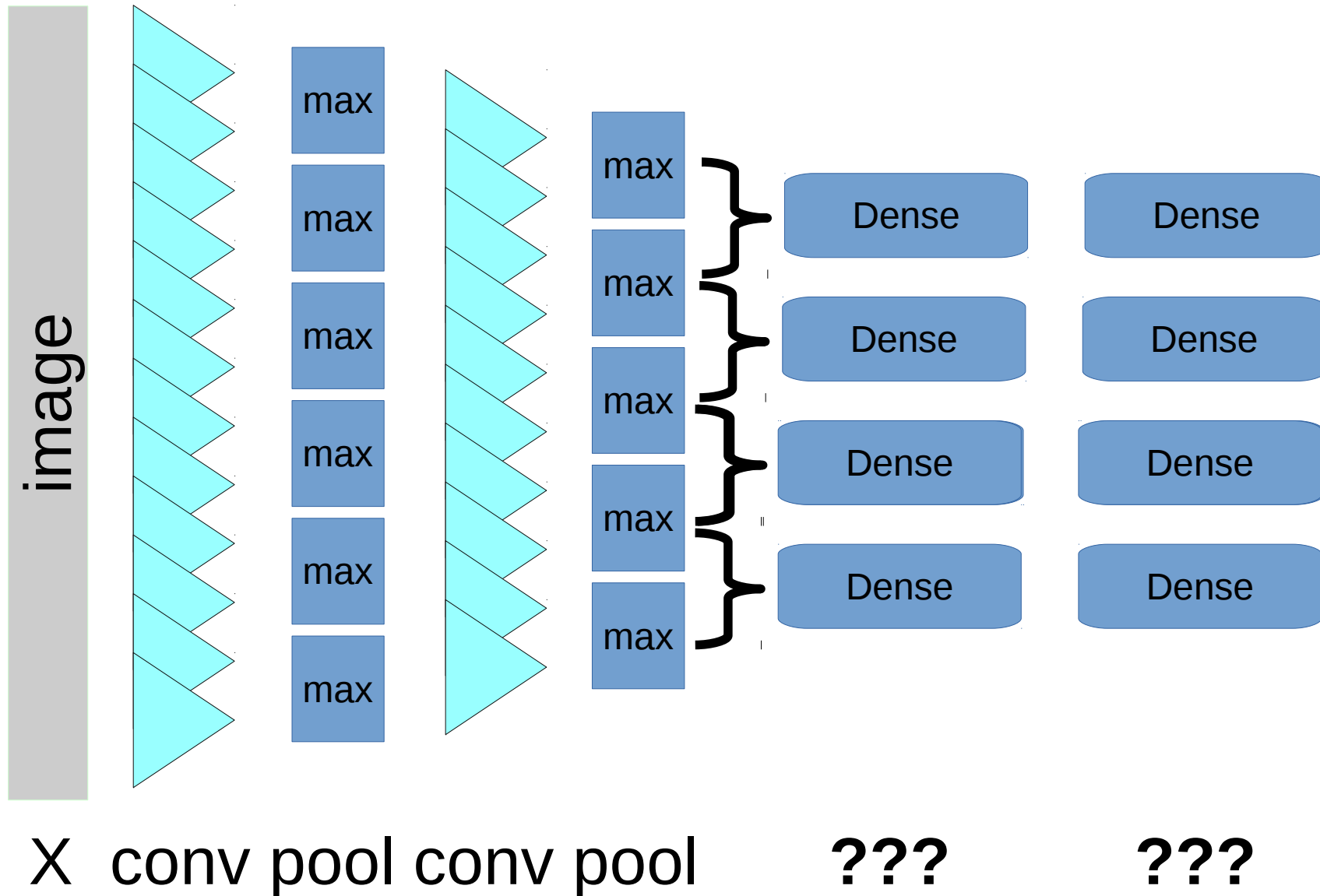
Convolutionalization



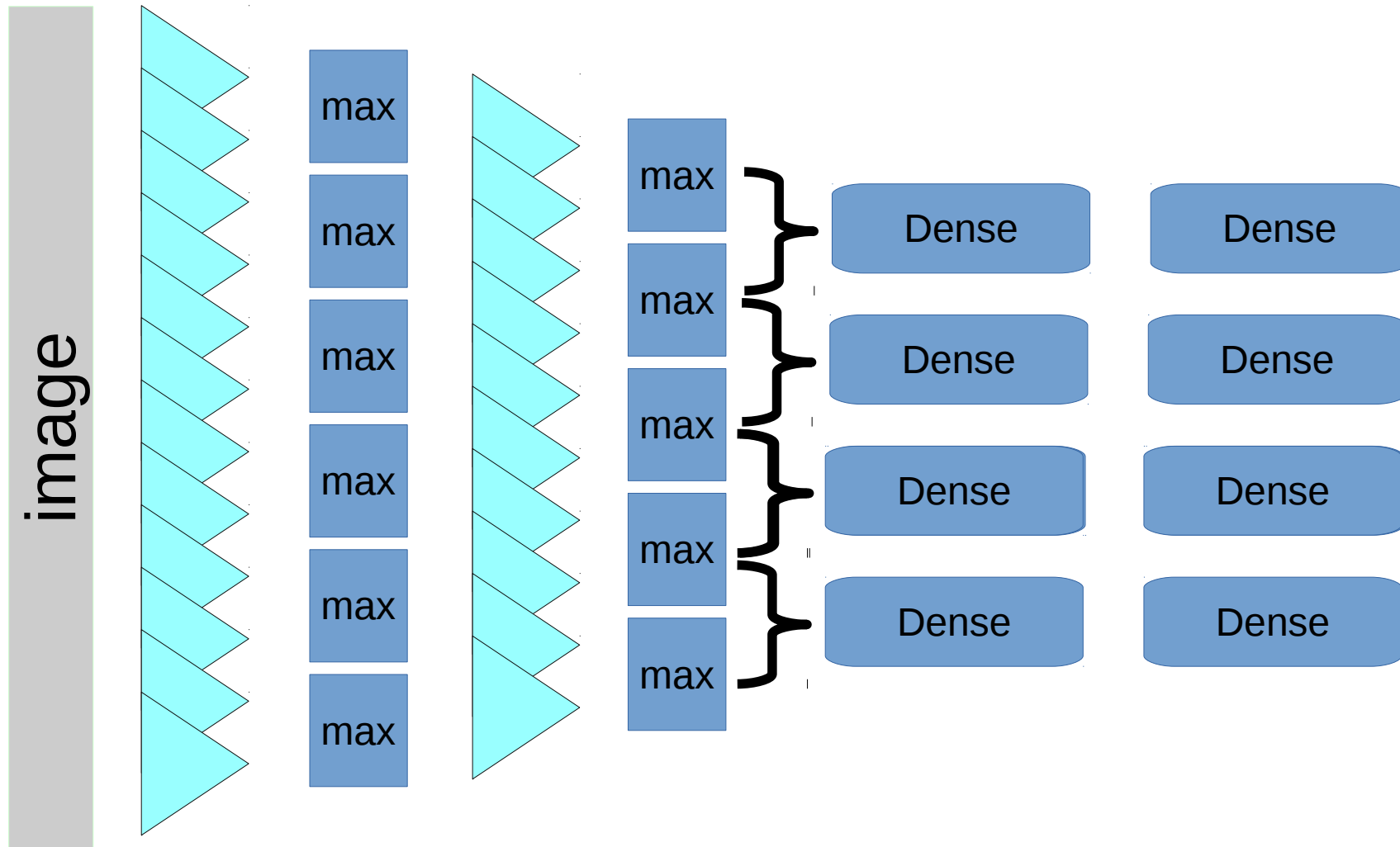
Convolutionalization



Convolutionalization



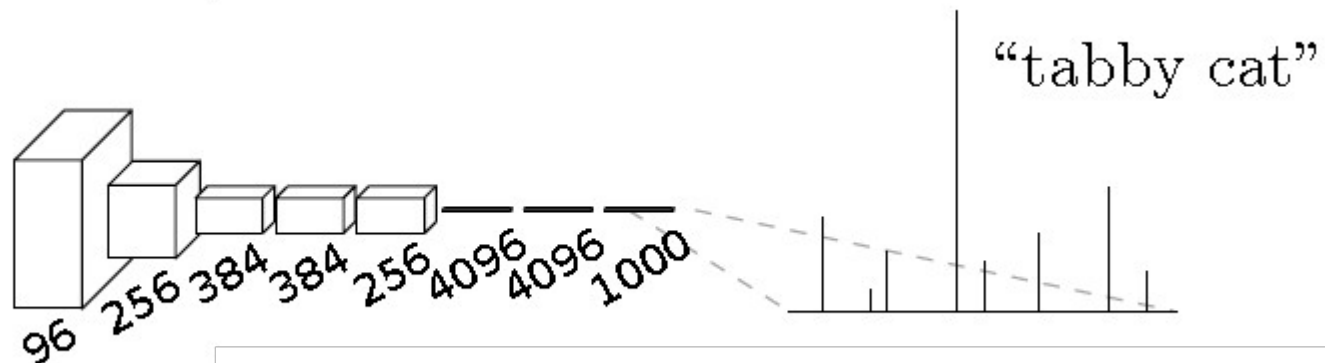
Convolutionalization



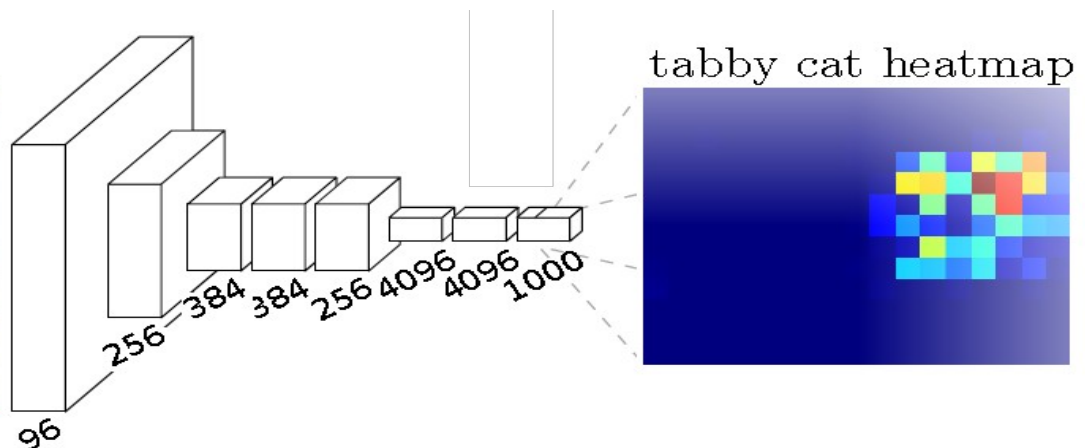
X conv pool conv pool conv2x2 conv1x1

Semantic segmentation

- Classification



- Segmentation/Detection



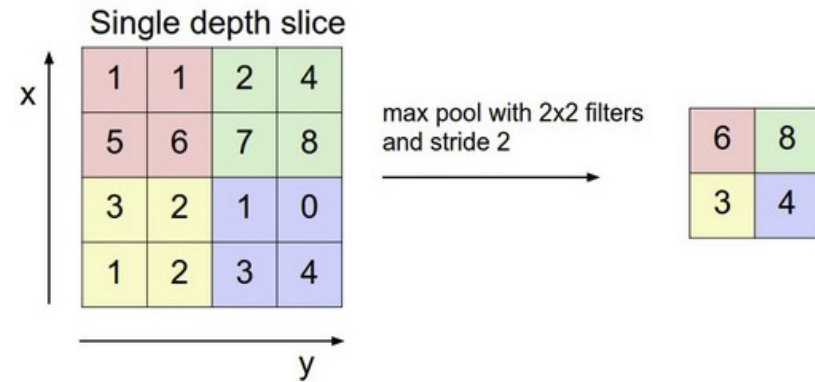
What if output must have same resolution as in the input image? Or even bigger?

How can we **increase** activation size?

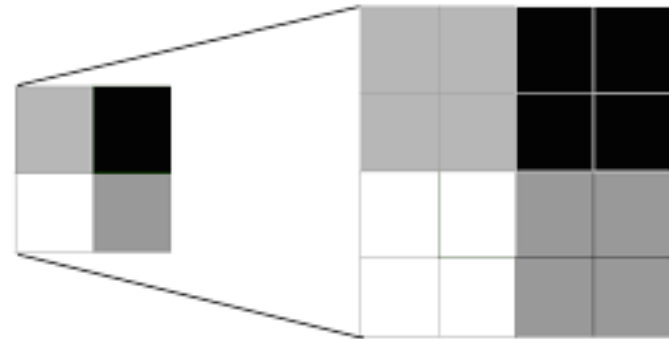
Ideas?

Upscale

- Pooling =
