

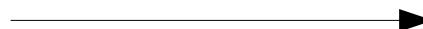
Deep Learning

Episode 3 fall'22

Transfer learning & more vision tasks

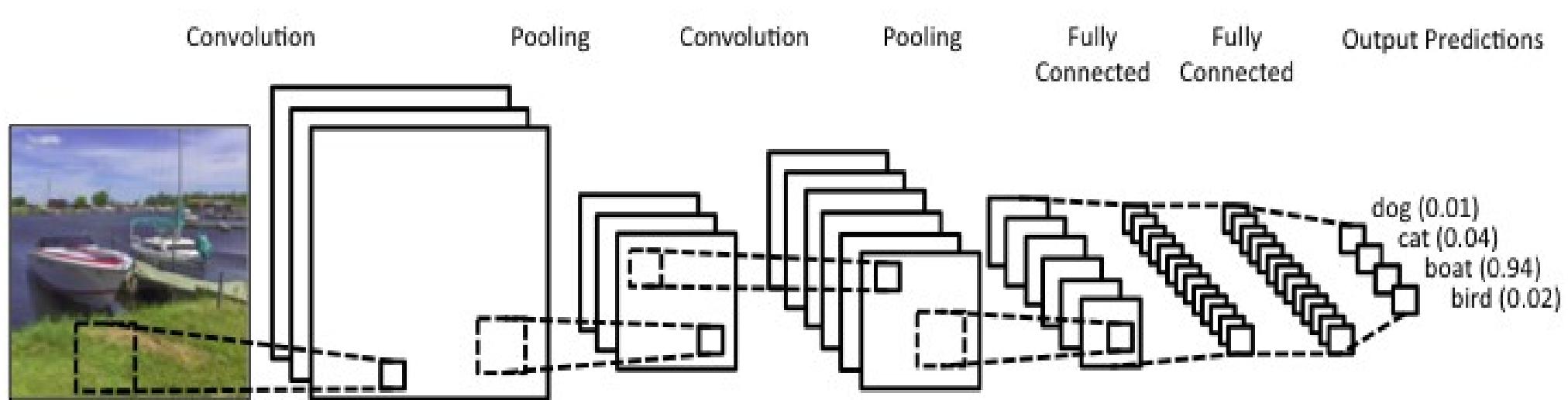


Previously...

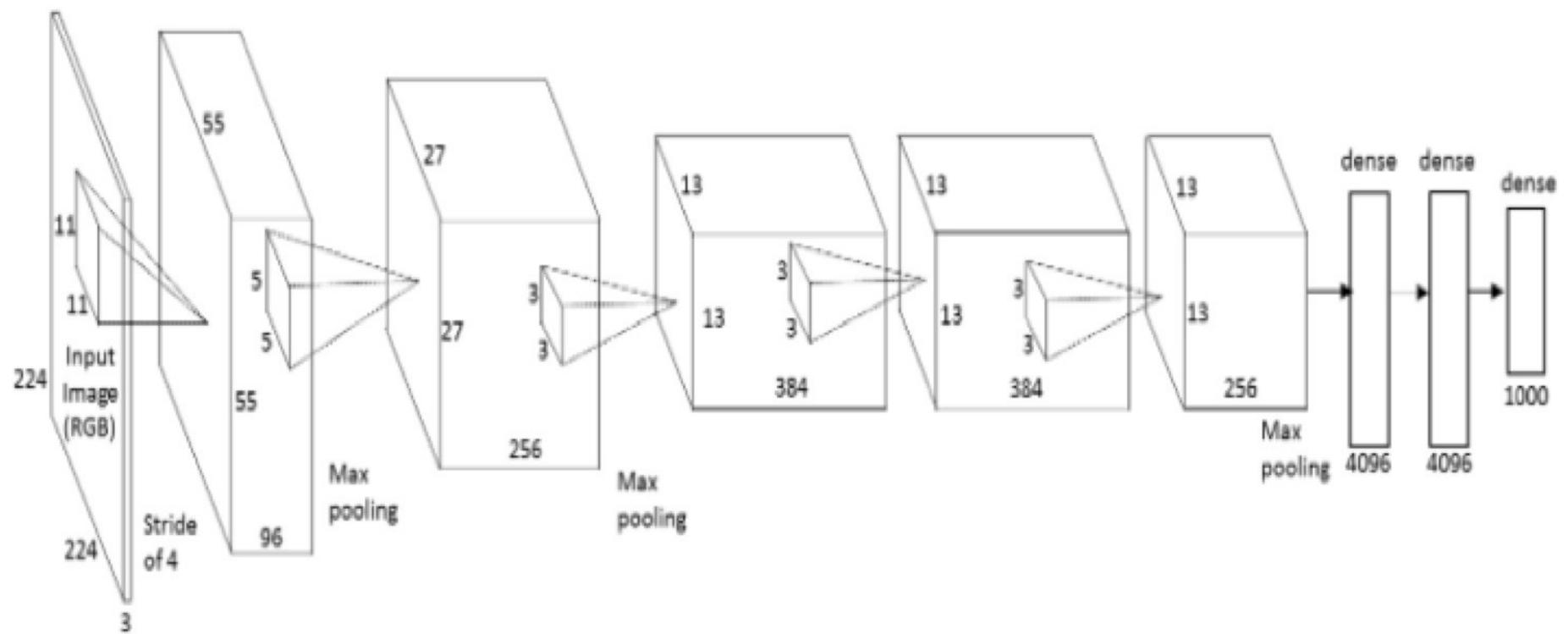


“Dog”

Previously...

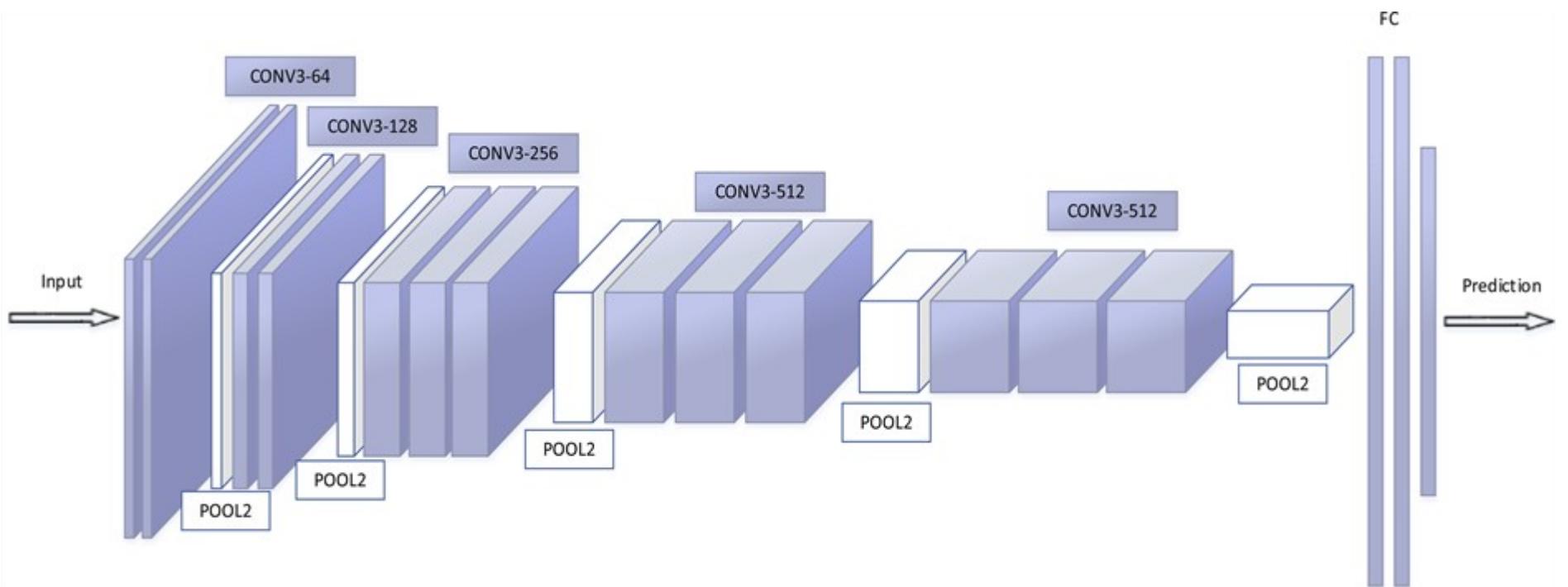


Alexnet



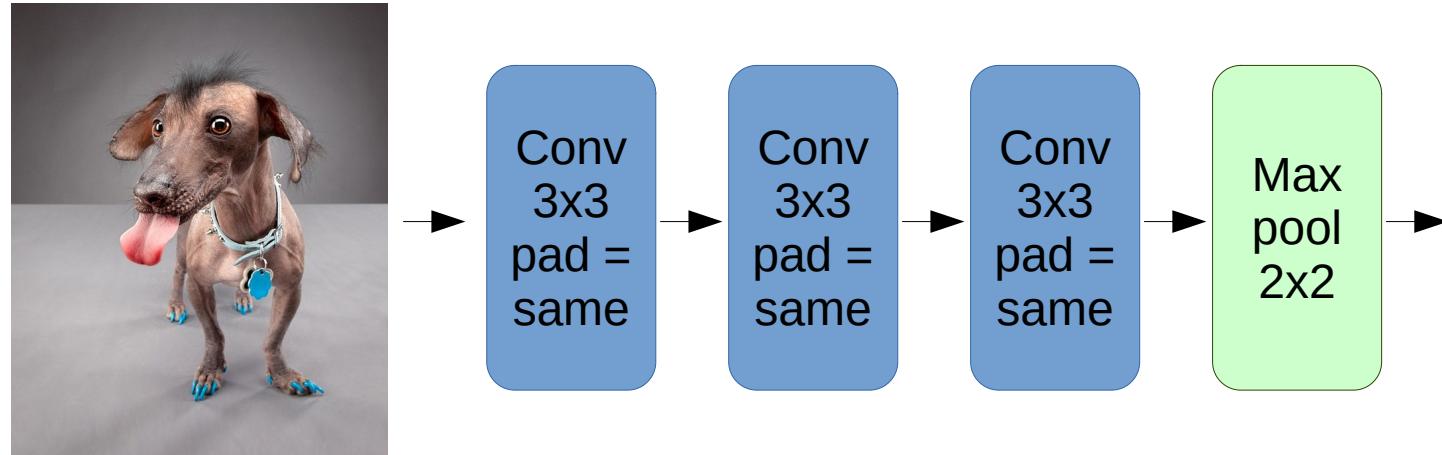
(Krizhevsky et al.)

VGG16

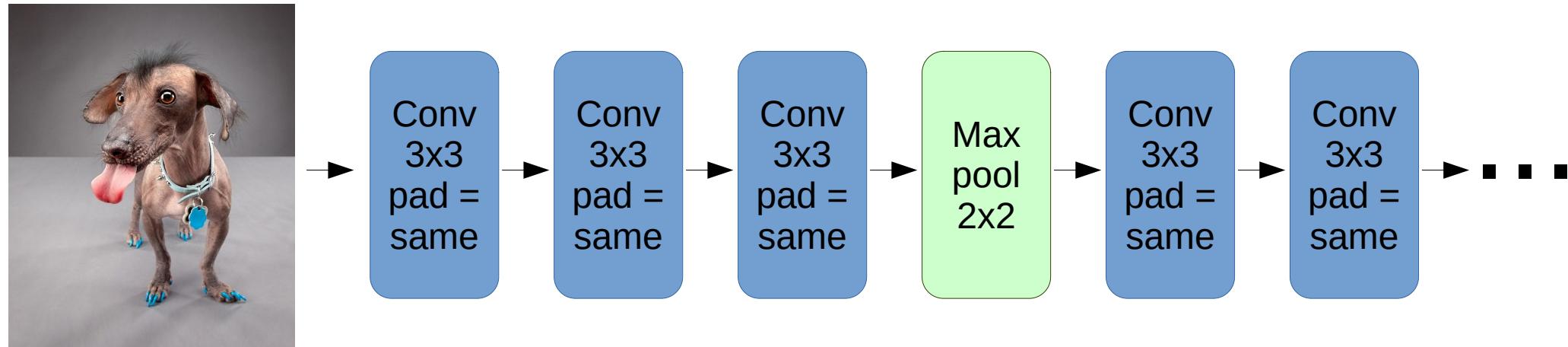


Karen Simonyan and Andrew Zisserman

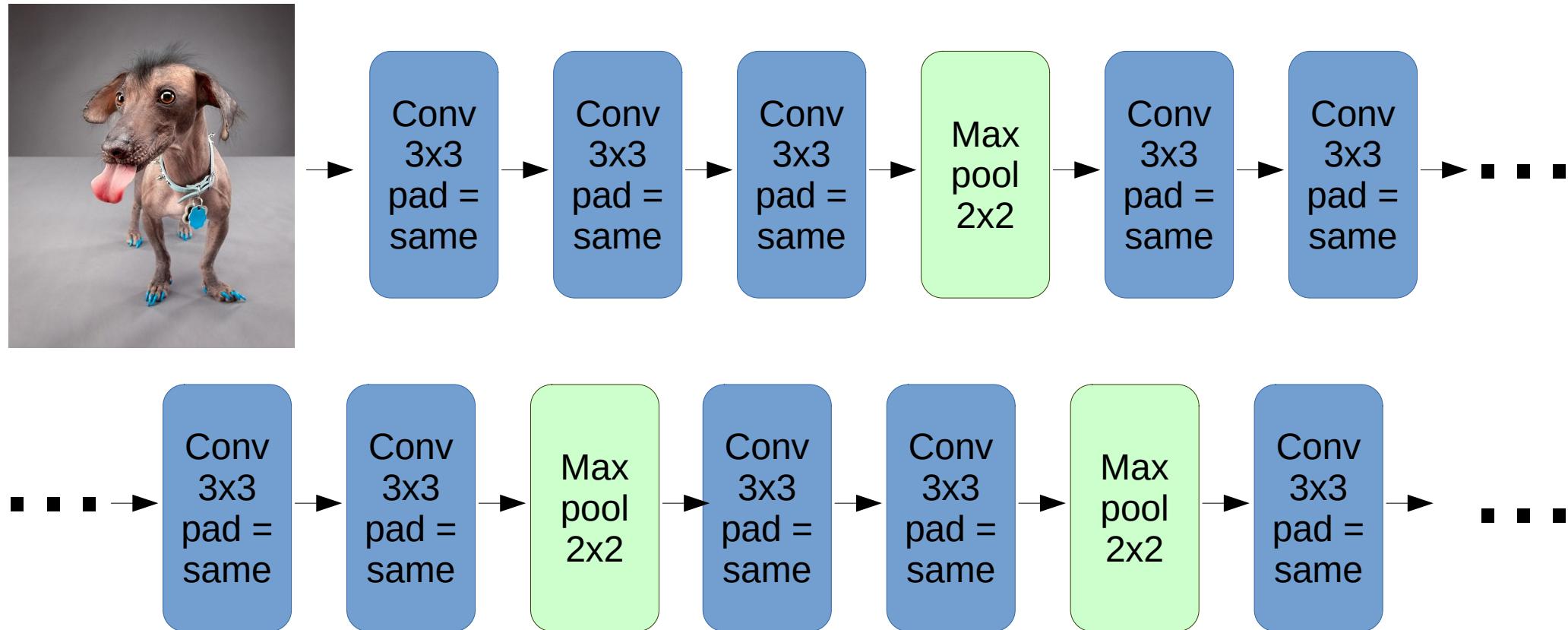
Previously...



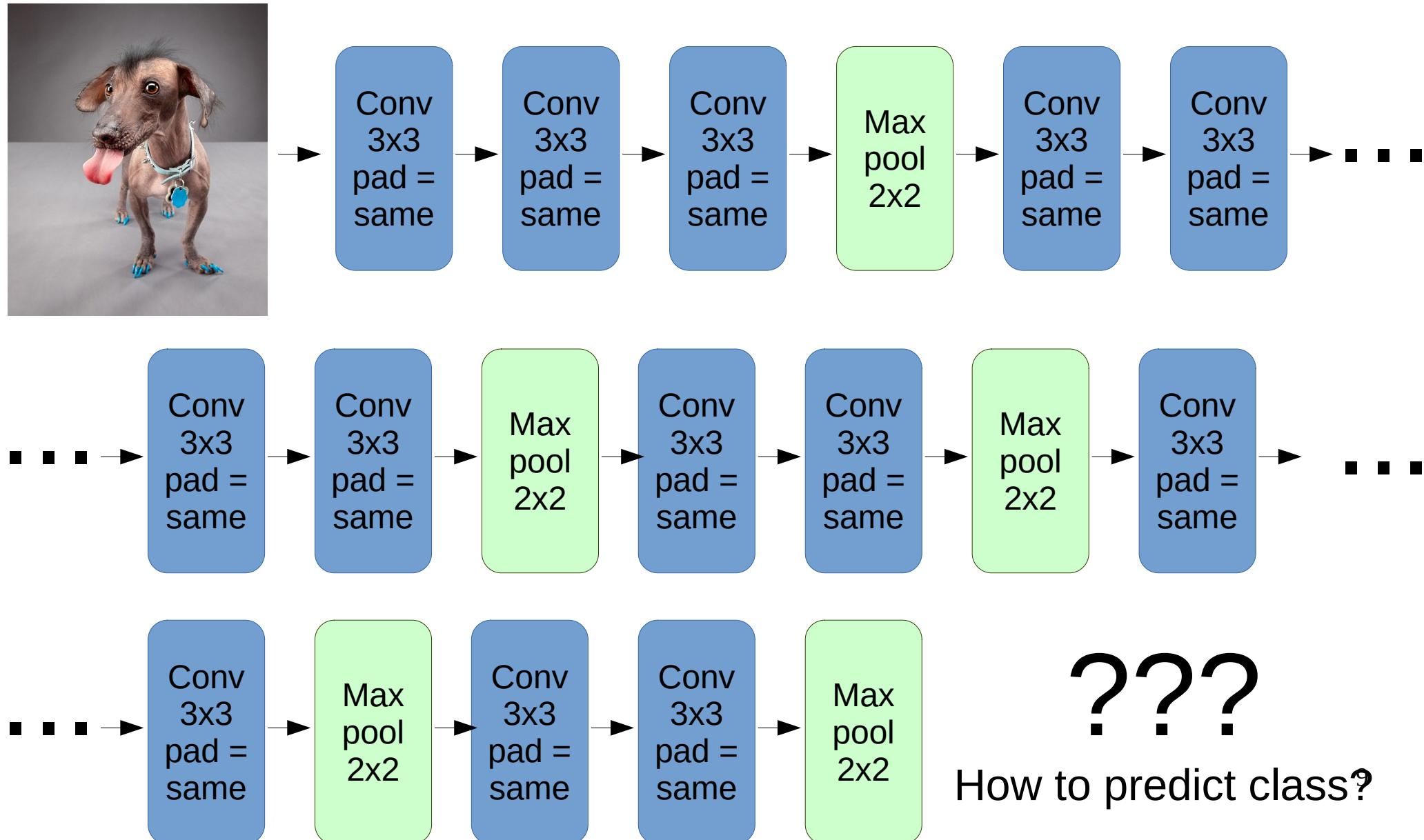
Previously...



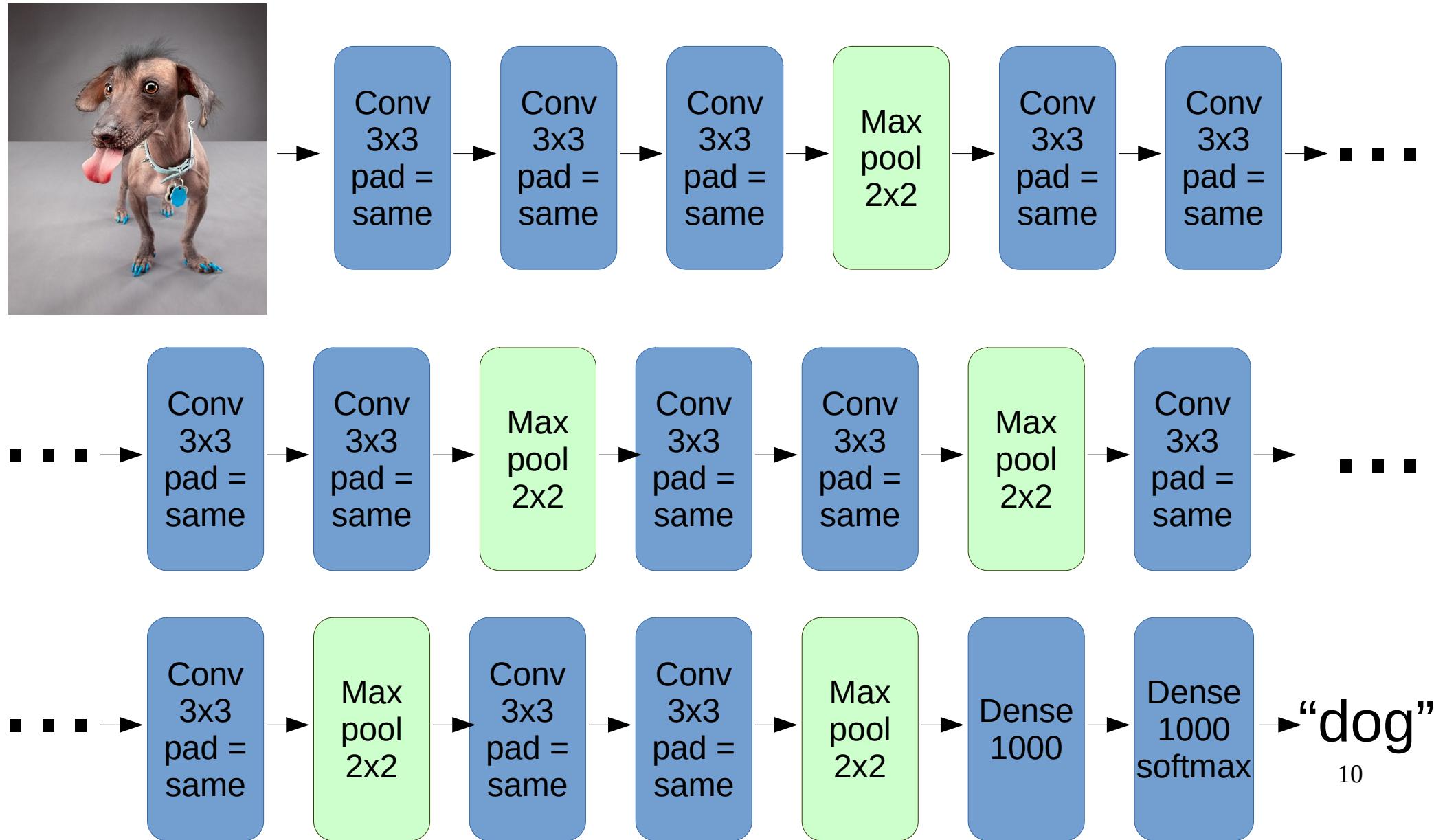
Previously...



