



convolutional_neural_

networks

(CNN)

stats403_deep_learning
spring_2025
lecture_4

4.1 more on regularization

regularization

- neural networks can approximate “any” continuous function
- a big number of combinations of parameters that we need to choose from
- need to restrict the capacity of a neural network

regularization

suppose the objective function $J(\theta; X, y)$ in training is

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$



penalty
term

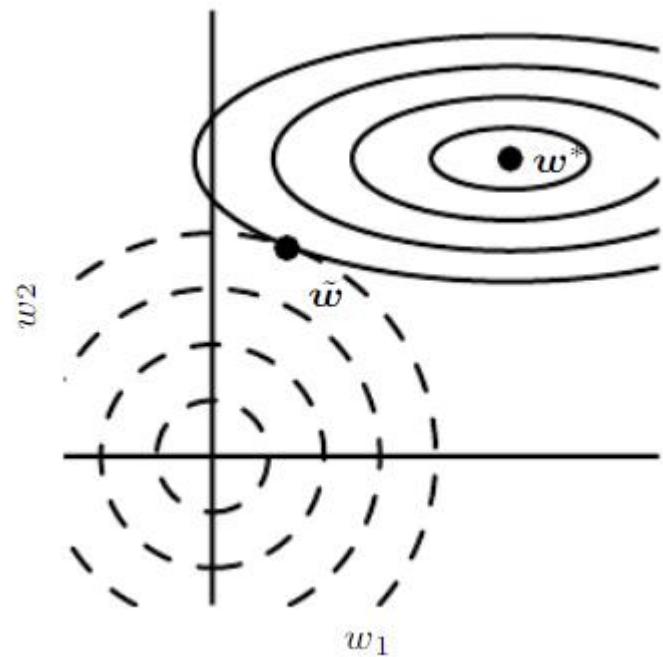
L2 regularization

ridge regression / Tikhonov regression

$$\Omega(\theta) = \|\boldsymbol{w}\|_2^2 / 2$$

$$\tilde{J}(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) = J(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) + \frac{\alpha}{2} \boldsymbol{w}^T \boldsymbol{w}$$

L2 regularization



question: why does the solution move more significantly along the direction of w_1 ?

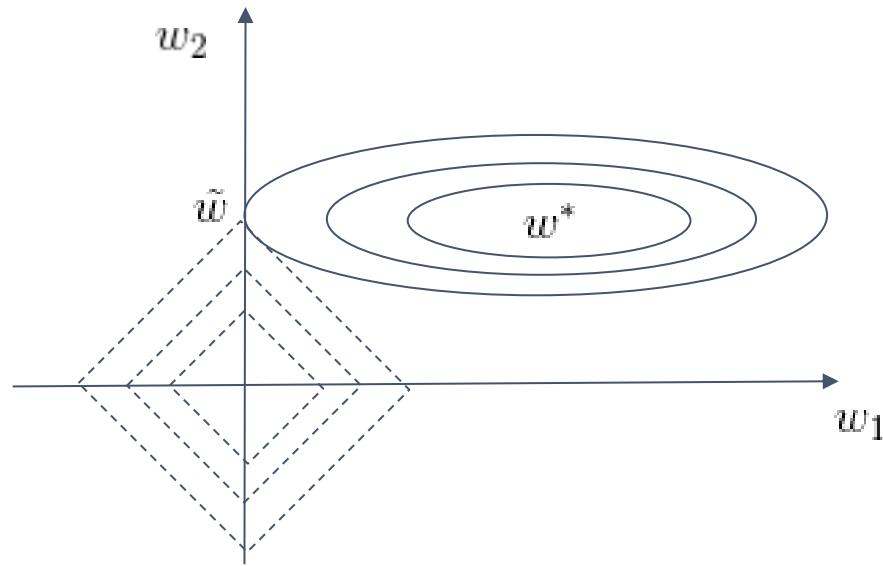
L1 regularization

$$\Omega(\theta) = \|w\|_1$$

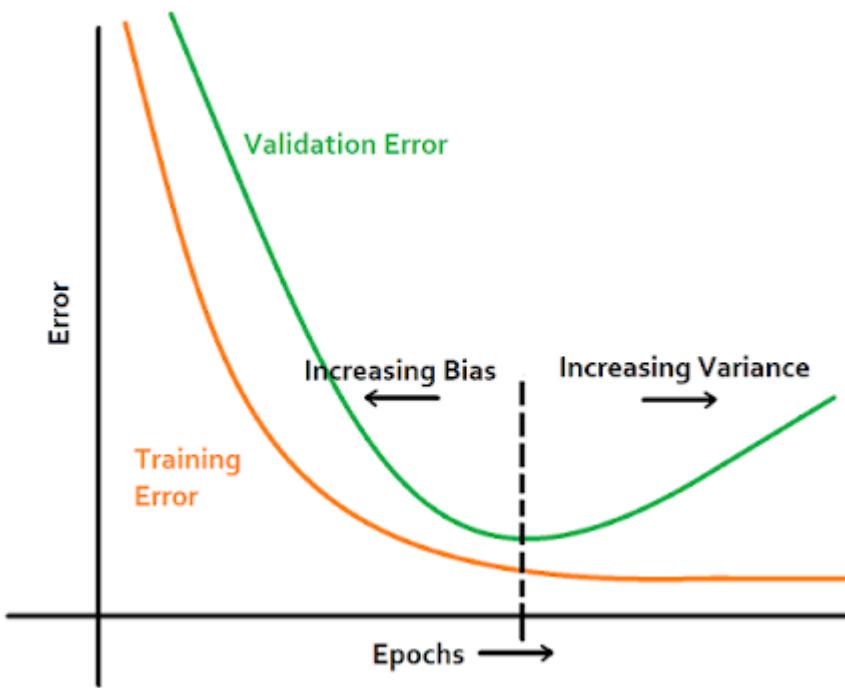
$$\tilde{J}(\theta; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \alpha \|w\|_1$$

L1 regularization

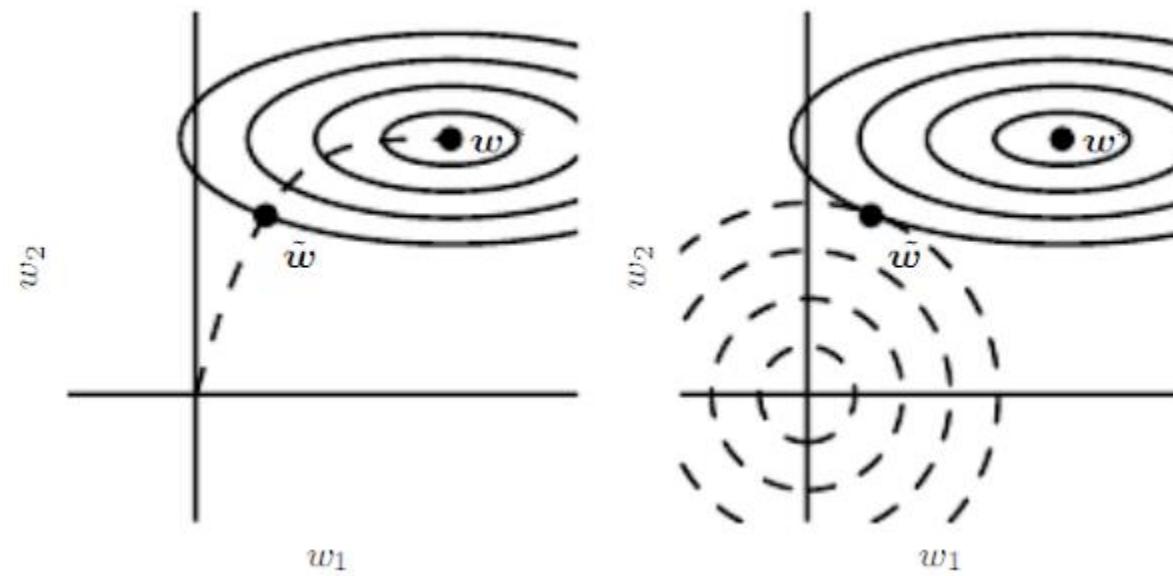
encourages sparsity of solution



early stopping



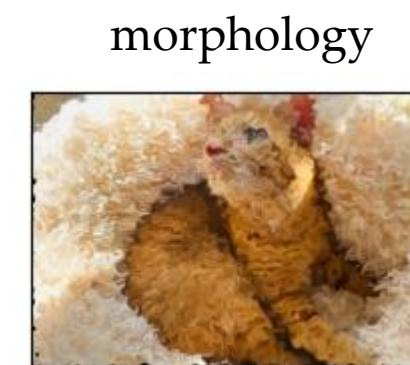
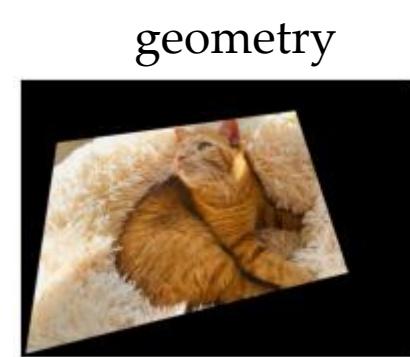
early stopping is similar to regularization



dataset augmentation

- One effective way to prevent overfitting is to use more data.
- In practice, we do not have an infinite amount of data.
- A practical solution is to generate data by ourselves.

dataset augmentation



dataset augmentation

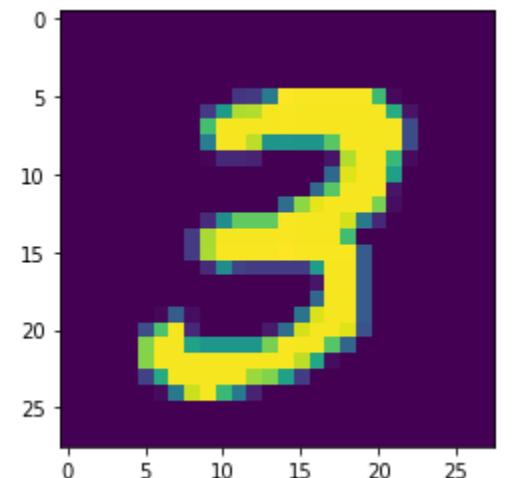
- mixup



label smoothing

- instead of one-hot labels, we use a smooth label such as

[0.01, 0.01, 0.01, 0.91, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01]



preprocess image data



preprocess image data

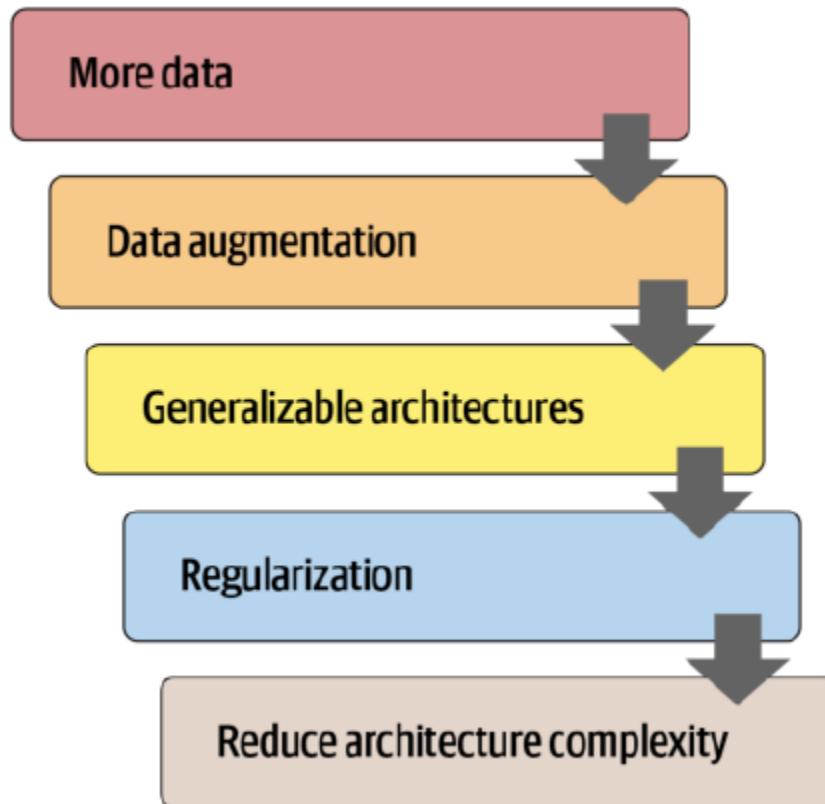
- normalizing:
 - before normalizing, the pixel values are from 0 to 255
 - need to apply a transformation so that each pixel has zero mean and unit variance

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	0	0	0	0	0	29	150	195	254	255	254	176	193	150	96	0	0	
2	0	0	0	48	166	224	253	253	234	196	253	253	253	253	233	0	0	
3	0	83	244	249	253	187	46	10	8	4	10	194	253	253	233	0	0	
4	0	107	253	253	230	48	0	0	0	0	0	192	253	253	156	0	0	
5	0	3	20	20	15	0	0	0	0	0	43	224	253	245	74	0	0	
6	0	0	0	0	0	0	0	0	0	0	249	253	245	126	0	0	0	
7	0	0	0	0	0	0	0	14	101	223	253	248	124	0	0	0	0	
8	0	0	0	0	0	11	166	239	253	253	253	187	30	0	0	0	0	
9	0	0	0	0	0	16	248	250	253	253	253	253	232	213	111	2	0	
10	0	0	0	0	0	0	0	43	98	98	208	253	253	253	187	22	0	

parameter sharing

- force sets of parameters to be equal
- reduces the number of parameters and therefore reduces overfitting
- typical example: Convolutional Neural Networks (CNN)

reduce overfitting



4.2 convolutional neural networks (CNN)

using NN for image data



using NN for image data

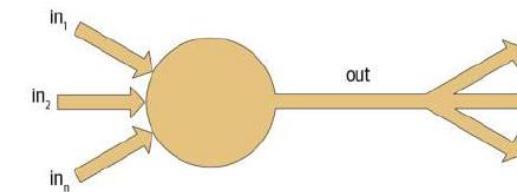
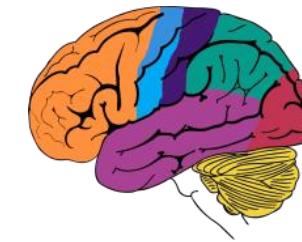


image as 2D function

- instead of flatten an image into a vector, we can treat an image as a function over a two-dimensional space $f(x, y)$
- that is, for each pixel position (x, y) , the function value is $f(x, y)$

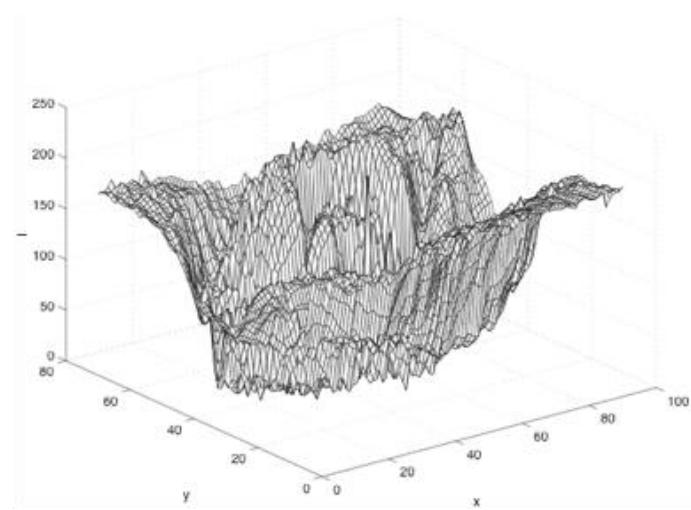
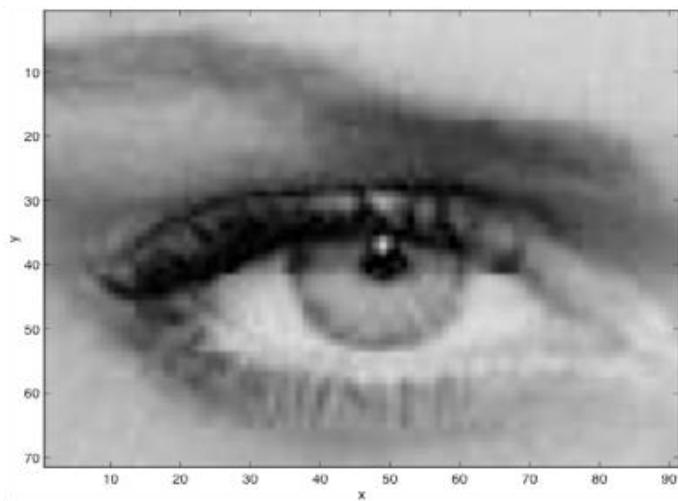


image as 2D function

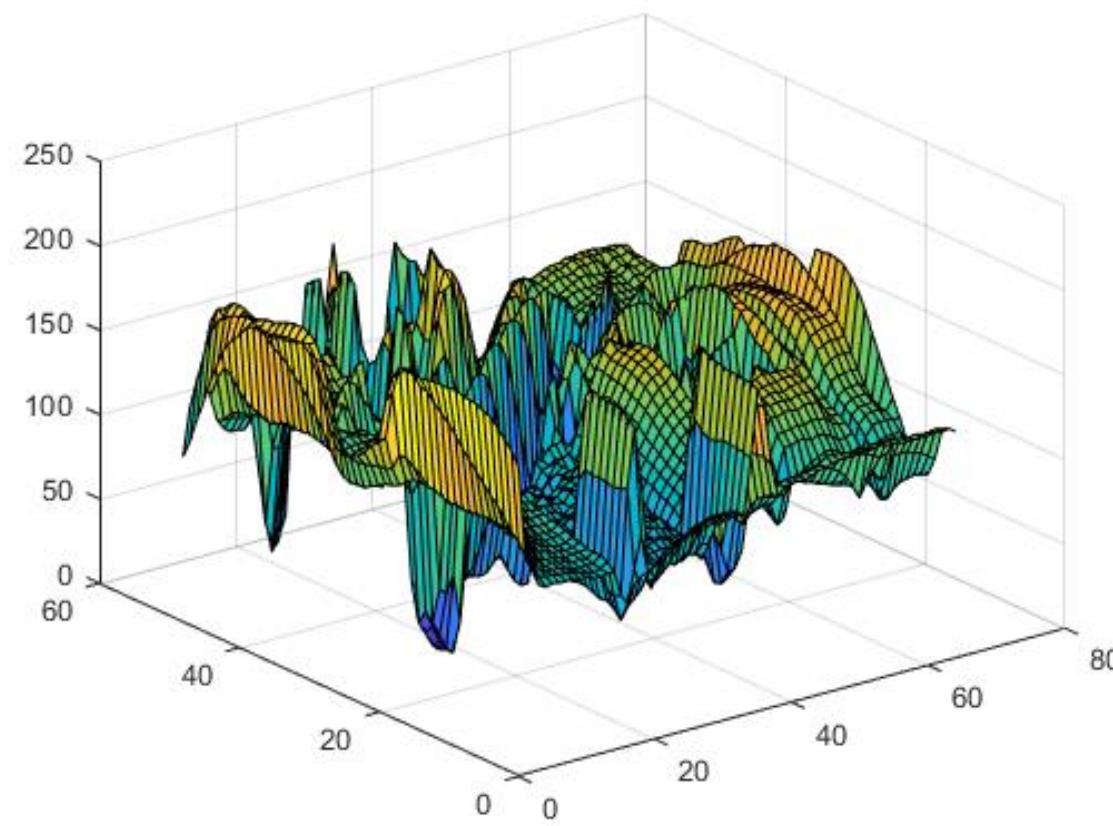
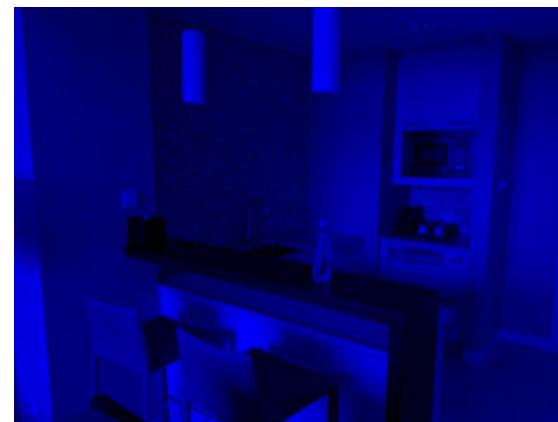
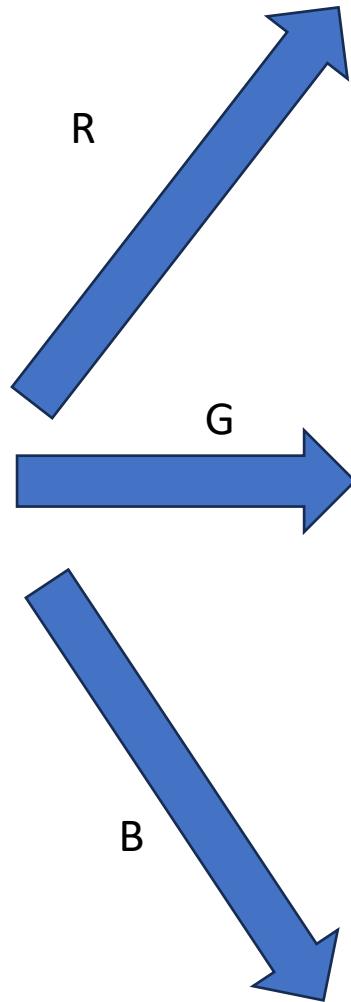
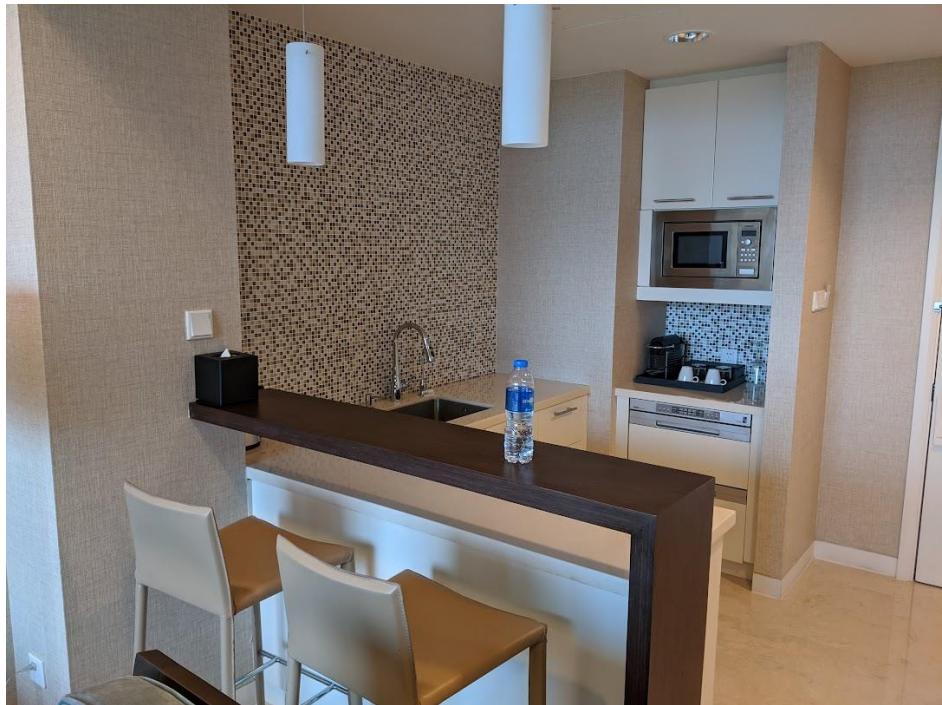


image as 2D function



RGB channels



convolution / neighborhood filtering

- in convolution (a.k.a. neighborhood filtering), a two-dimensional function $h(x, y)$, called a **kernel**, acts on the input image $f(x, y)$, to produce an output $g(x, y)$