 - take $(n*m)$ region,
 - output one / channel

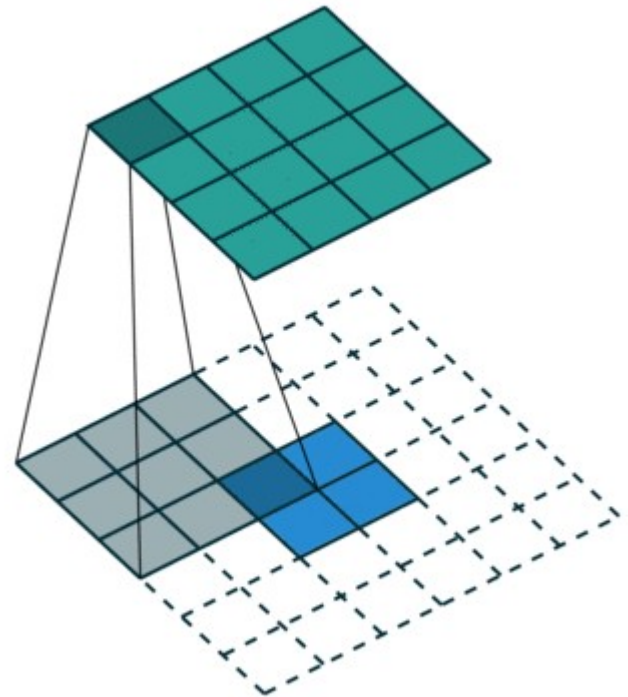


- Upscaling / Unpooling / Upsampling
 - Take one value
 - Output $(n*m)$ region
 - Several strategies



Deconvolution

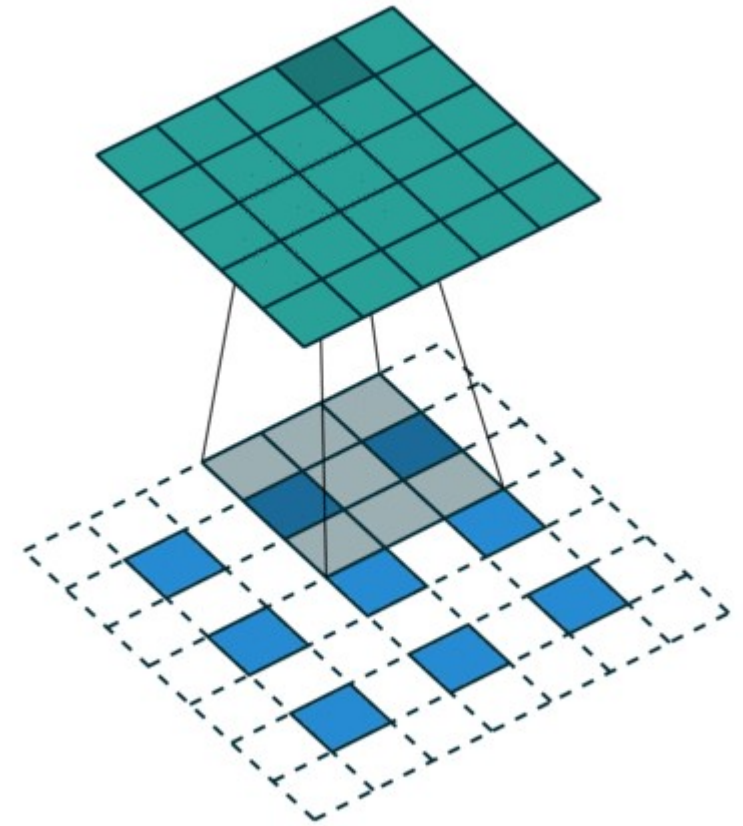
- The “inverse” of convolution
 - = transposed convolution
 - ~ wide convolution
- `deconv(conv(img))`
- preserves image shape
- Stride ~ upscale



deconvolution

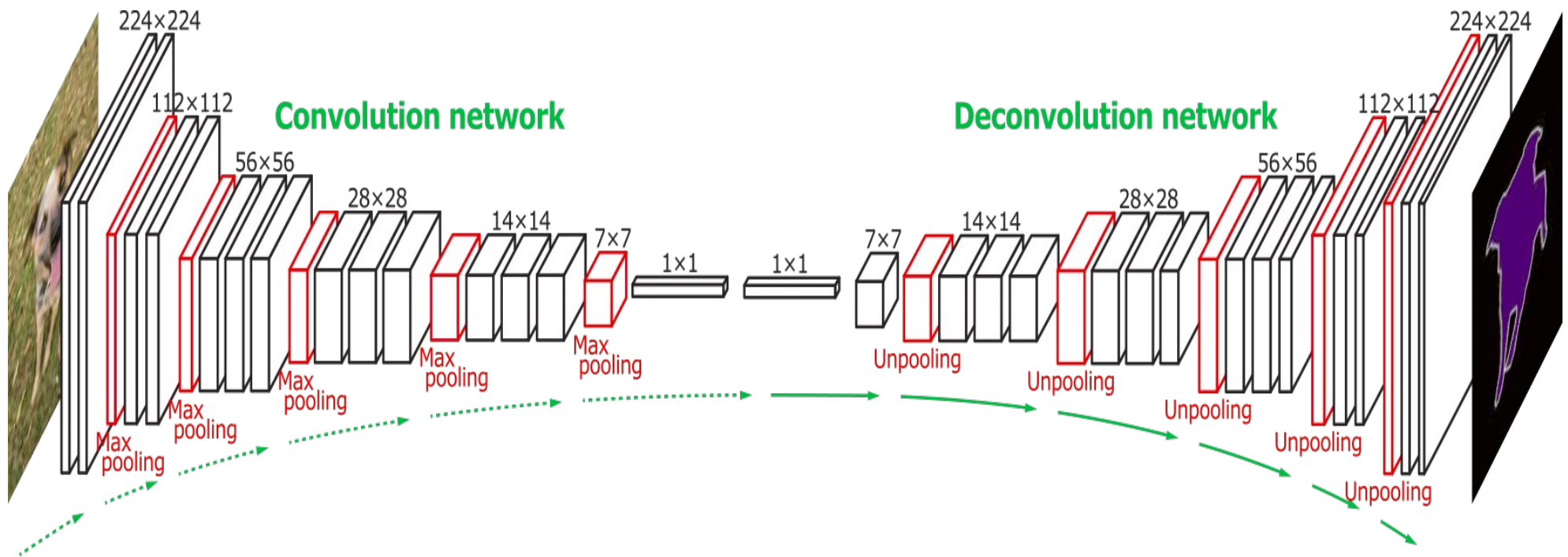
Deconvolution

- The “inverse” of convolution
 - = transposed convolution
 - ~ wide convolution
- `deconv(conv(img))`
- preserves image shape
- Stride ~ upscale