Previously...



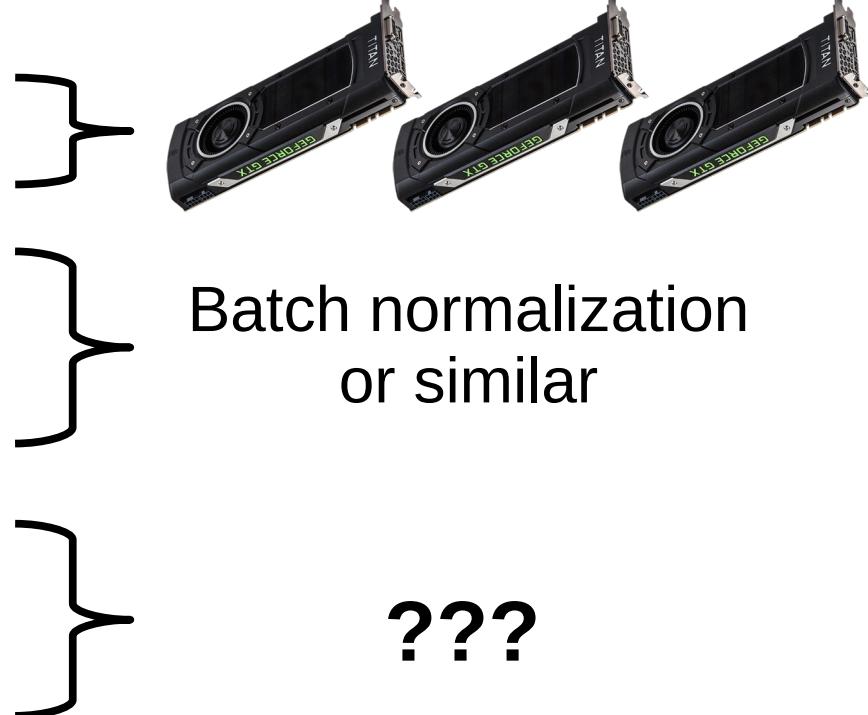
Previously...



Previously...

What you sign for if you stack 1000 layers:

- MemoryError(0x...)
- Activations can vanish
- Activations can explode
- Gradients can vanish
- Gradients can explode



Large-scale recognition



CIFAR



- 60k images
 - 10 classes (left)
 - 32x32 RGB
 - subclasses
- | | |
|-----------------|---|
| aquatic mammals | beaver, dolphin, otter, seal, whale |
| fish | aquarium fish, flatfish, ray, shark, trout |
| flowers | orchids, poppies, roses, sunflowers, tulips |
| food containers | bottles, bowls, cans, cups, plates |

Pfft... weakling

IMAGENET Large Scale Visual Recognition Challenge (ILSVRC) 2010-2014

200 object classes

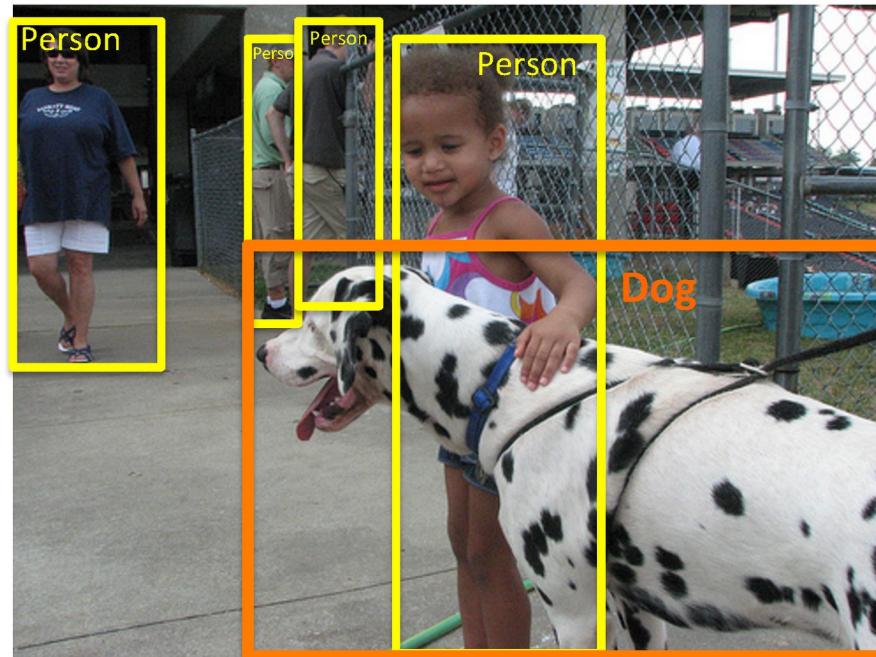
1000 object classes

456,567 images

1,431,167 images

DET

CLS-LOC

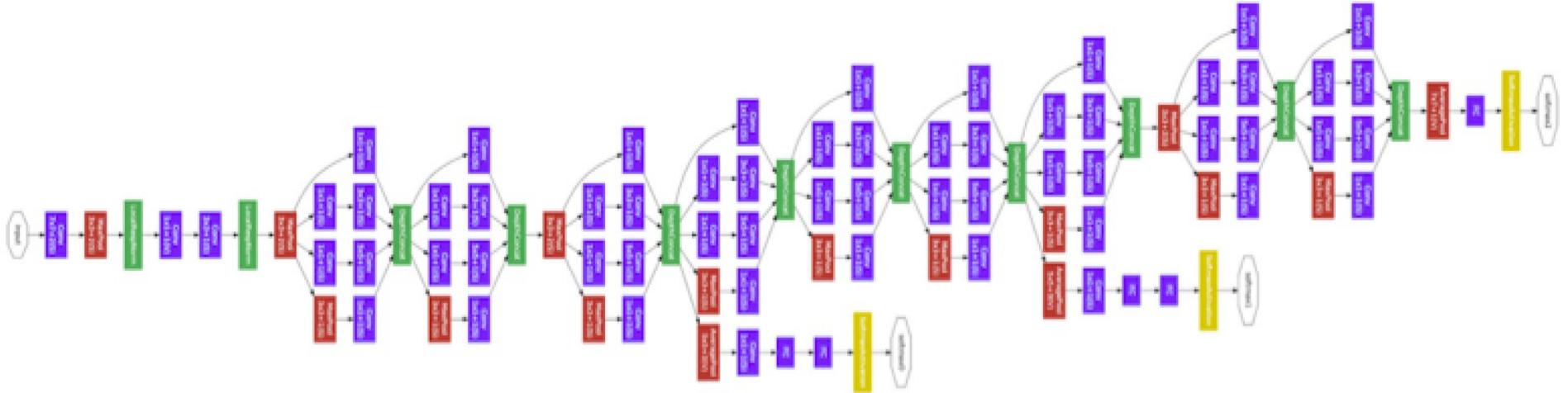


<http://image-net.org/challenges/LSVRC/>

Variety of object classes in ILSVRC

	DET	CLS-LOC				
birds						
bottles						
cars						

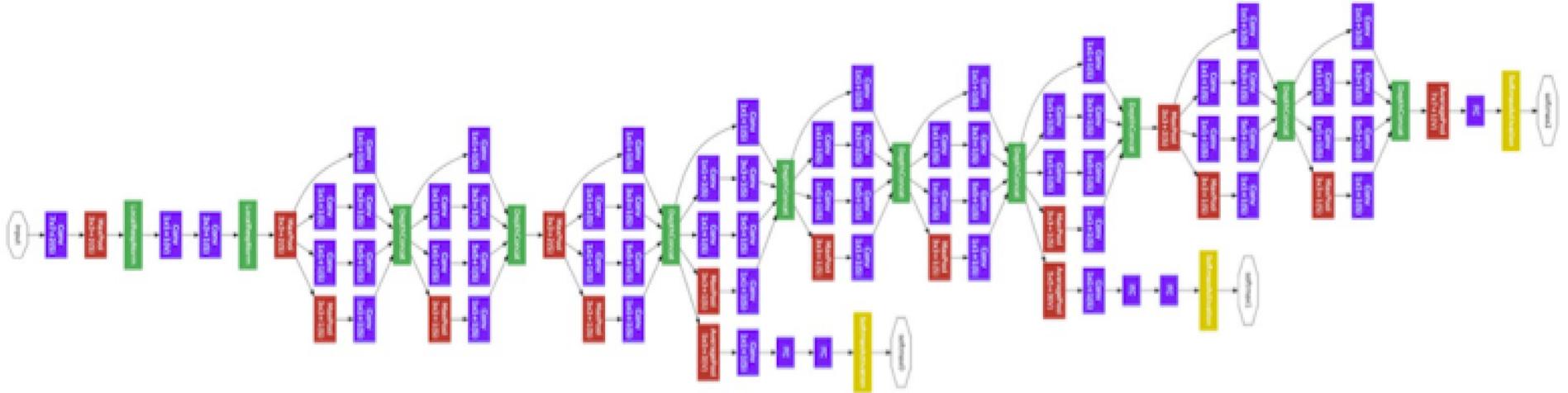
Recap: Inception



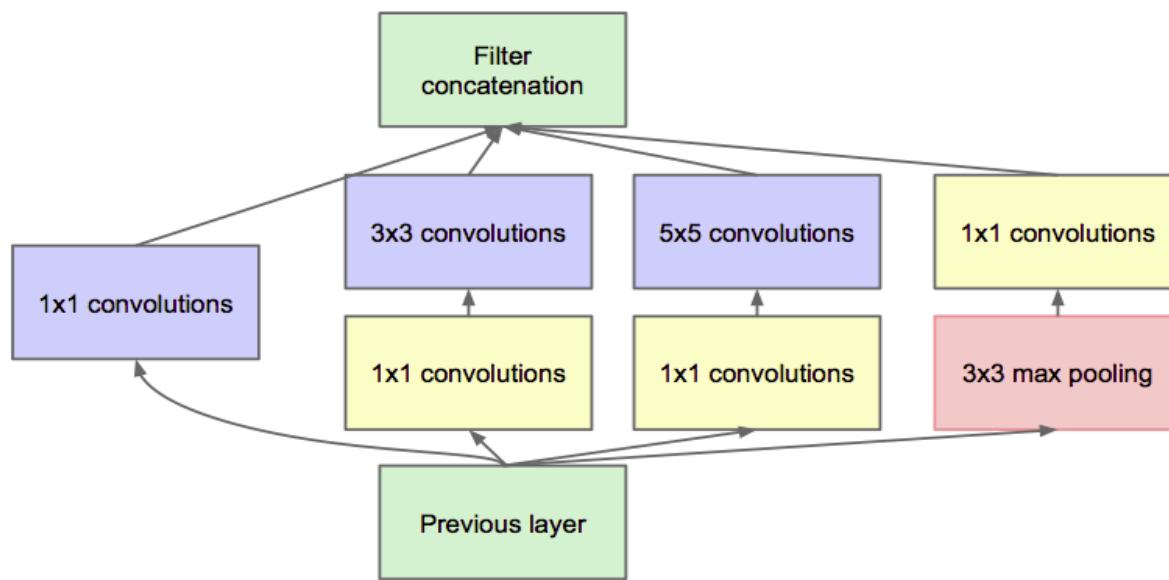
Convolution
Pooling
Softmax
Other

It is not a moon. It is a space station (c)

Recap: Inception



“Inception module”

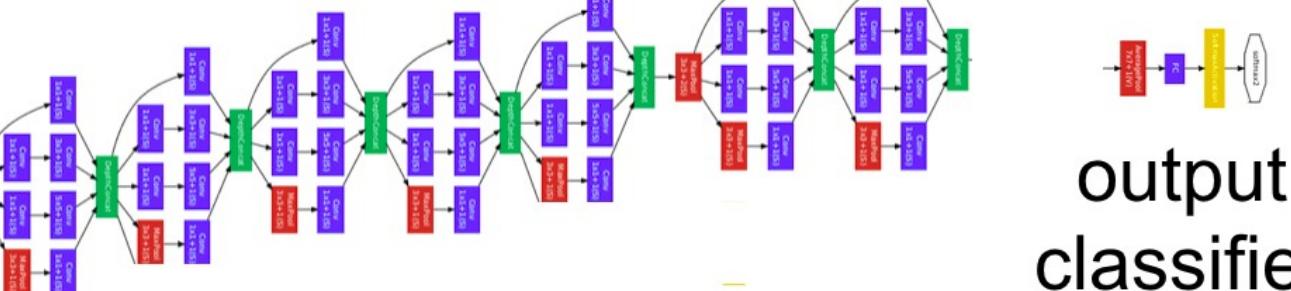
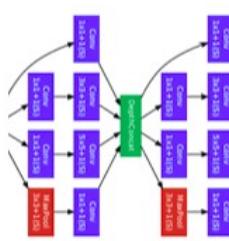


Convolution
Pooling
Softmax
Other

Recap: Inception

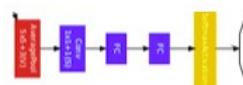
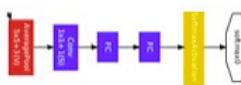
inception modules

stem

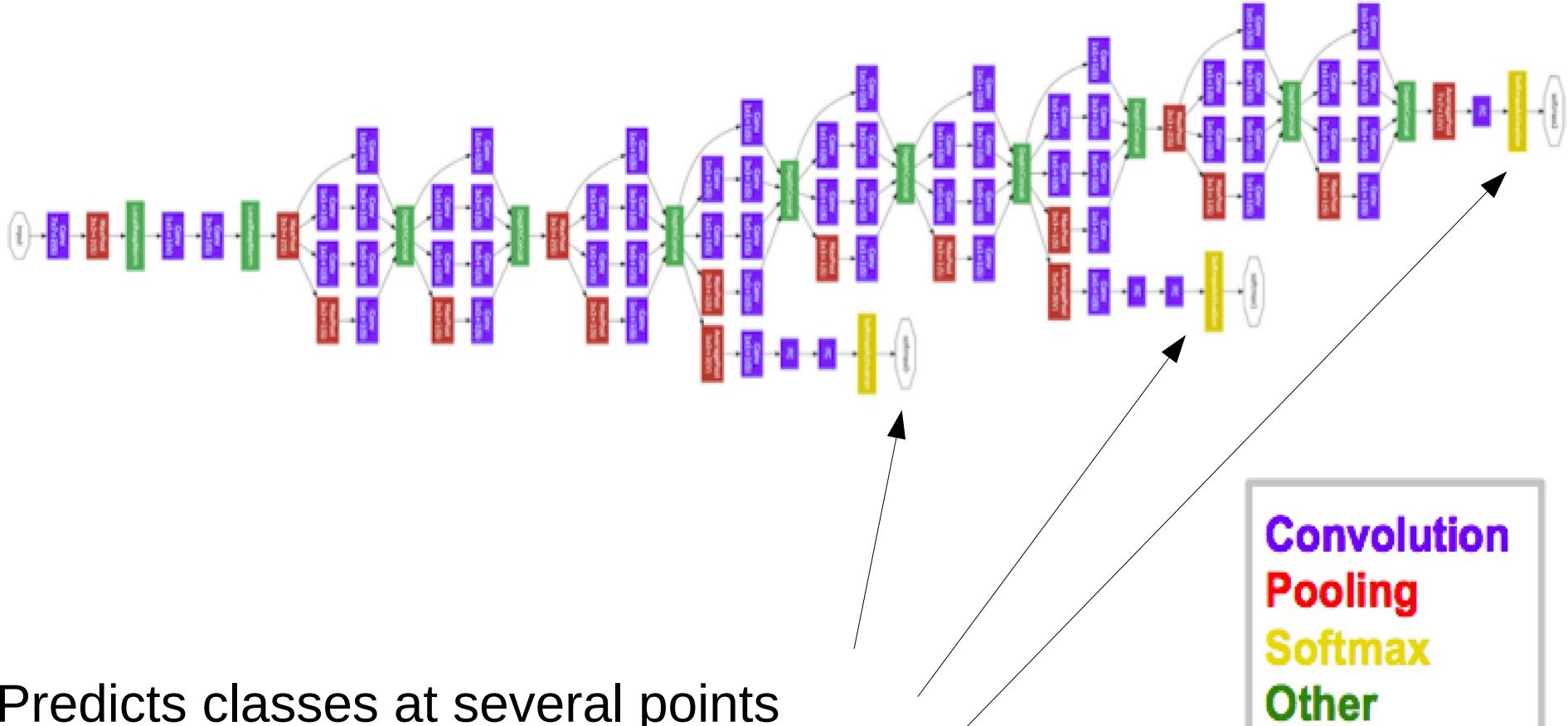


output
classifier

auxiliary classifiers



Recap: Inception

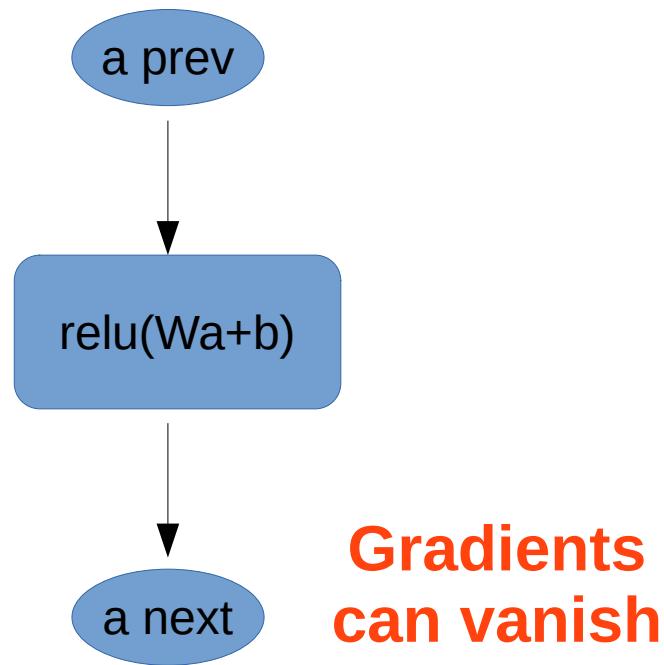


Predicts classes at several points
to make a shorter route for gradients

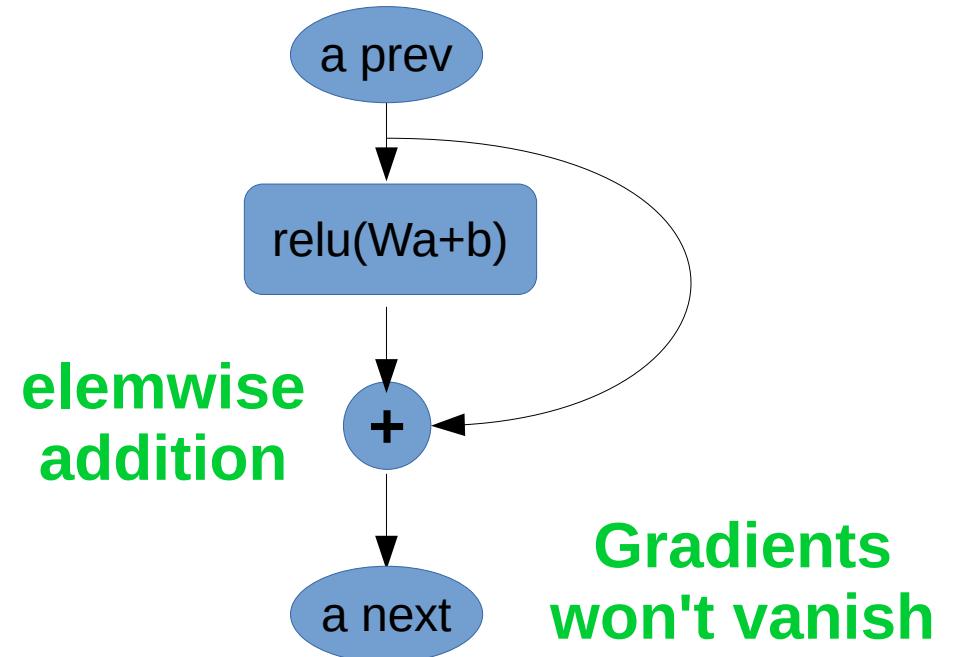
Recap: Residual Layer

Idea: let's create a shortcut for gradients

Normal layer



Residual layer



$$f_{w,b}(x) = \text{relu}(W \cdot a + b)$$

$$\nabla f_{w,b}(x) = \nabla \text{relu}(W \cdot a + b)$$

$$f_{w,b}(x) = \text{relu}(W \cdot a + b) + X$$

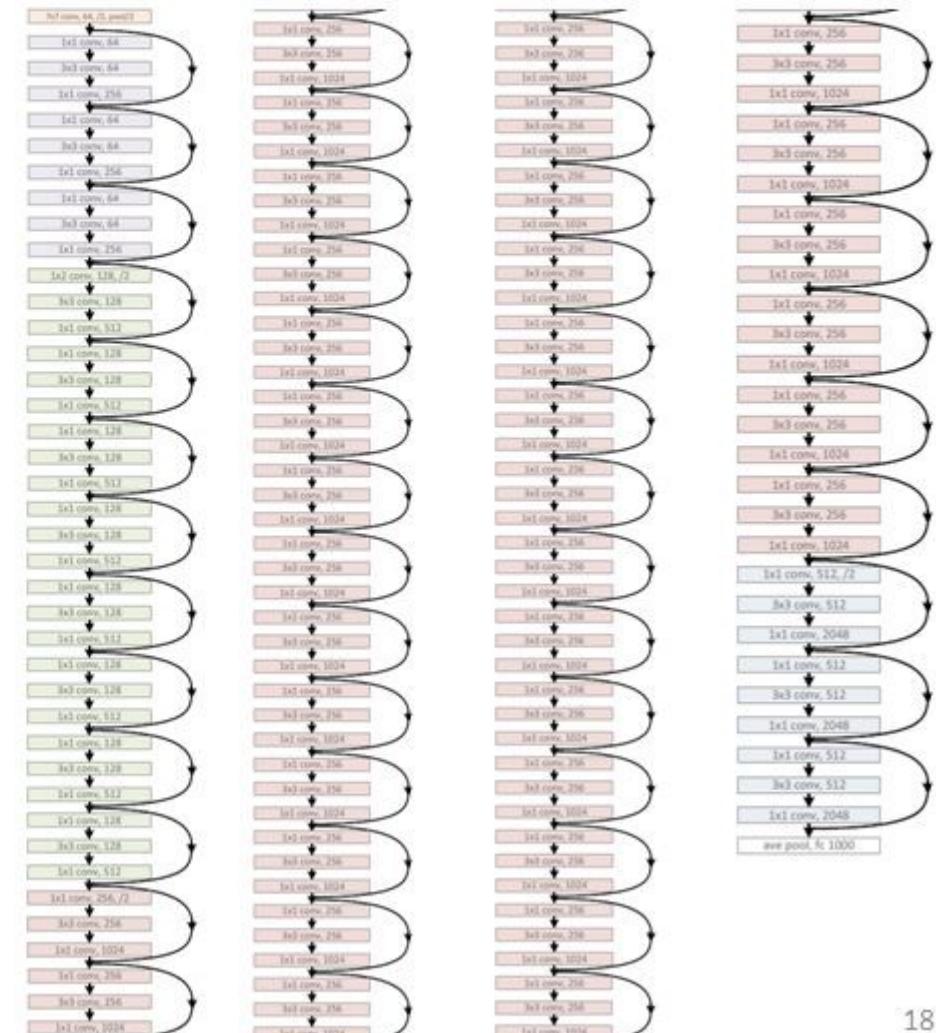
$$\nabla f_{w,b}(x) = \nabla \text{relu}(W \cdot a + b) + \vec{1}$$