45	60	98	127	132	133	137	133
46	65	98	123	126	128	131	133
47	65	96	115	119	123	135	137
47	63	91	107	113	122	138	134
50	59	80	97	110	123	133	134
49	53	68	83	97	113	128	133
50	50	58	70	84	102	116	126
50	50	52	58	69	86	101	120

$f(x,y)$

$$\begin{array}{|c|c|c|} \hline 0.1 & 0.1 & 0.1 \\ \hline 0.1 & 0.2 & 0.1 \\ \hline 0.1 & 0.1 & 0.1 \\ \hline \end{array}$$

*

69	95	116	125	129	132
68	92	110	120	126	132
66	86	104	114	124	132
62	78	94	108	120	129
57	69	83	98	112	124
53	60	71	85	100	114

$h(x,y)$

$g(x,y)$

convolution / neighborhood filtering

2	7	1	8
2	8	1	8
2	8	4	5
9	0	4	5

$$\begin{matrix} 3 & 1 \\ 4 & 5 \end{matrix} * \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \end{matrix} = \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \end{matrix}$$

2	7	1	8
2	8	1	8
2	8	4	5
9	0	4	5

$$\begin{matrix} 3 & 1 \\ 4 & 5 \end{matrix} * \begin{matrix} 2 & 7 & 1 & 8 \\ 2 & 8 & 1 & 8 \\ 2 & 8 & 4 & 5 \\ 9 & 0 & 4 & 5 \end{matrix} = \begin{matrix} 61 & & & \\ & & & \\ & & & \\ & & & \end{matrix}$$

2	7	1	8
2	8	1	8
2	8	4	5
9	0	4	5

$$\begin{matrix} 3 & 1 \\ 4 & 5 \end{matrix} * \begin{matrix} 2 & 7 & 1 & 8 \\ 2 & 8 & 1 & 8 \\ 2 & 8 & 4 & 5 \\ 9 & 0 & 4 & 5 \end{matrix} = \begin{matrix} 61 & 59 & & \\ & & & \\ & & & \\ & & & \end{matrix}$$

convolution / neighborhood filtering

original



blurring



sharpening



smoothing



convolution / neighborhood filtering

$$\frac{1}{K^2}$$

1	1	...	1
1	1	...	1
\vdots	\vdots	1	\vdots
1	1	...	1

$$\frac{1}{16}$$

1	2	1
2	4	2
1	2	1

$$\frac{1}{256}$$

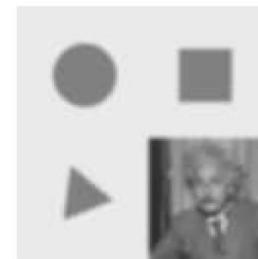
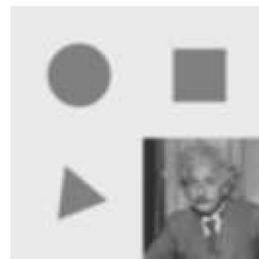
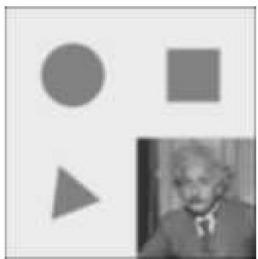
1	4	6	4	1
4	16	24	16	4
6	24	36	24	6
4	16	24	16	4
1	4	6	4	1

$$\frac{1}{8}$$

-1	0	1
-2	0	2
-1	0	1

$$\frac{1}{4}$$

1	-2	1
-2	4	-2
1	-2	1



convolution formulae

$$f * h(t) = \int f(t-u)h(u) du = \int f(u)h(t-u) du$$

↑ ↑
input kernel

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n) = \sum_m \sum_n I(i - m, j - n)K(m, n)$$

$$S(i, j) = \sum_m \sum_n I(i + m - 1, j + n - 1)K(m, n) \quad (\text{cross-correlation form})$$

convolution as matrix multiplication

$$\mathbf{x} * \mathbf{h} = \mathbf{Ax} = \begin{bmatrix} h_m & h_{m-1} & \cdots & h_1 & 0 & 0 & \cdots & 0 \\ 0 & h_m & h_{m-1} & \cdots & h_1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & & & \ddots & \vdots & \\ 0 & \cdots & 0 & h_m & h_{m-1} & \cdots & h_1 & 0 \\ 0 & \cdots & 0 & 0 & h_m & h_{m-1} & \cdots & h_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

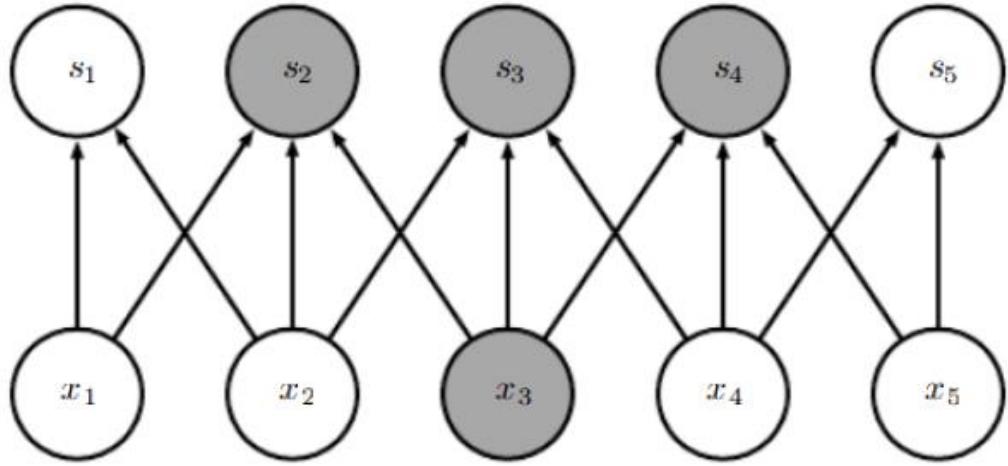
CNN

- idea: replacing linear/dense/fully-connected layers with convolution

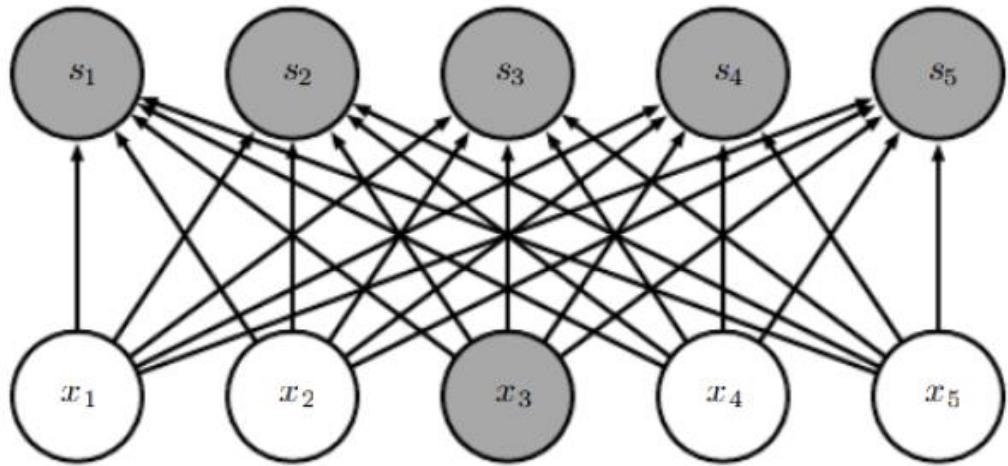
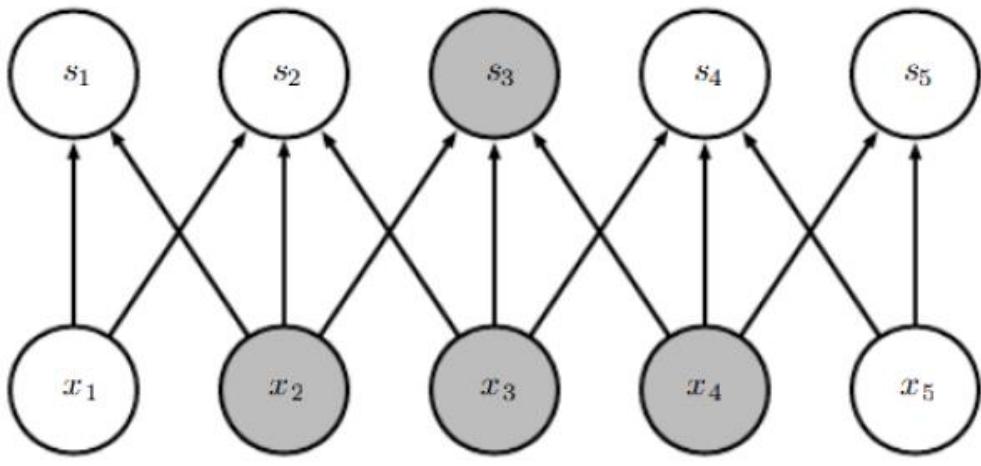
motivation of CNN

- sparse interactions
- parameter sharing
- equivariant representations

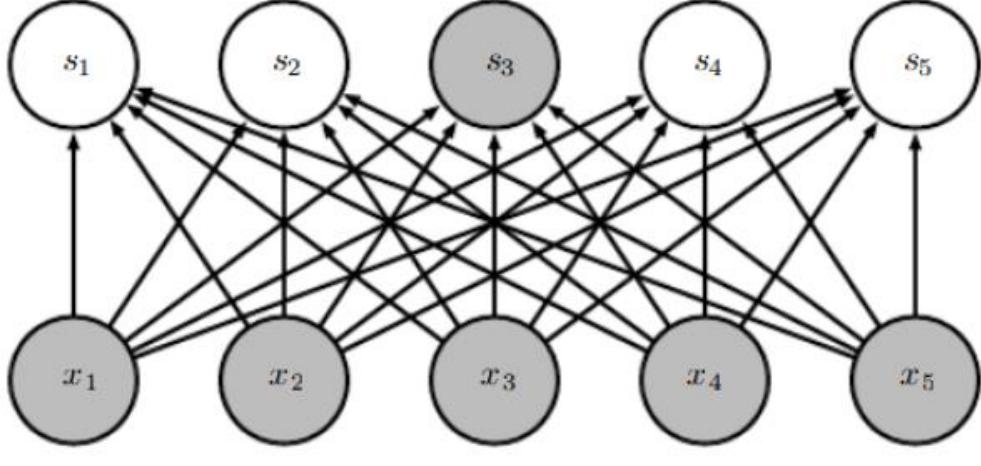
sparse connection



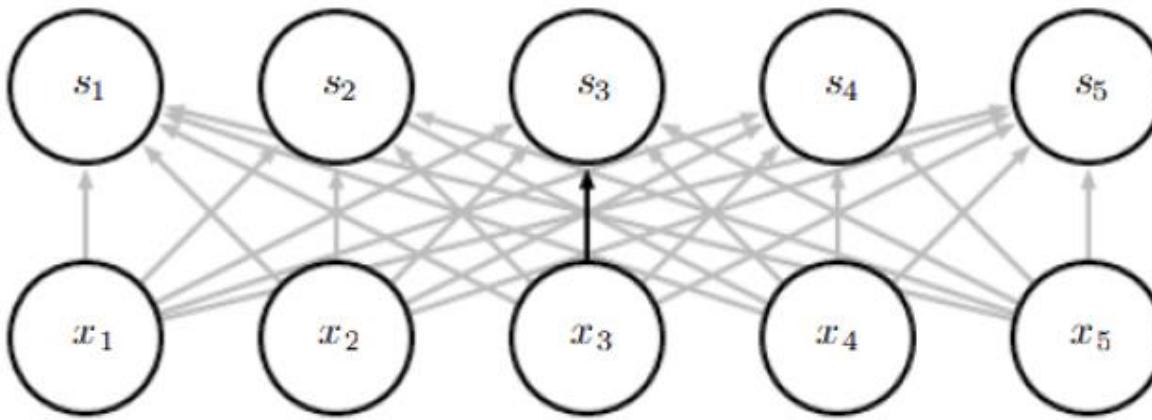
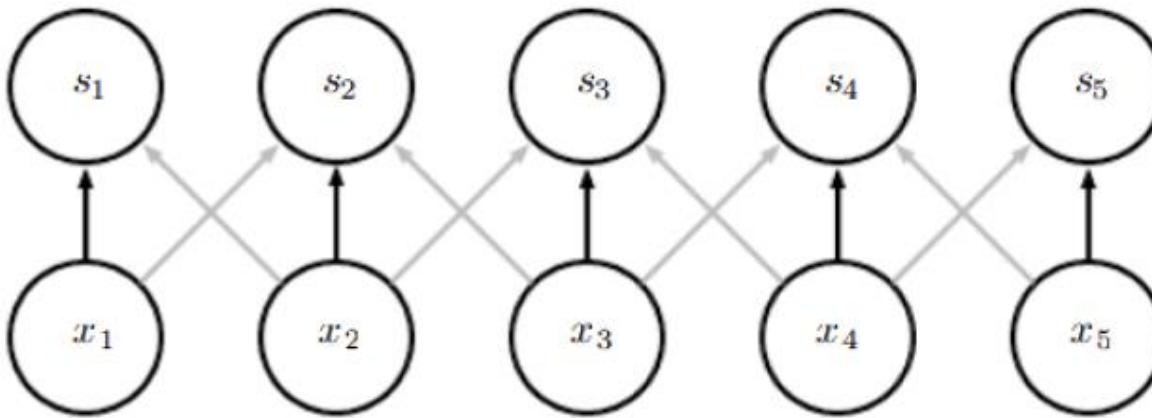
CNN



MLP



parameter sharing



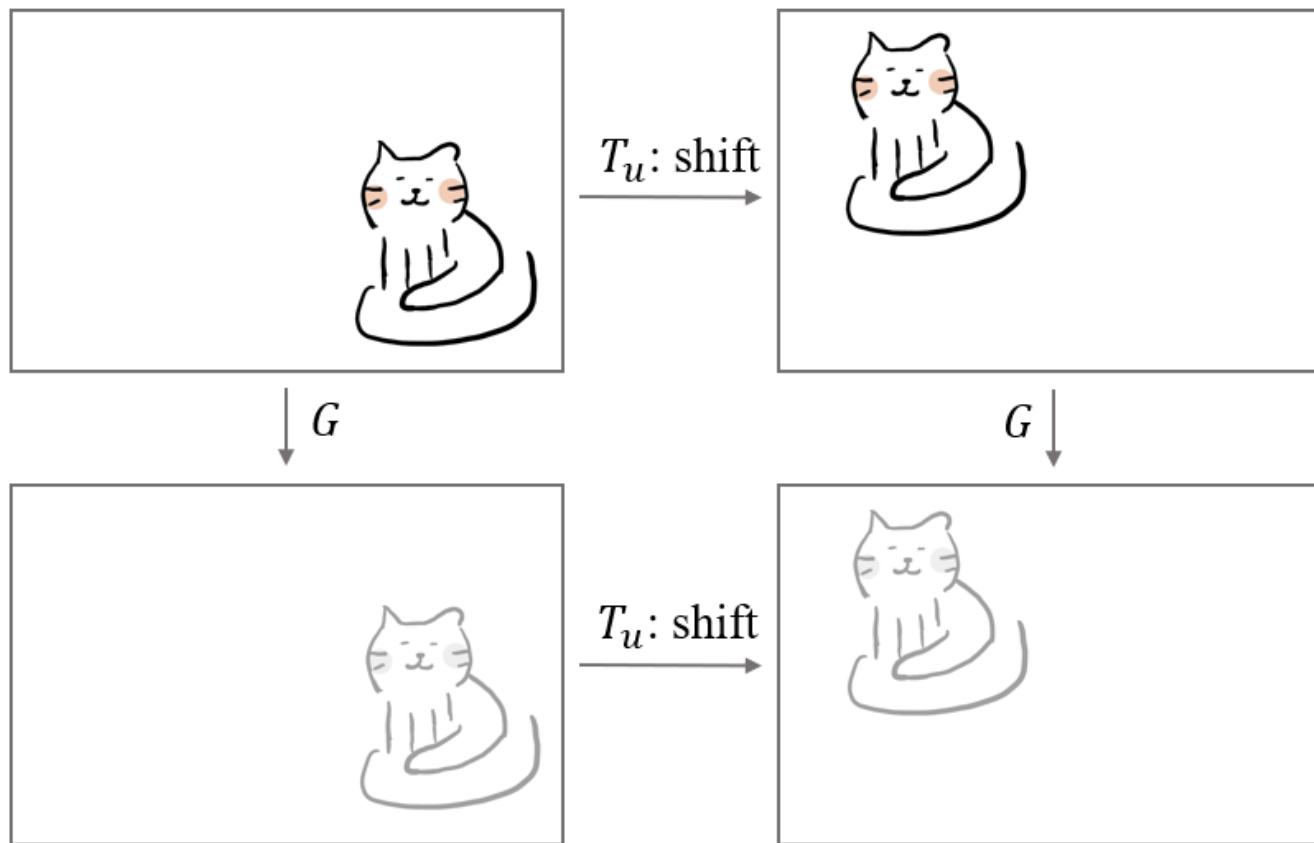
equivariance

- consider a translated/shifted version of the input:

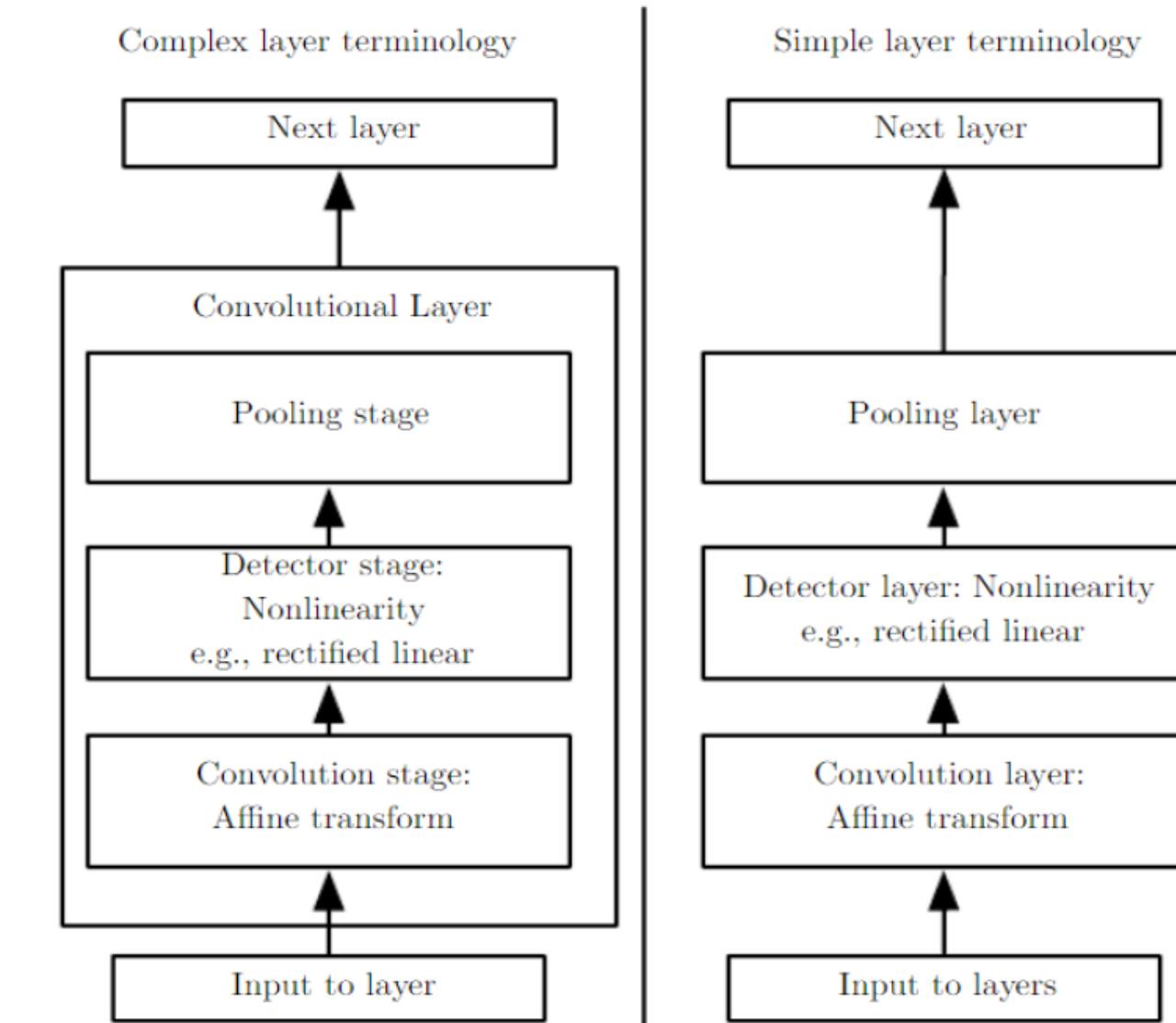
$$f_\tau(t) = f(t - \tau)$$

- what is $f_\tau * g$?