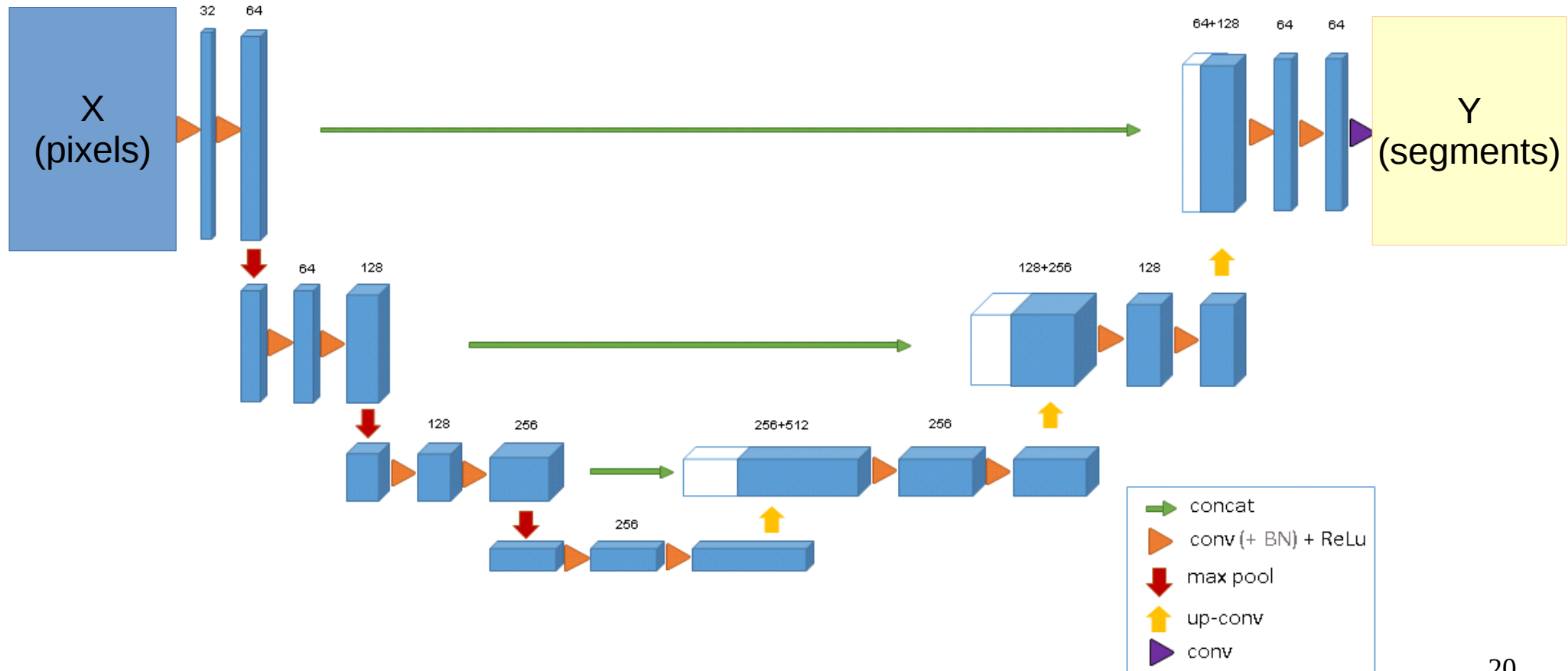
strided deconvolution

Fully-convolutional

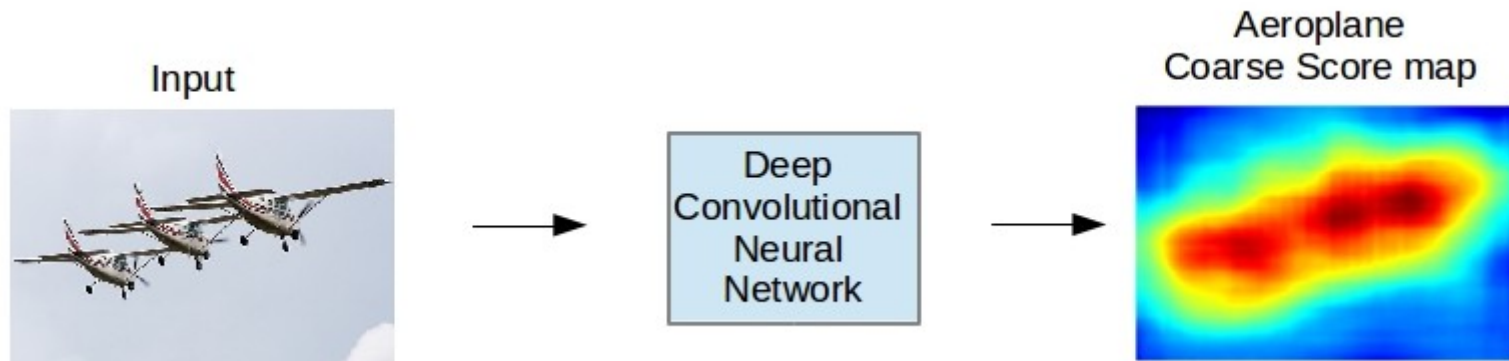


The U-net

- Add connections between layers of similar abstraction
- The idea is similar to gaussian/laplacian pyramids



Raw network output is usually rather blurred

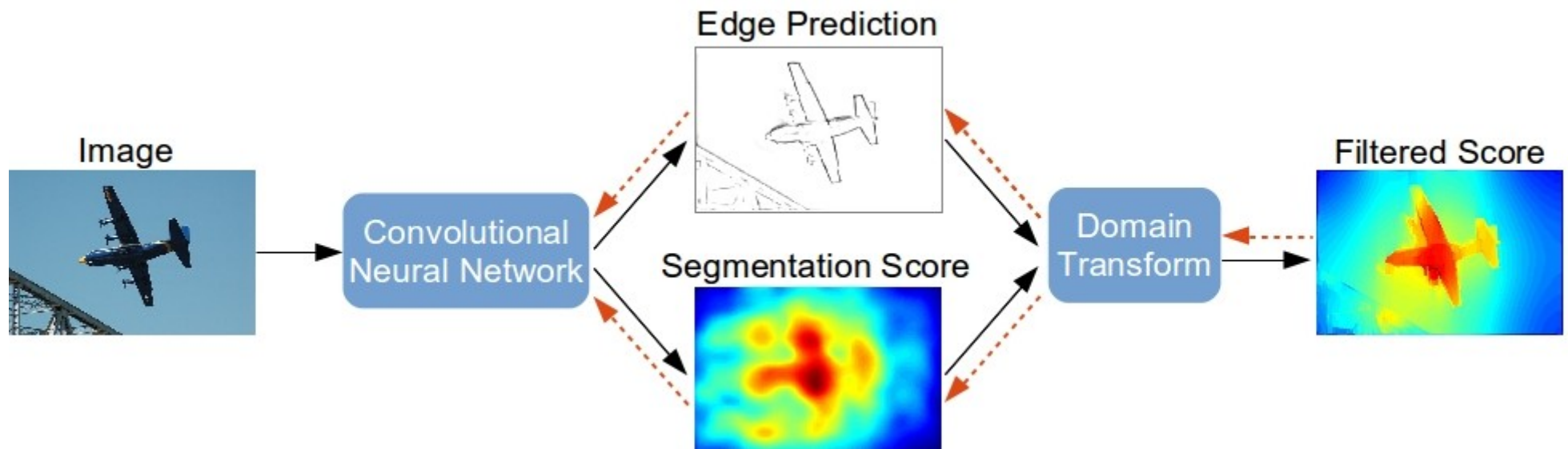


How can we get sharper edges?