Recap: Resnet-152

Each layer learns to “modify”
previous activation

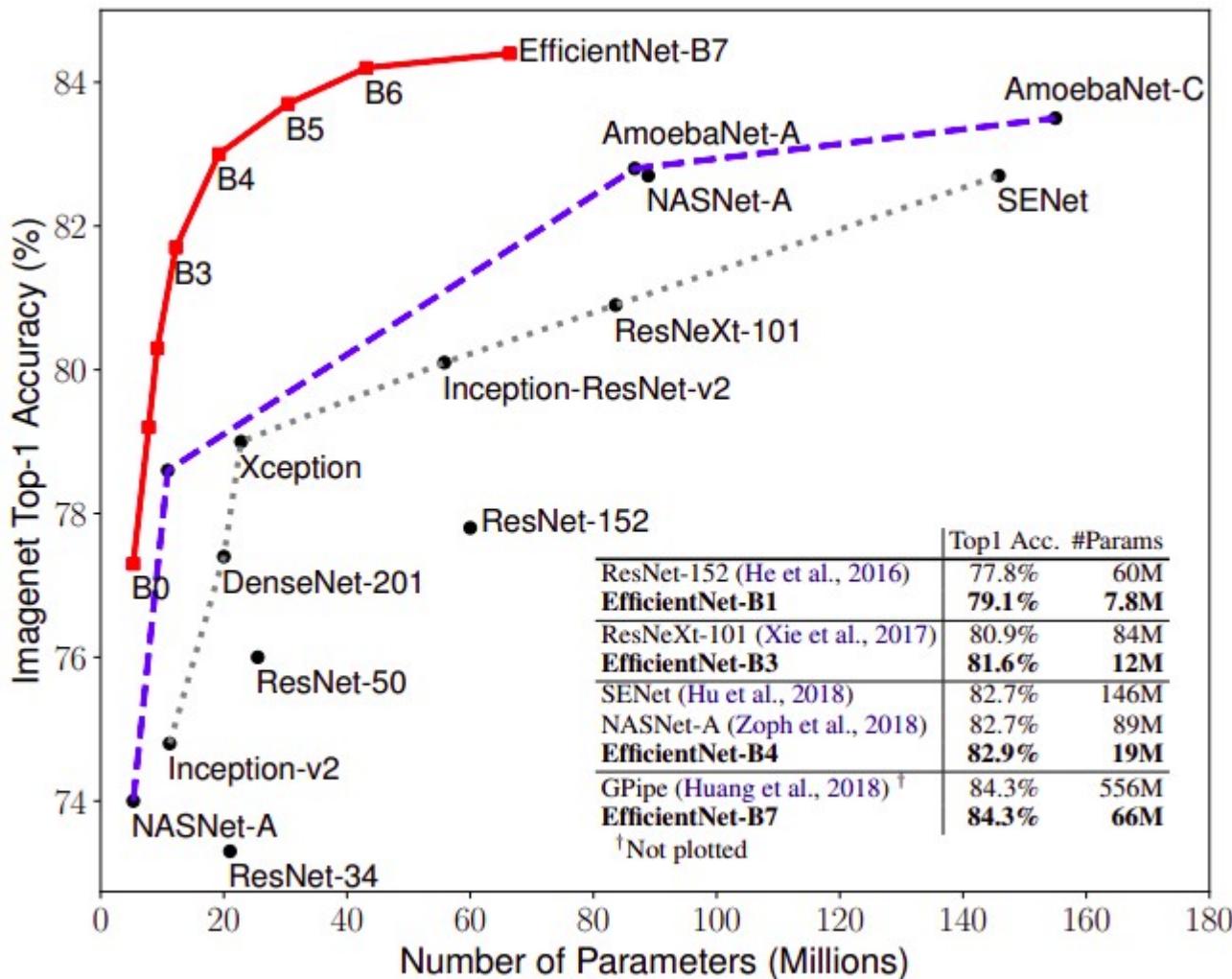
Much less parameters than
vgg16/vgg19

Can be used for other tasks:
segmentation, detection,
captioning
(just like vgg/inception/...)



Recap: EfficientNet

<https://arxiv.org/abs/1905.11946>



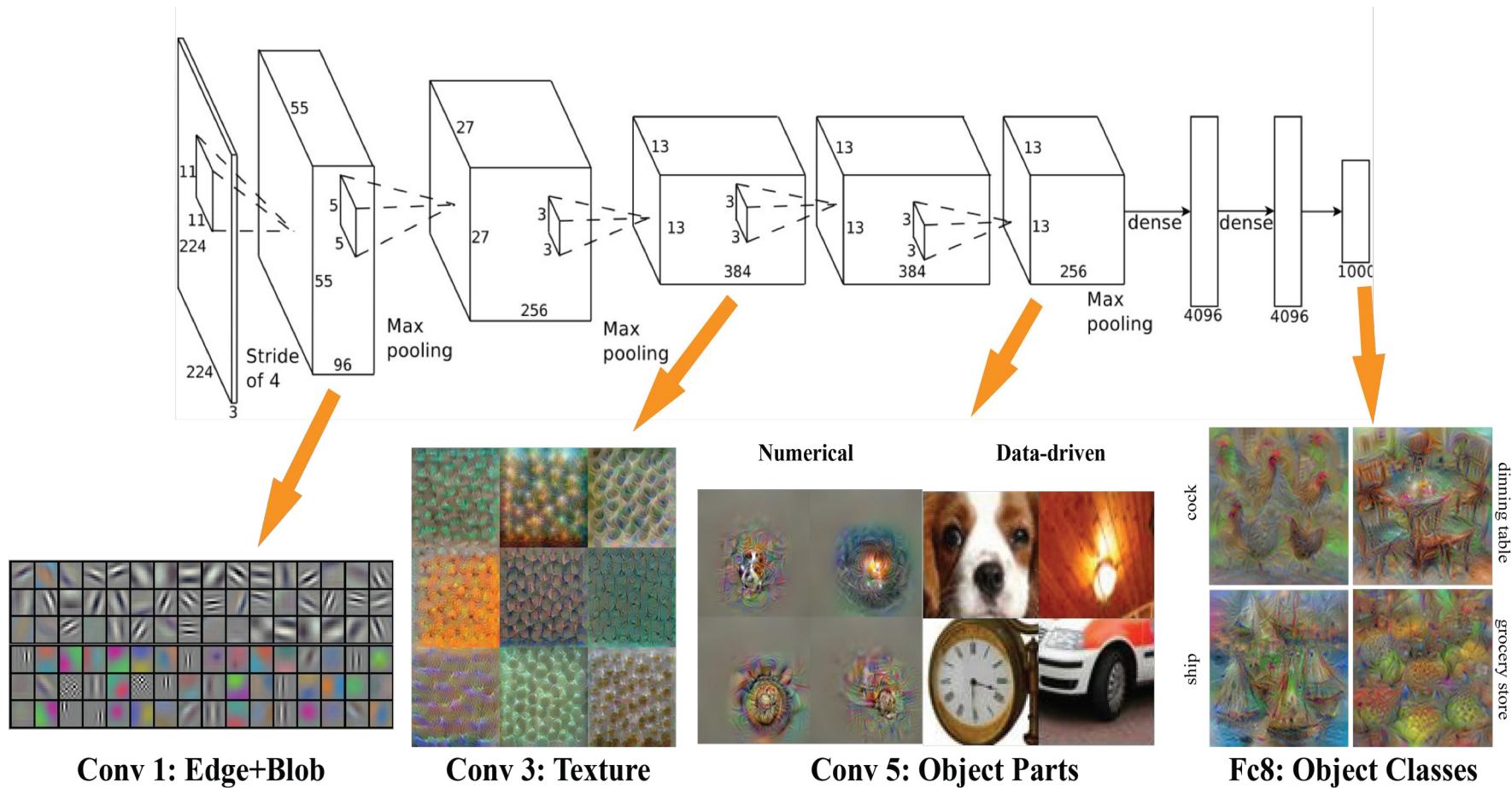
The only problem: you only have 5k images
e.g. transport vs military aircrafts

Ideas?

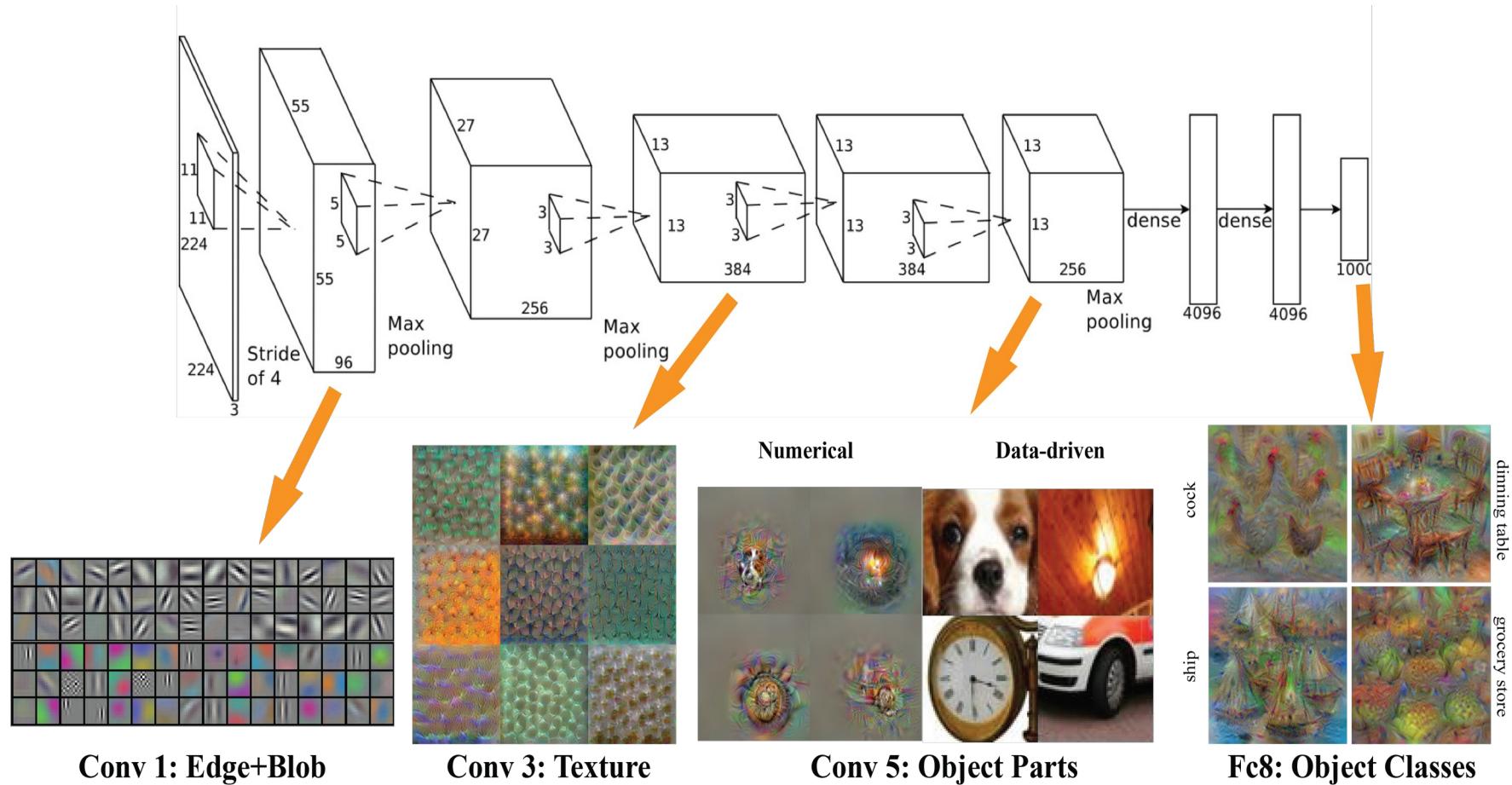
Traditional solutions

- Regularize really hard
- Super small network
- Data augmentation
- Whatever

Feature learning

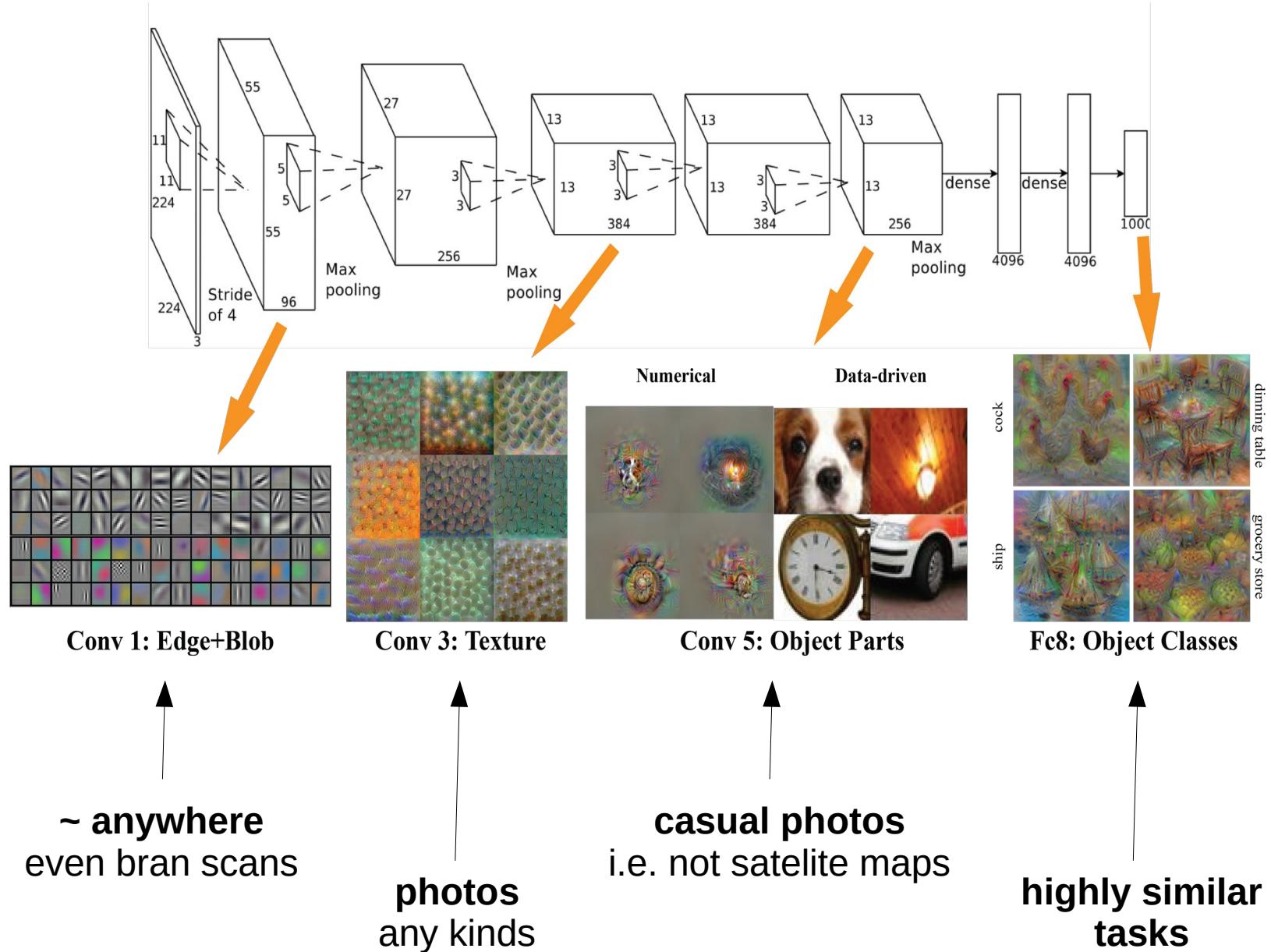


Feature learning



Idea: let's pre-train network on a larger dataset

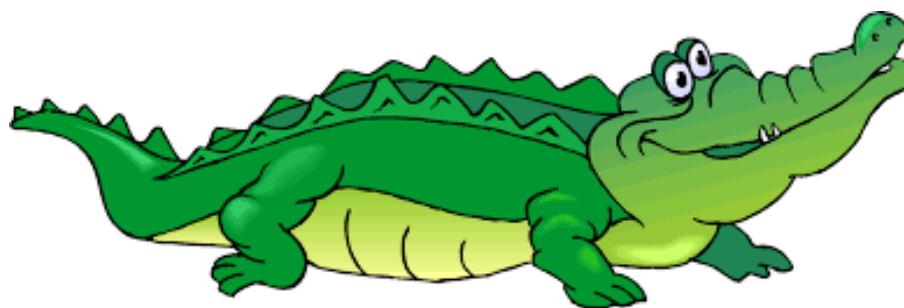
Can we use them?