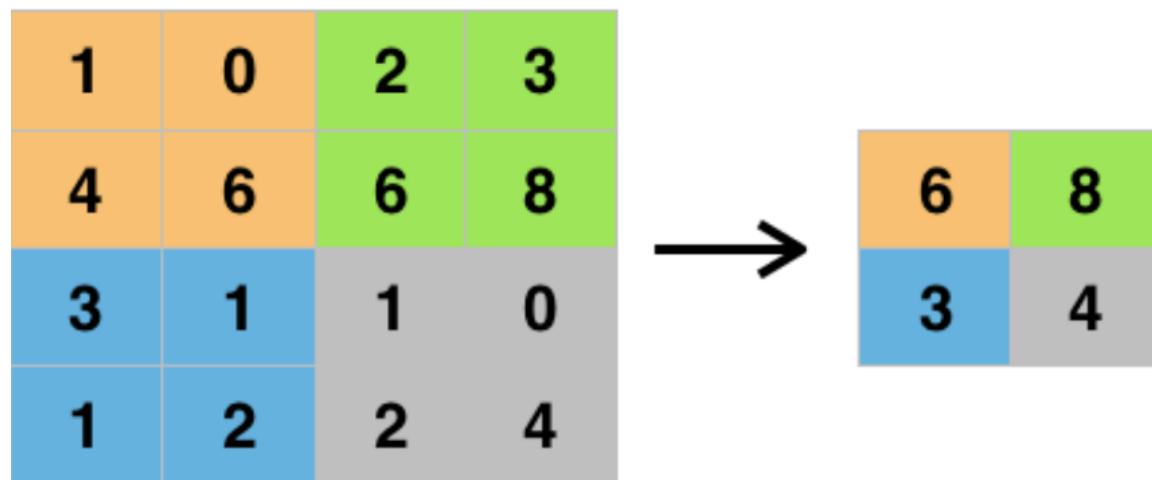
equivariance



structure

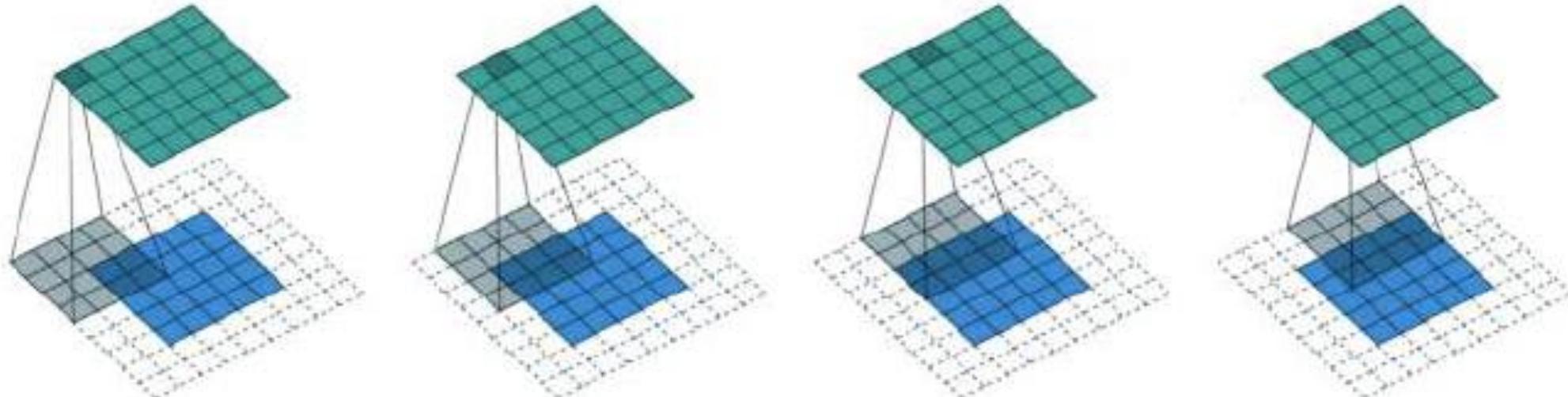


pooling



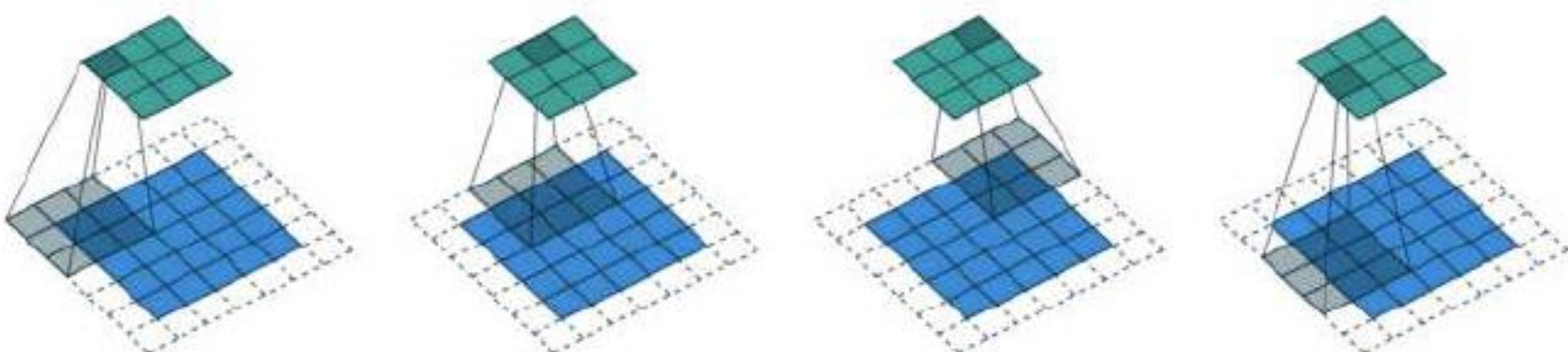
padding

a 4x4 kernel with 5x5 input and 2 pixels of padding



stride

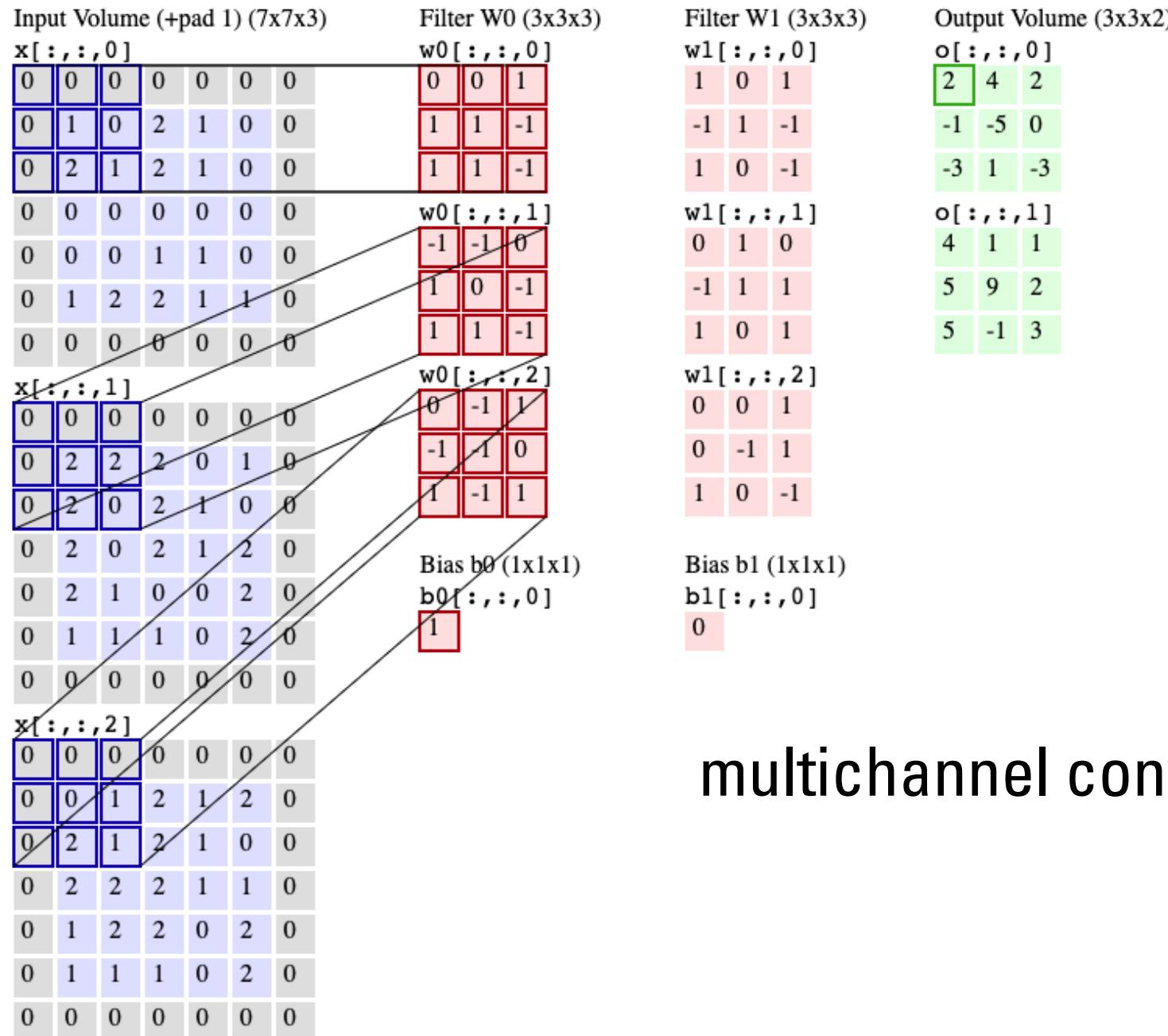
A 3x3 kernel with 5x5 input, stride-2 convolution, and 1 pixel of padding



multichannel convolution

- in general, the input and output have multiple channels (e.g., the original input image may have three RGB channels)
- 4D kernel: $\mathbf{K}(i, j, k, l)$
 - i : index for output channels
 - j : index for input channels
- suppose an input is $\mathbf{V}(j, k, l)$, then the output is:

$$\mathbf{Z}(i, k, l) = \sum_{j,m,n} \mathbf{V}(j, k + m - 1, l + n - 1) \mathbf{K}(i, j, m, n)$$



multichannel convolution

Input Volume (+pad 1) (7x7x3)

 $x[:, :, 0]$

0	0	0	0	0	0	0
0	1	0	2	1	0	0
0	2	1	2	1	0	0

Filter W0 (3x3x3)

 $w0[:, :, 0]$

0	0	1
1	1	-1
1	1	-1

Filter W1 (3x3x3)

 $w1[:, :, 0]$

1	0	1
-1	1	-1
1	0	-1

Output Volume (3x3x2)

 $o[:, :, 0]$

2	4	2
-1	-5	0
-3	1	-3

 $o[:, :, 1]$

4	1	1
5	9	2
5	-1	3

 $x[:, :, 1]$

0	0	0	0	0	0	0
0	2	2	2	0	1	0
0	2	0	2	1	0	0

w0[:, :, 1]

 $w0[:, :, 1]$

0	-1	0
-1	-1	0
1	1	-1

w1[:, :, 1]

 $w1[:, :, 1]$

0	1	0
-1	1	1
1	0	1

 $x[:, :, 2]$

0	2	0	2	1	2	0
0	2	1	0	0	2	0
0	1	1	1	0	2	0

w0[:, :, 2]

 $w0[:, :, 2]$

0	1	1
-1	-1	0
1	-1	1

w1[:, :, 2]

 $w1[:, :, 2]$

0	0	1
0	-1	1
1	0	-1

Bias b0 (1x1x1)

 $b0[:, :, 0]$

1

Bias b1 (1x1x1)

 $b1[:, :, 0]$

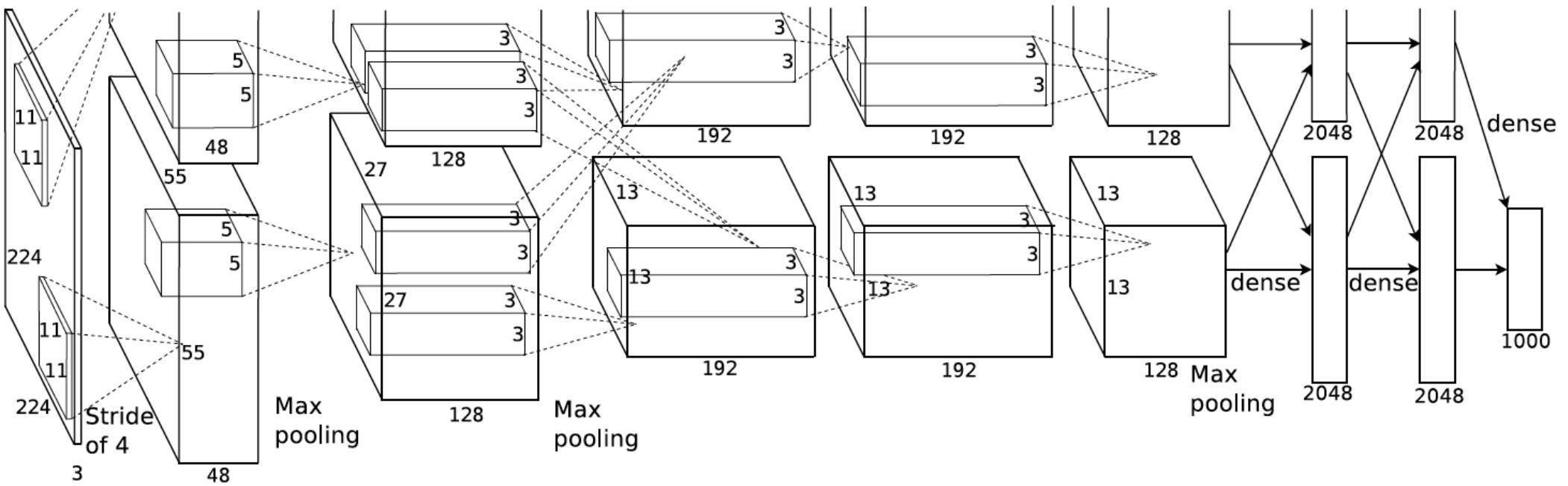
0

multichannel convolution

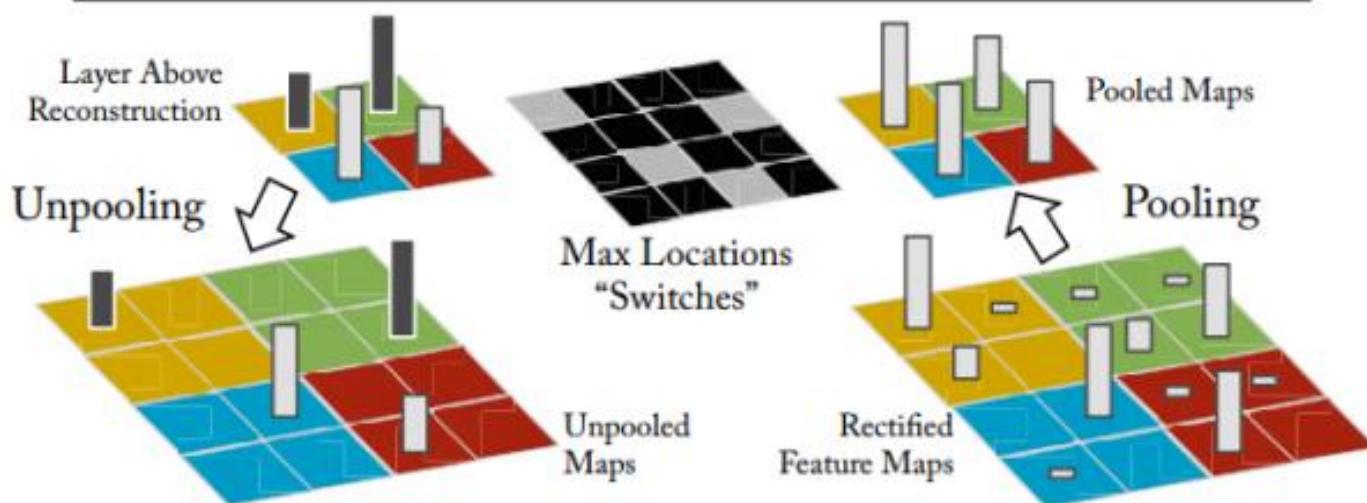
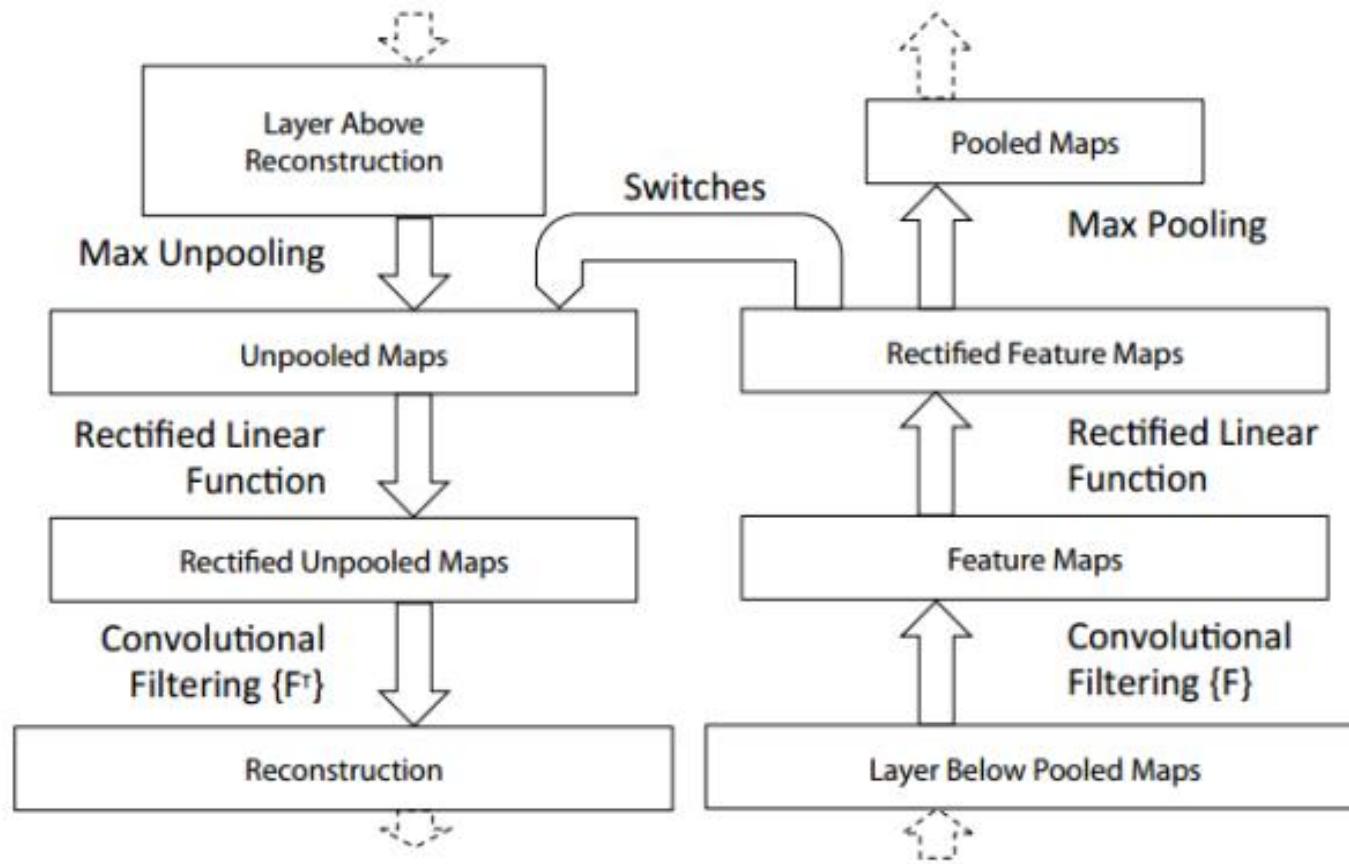
example