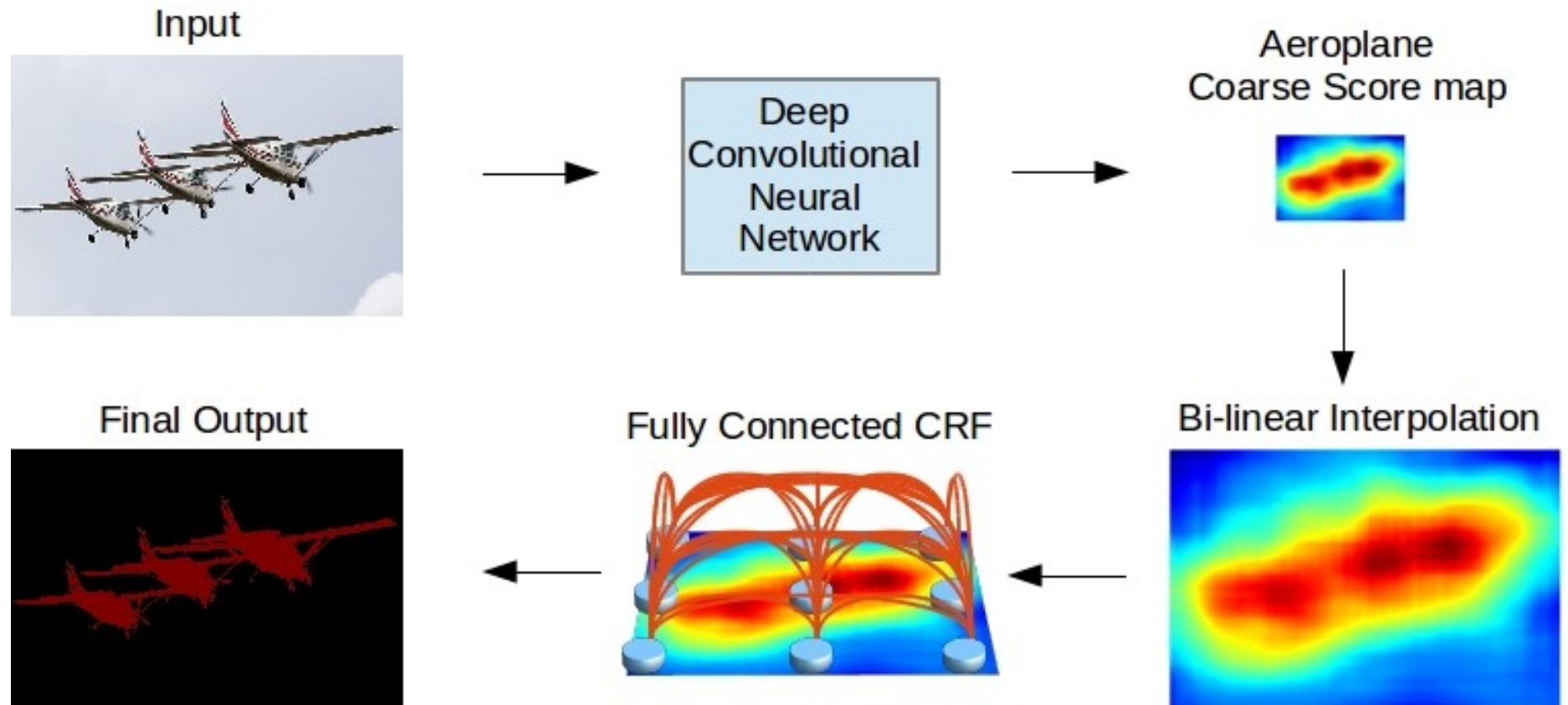
Segmentation: Edge prediction

- Train both class predictor and edge predictor



Segmentation: Graph models

- Tldr incentivize predicting same class over same texture / different classes over edges

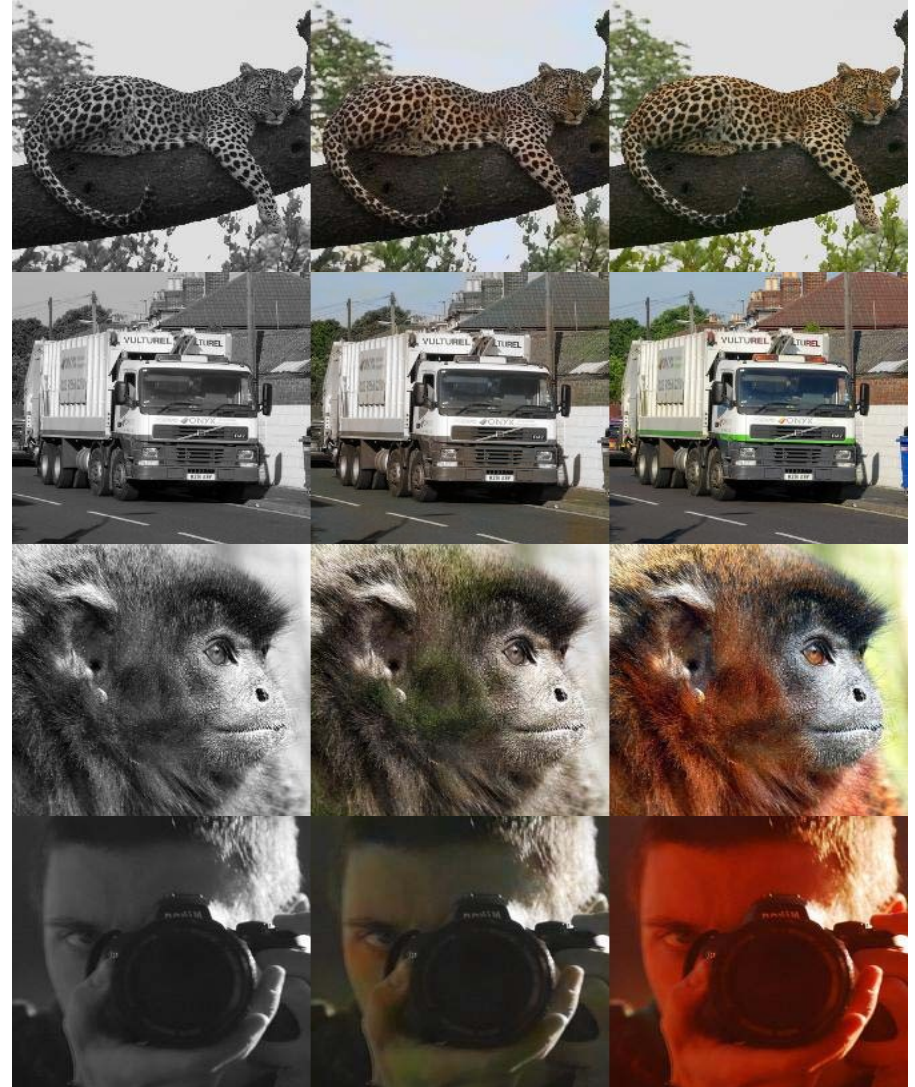




More on that in your Graphical Models course

Image coloring

- Input:
 - gray/sepia image
 - mb neighboring frames
- Output:
 - RGB image
 - Same size as input
- Ideas?
 - Data?
 - Loss?
 - Architecture?



Validation X / prediction / y

Image coloring: hypercolumns

- Idea: upscale all layers to image size and pixel-wise
- Predict 2 UV channels • Grayscale & UV \rightarrow RGB