Pre-training

- 1. Train a network on large dataset

cifar X

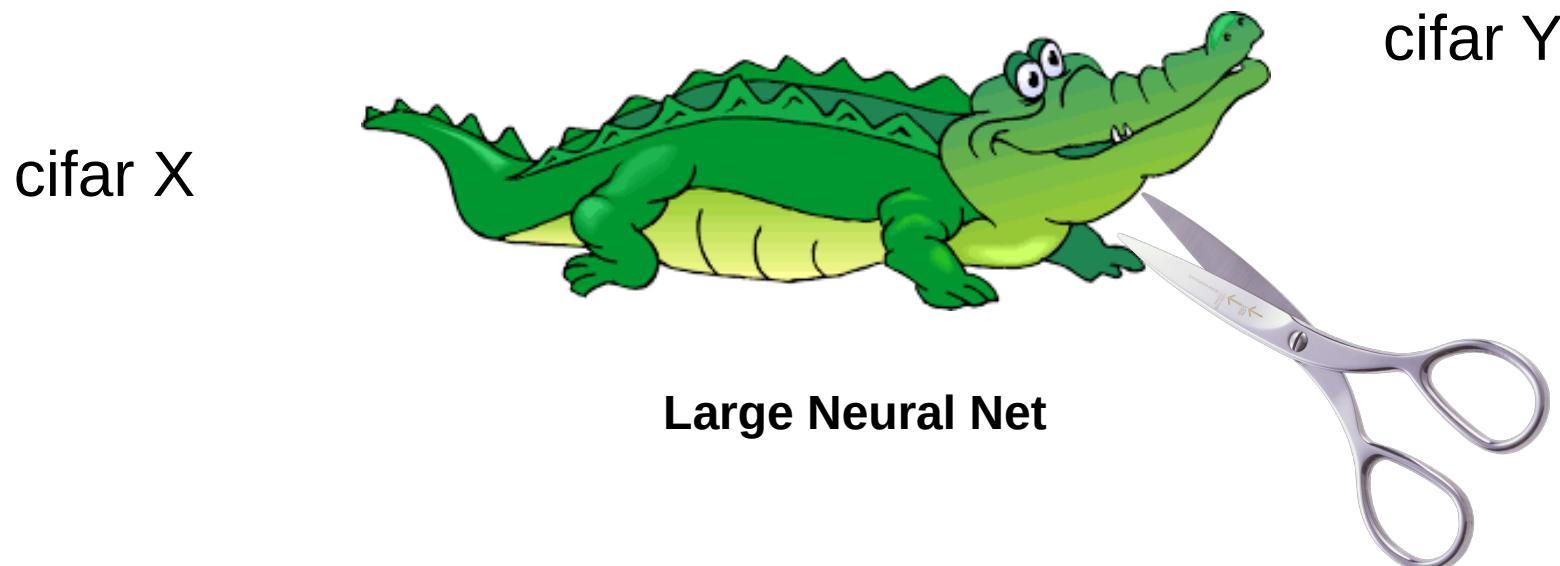


cifar Y

Large Neural Net

Pre-training

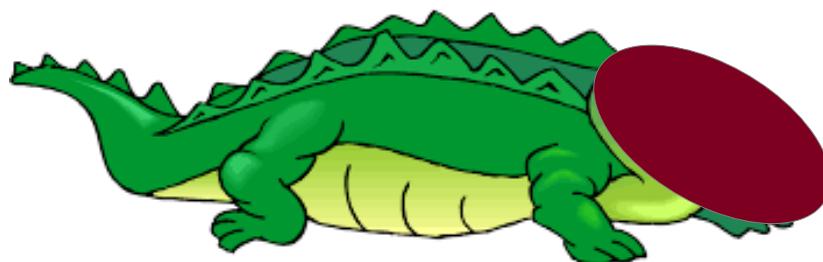
- 1. Train a network on large dataset
- 2. Take some intermediate layer



Pre-training

- 1. Train a network on large dataset
- 2. Take some intermediate layer

cifar X

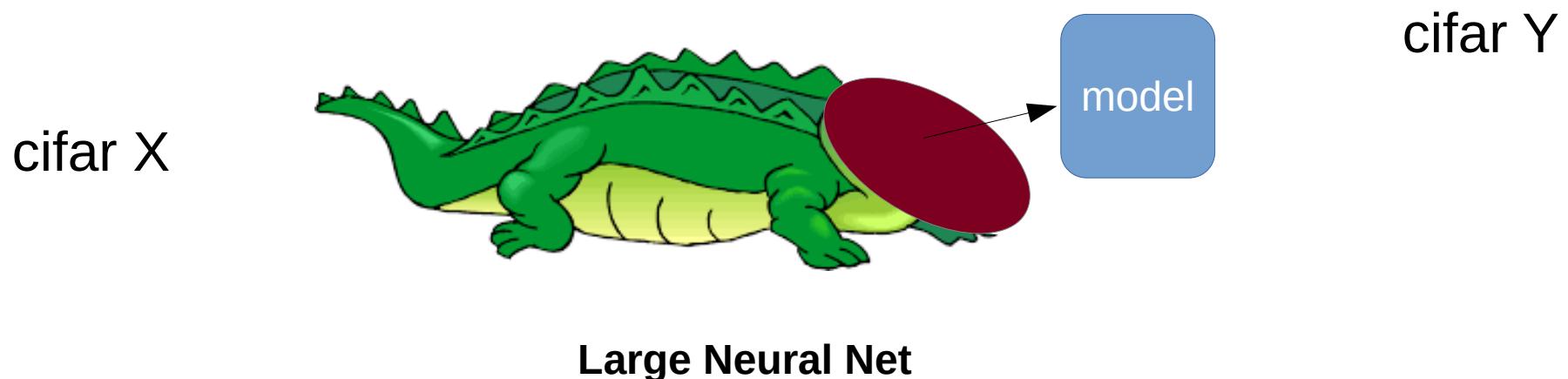


cifar Y

Large Neural Net

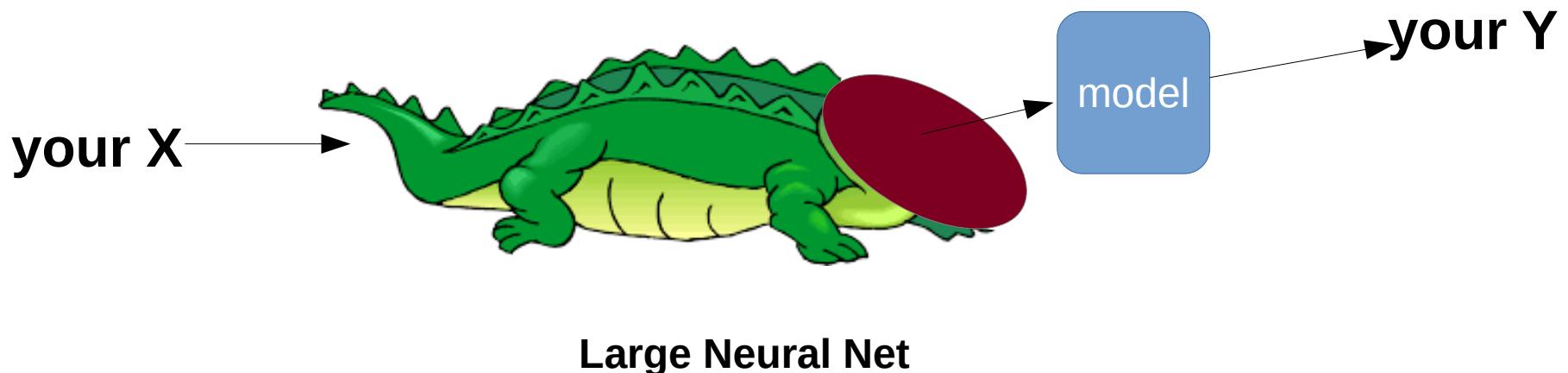
Pre-training

- 1. Train a network on large dataset
- 2. Take some intermediate layer
- 3. Build model on top of it



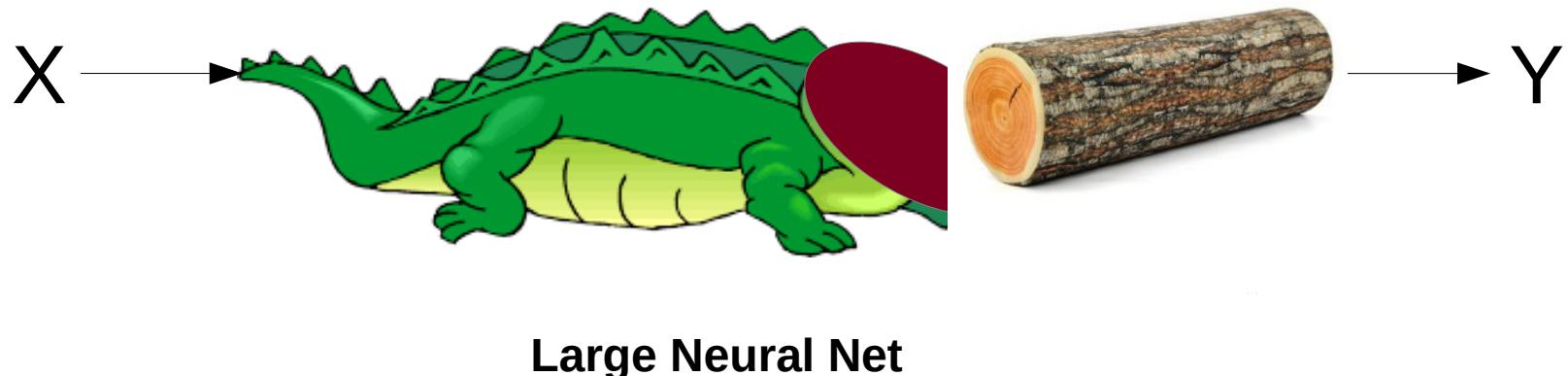
Pre-training

- 1. Train a network on large dataset
- 2. Take some intermediate layer
- 3. Build model on top of it
- 4. Train model for your objective



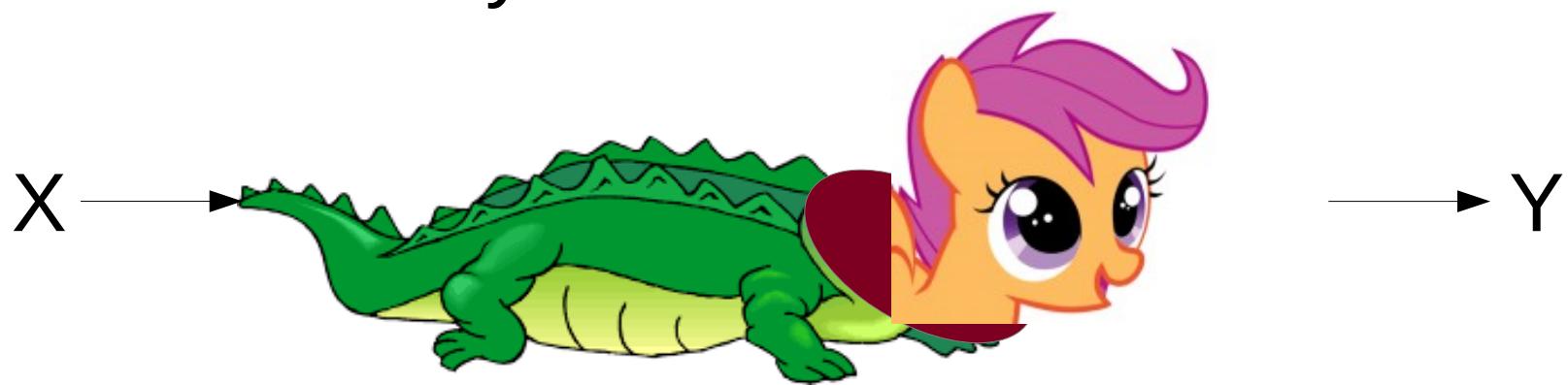
Reusing features

- Chop off “head”
- use “neck” as feature extractor
- Train ANY classifier
 - even random forest will do



Fine-tuning

- Chop off “head”
- “freeze” body (consider constant)
- Build new neural network in it's place
- Train “head” only for several iterations
- Un-freeze body and train full network



Large Neural Net

**Q: What if we only have
1-5 images for each class?**

Few-shot learning with KNN

Idea: run KNN over frozen pre-trained model

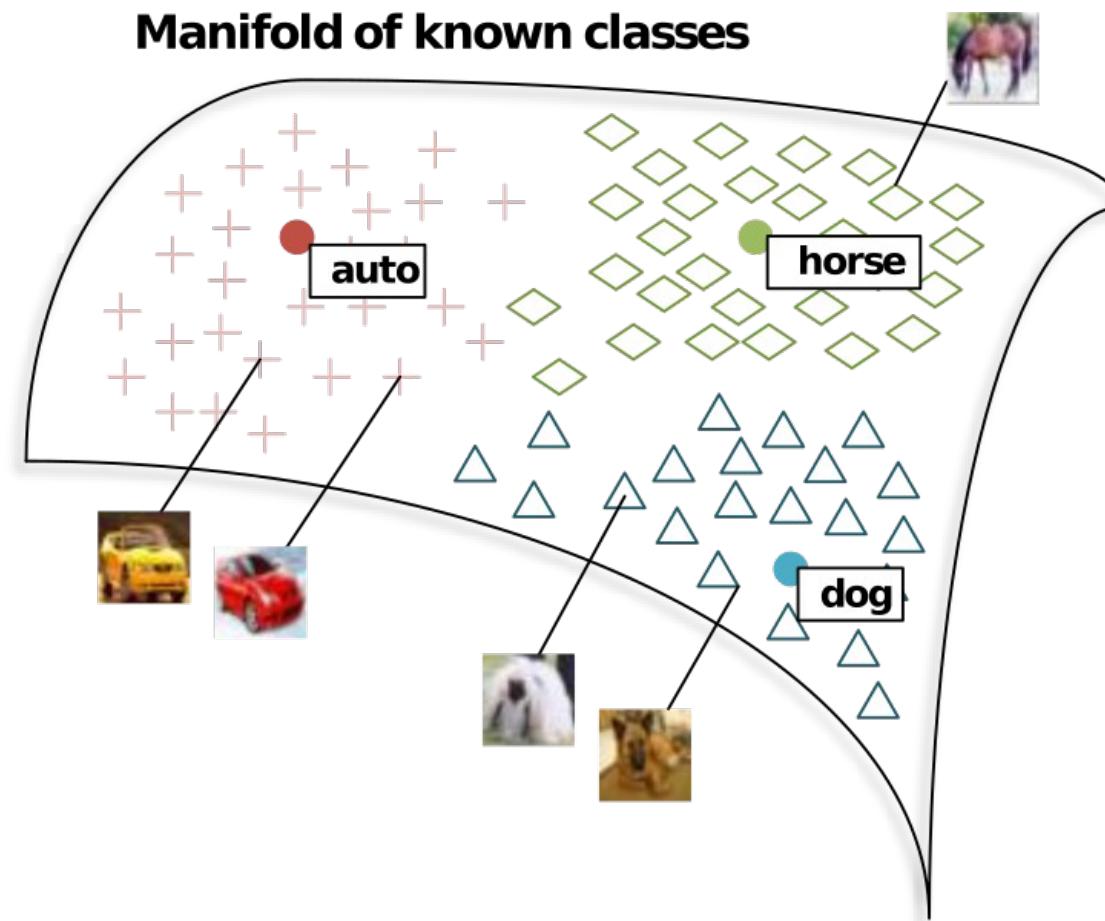
- ● ● - training images, color = class



Few-shot learning with KNN

Idea: run KNN over frozen pre-trained model

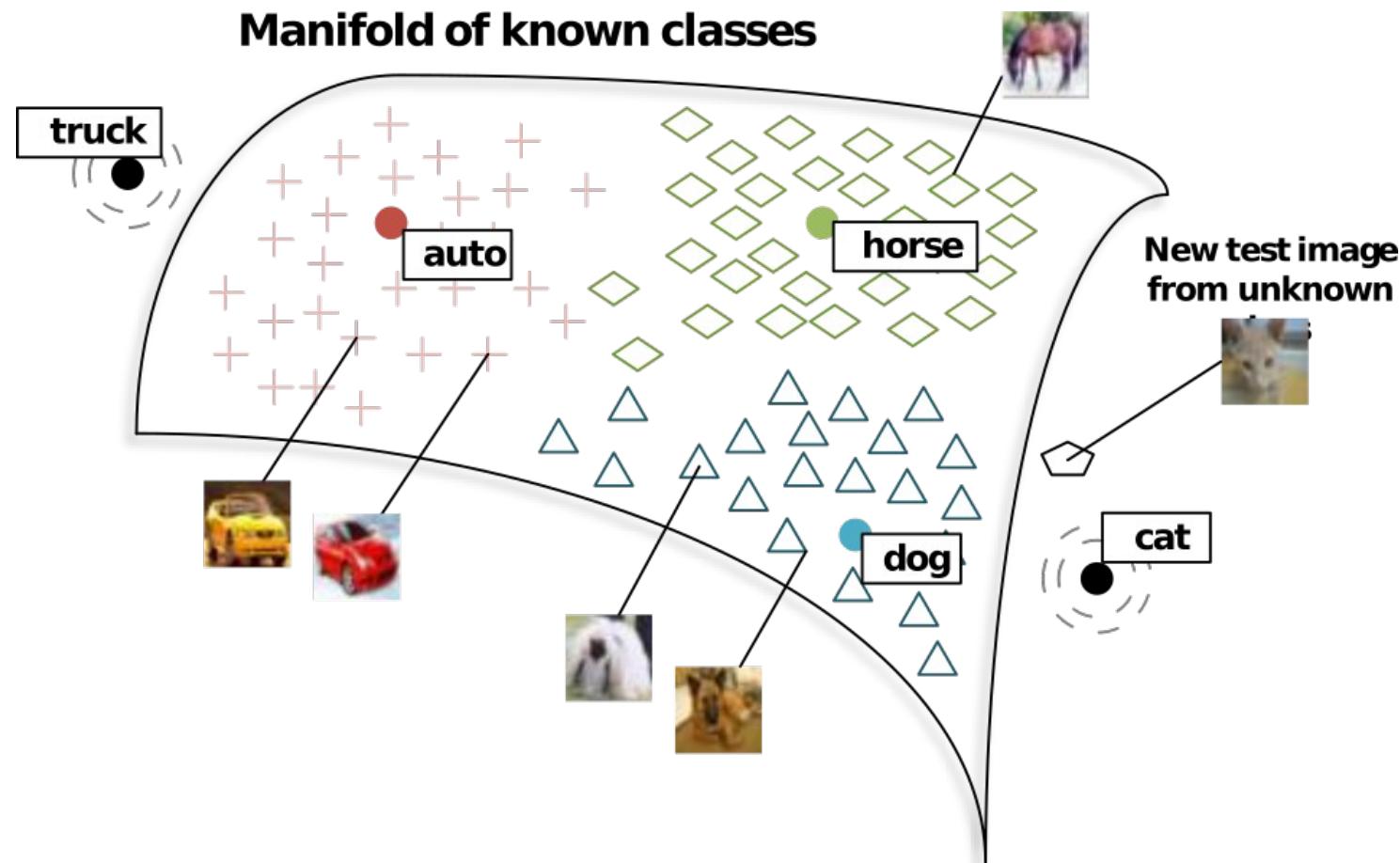
- ● ● - training images, color = class



Few-shot learning with KNN

Idea: run KNN over frozen pre-trained model

- ● ● - training images, color = class



**What if we do not have
any labeled dataset?**

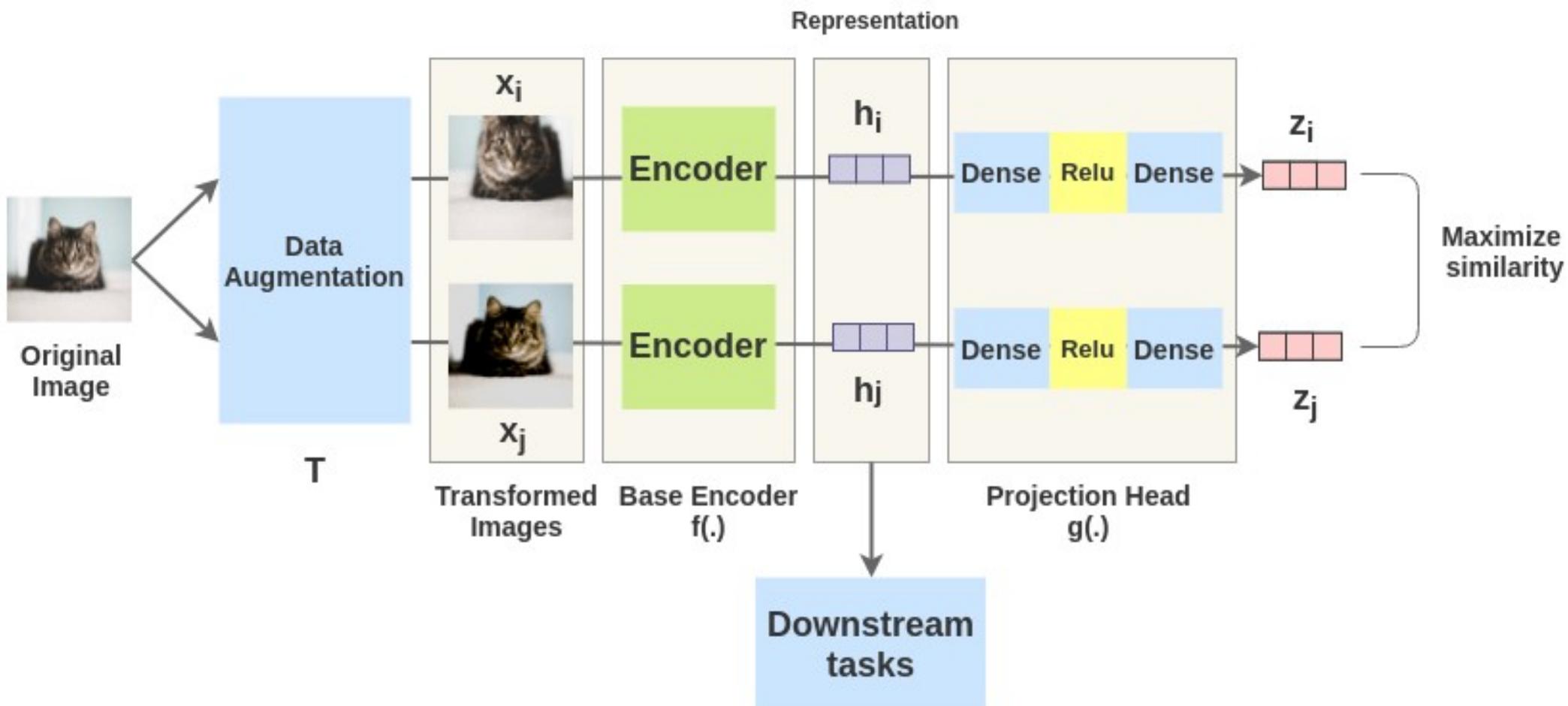
**What if we do not have
any ~~labeled~~ dataset?**

Can we learn features from just images?