Layer (type)	Output Shape	Param #	Trainable
Conv2d	64 x 4 x 14 x 14	40	True
ReLU	64 x 4 x 14 x 14	0	False
Conv2d	64 x 8 x 7 x 7	296	True
ReLU	64 x 8 x 7 x 7	0	False
Conv2d	64 x 16 x 4 x 4	1,168	True
ReLU	64 x 16 x 4 x 4	0	False
Conv2d	64 x 32 x 2 x 2	4,640	True
ReLU	64 x 32 x 2 x 2	0	False
Conv2d	64 x 2 x 1 x 1	578	True
Flatten	64 x 2	0	False

AlexNet



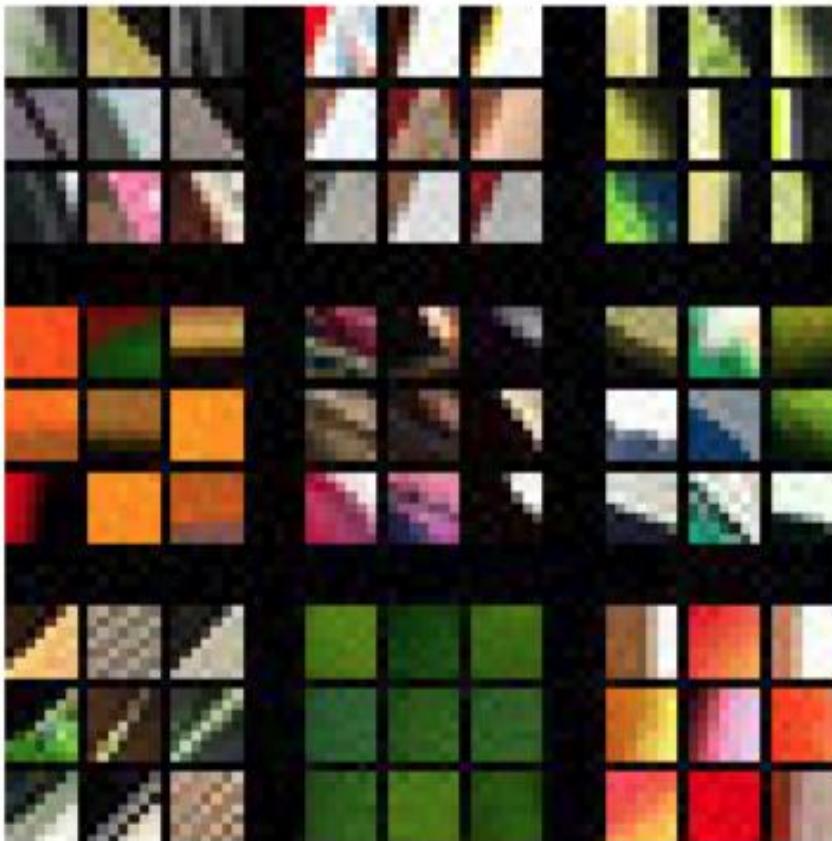
what do convolution layers learn



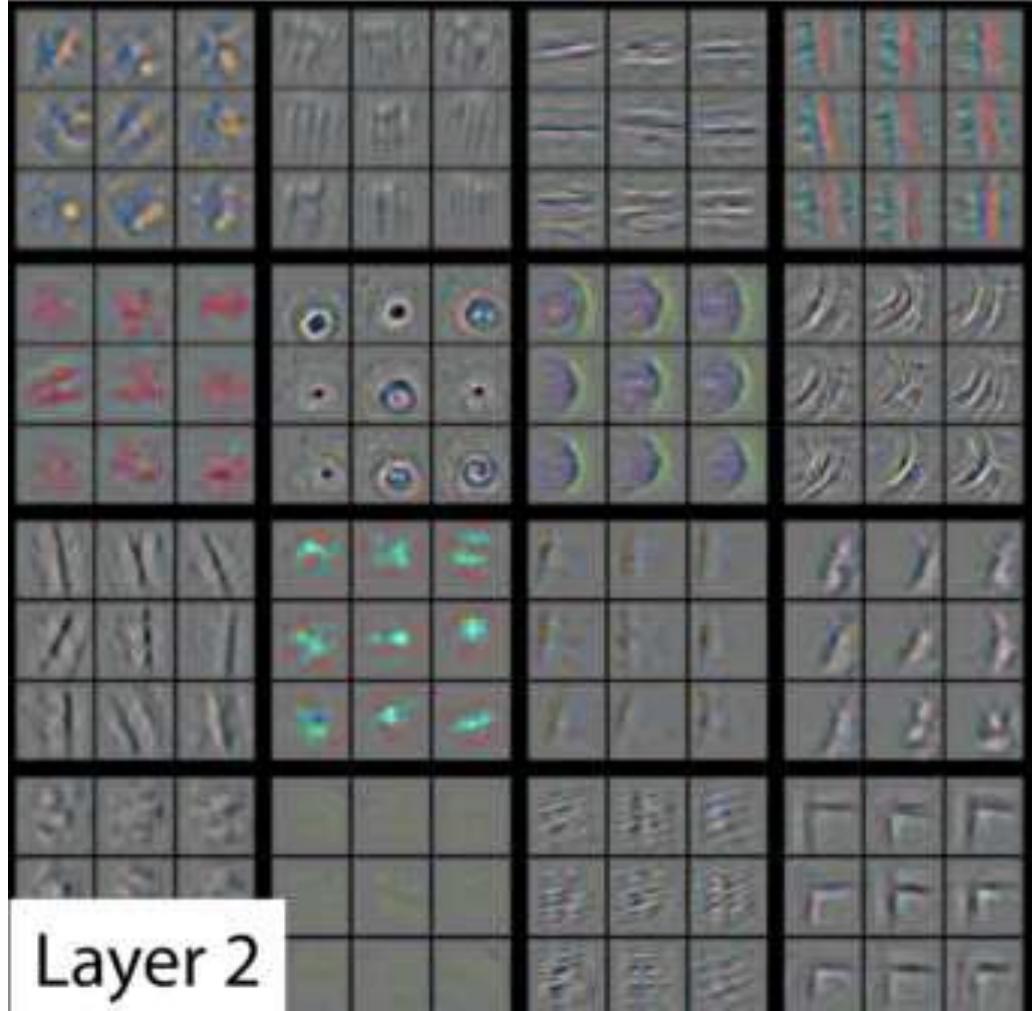
Visualizing and Understanding
Convolutional Networks
(ECCV 2014)



Layer 1



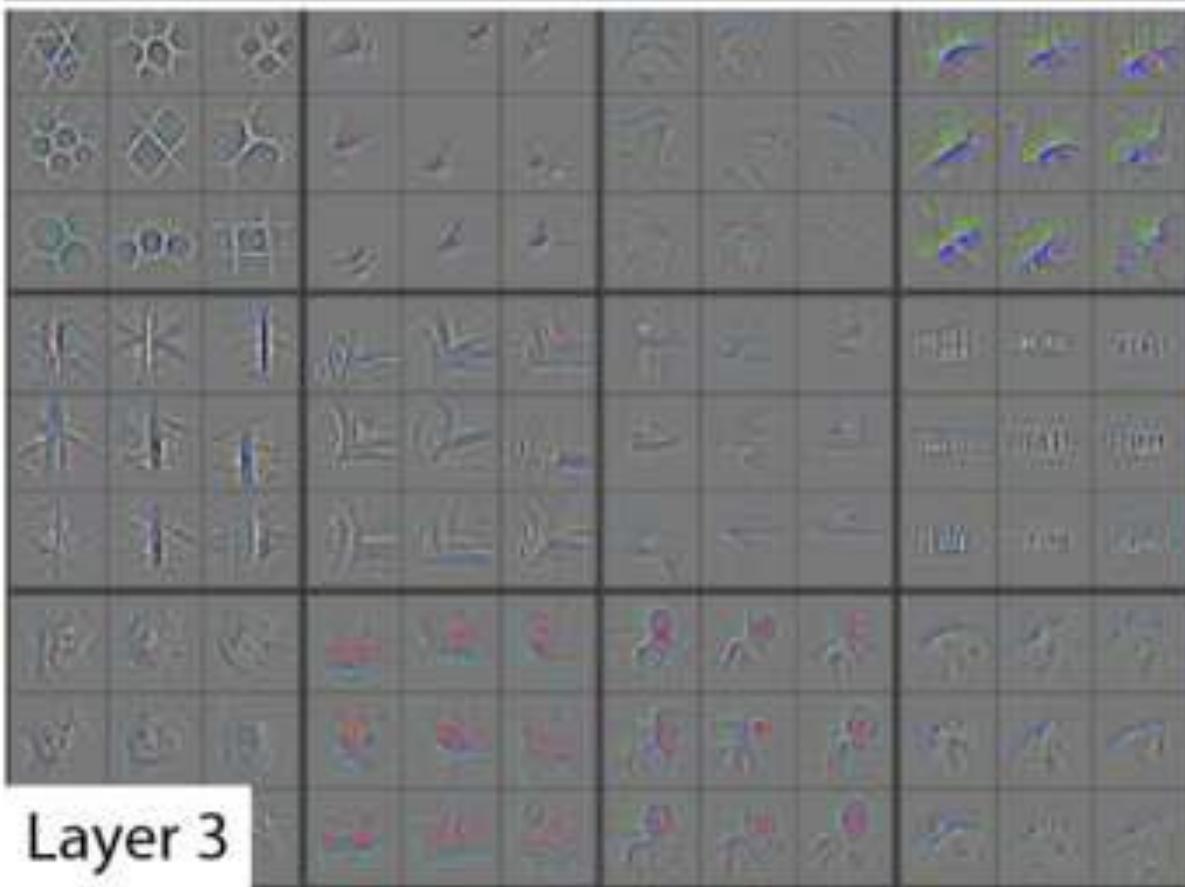
in layer 1, what we can see is that the model has discovered weights that represent diagonal, horizontal, and vertical edges, as well as various gradients



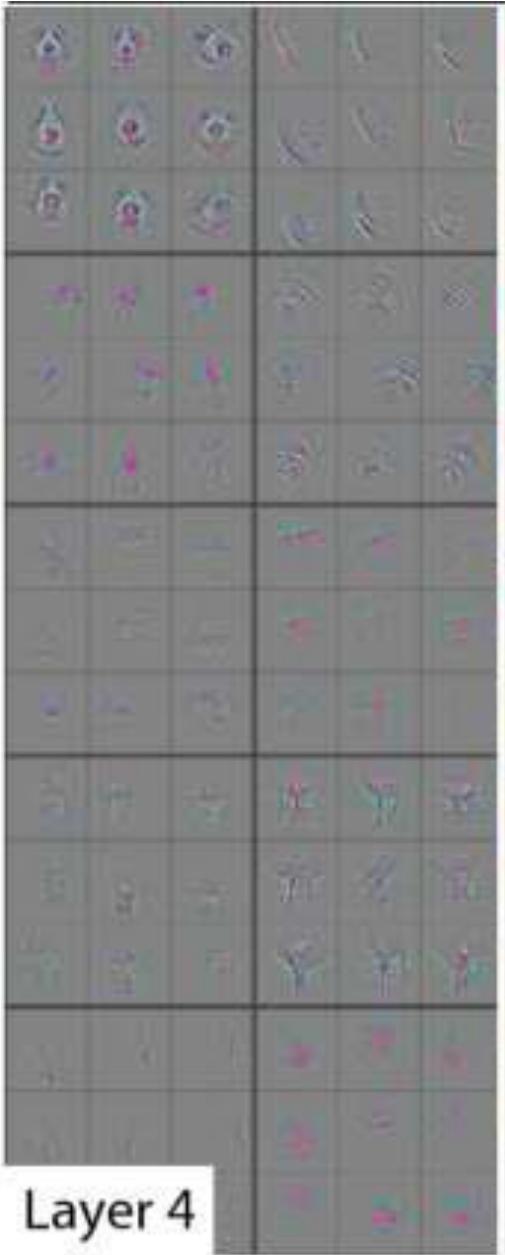
Layer 2



in layer 2, we can see that the model has learned to create feature detectors that look for corners, repeating lines, circles, and other simple patterns



as it goes deeper, more higher-level concepts



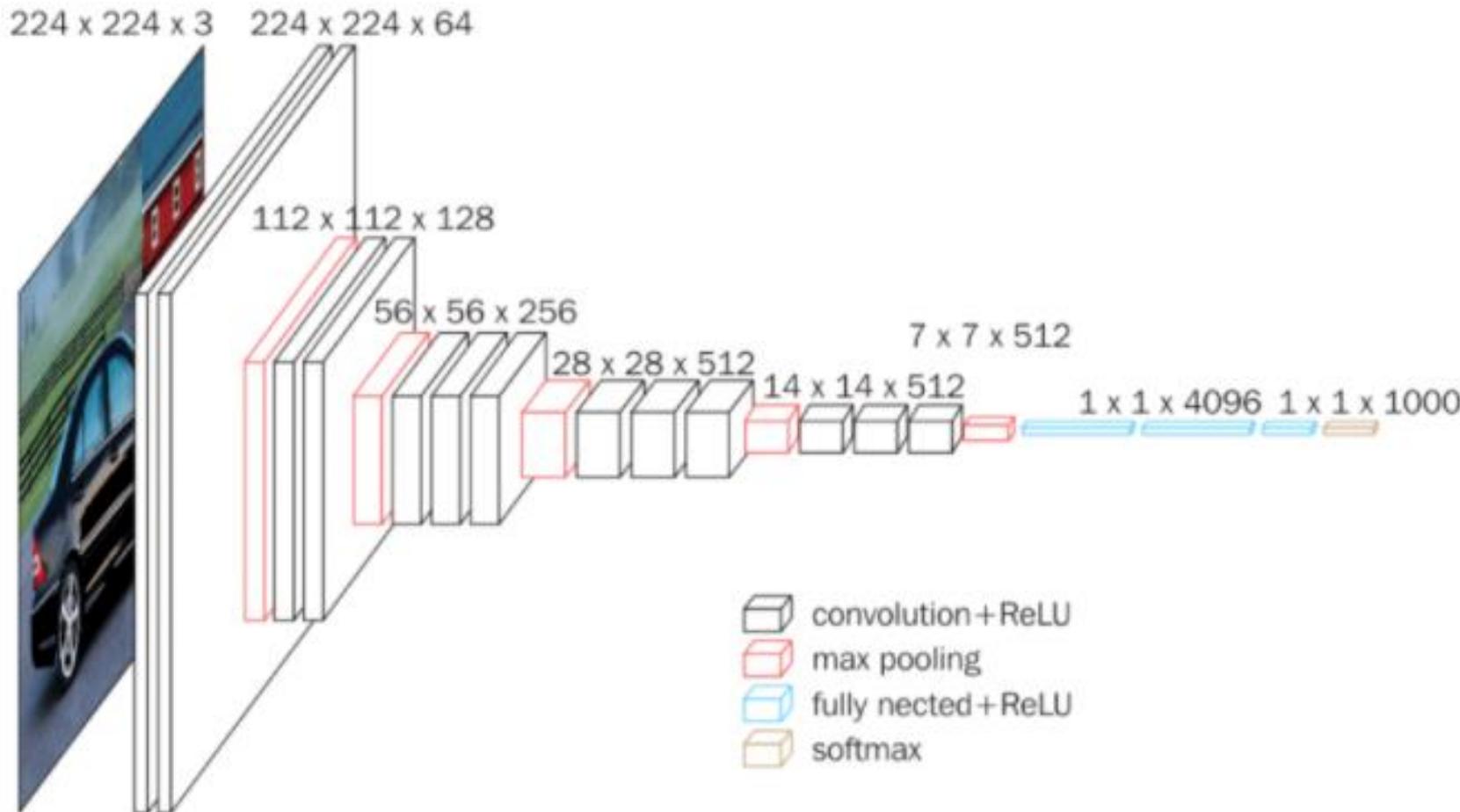
datasets

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

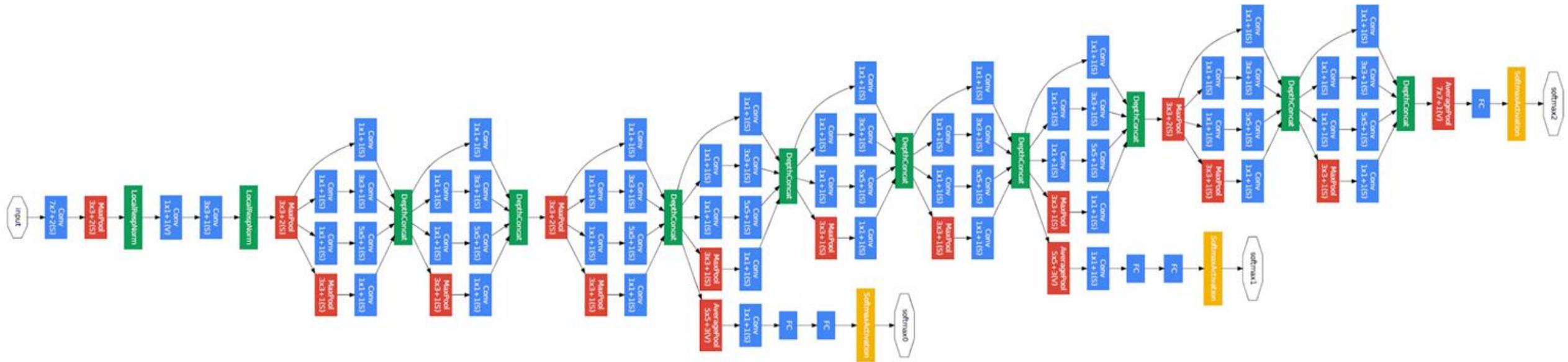


VGG16



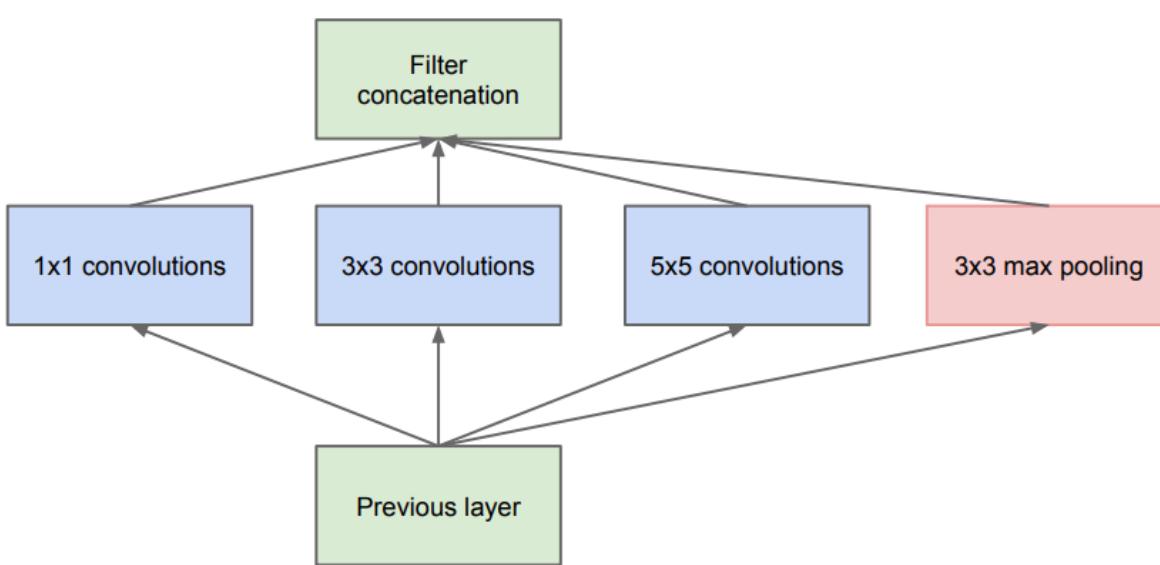
Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. In *ICLR 2015*.