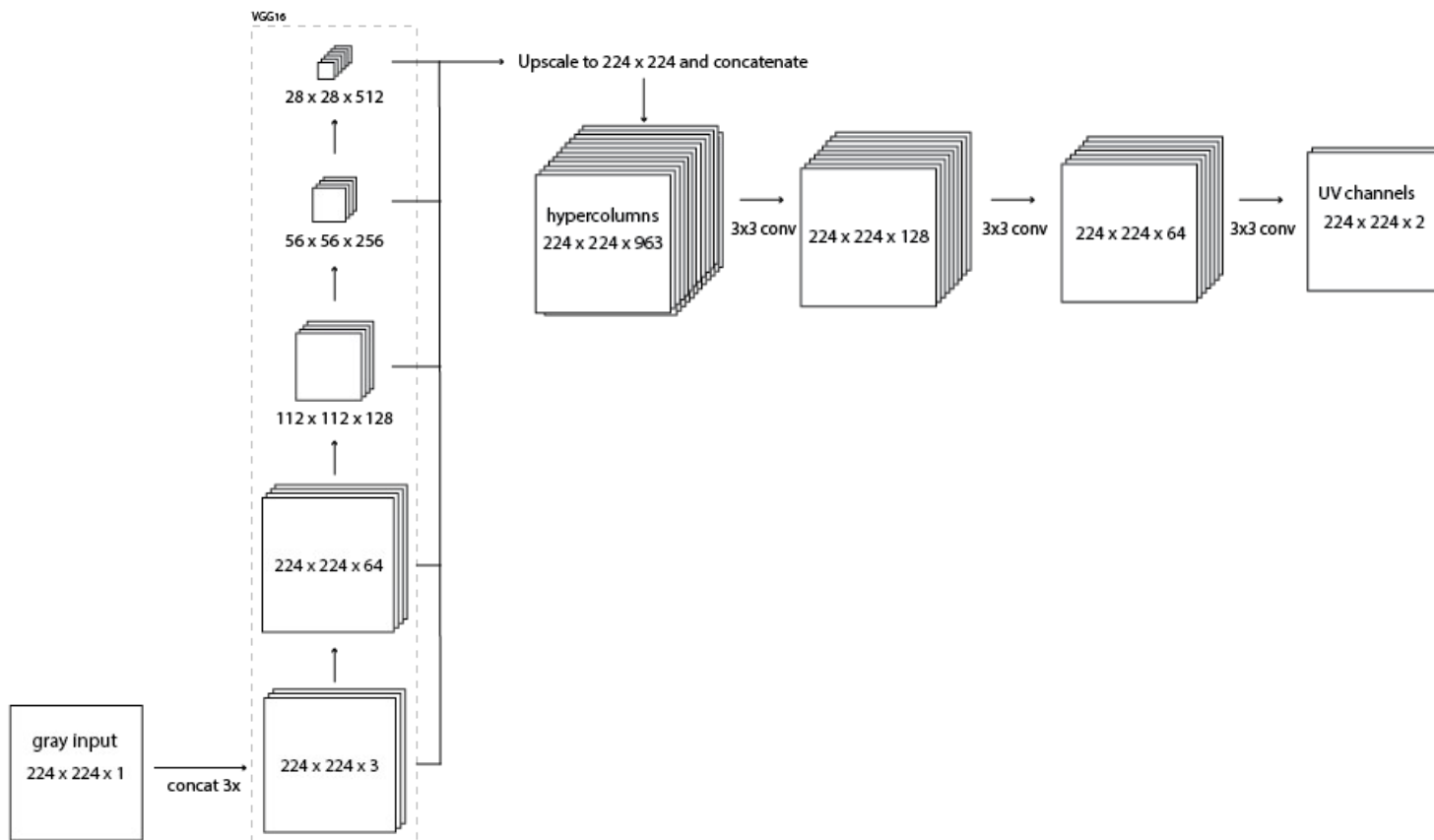
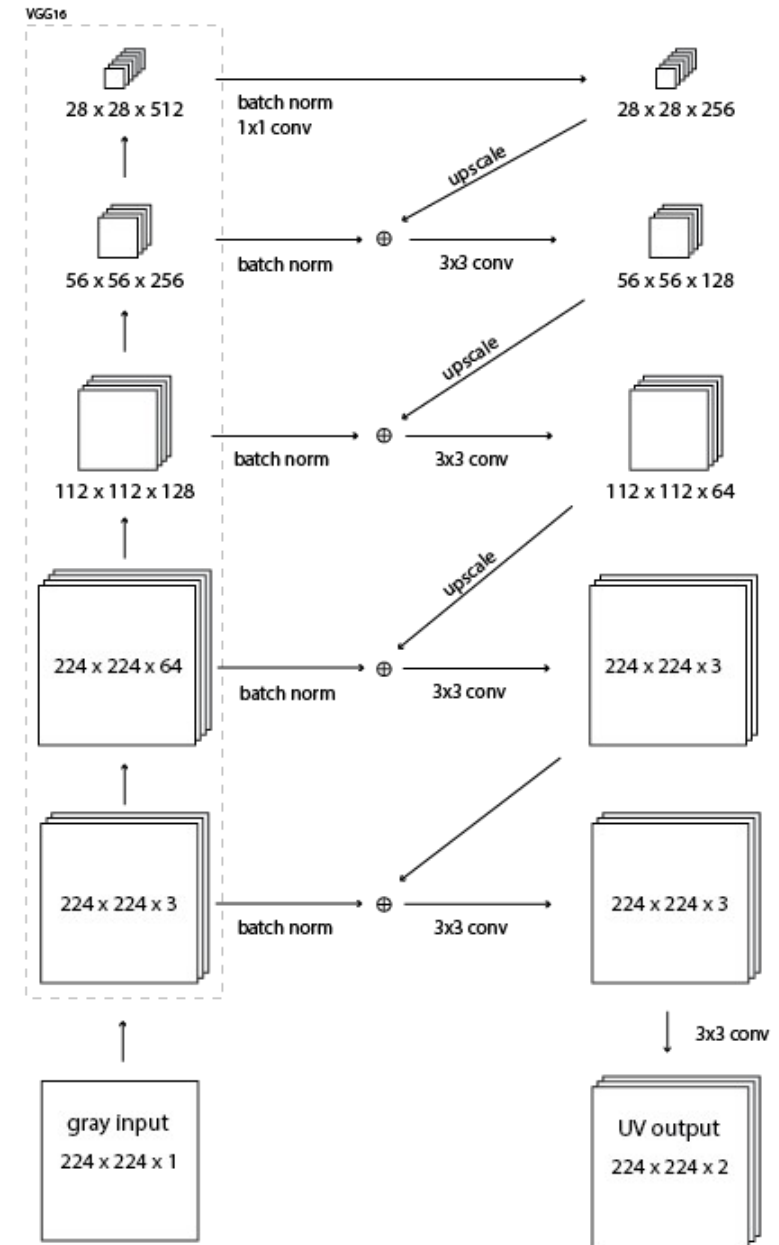
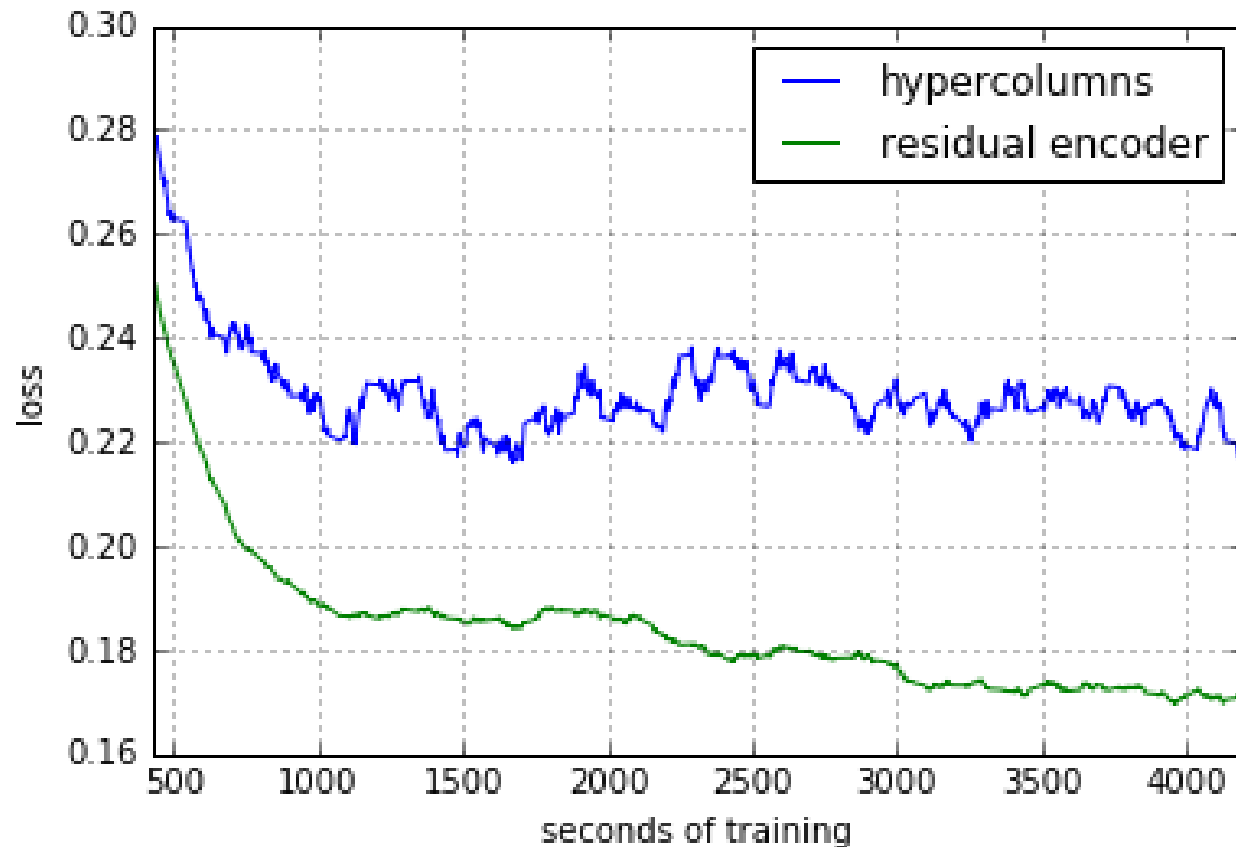


Image coloring: residual

- Similar to U-net
- Add instead of concat
 - Like ResNet
- Consumes less memory
 - No $W \times H \times 1k$ blob

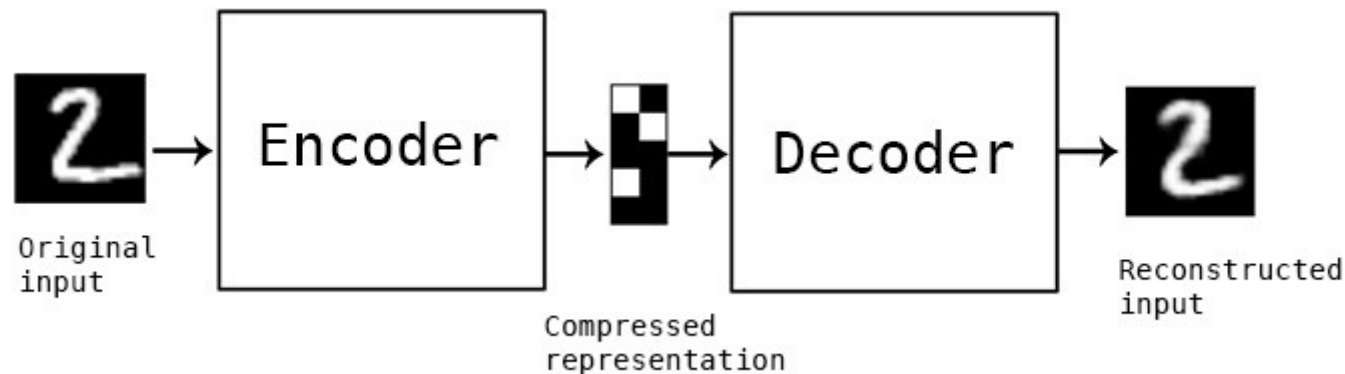


Comparing these 2



Autoencoders, briefly

- image \rightarrow code \rightarrow image
- Previous tricks with architecture still valid
- A few super cool task-specific hacks