Self-supervised learning

- SimCLR: <https://arxiv.org/abs/2002.05709>

SimCLR Framework



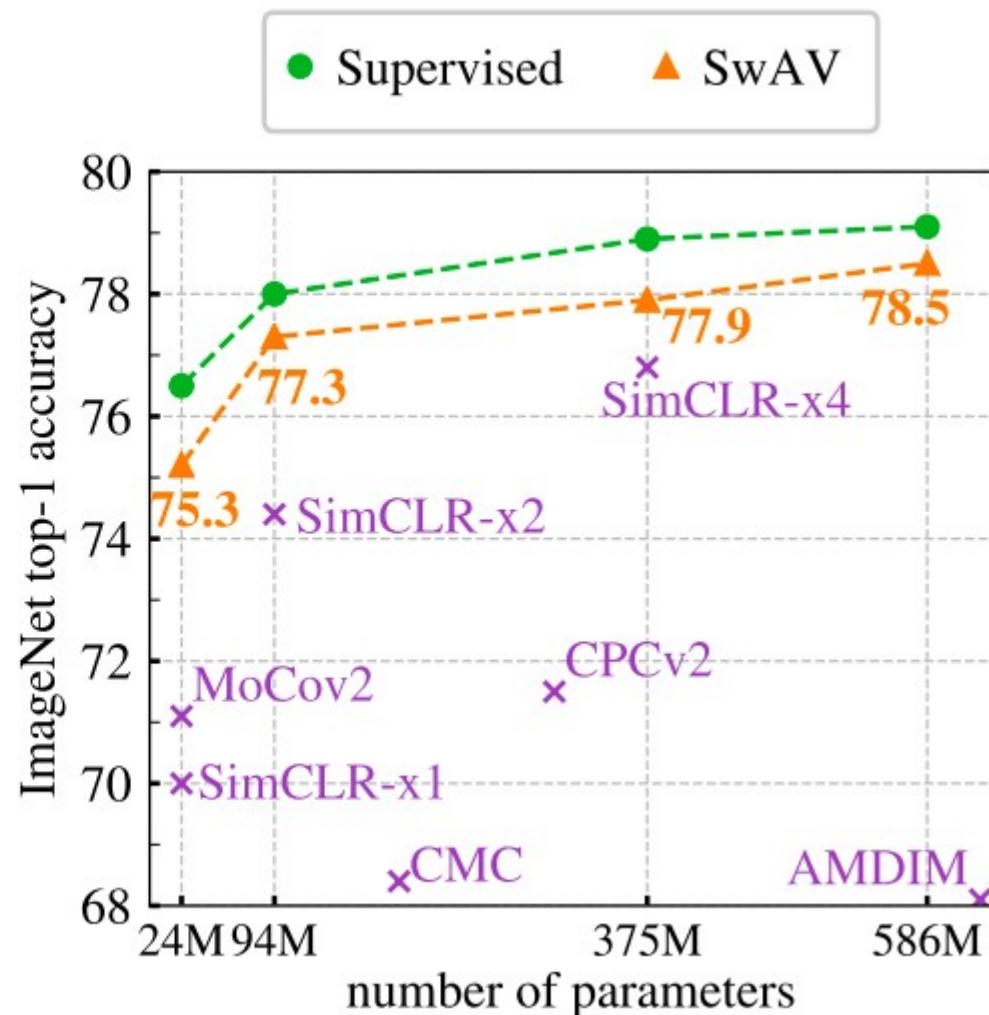
Self-supervised learning

- SwAV: <https://arxiv.org/abs/2006.09882>

Self-supervised learning

- ImageNet top-1 accuracy: getting close!

Method	Arch.	Param.	Top1
Supervised	R50	24	76.5
Colorization [65]	R50	24	39.6
Jigsaw [46]	R50	24	45.7
NPID [58]	R50	24	54.0
BigBiGAN [15]	R50	24	56.6
LA [68]	R50	24	58.8
NPID++ [44]	R50	24	59.0
MoCo [24]	R50	24	60.6
SeLa [2]	R50	24	61.5
PIRL [44]	R50	24	63.6
CPC v2 [28]	R50	24	63.8
PCL [37]	R50	24	65.9
SimCLR [10]	R50	24	70.0
MoCov2 [11]	R50	24	71.1
SwAV	R50	24	75.3

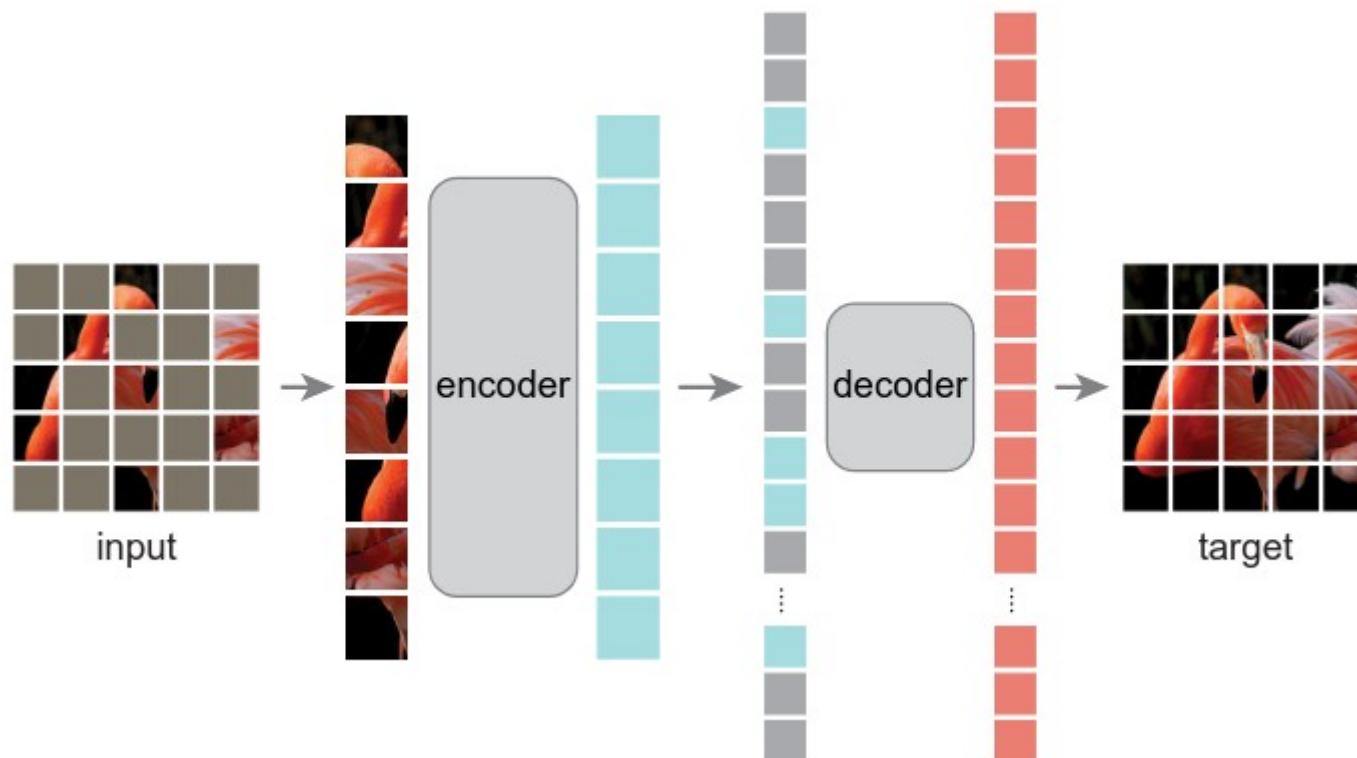


Masked Image Modeling

<https://arxiv.org/abs/2111.09886>

<https://arxiv.org/abs/2111.06377>

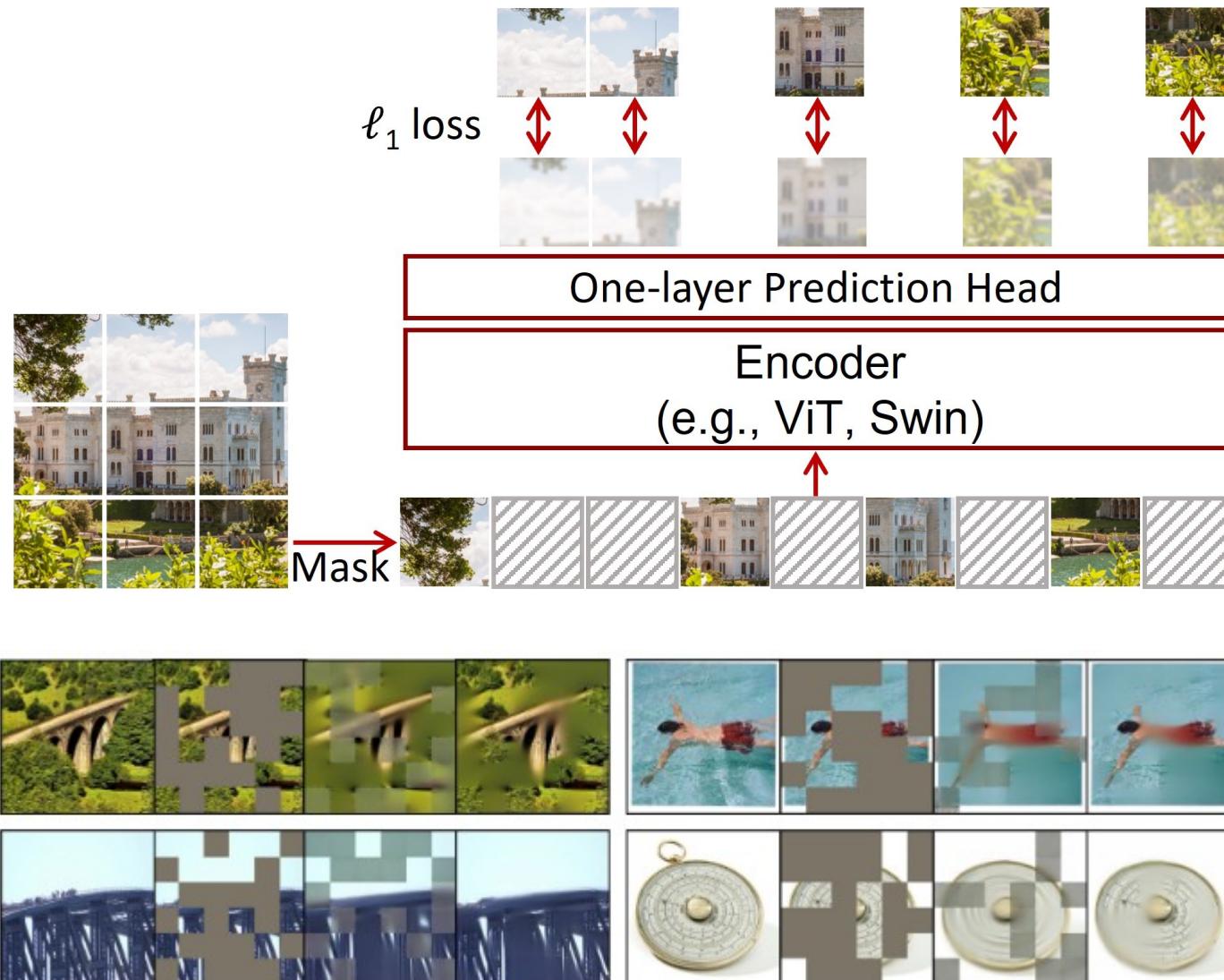
- Remove random tiles, try to reconstruct them!
reuse the pre-trained network for other tasks



Masked Image Modeling

<https://arxiv.org/abs/2111.09886>

<https://arxiv.org/abs/2111.06377>

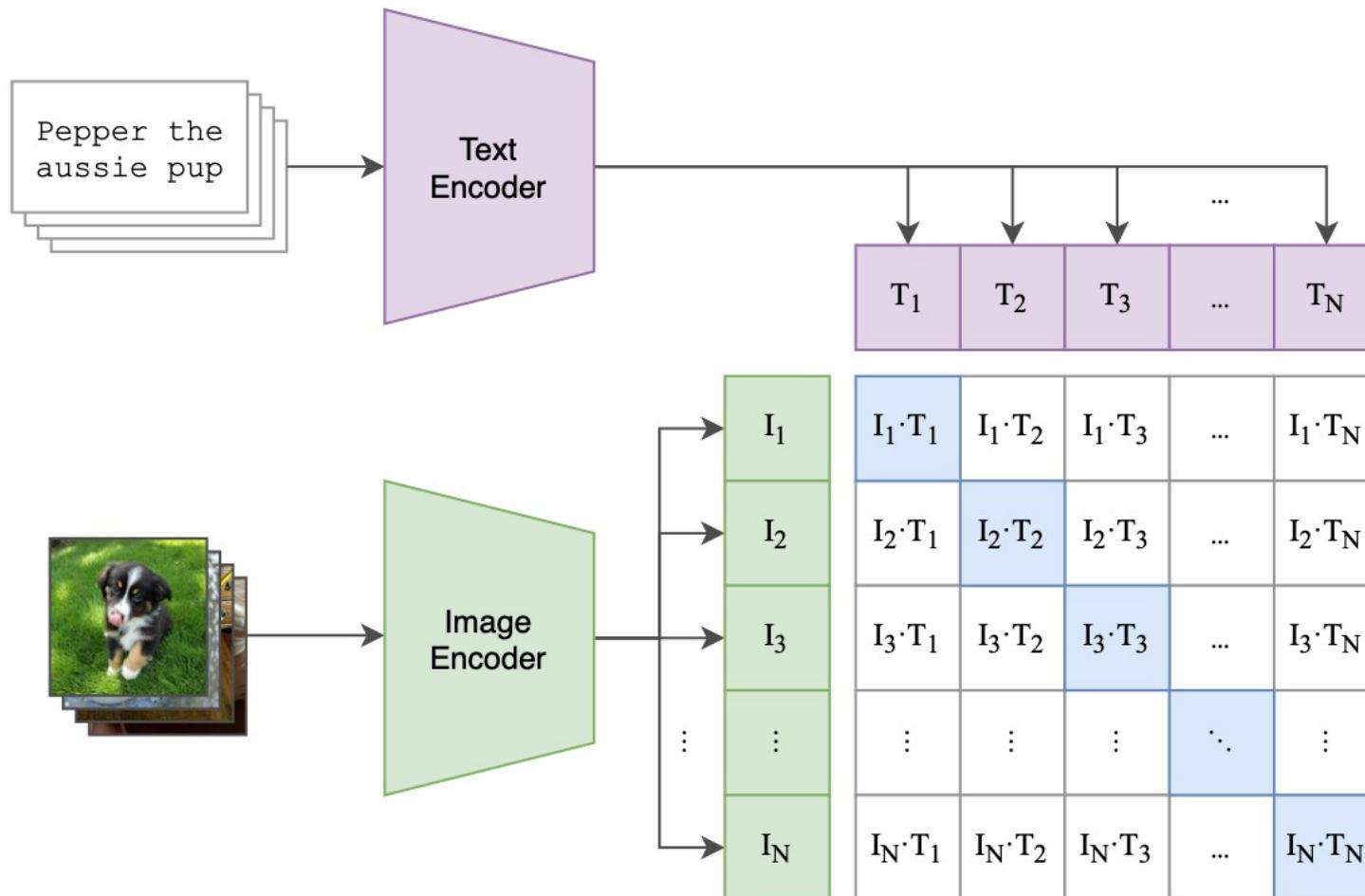


Side-quest: CLIP

<https://arxiv.org/abs/2103.00020>

Data: a large dataset of images with text captions

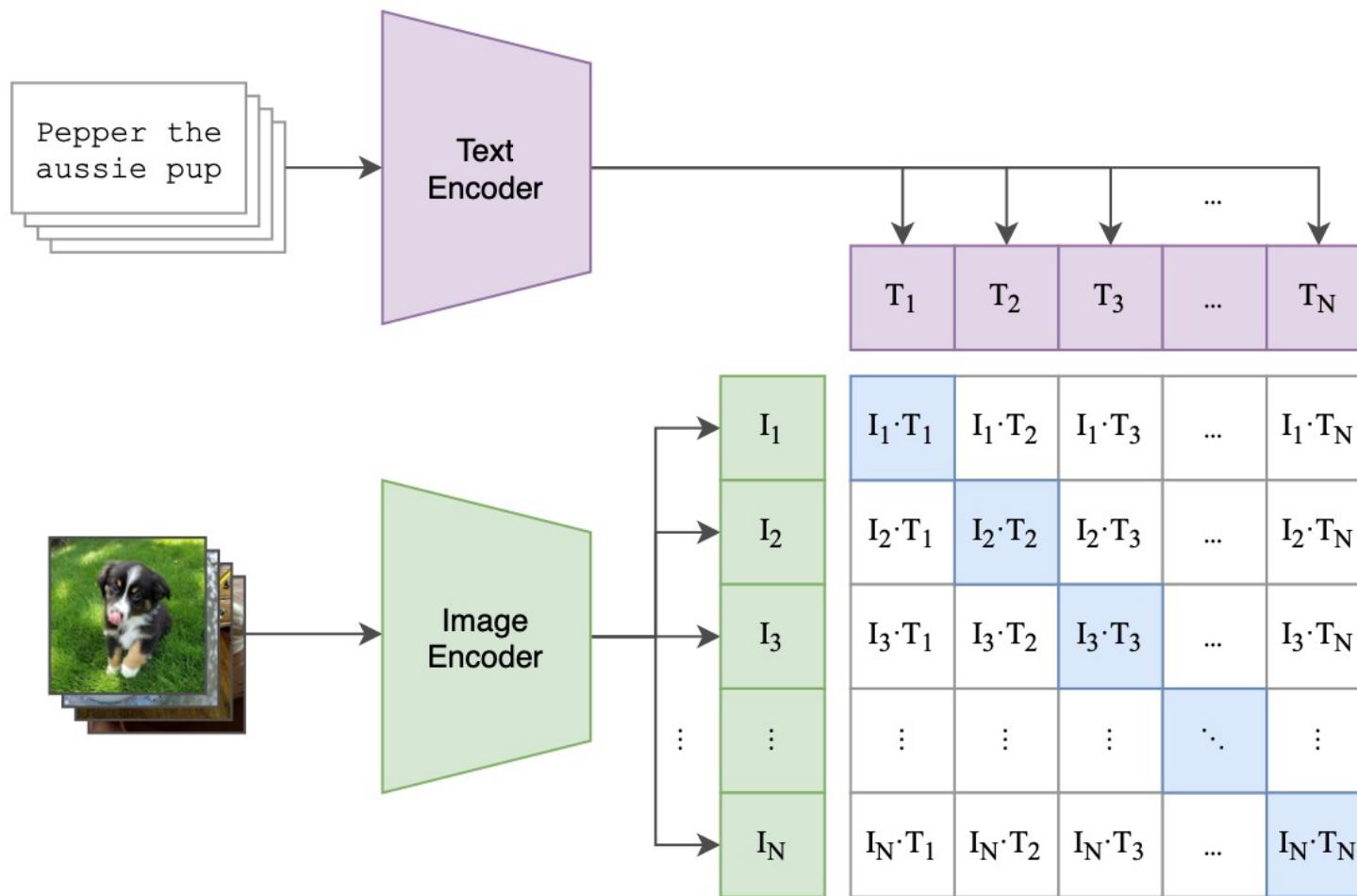
Objective: image and it's caption must have close vectors
different images/captions should not have close vectors



Side-quest: CLIP

<https://arxiv.org/abs/2103.00020>

Data: a large dataset of images with text captions
Loss function: contrastive, similar to SimCLR



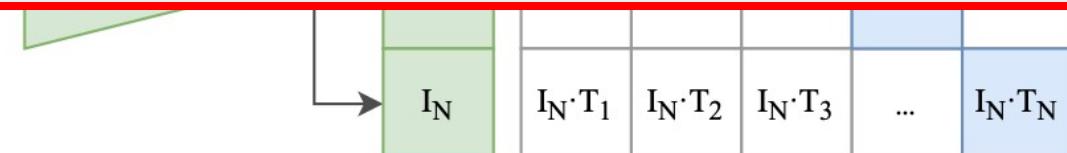
Side-quest: CLIP

<https://arxiv.org/abs/2103.00020>

Data: a large dataset of images with text captions

Loss function: contrastive, similar to SimCLR

More on CLIP later in the course!



What's the best method?

There are (small) improvements every year...

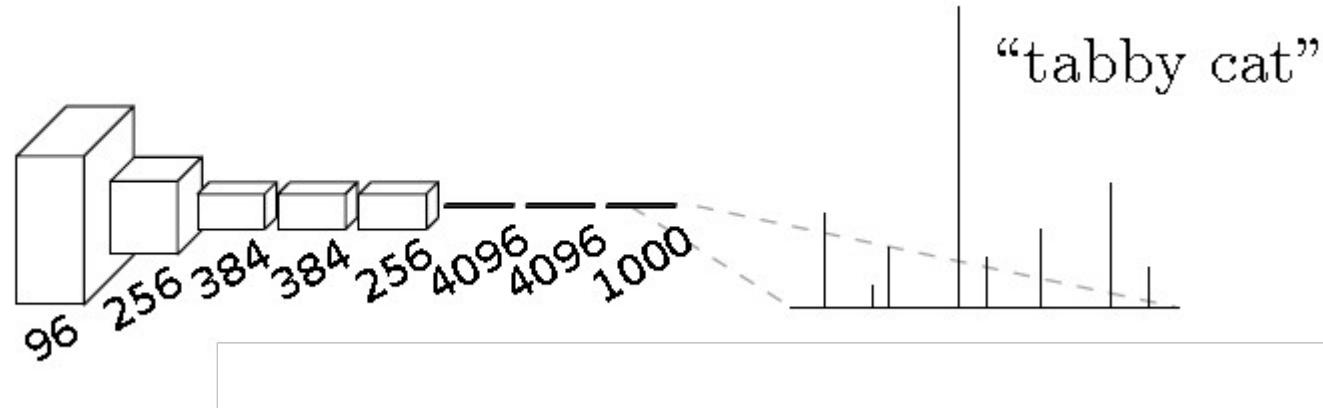
Let's go look at paperswithcode :)

Alright, we did it!

What else can we do besides image classification?

Semantic segmentation

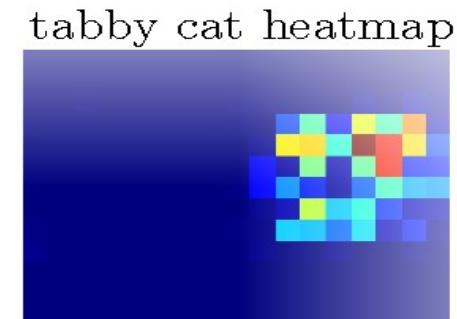
- Classification



- Segmentation/Detection



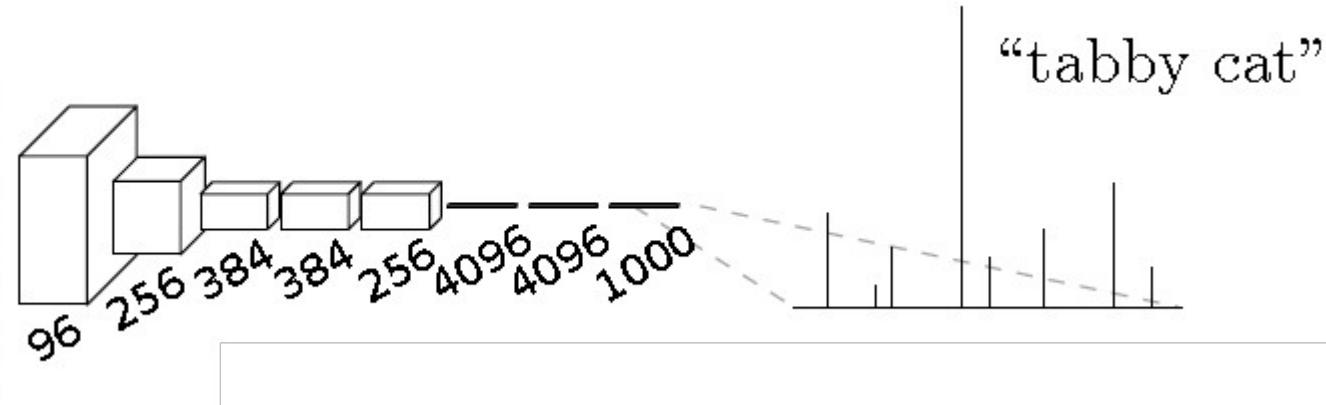
???
Ideas on how to
predict that?