GoogleNet

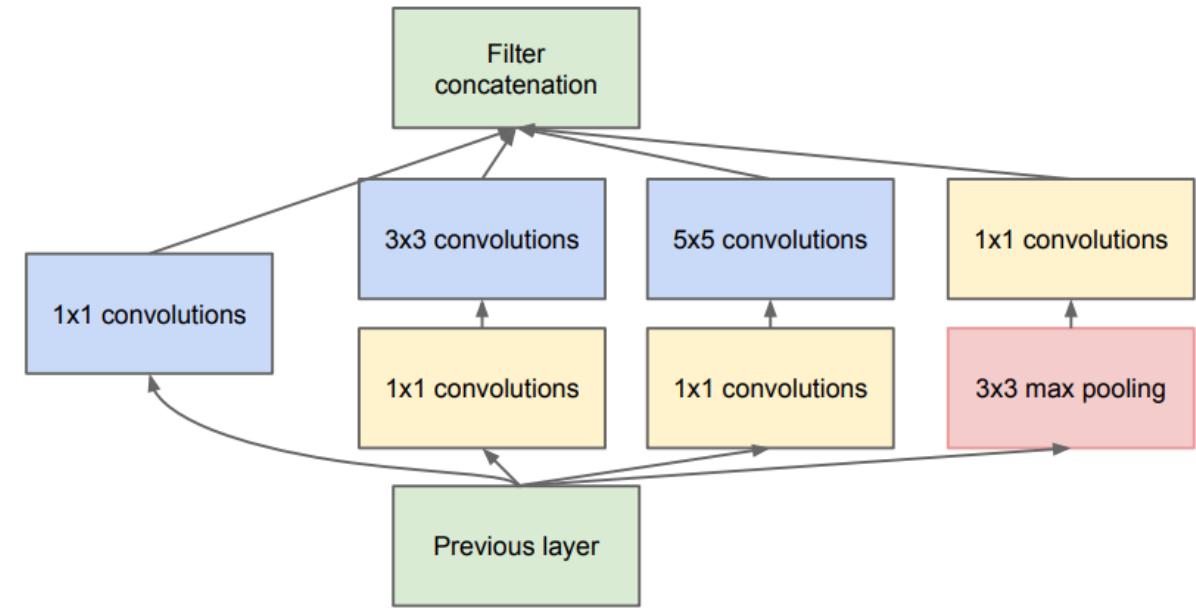


Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions. In *CVPR* (pp. 1-9).

GoogleNet: inception modules



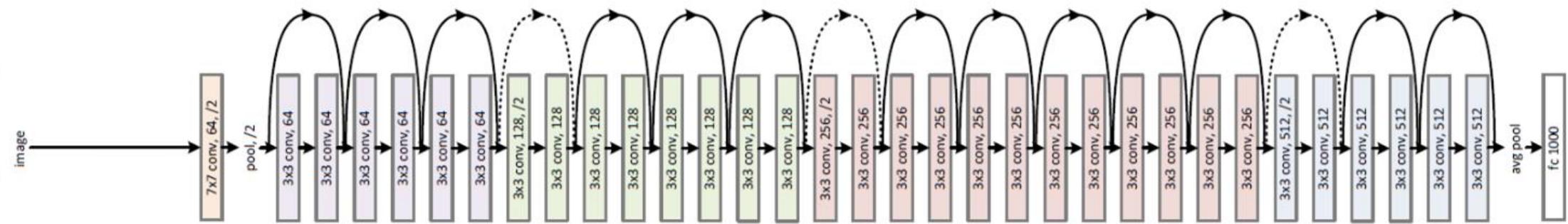
(a) Inception module, naïve version



(b) Inception module with dimensionality reduction

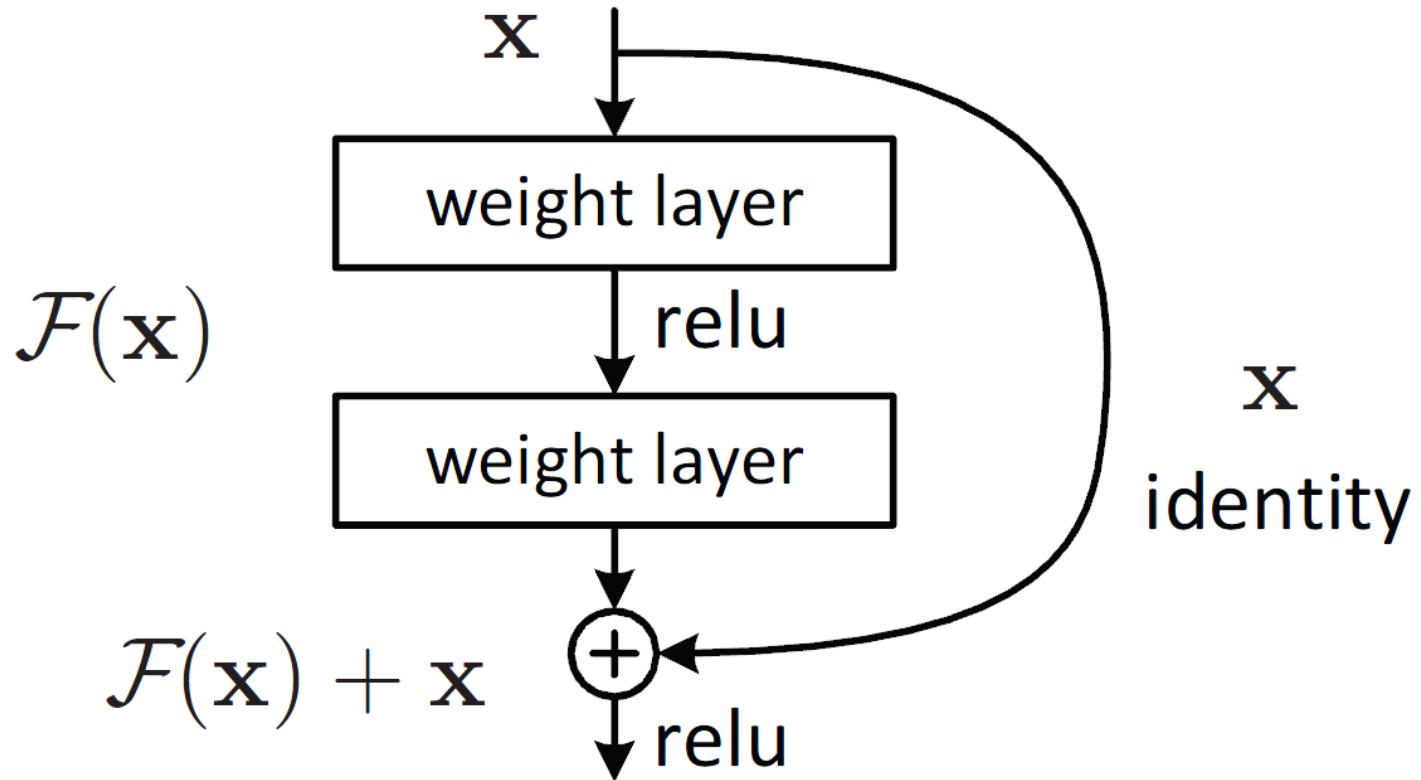
ResNet

34-layer residual



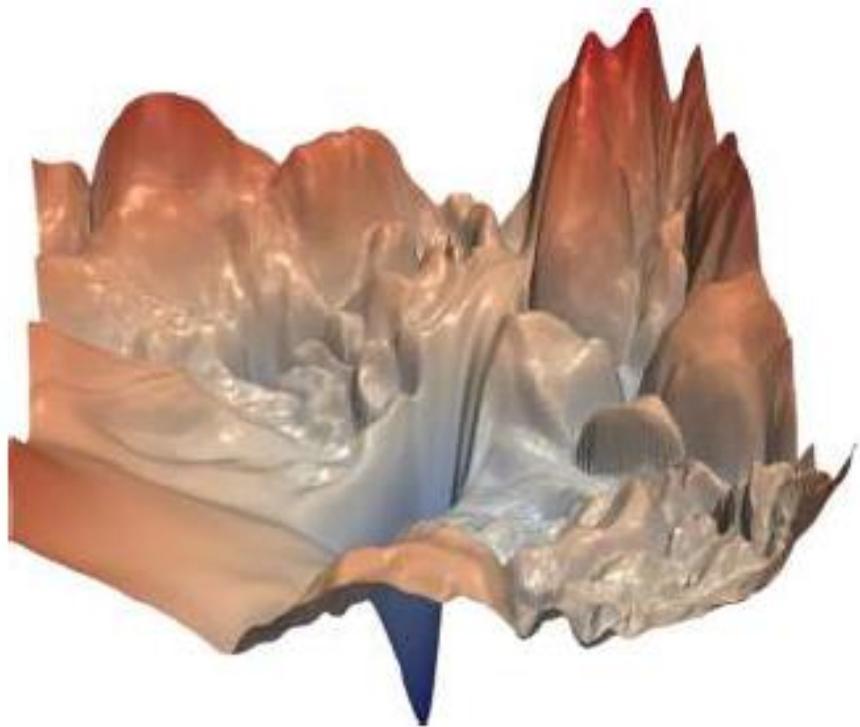
He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *CVPR* (pp. 770-778).

ResNet

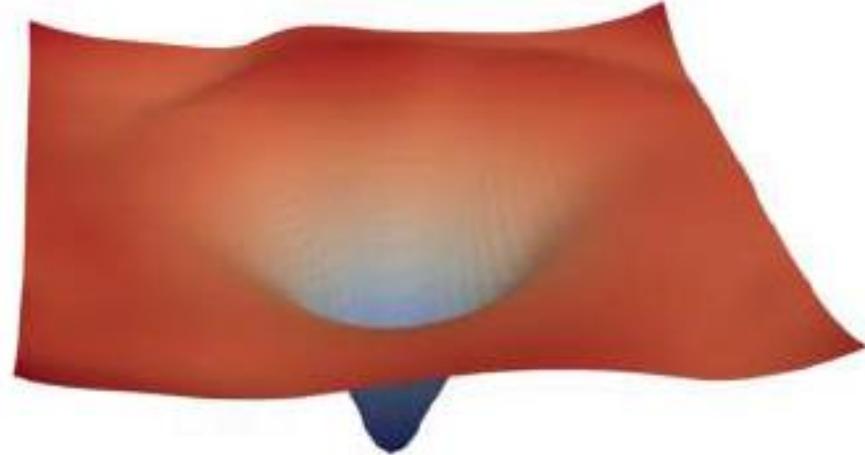


ResNet

- landscape of the loss functions



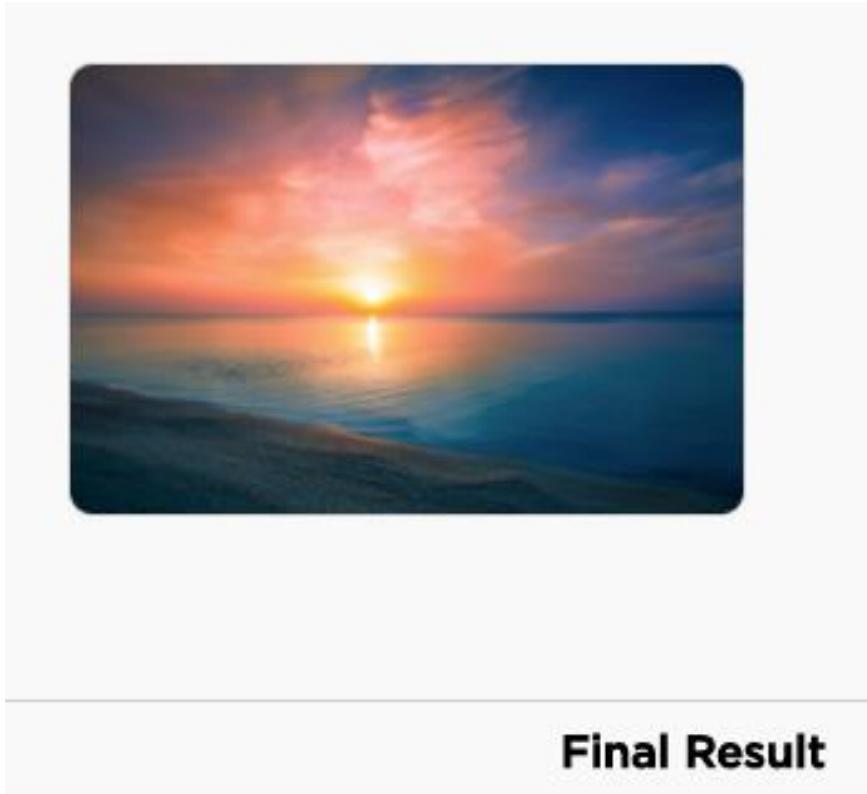
regular CNN



ResNet

4.3 CNN in classification

scene recognition



application: smart album, auxiliary photography, ...

scene: landscape

confidence: 99.81%

scene: sky

confidence: 83.85%

scene: sunrise

confidence: 99.83%

scene: lake

confidence: 64.05%

scene: sea

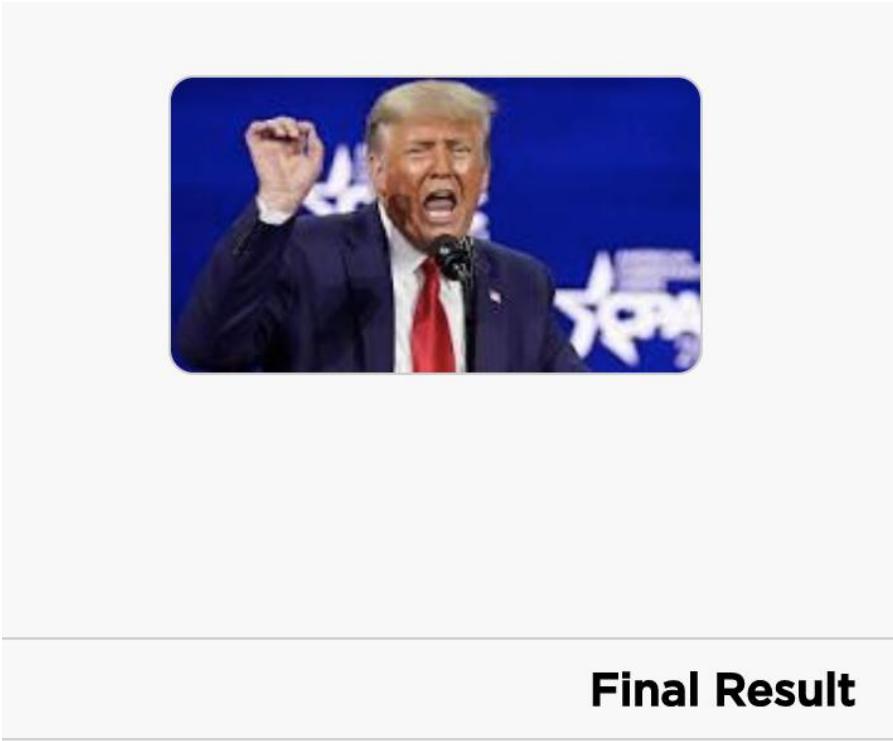
confidence: 73.72%

Places365 dataset

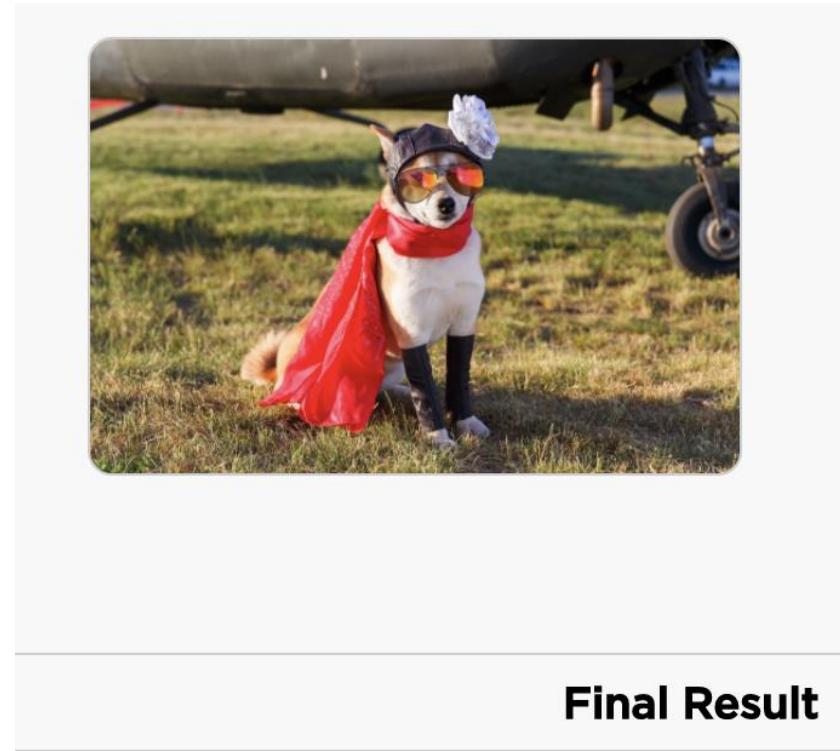


<https://github.com/CSAILVision/places365>

photo evaluation



images quality: 4.593186331097968



images quality: 5.335587732726708

AVA (Aesthetic Visual Analysis) dataset

TITLE: Skyscape

Description:

Make the sky the subject
of your photo this week.

Stats

Voting Dates:

13/07/2010 - 19/07/2010

Numbers & Statistics:

Submissions: 136

Disqualifications: 1

Votes: 16,009

Comments: 595

Average Score: 5.64014

1st place with an
average vote of **7.4831**



4th place with an
average vote of **6.8547**



7th..

2nd place with an
average vote of **7.0328**



5th place with an
average vote of **6.7073**



8th..

3rd place with an
average vote of **6.9333**



6th place with an
average vote of **6.6667**

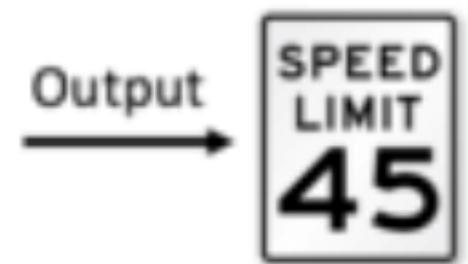


9th..



robustness

- view angles
- sizes
- distortion
- lighting
- cover
- ...



robustness



x
“panda”
57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

$=$



$x +$
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

face recognition

False
Positives



frontal faces

False
Negatives



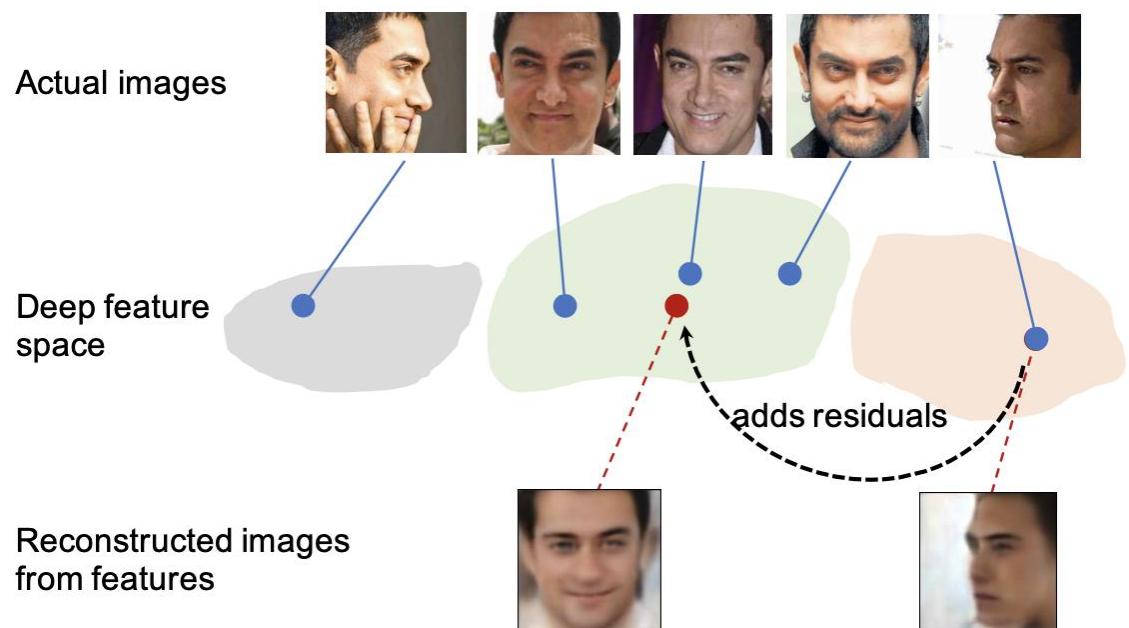
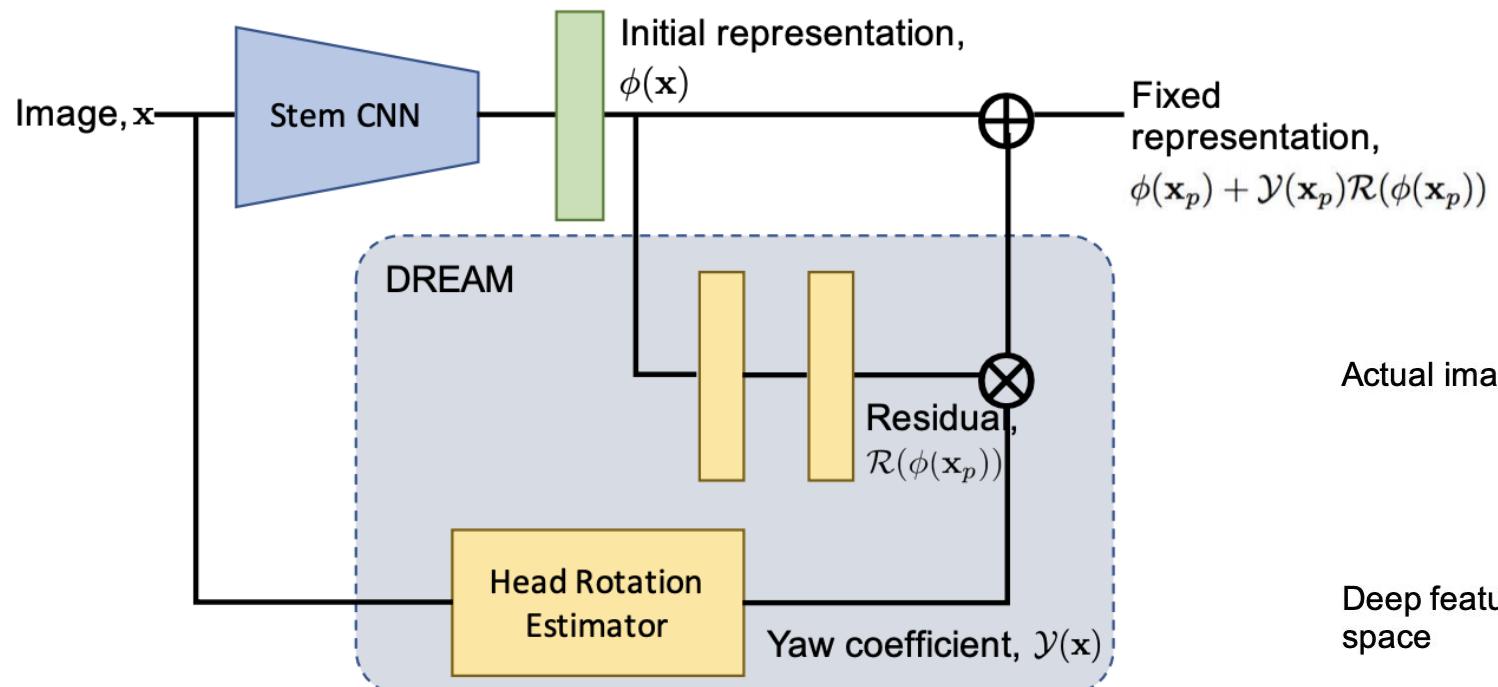
profile faces