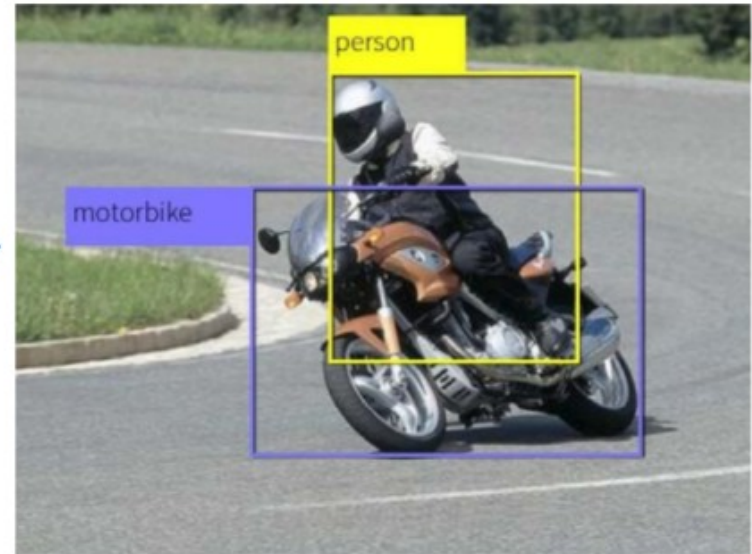


More on that next week

Bounding Box Regression



Input image

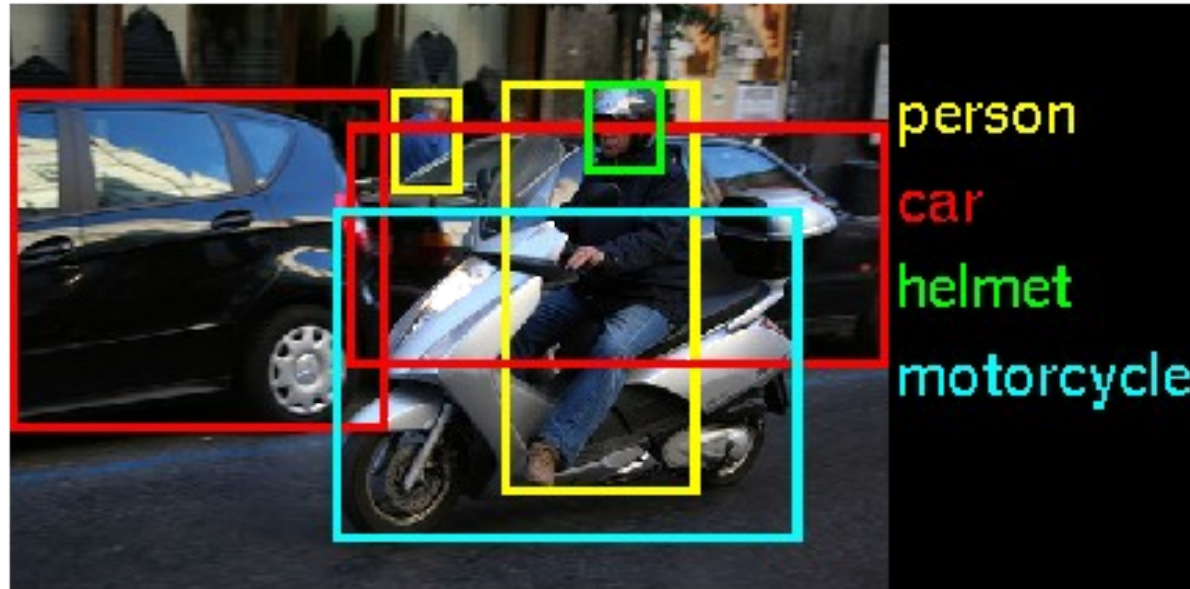


Desired output

- Predict bounding rectangle for the object
- How do we deal with a single square object out of 100 possible objects?

Ideas?

Bounding Box Regression

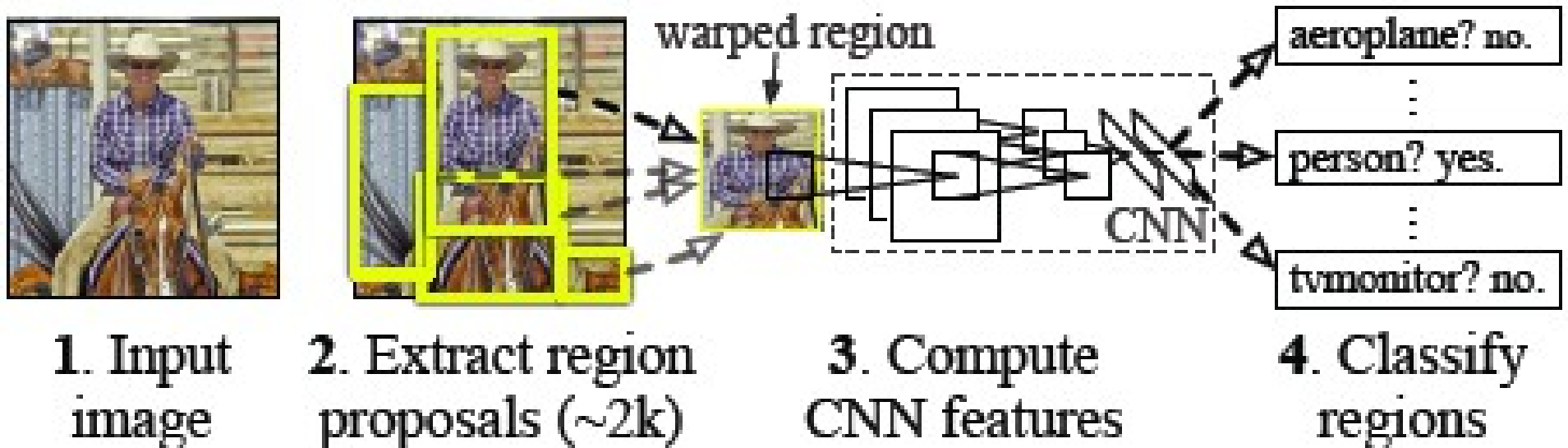


- Multiple squares?
- Non-square shape?
- Arbitrary rotation?

RCNN

- Generate candidates
- Adjust region shapes
- Apply classifier to each region

R-CNN: *Regions with CNN features*



Bounding Box Adjustment

- Learn to predict how to adjust bounding box
(dx, dy, dw, dh)
- Train on random candidates / RCNN output
- Iteratively adjust until good enough
- Learn the adjustment **policy**

Training image regions



Cached region features



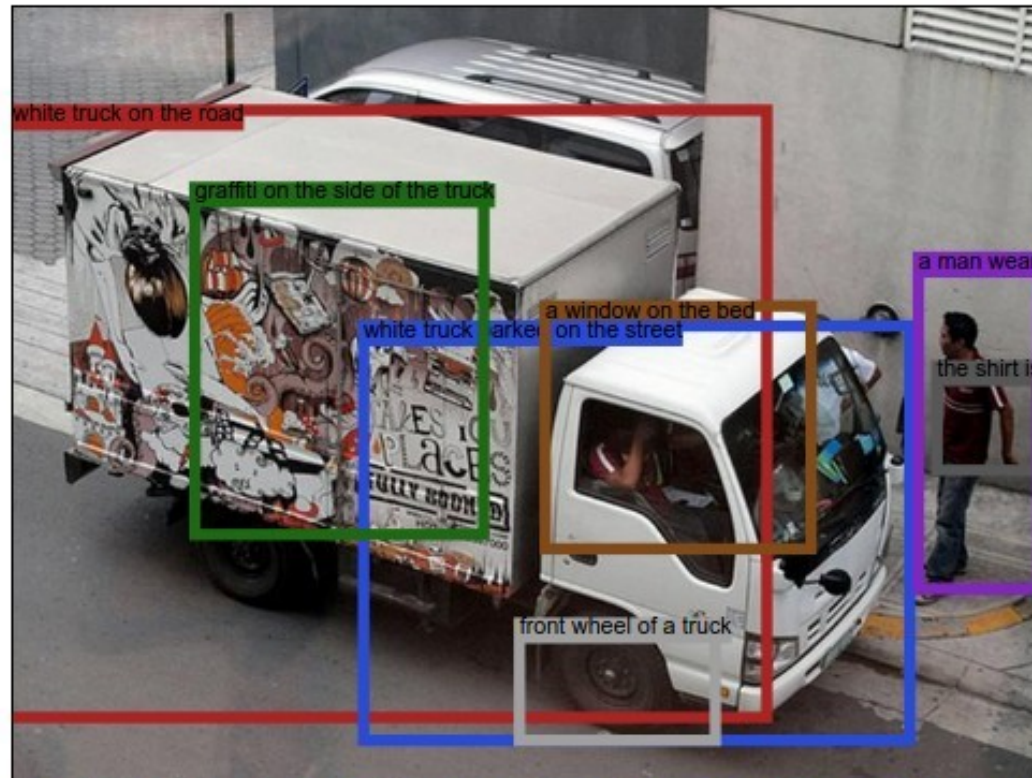
Regression targets
(dx, dy, dw, dh)
Normalized coordinates