- Pics: V. Lempitsky

Semantic segmentation

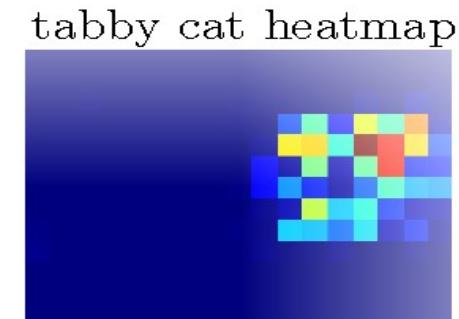
- Classification



- Segmentation/Detection (naïve)



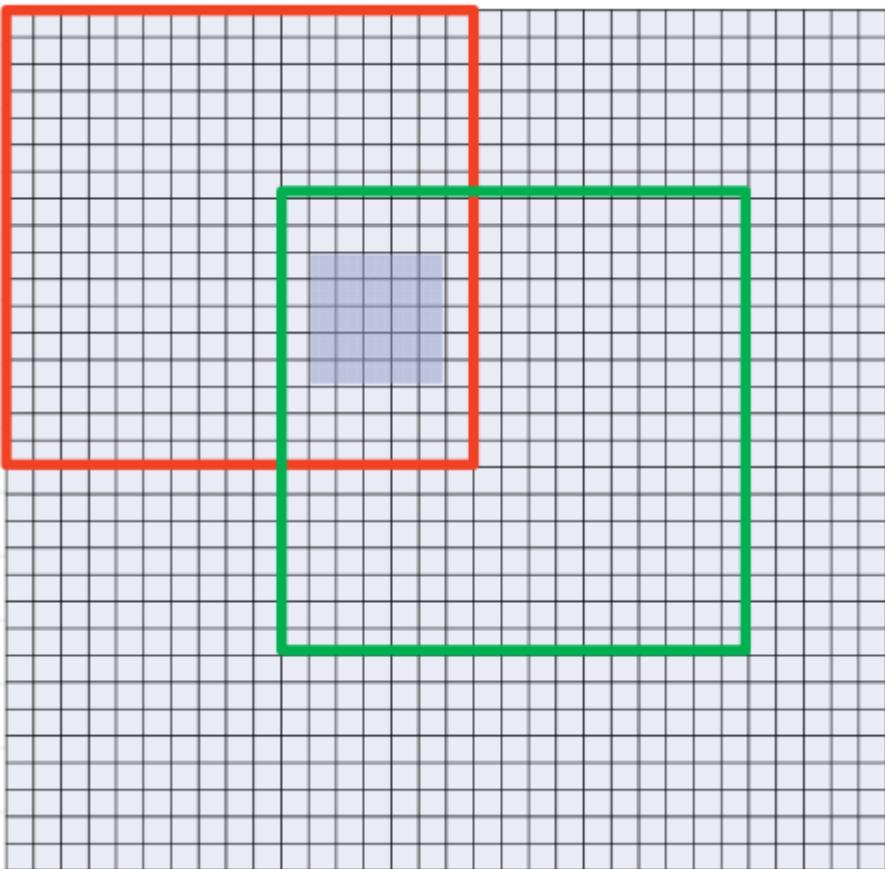
Apply to each spot



"convolution with whole
NN as a filter"

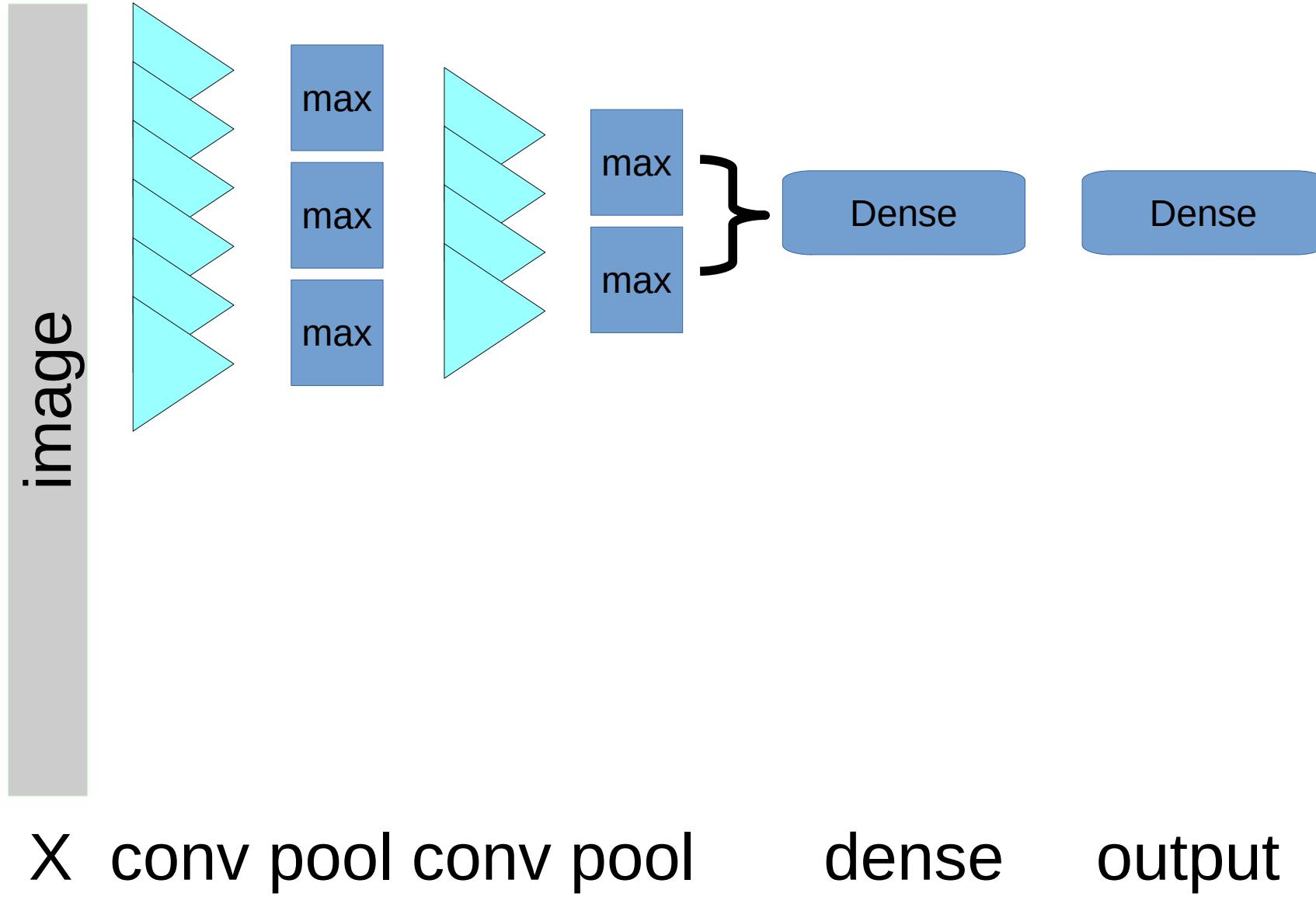
- Pics: V. Lempitsky

Convolutionalization

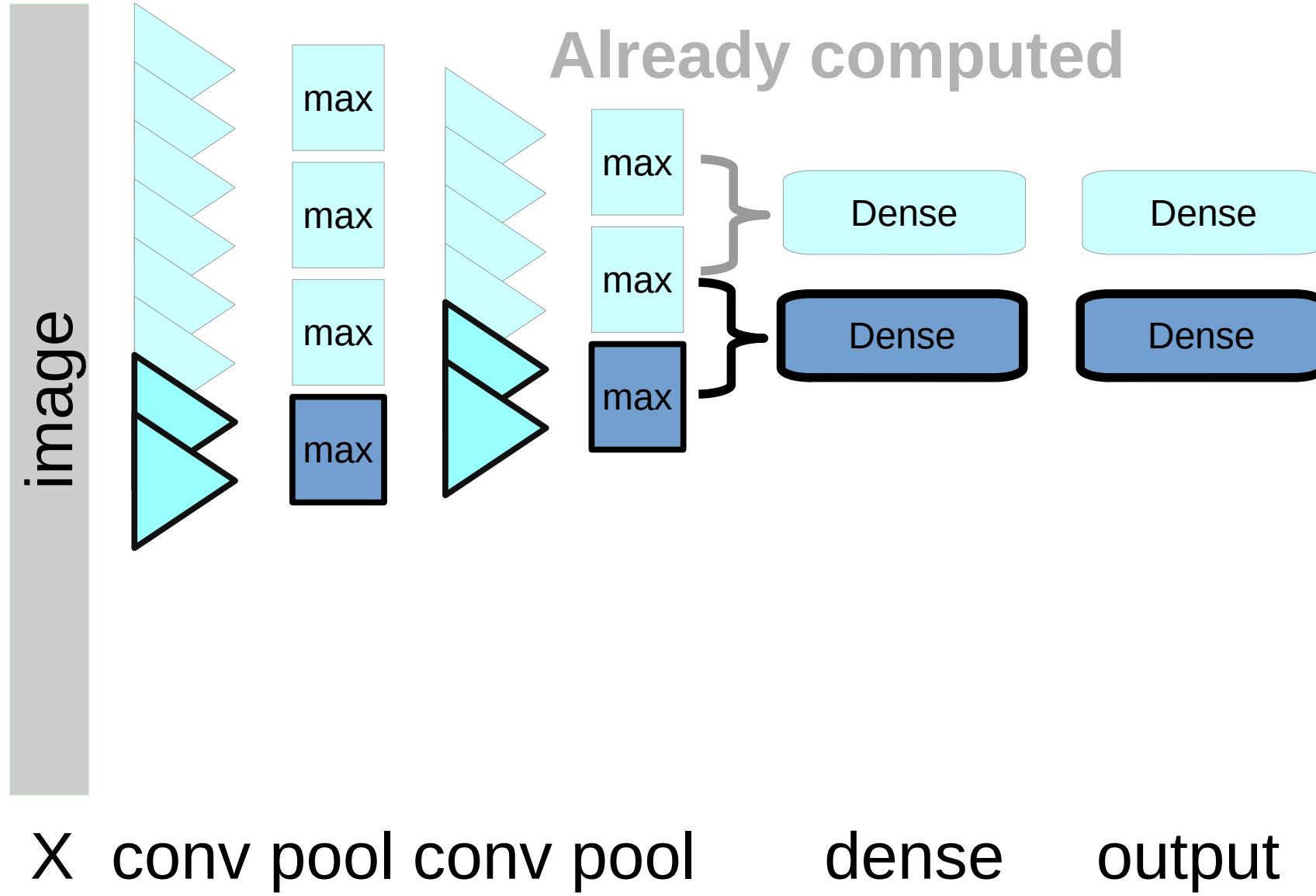


- Idea: if you apply the NN to neighboring patches, you'll have to compute same convolutions again
- Instead let us precompute convolutions for the whole image in advance.

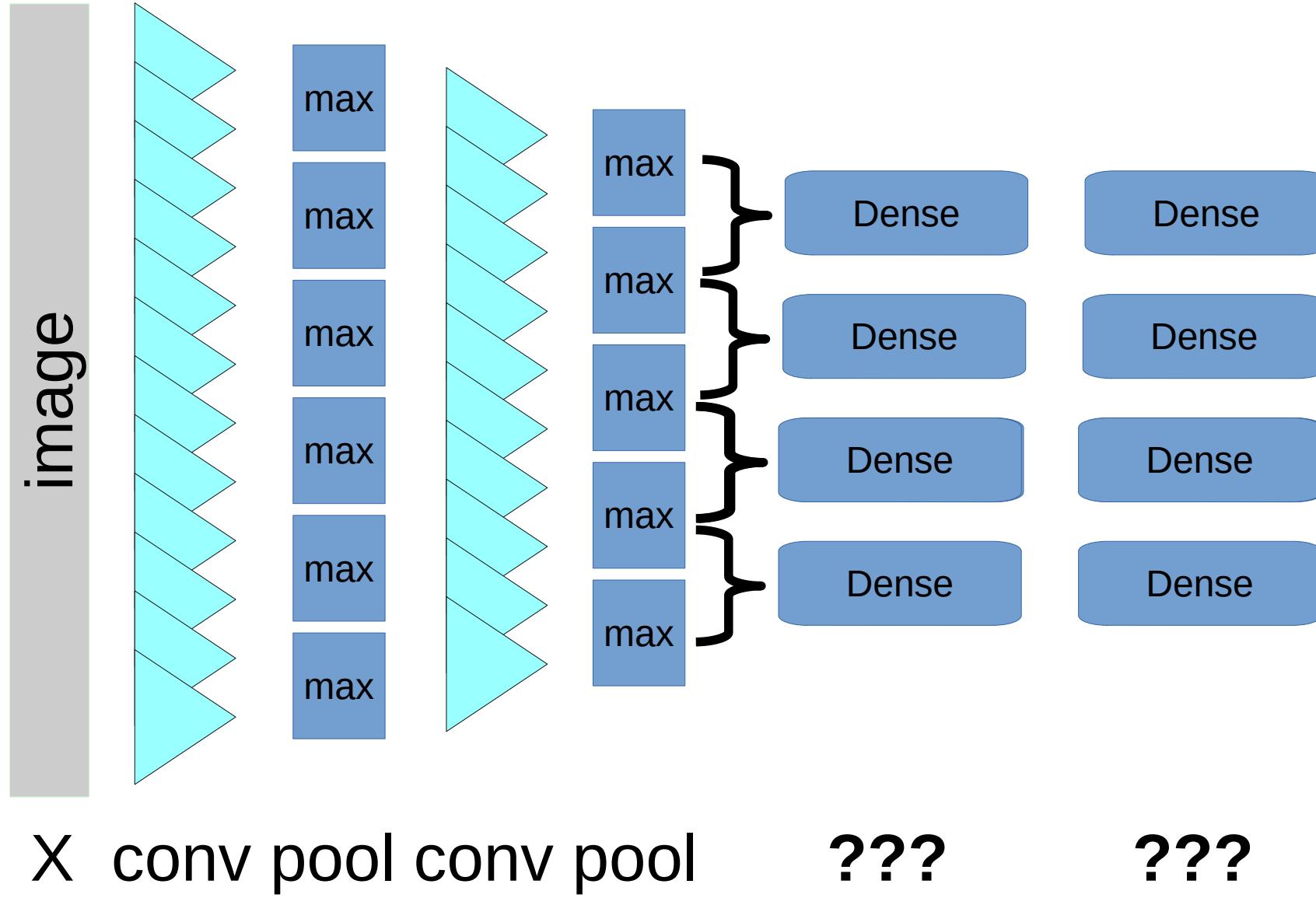
Convolutionalization



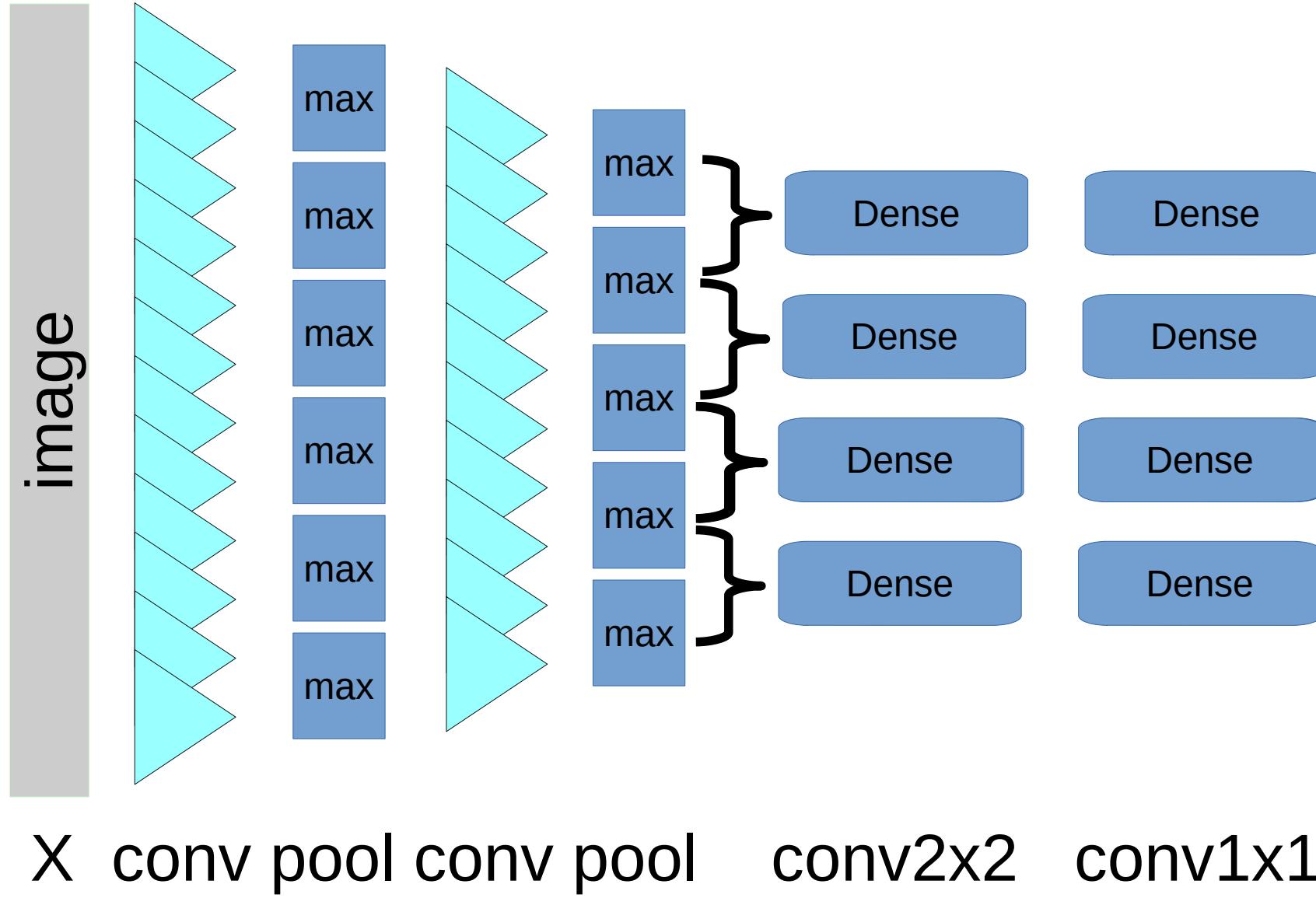
Convolutionalization



Convolutionalization

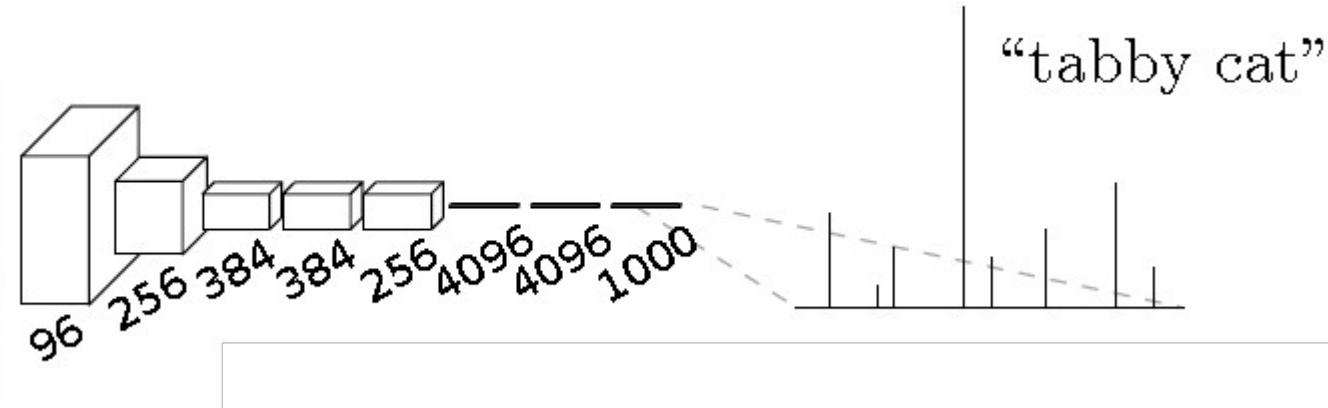


Convolutionalization

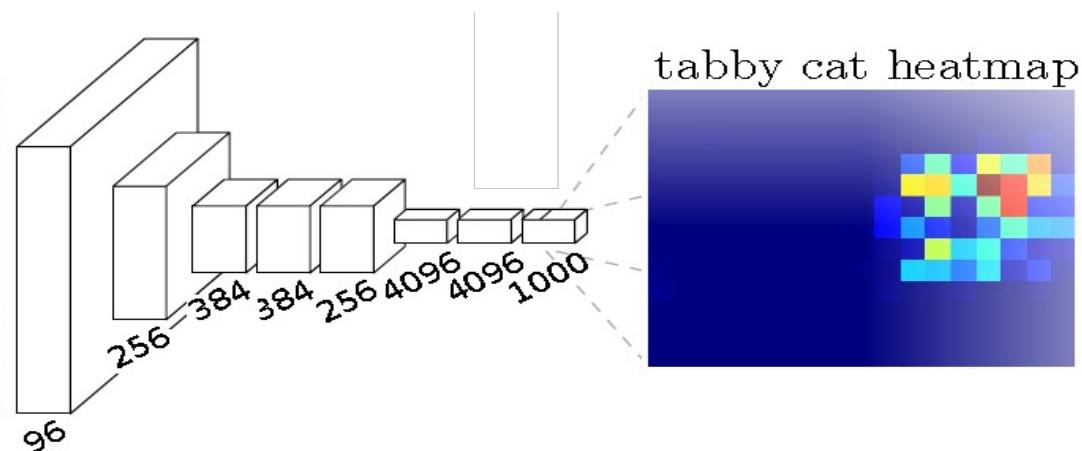


Semantic segmentation

- Classification



- Segmentation/Detection



- Pics: V. Lempitsky

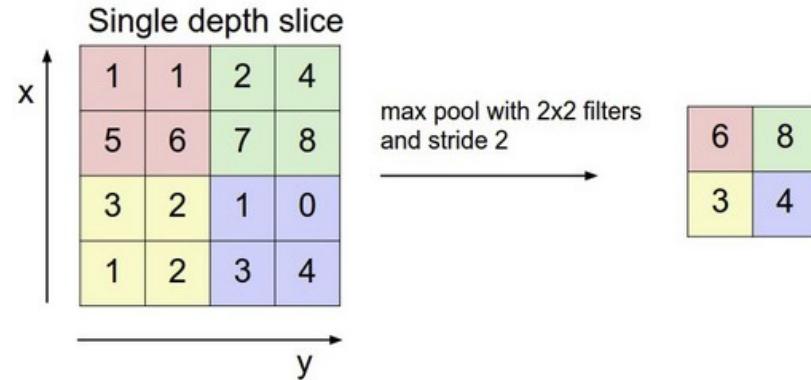
What if output must have same resolution as in
the input image? Or even bigger?