face recognition

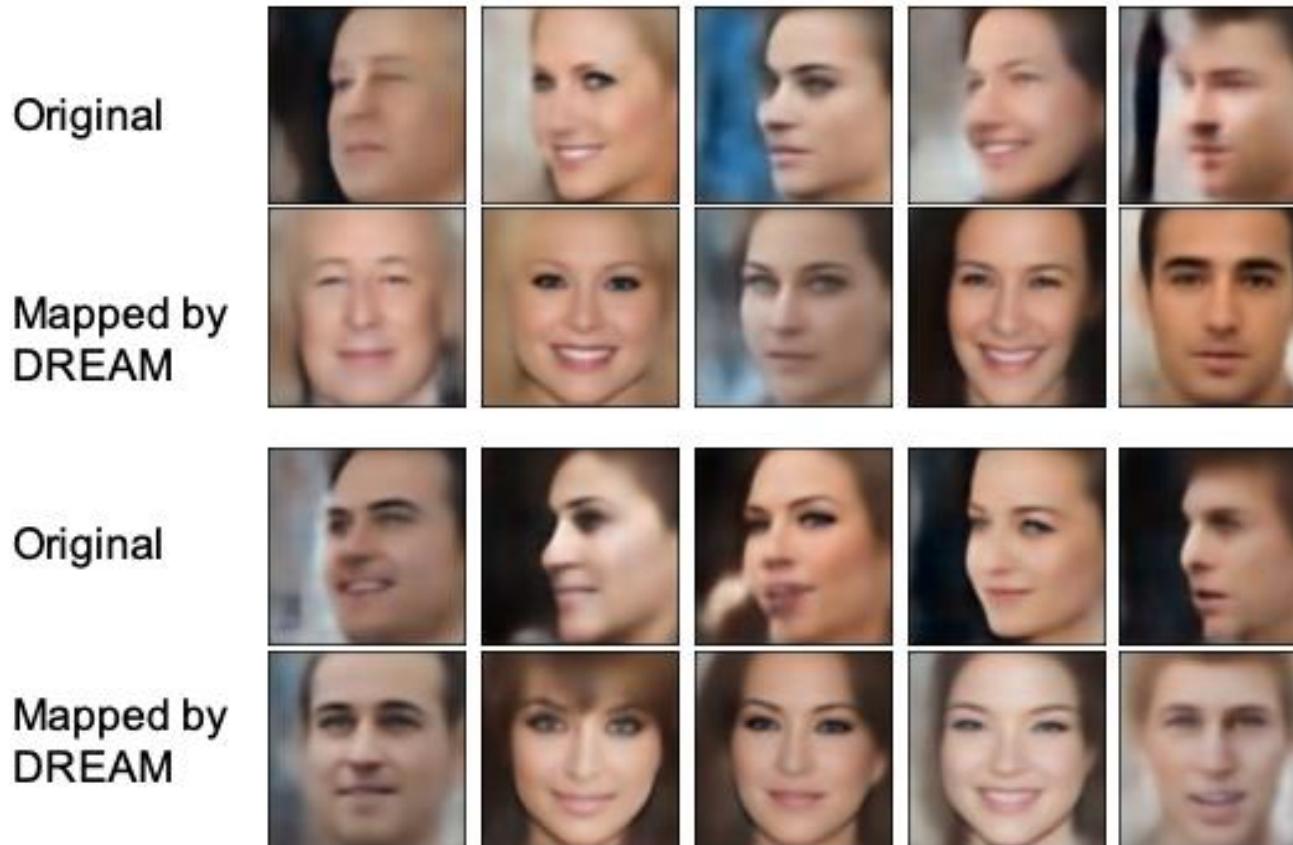
why does face recognition work poorly on profile faces?

- the generalization power of deep models is usually proportional to the training data size
- given an **uneven distribution of profile and frontal faces in the dataset**, deeply learned features tend to bias on distinguishing frontal faces rather than profile faces

face recognition



face recognition

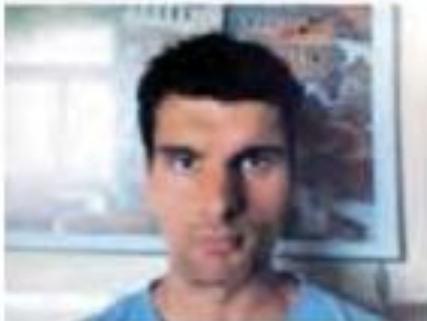


face spoof detection

which of the three is also a genuine face image?



(a)



(b)



(c)



(d)

genuine face image

face spoof detection

which of the three is also a genuine face image?



(b)



(a)



(c)



(d)

print attack:
photo of a printed photo

replay attack:
photo of a tablet screen

a person wearing a mask

face spoof detection

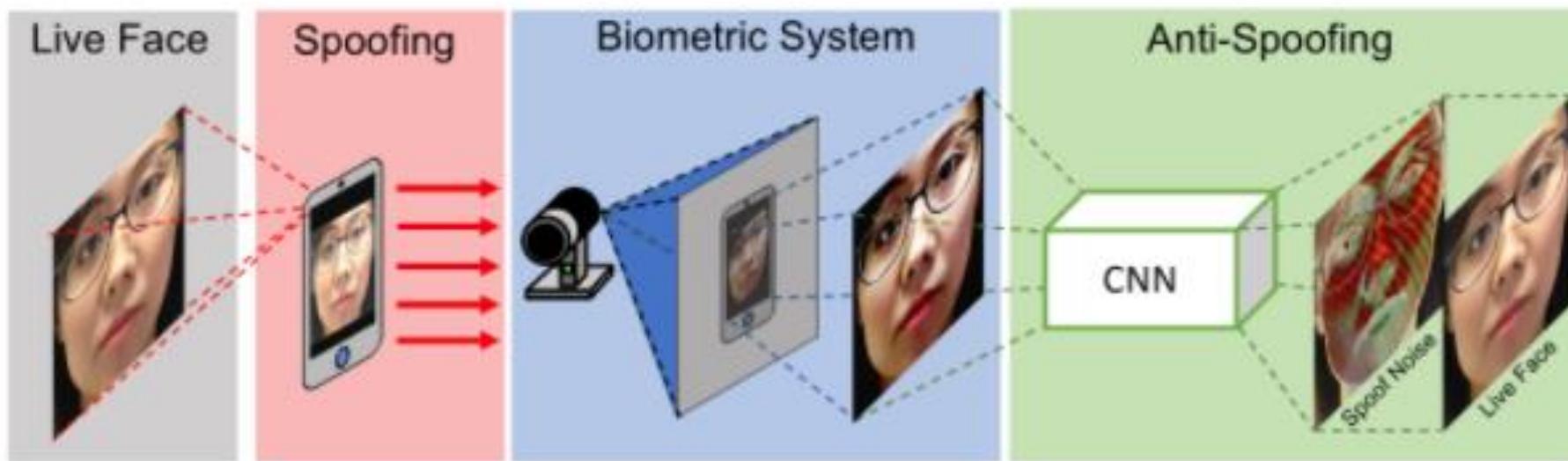


image credit: Jourabloo et al. (2018), Face De-Spoofing

face spoof detection

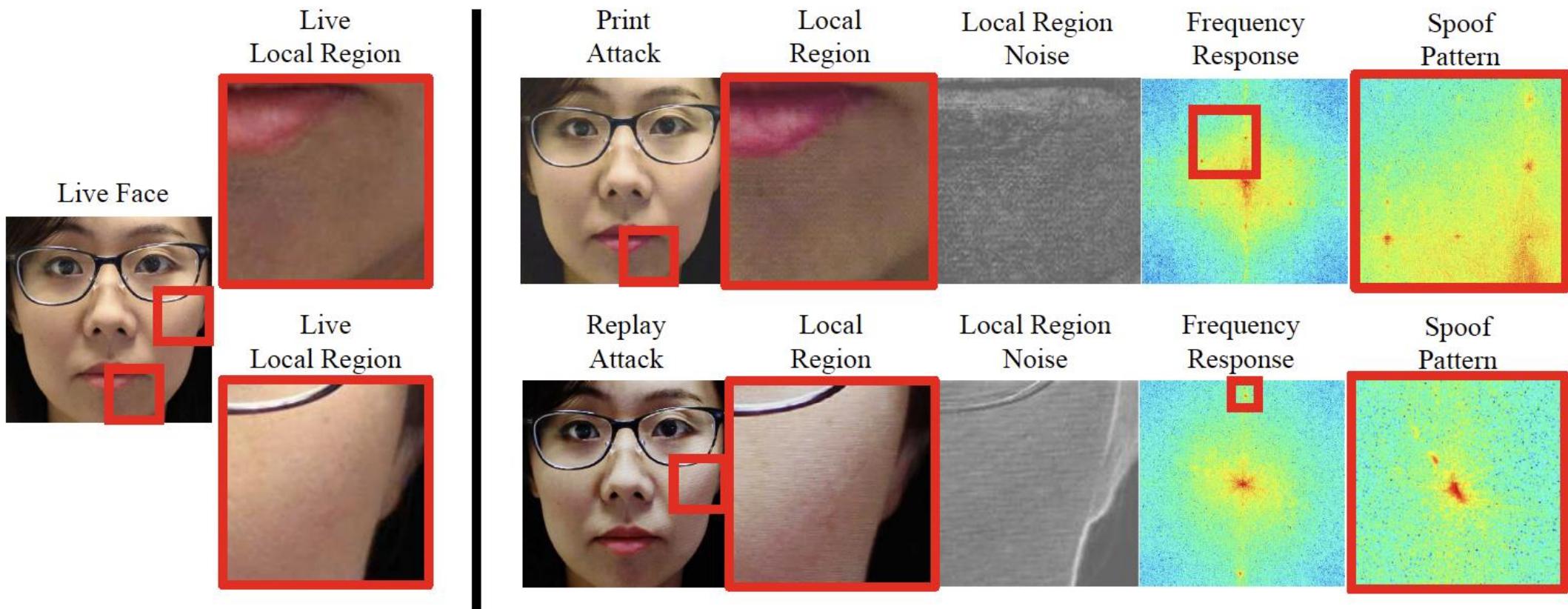


image credit: Jourabloo et al. (2018), Face De-Spoofing

face spoof detection

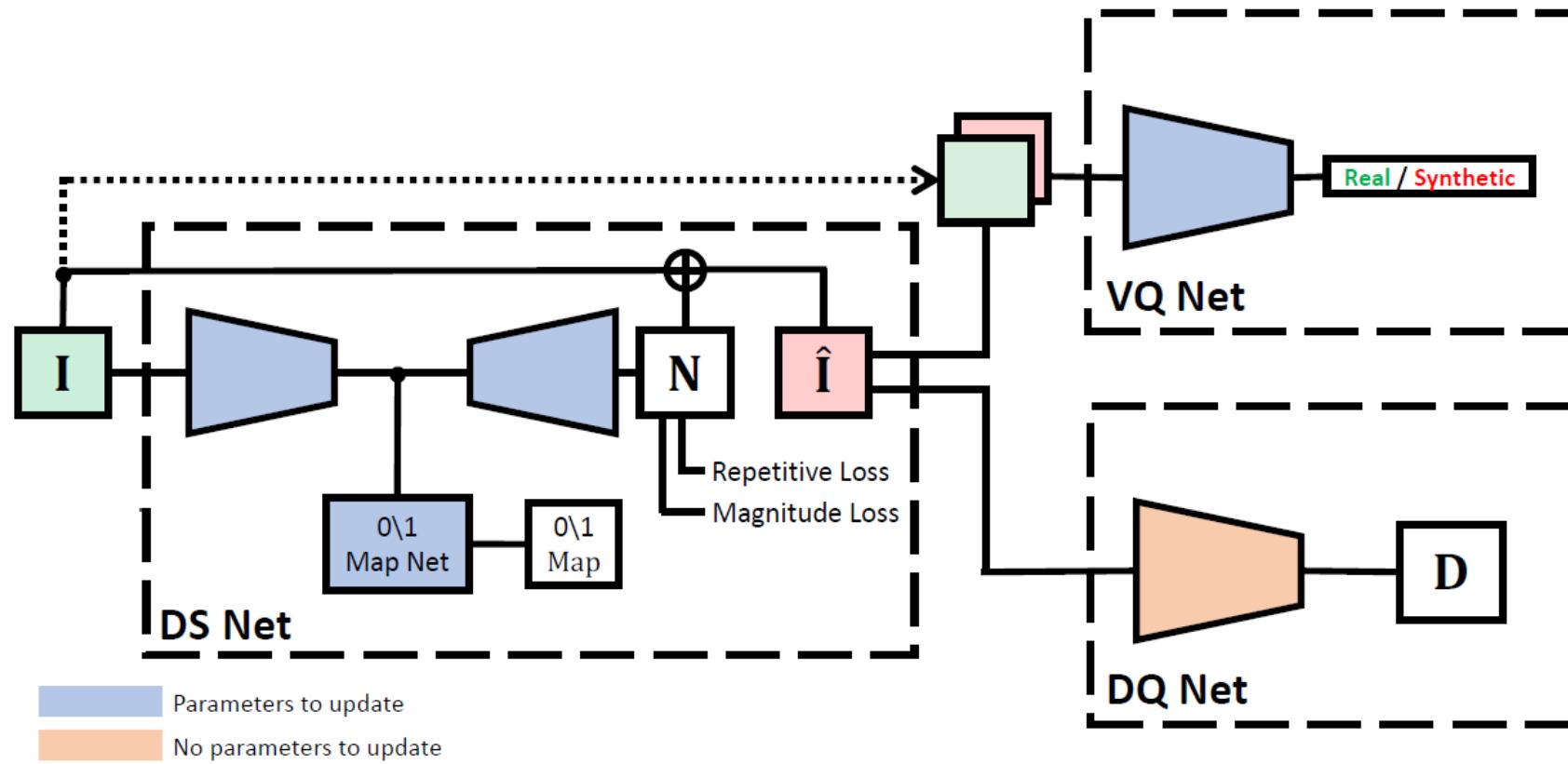
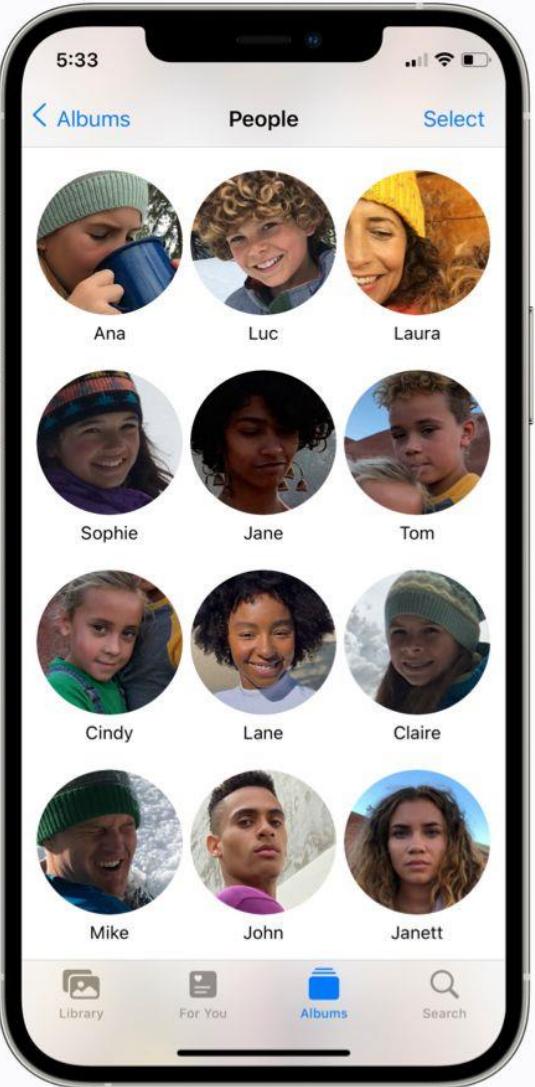