(0, 0, 0, 0)
Proposal is good

(.25, 0, 0, 0)
Proposal too
far to left

(0, 0, -0.125, 0)
Proposal too
wide

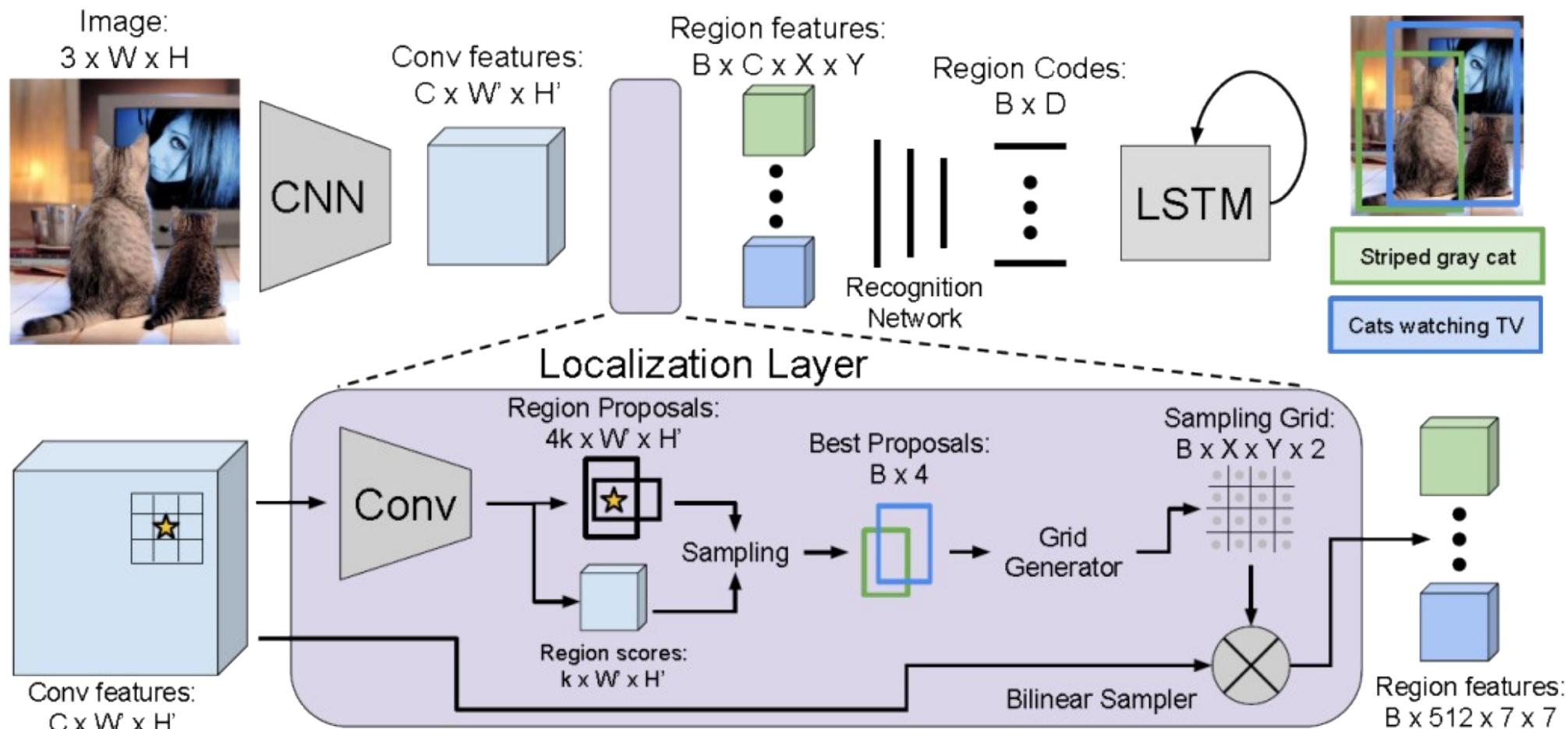
Image description / captioning



a white truck on the road. white truck parked on the street. the shirt is red. graffiti on the side of the truck. a window on the bed. a man wearing a black shirt. front wheel of a truck.

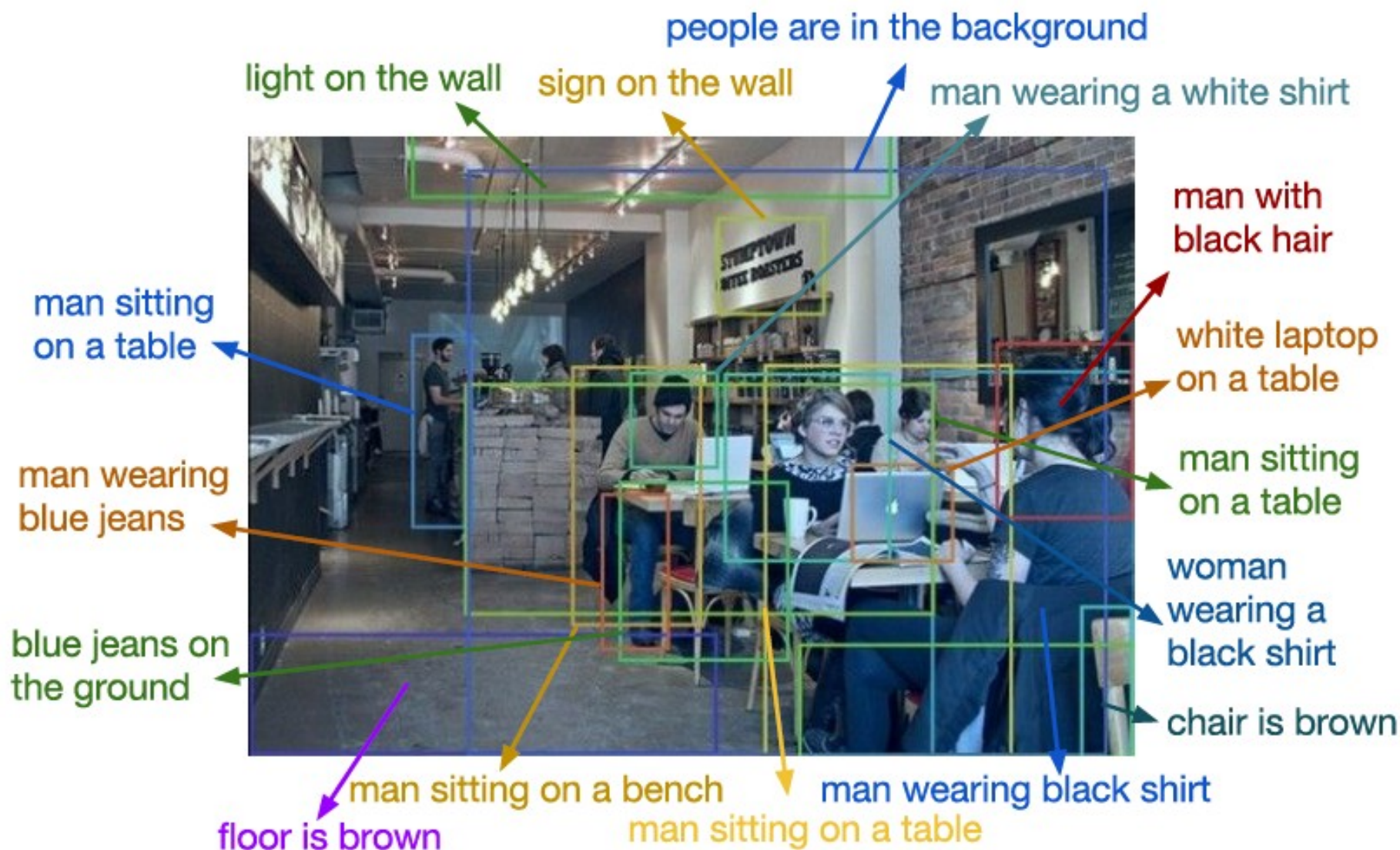
How could we approach this one?

Image description / captioning



Combines EVERYTHING we learned so far

Image description / captioning

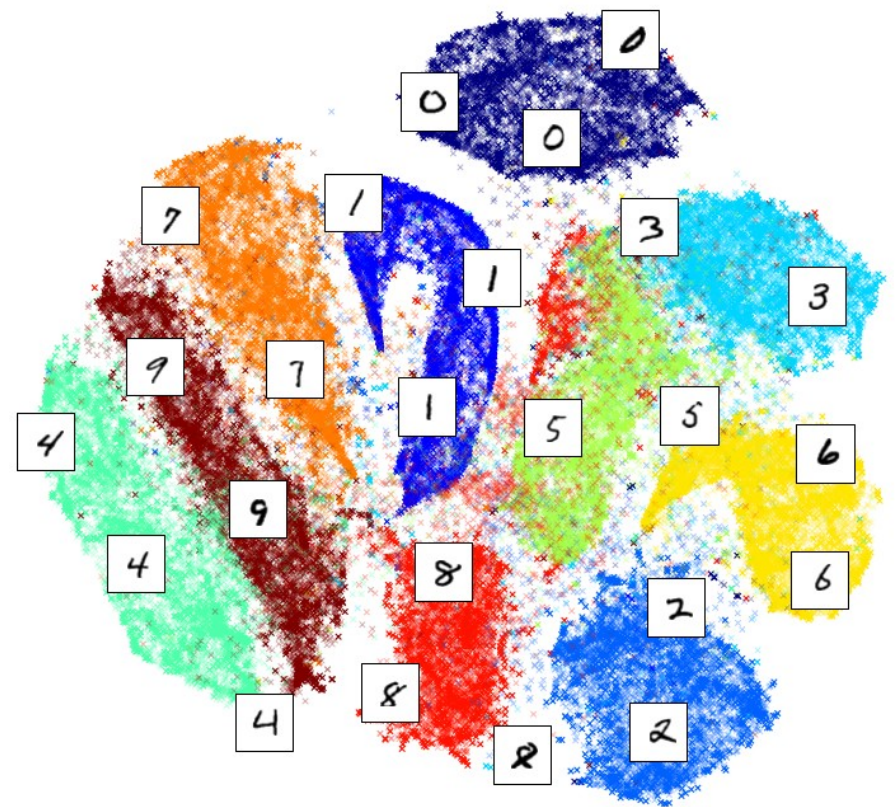


Representation learning

- Layer activations = representations
- Word2vec = word representation
- CNN activation = image2vec
- Image captioning = common space for images and text