How can we **increase** activation size?

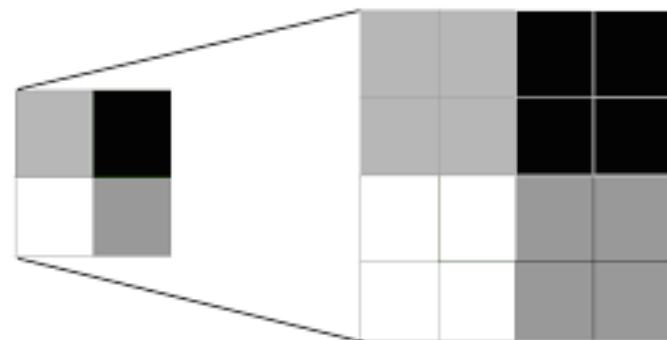
Ideas?

Upscale

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

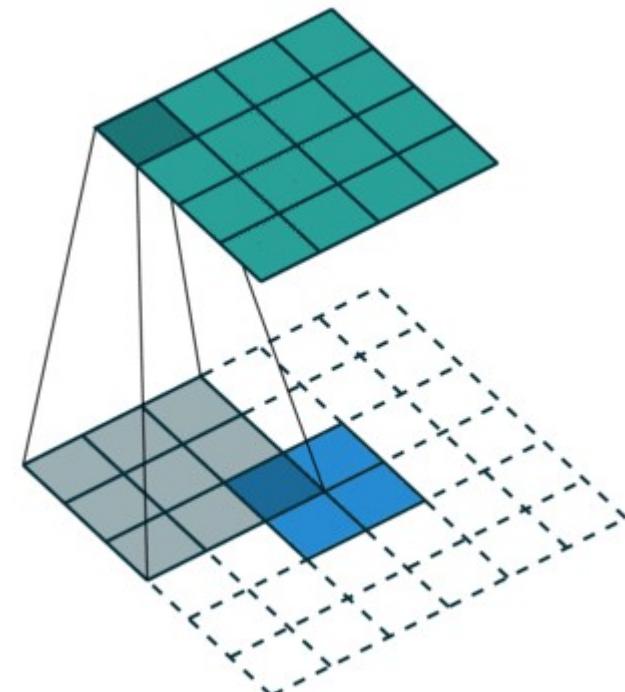


- Upscaling / Unpooling / Upsampling
 - Take one value
 - Output $(n*m)$ region
 - Several strategies



Deconvolution

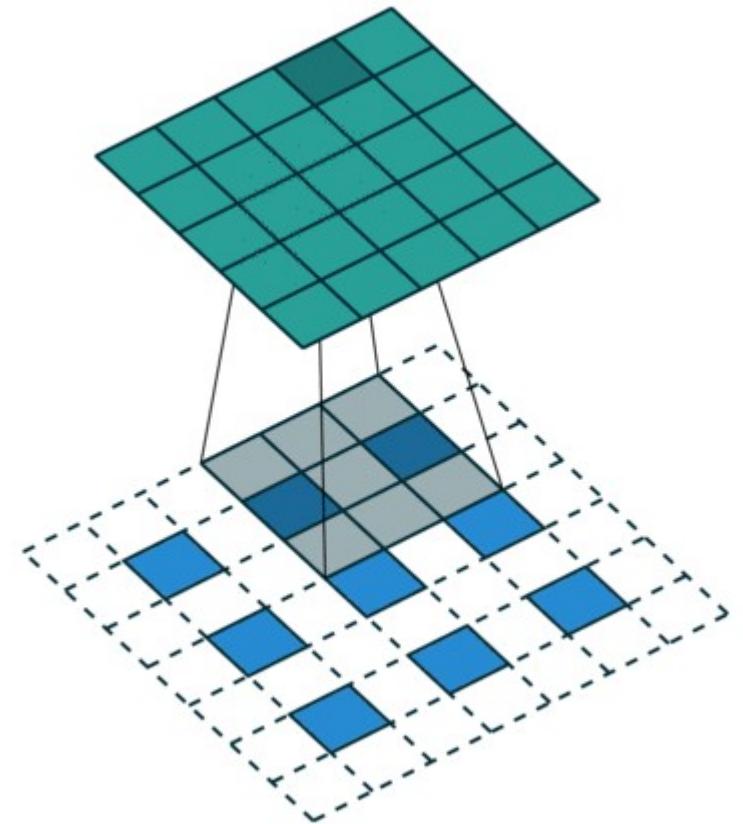
- The “inverse” of convolution
 - = transposed convolution
 - ~ wide convolution
- `deconv(conv(img))`
- preserves image shape
- Stride ~ upscale



deconvolution

Deconvolution

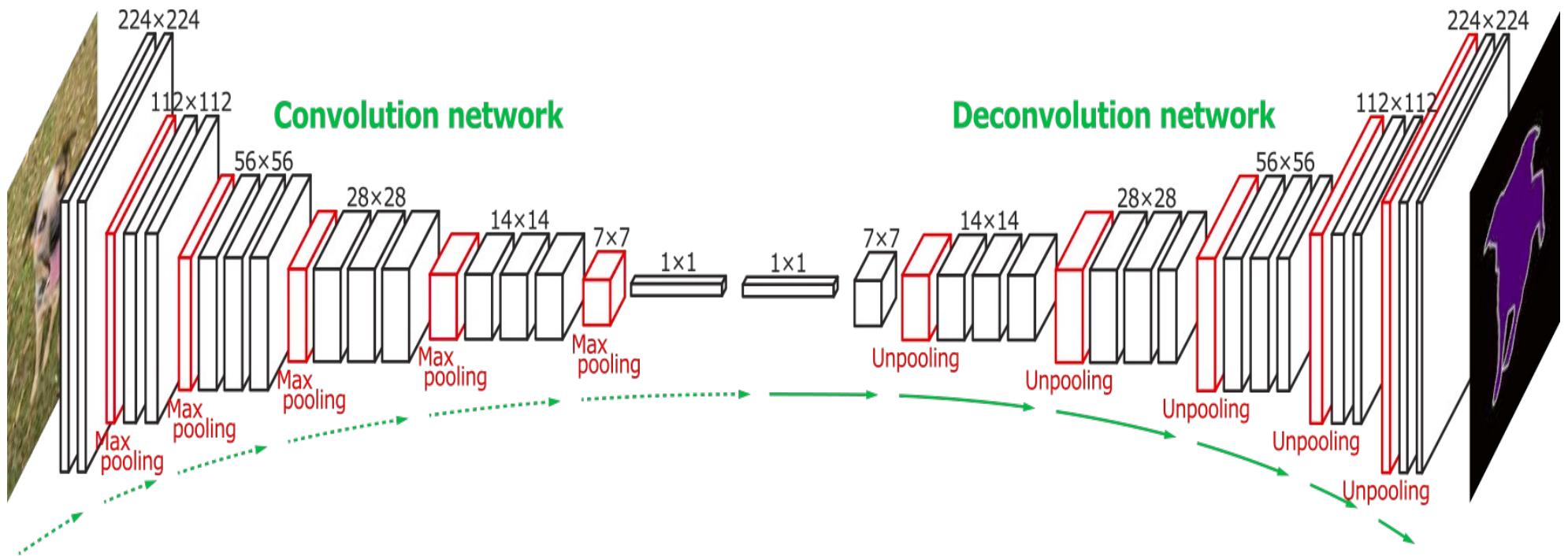
- The “inverse” of convolution
 - = transposed convolution
 - ~ wide convolution
- `deconv(conv(img))`
- preserves image shape
- Stride ~ upscale



strided deconvolution

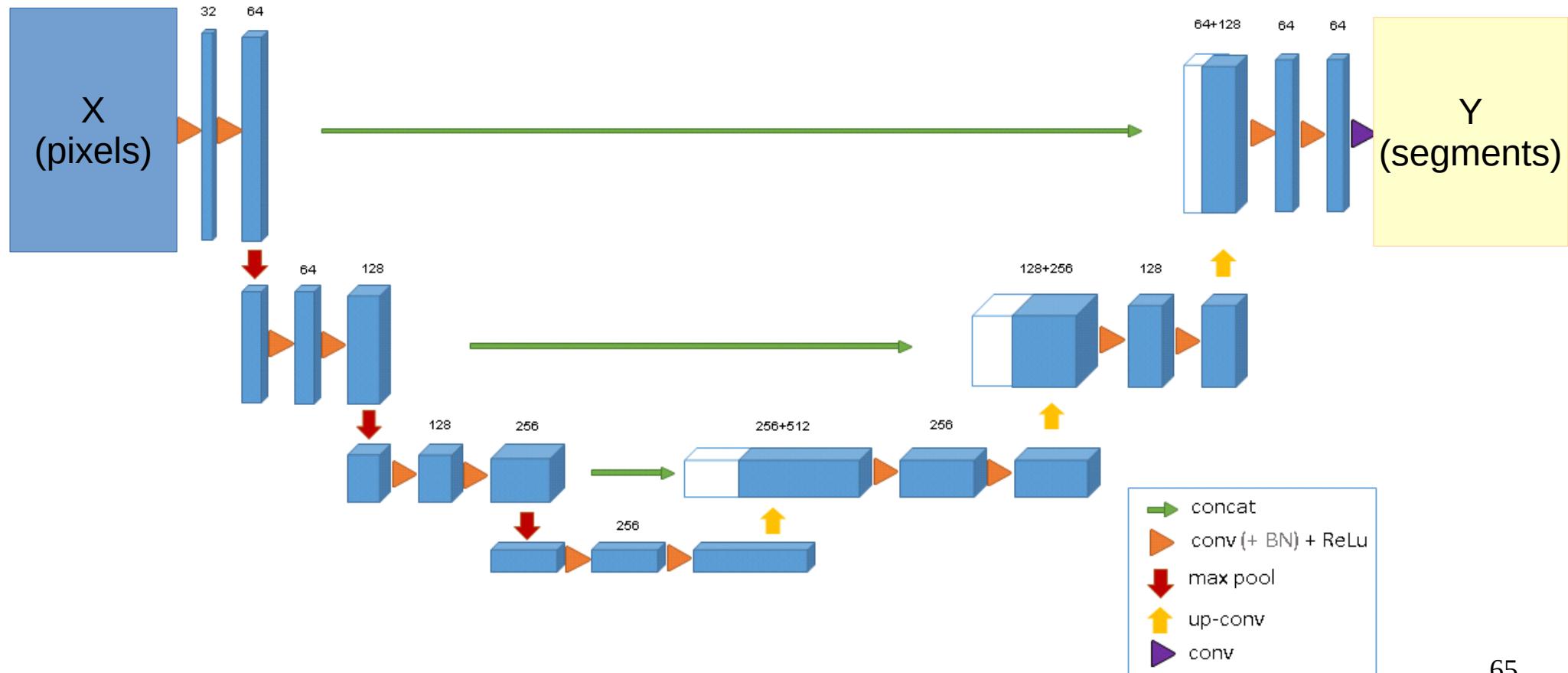
63

Fully-convolutional

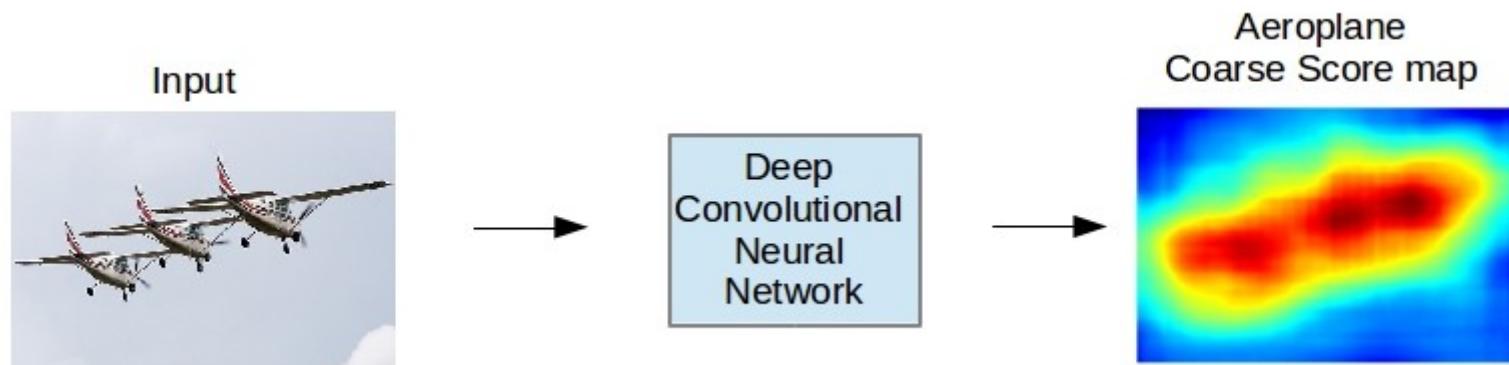


The U-net

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



Raw network output is usually rather blurred

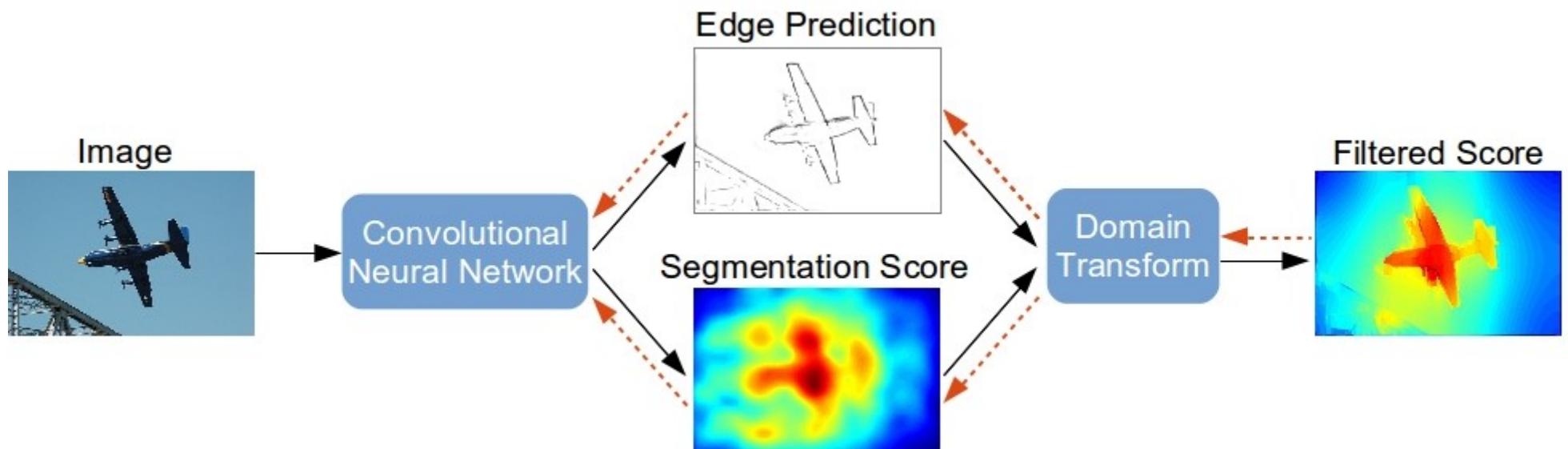


How can we get sharper edges?



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

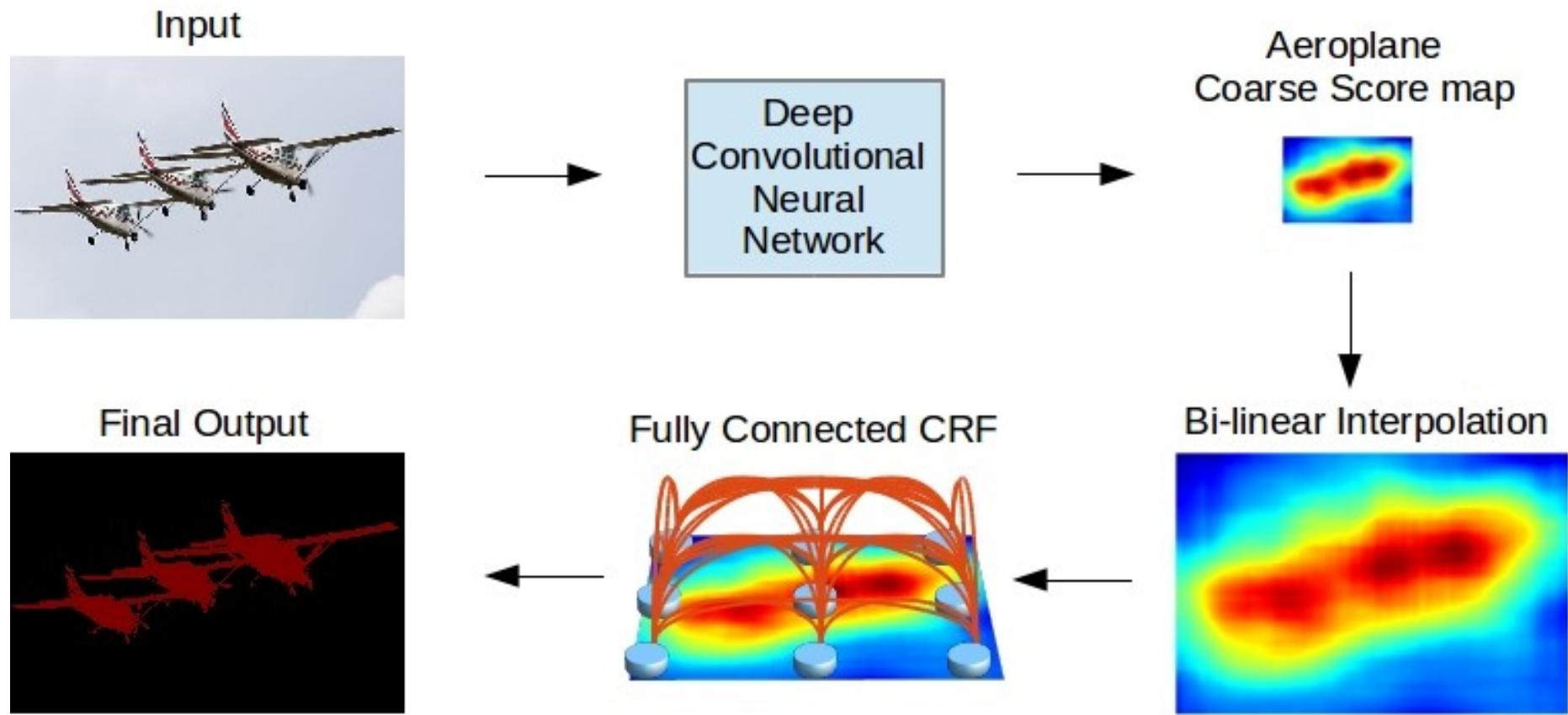
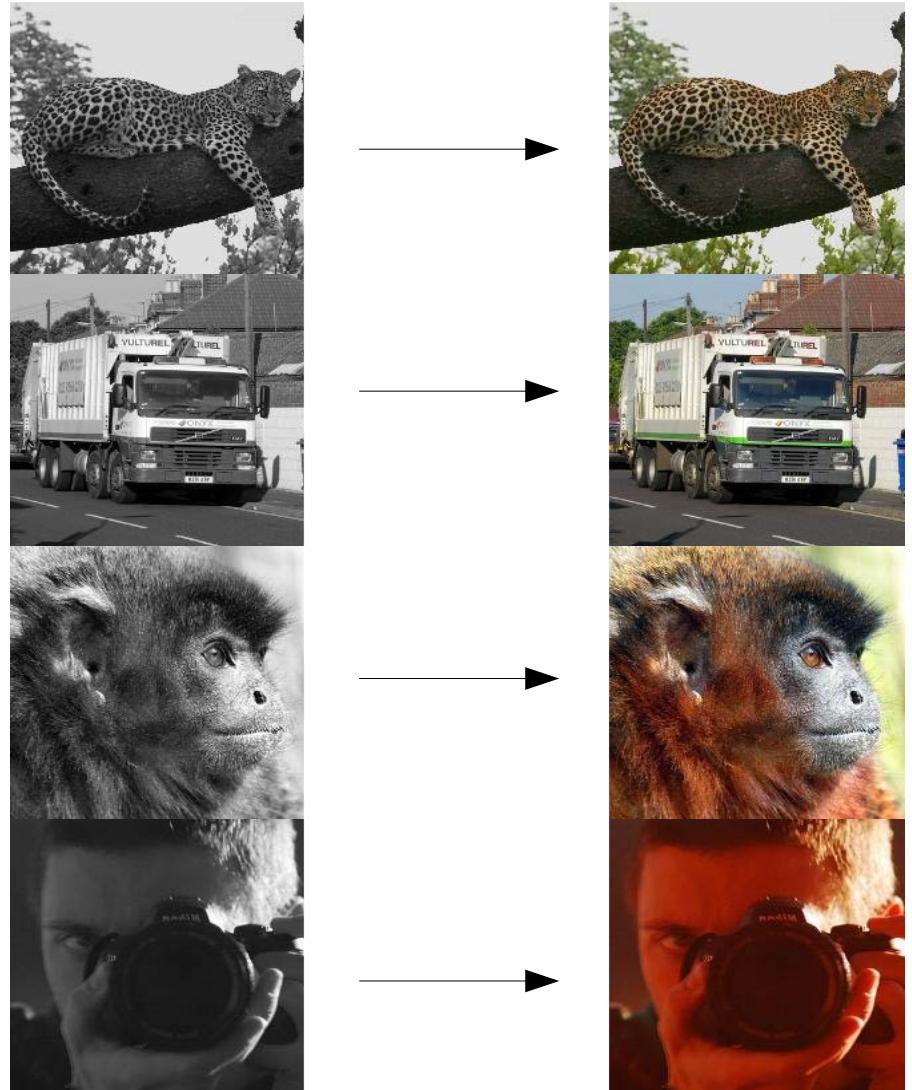