image credit: Jourabloo et al. (2018), Face De-Spoofing



A

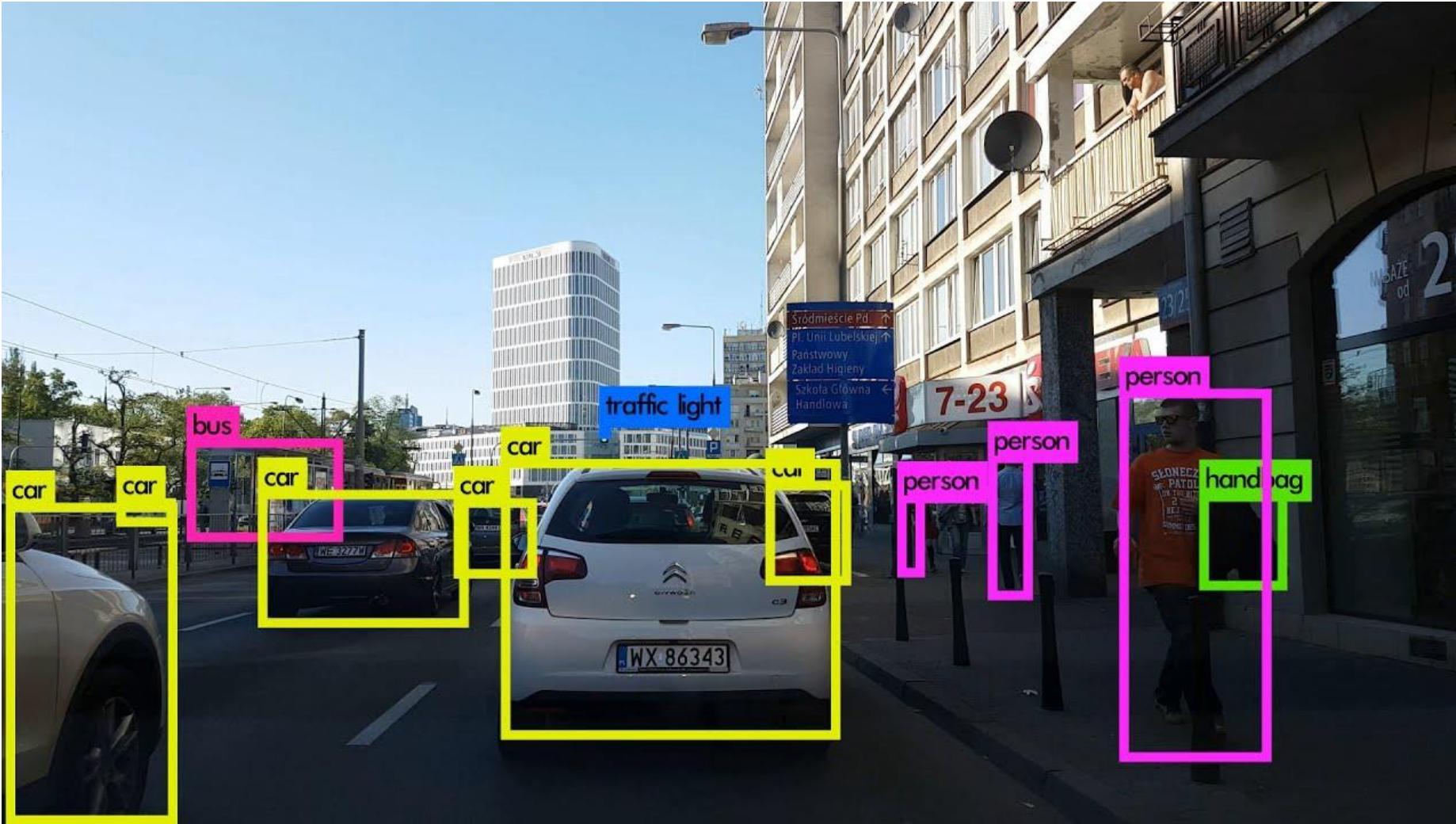
B

C

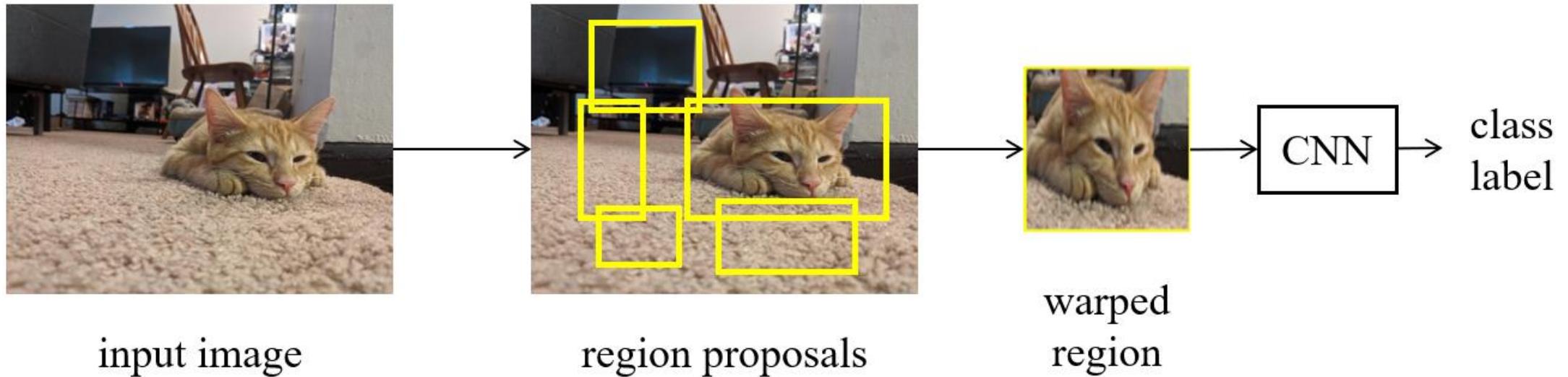
D

4.4 CNN in object detection

object detection



R-CNN



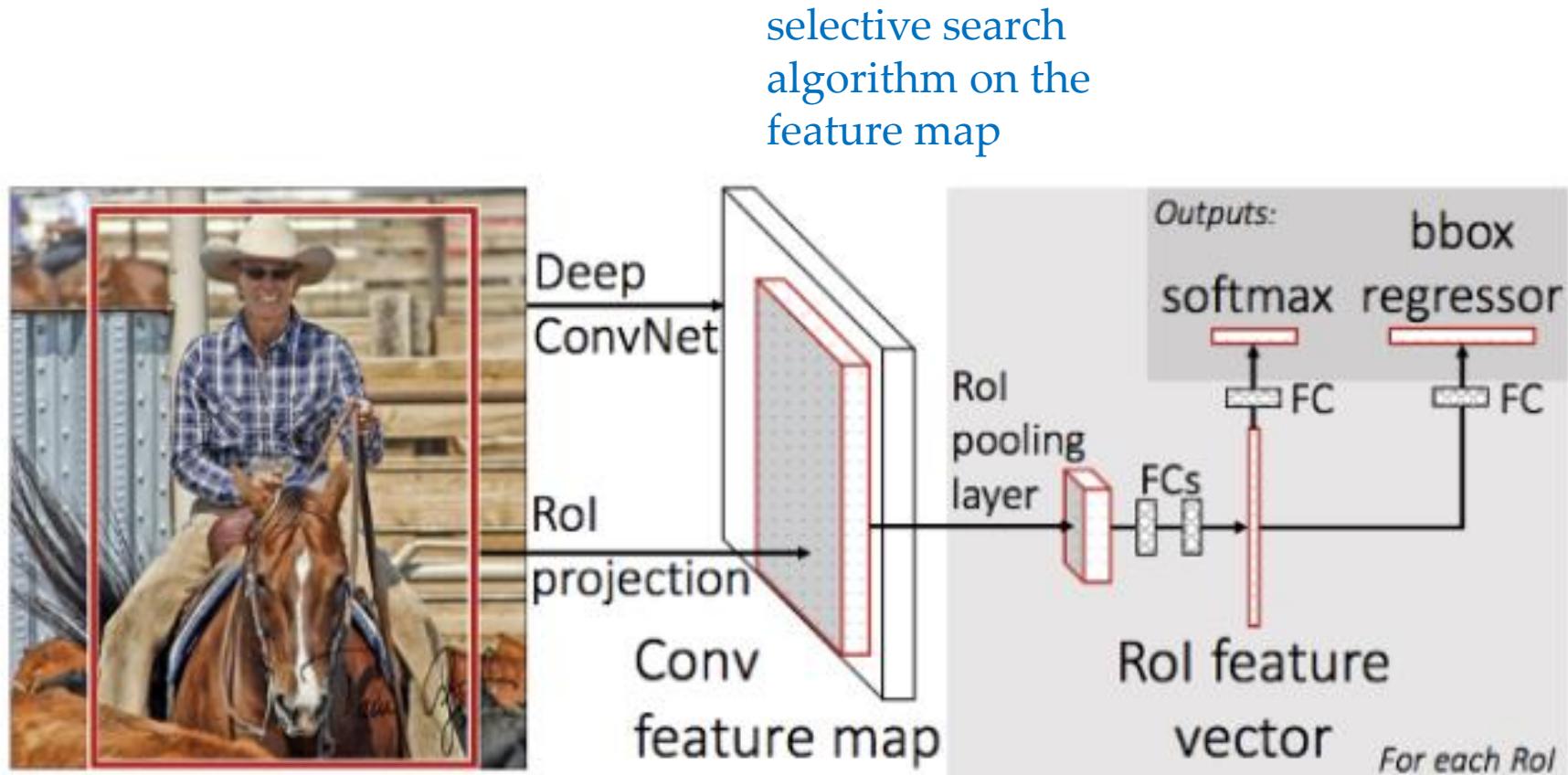
R-CNN

- regions are proposed via “selective search”, a fixed algorithm (no learning involved):
 1. generate initial sub-segmentation (many candidates)
 2. use greedy algorithm to recursively combine similar regions
 3. produce the final candidate region proposals

R-CNN

- relatively slow in training (2000 proposed regions)
- cannot be implemented real time (~47 seconds for each test image)
- the selective search algorithm is a fixed algorithm involves no learning

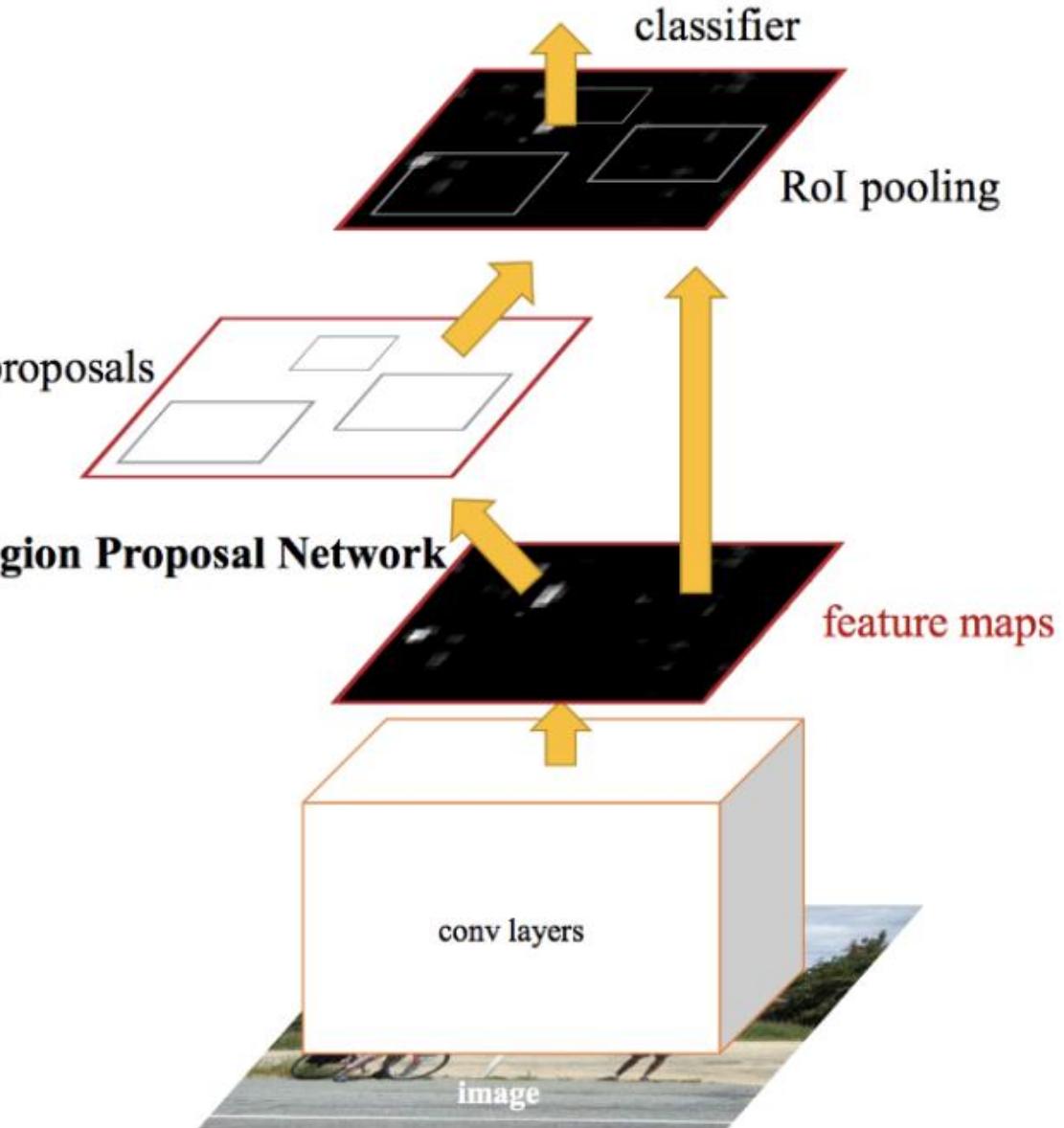
fast R-CNN



the input image, instead of the proposed regions,
is fed to the CNN to generate a convolutional
feature map

faster R-CNN

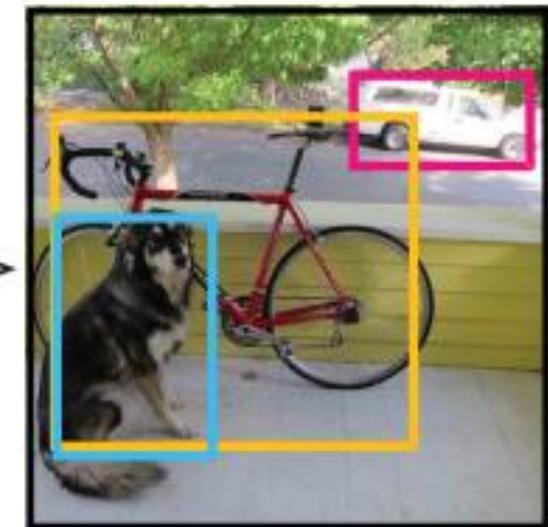
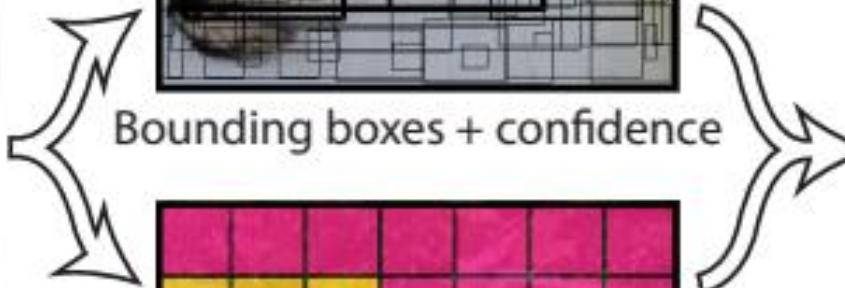
no selective search algorithm used on the feature map



A single neural network predicts bounding boxes and class probabilities directly from full images in [one evaluation](#).

YOLO

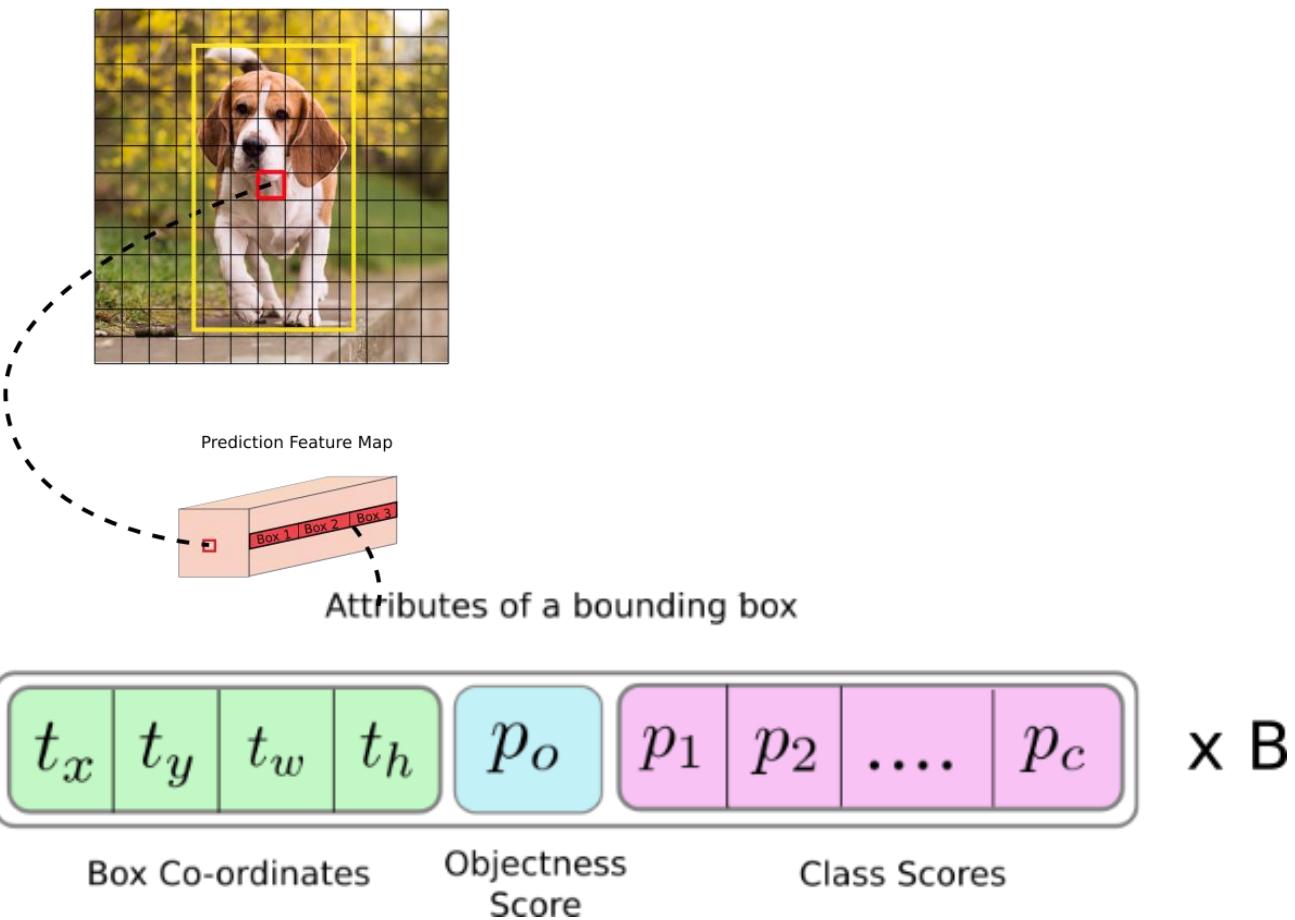
- YOLO = You Only Look Once



The image is divided into an $S \times S$ grid and for each grid cell predicts **B bounding boxes**, **confidence** for those boxes, and **C class probabilities**. These predictions are encoded as an $S \times S \times (B \times 5 + C)$

Class probability map

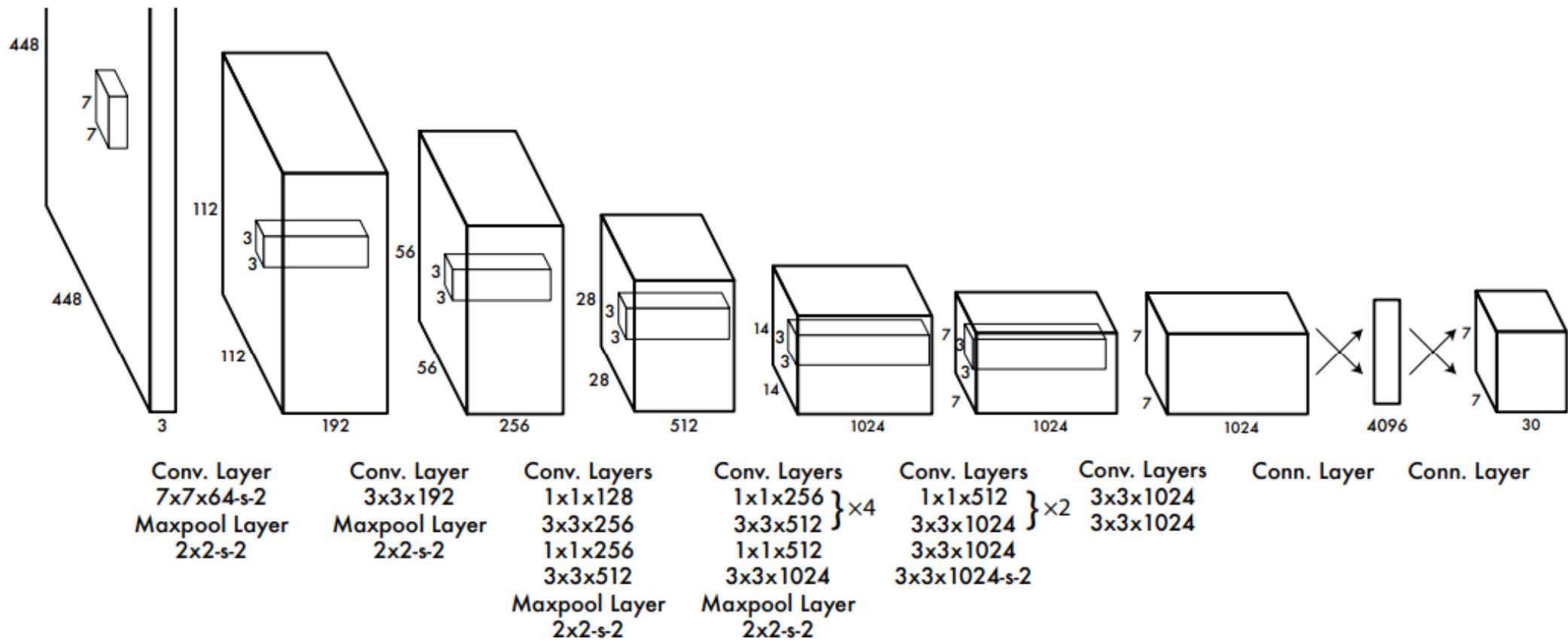
YOLO



The image is divided into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B \times 5 + C)$

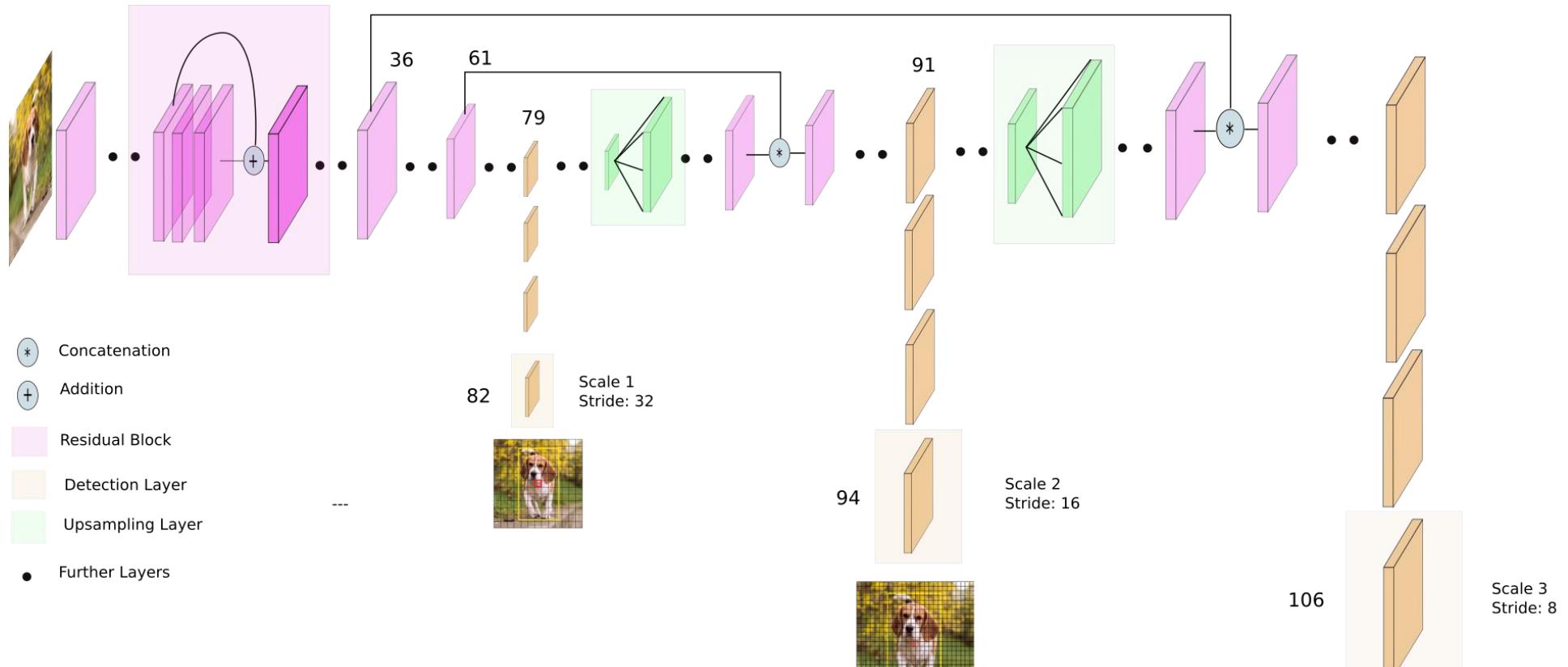
YOLO

too big, too strong, too fast, too good!!!