MNIST dataset

Two-dimensional embedding of 70,000 handwritten digits with t-SNE

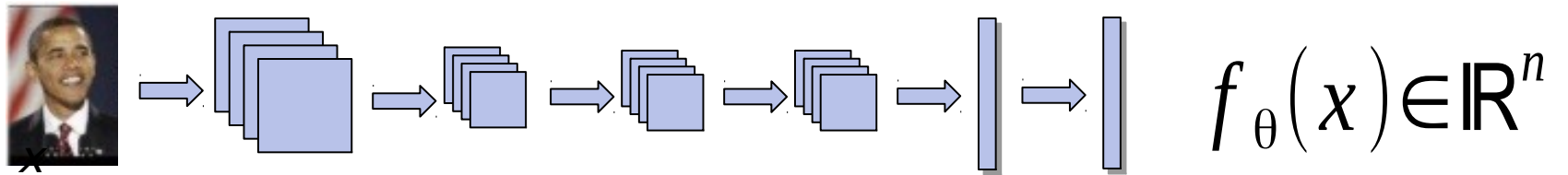


Metric learning

- Deliberately train representations that satisfy some properties
- **Word2vec:** words that appear in similar context should have similar vectors;
- **Image verification:** photos of the same person should have similar vectors, different persons = farther vectors;
- **Image captioning (alternative):** image must be close to it's correct descriptions, far from incorrect ones;
- One way: close \sim cosine, far $\sim 1 - \text{cosine}$

Image verification

- Image2vec



- We want:

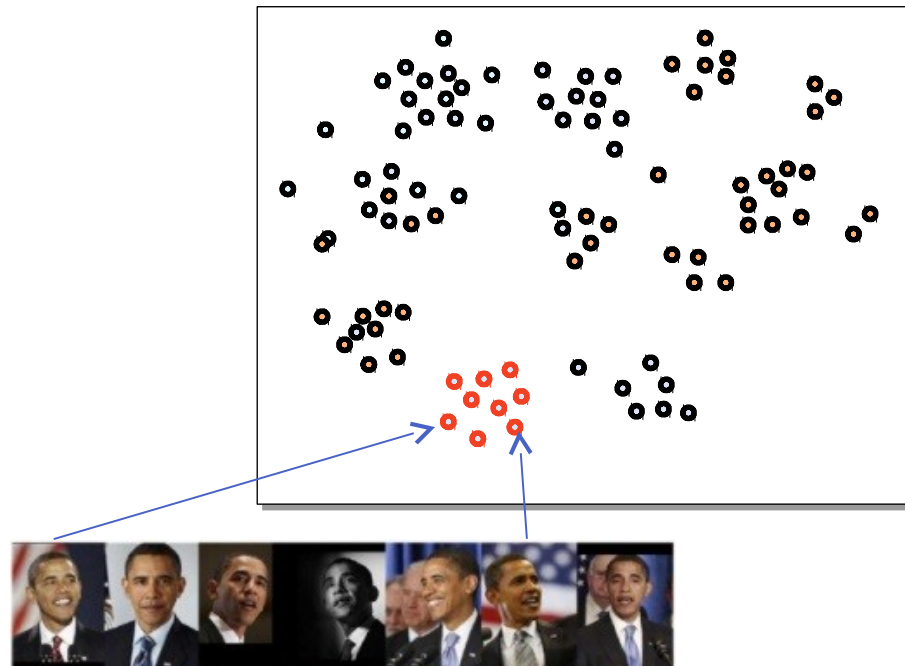
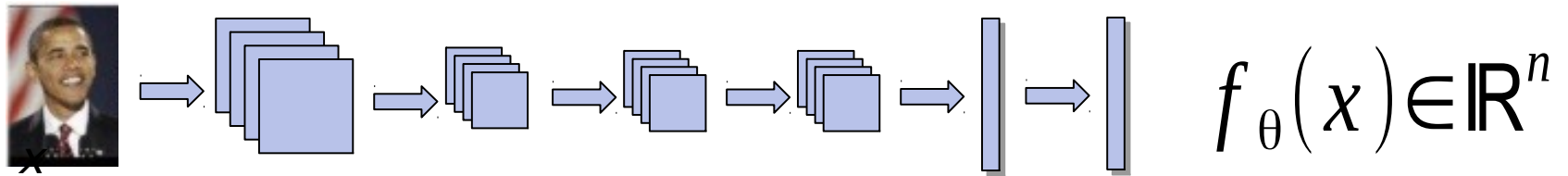


Image verification

- Image2vec



- Similarity metric(example)

$$M(x_1, x_2) = \cos(f_{\theta}(x_1), f_{\theta}(x_2))$$

- x_1, x_2 - pair of images, $T = \# \text{same person}$

Loss:

$$L = T \cdot M(x_1, x_2) + [1 - T] \cdot [1 - M(x_1, x_2)]$$

Case Study: Image search

- Naive approach: (image,query) \rightarrow relevance
- Problem: need $O(\text{images} \times \text{queries})$ runs
- Main objective = top 5~10 candidates



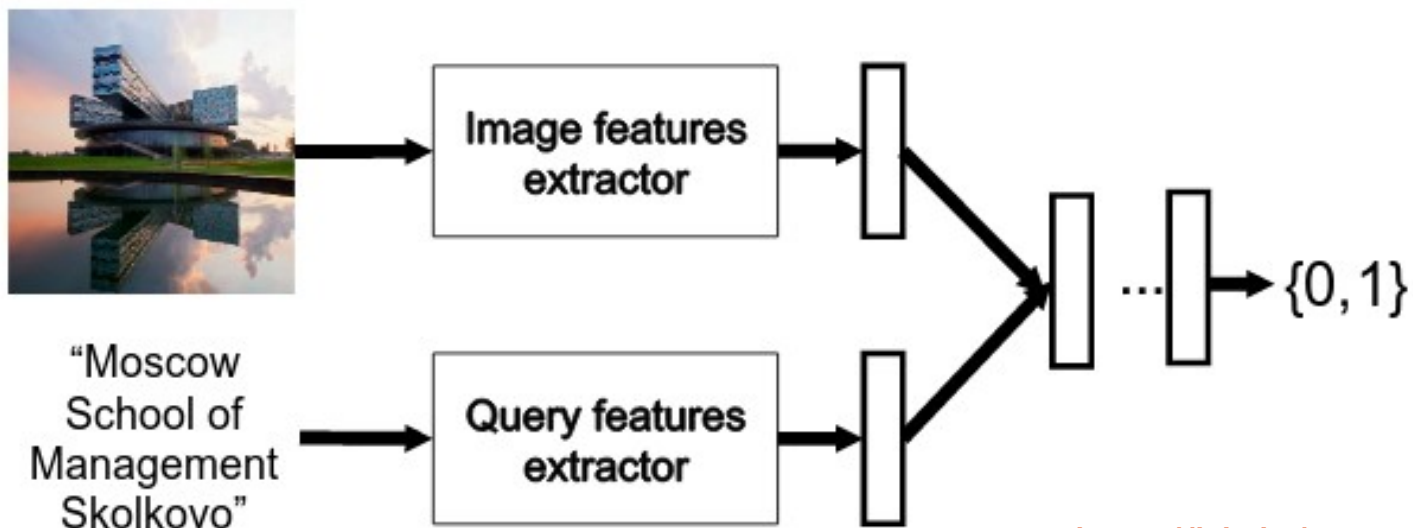
Case Study: Image search

Triples (query, good image, bad image)

$$(q, i^{good}, i^{bad})$$

Pairwise hinge loss

$$L = \max(0, \delta - M(q, i^{good}) + M(q, i^{bad}))$$



Hard negatives

- Baseline: pick negative images at random
 - Most cases are far too simple
- Idea: deliberately pick hardest negative images
 - Hardest = highest score with current NN

racehorse



Hard negatives

- Baseline: pick negative images at random
 - Most cases are far too simple
- Idea: deliberately pick hardest negative images
 - Hardest = highest score with current NN

racehorse



Similarity

0.65



← Pick this one

0.3



Speed up:

- Precompute all query vectors $O(n \text{ queries})$
- Precompute all image vectors $O(n \text{ images})$
- Compute cosine in prediction time
- Locally Sensitive Hashing

...



Google
Yandex

Similar problems

- Other Information Retrieval problem
- Recommendation
- Banner ads
- Classification with a LOT of classes



Brace yourselves