Image coloring

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



Dataset

Image coloring: hypercolumns

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

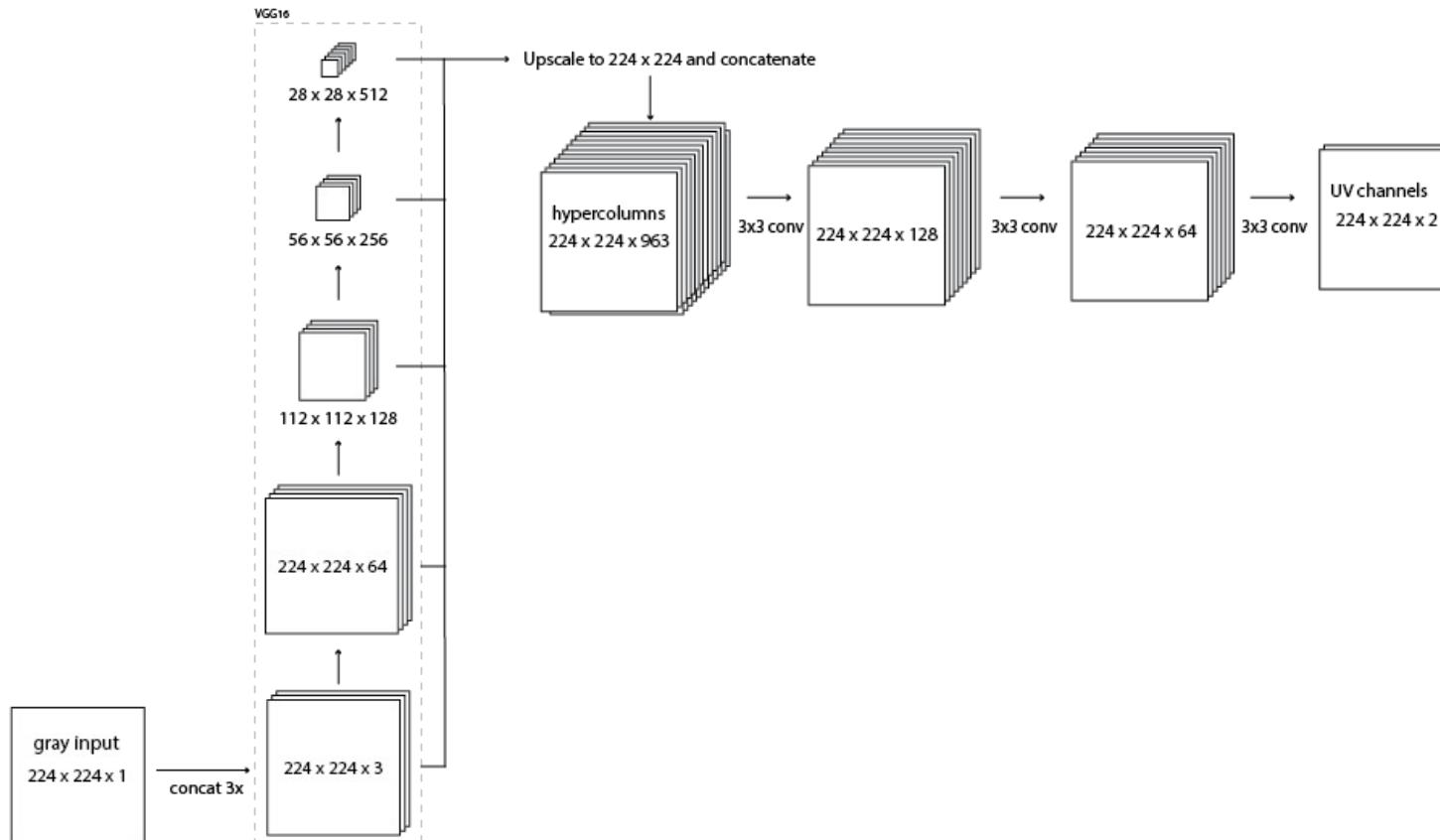


Image coloring: residual

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

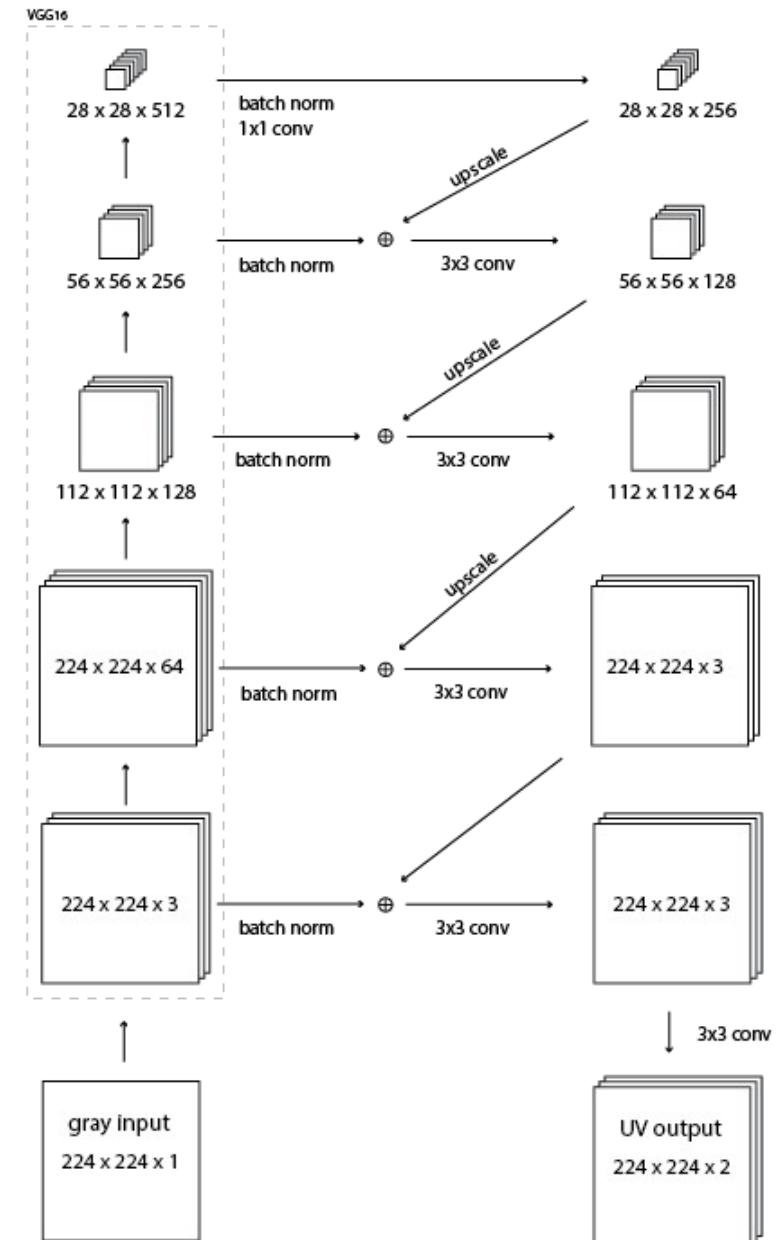


Image coloring: examples

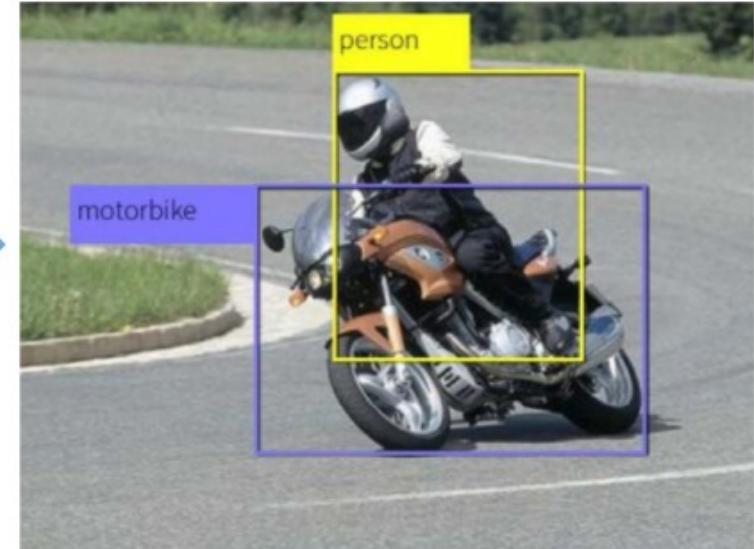


Input/prediction/reference

Bounding Box Regression



Input image

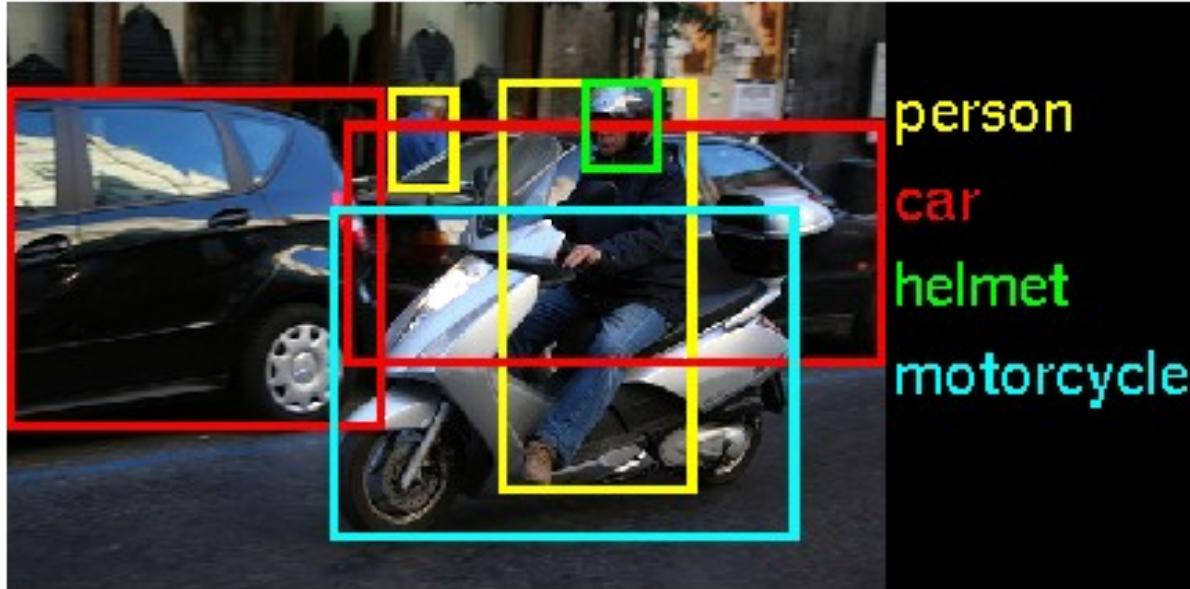


Desired output

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

Ideas?

Bounding Box Regression

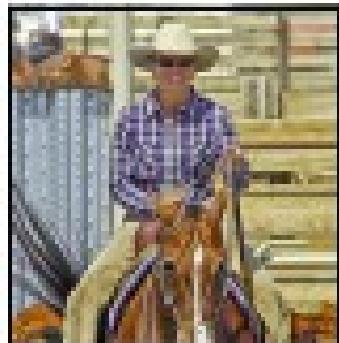


- Multiple squares?
- Non-square shape?
- Arbitrary rotation?

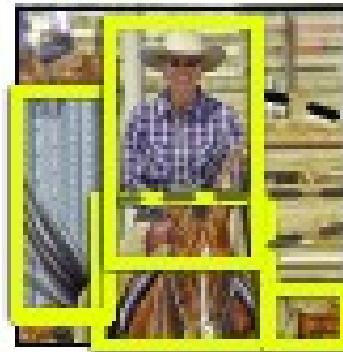
RCNN

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

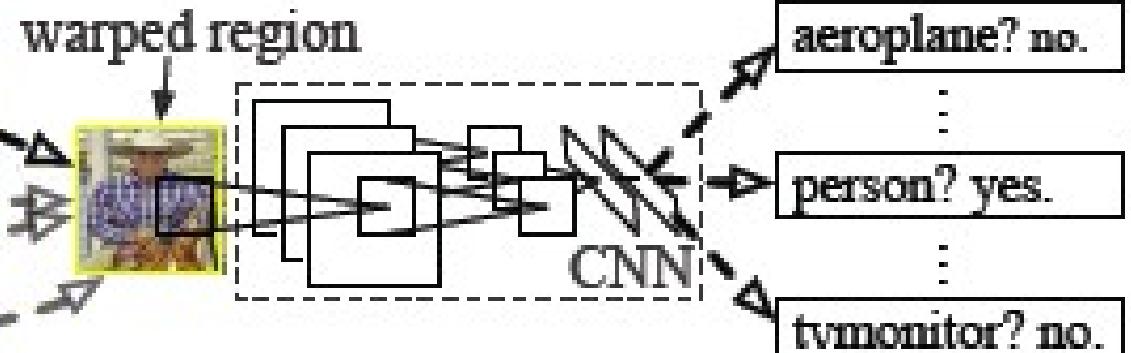
R-CNN: *Regions with CNN features*



1. Input image



2. Extract region proposals (~2k)

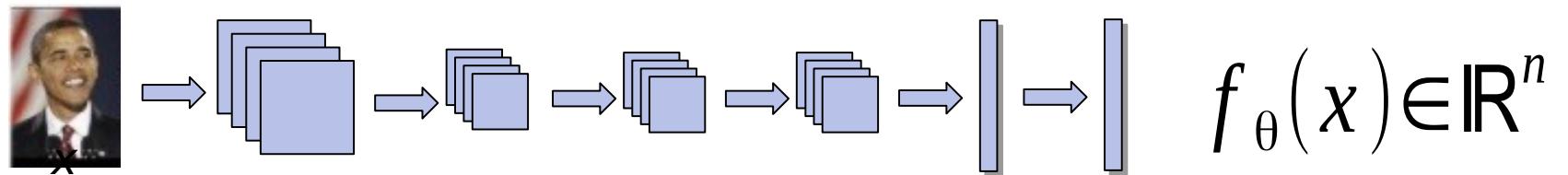


3. Compute CNN features

4. Classify regions

Image verification

- Image2vec



- We want:

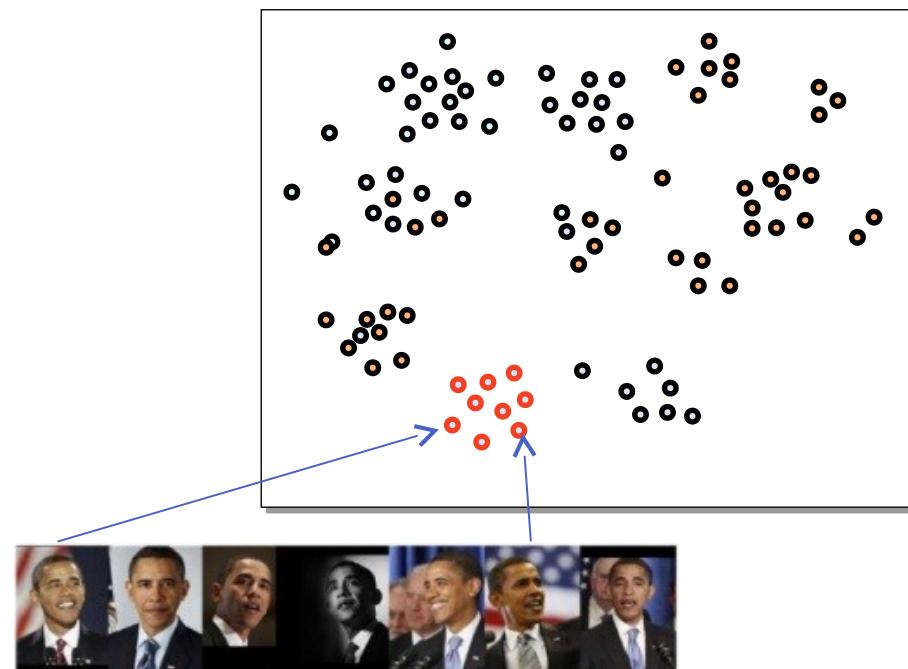
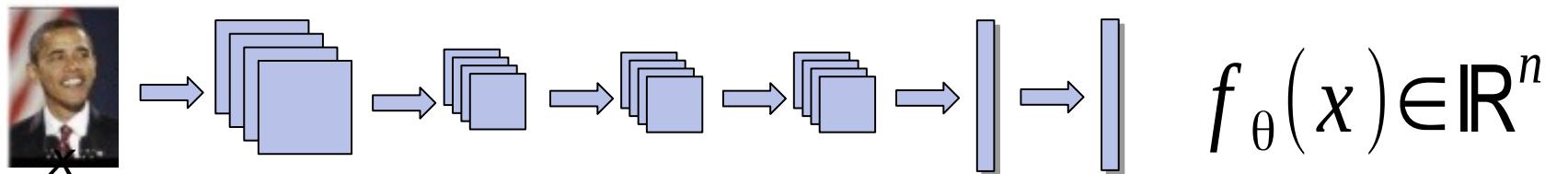


Image verification

- 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:

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

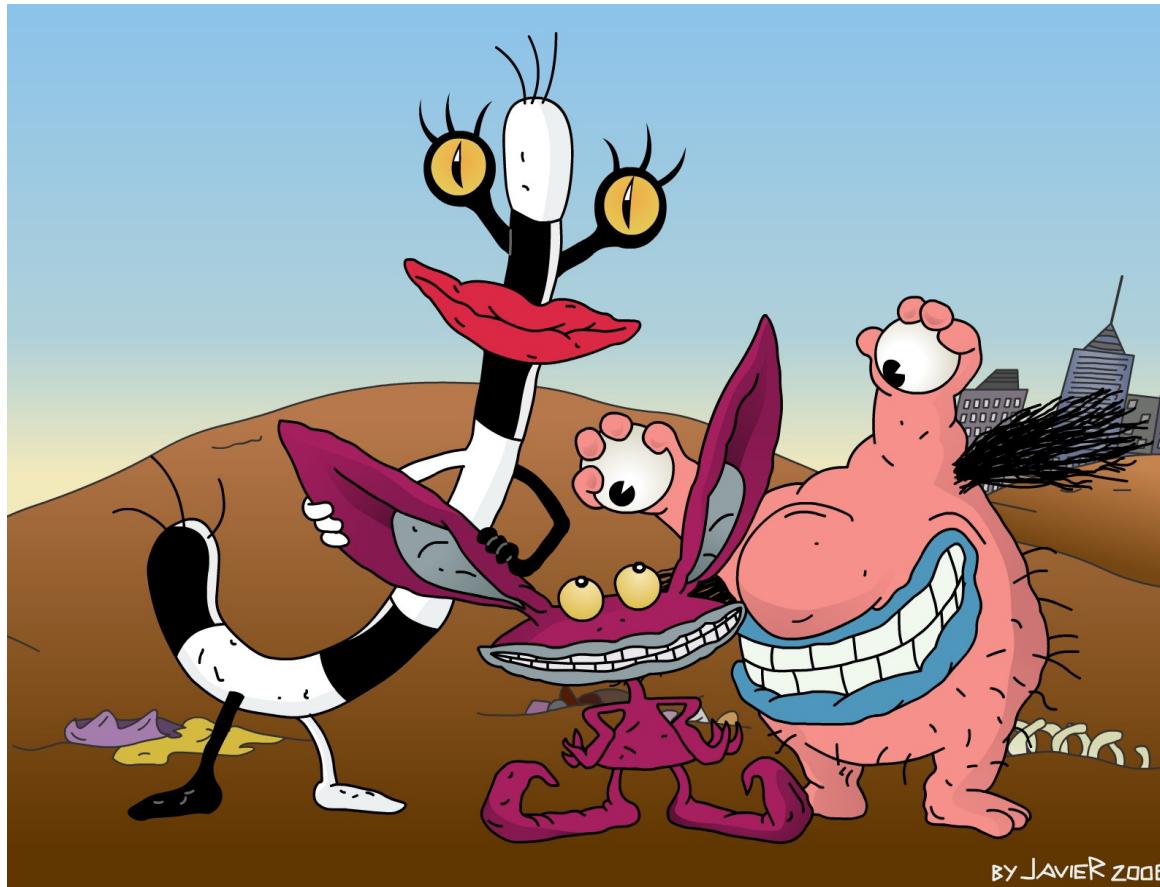
Speed up:

- 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



by JAVIER 2006

Model zoo in today's seminar