YOLOv3

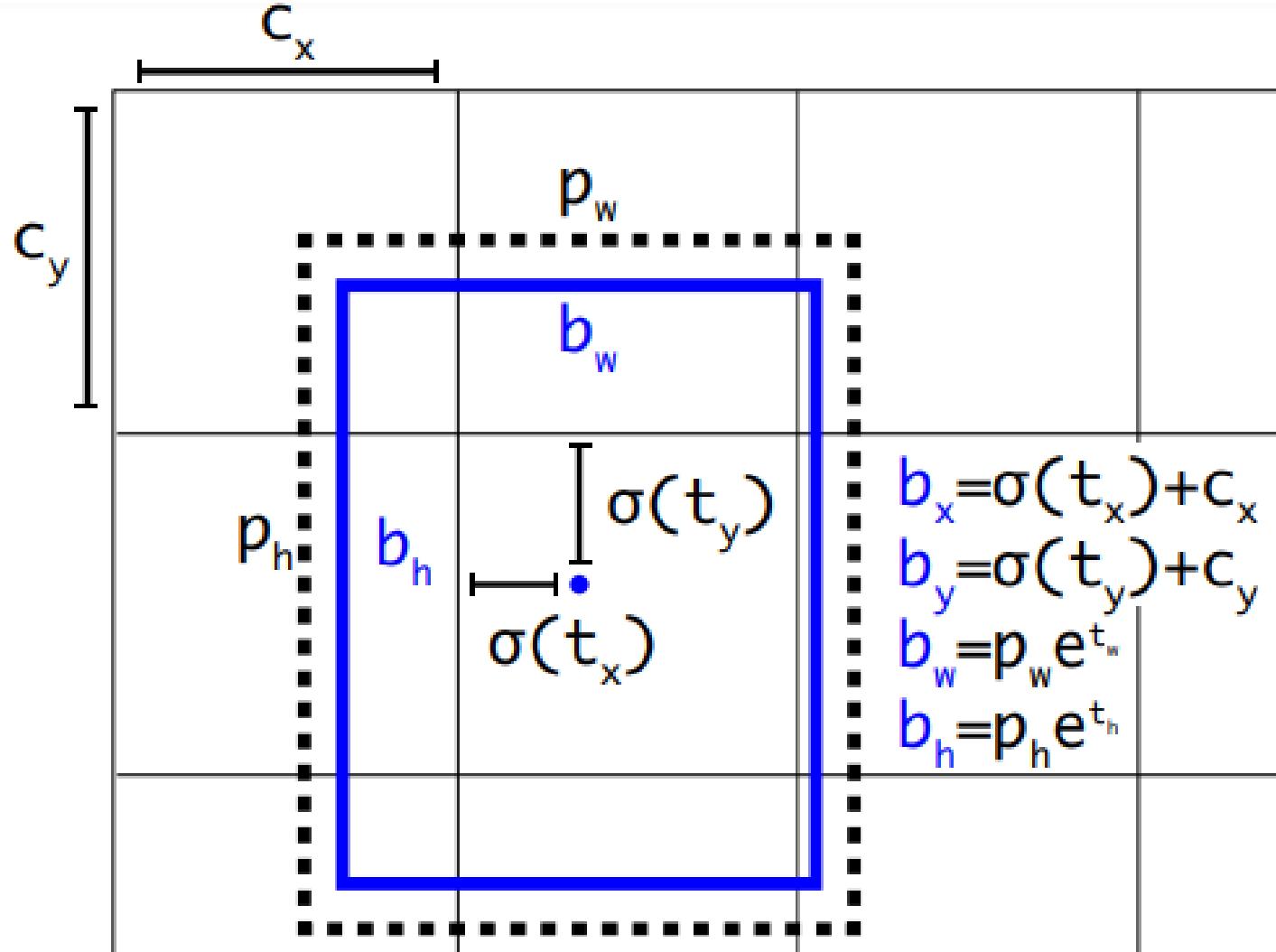
fully convolutional network: only convolutional layers



YOLO v3 network Architecture

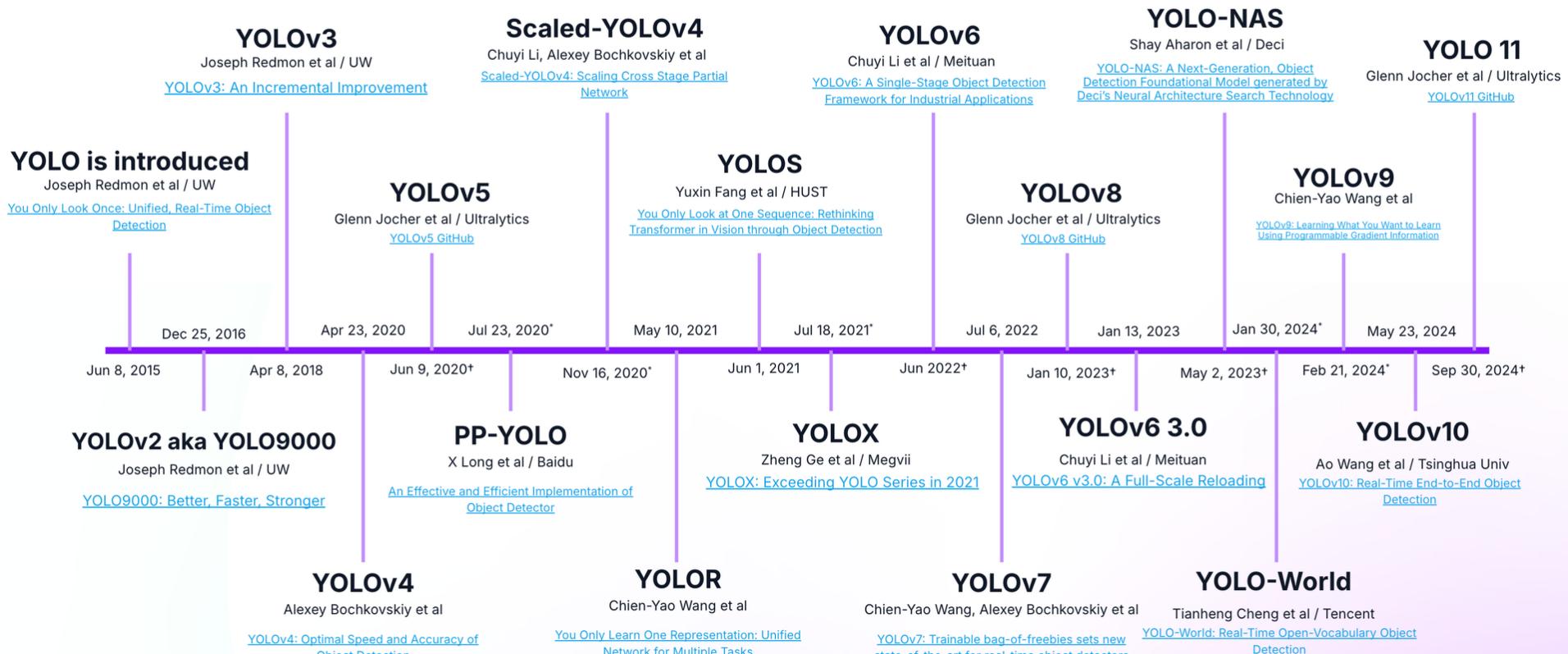
YOLOv3

predicts bounding boxes using dimension clusters as anchor boxes



genesis of YOLO

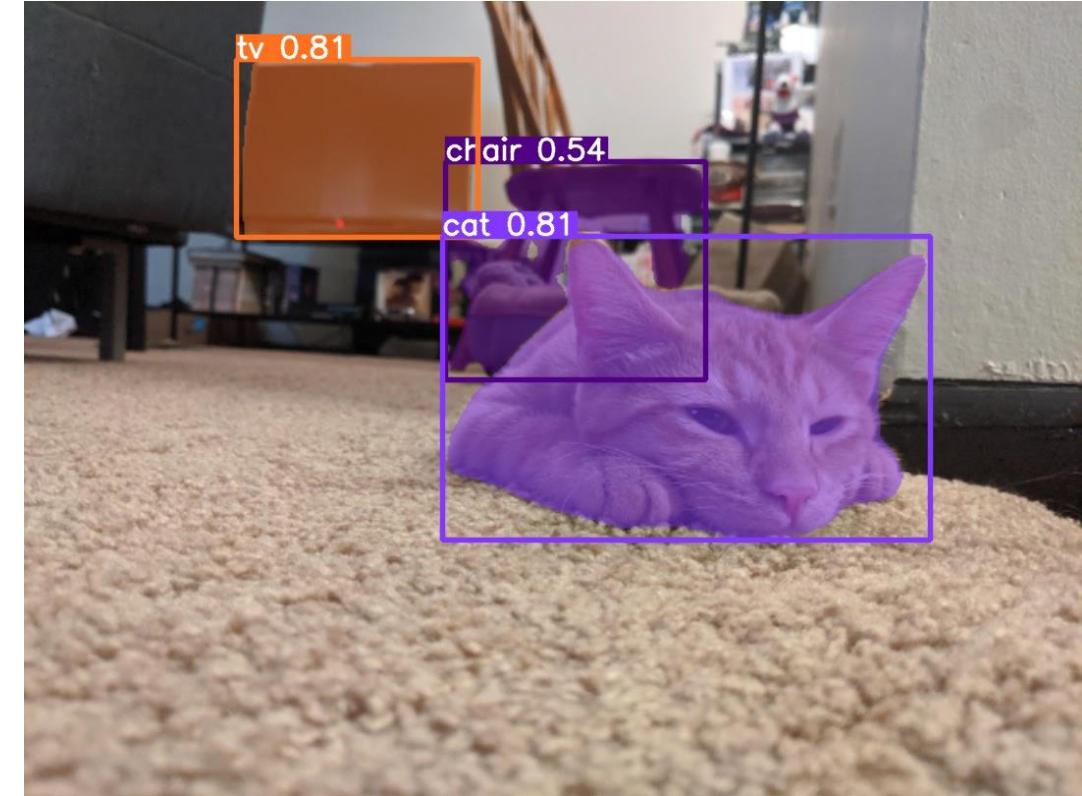
YOLO Greatest Hits

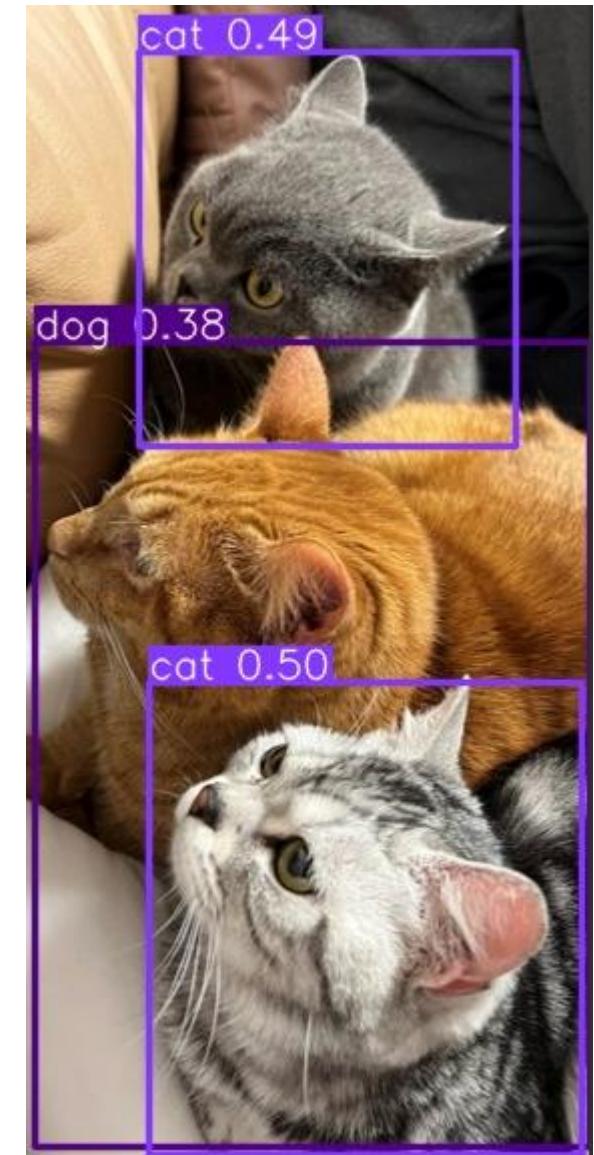


^{*}Denotes paper updated after first publication date

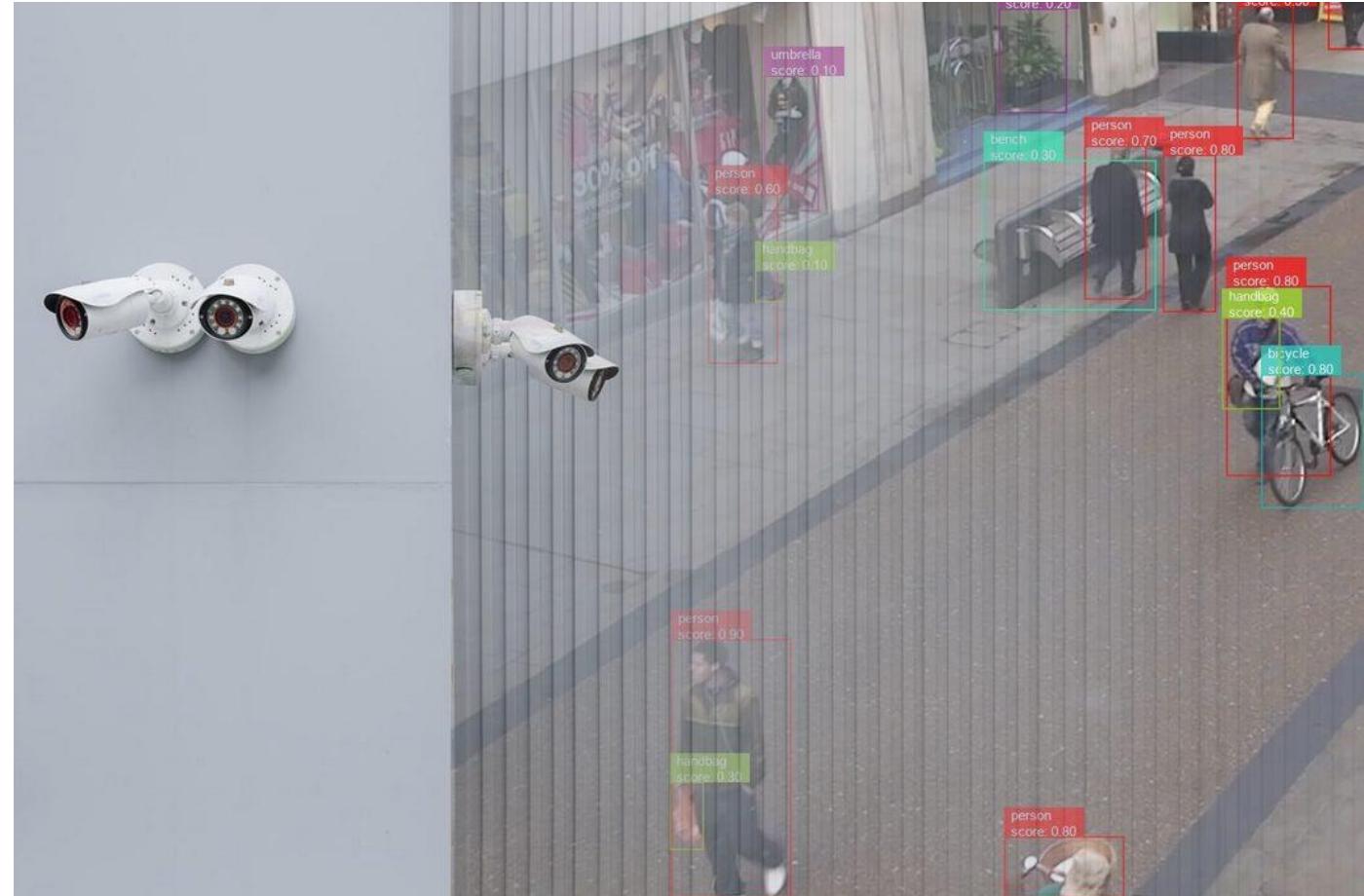
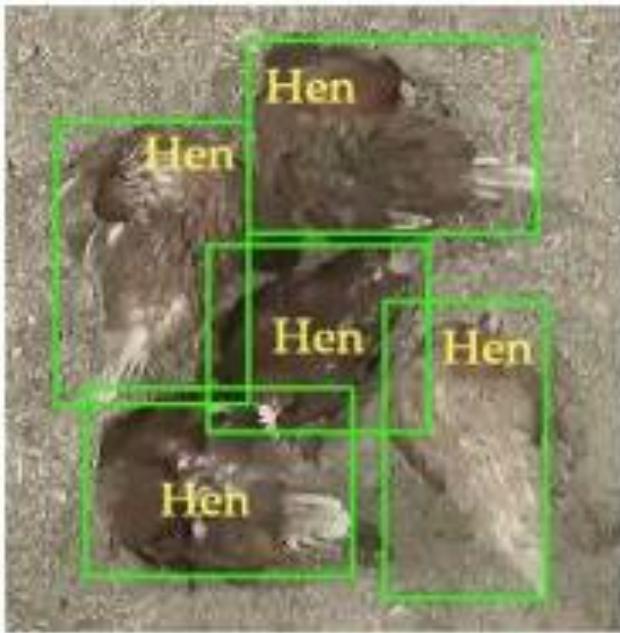
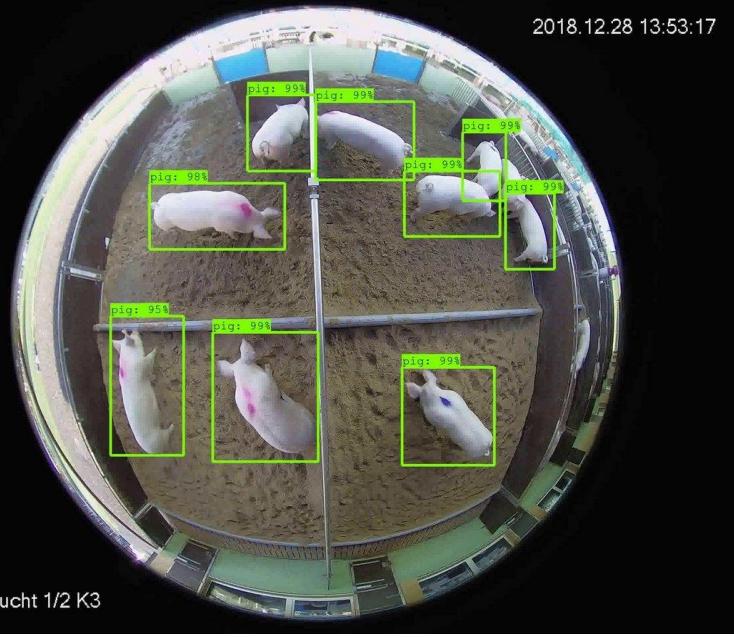
[†]Denotes repository predates paper publication date

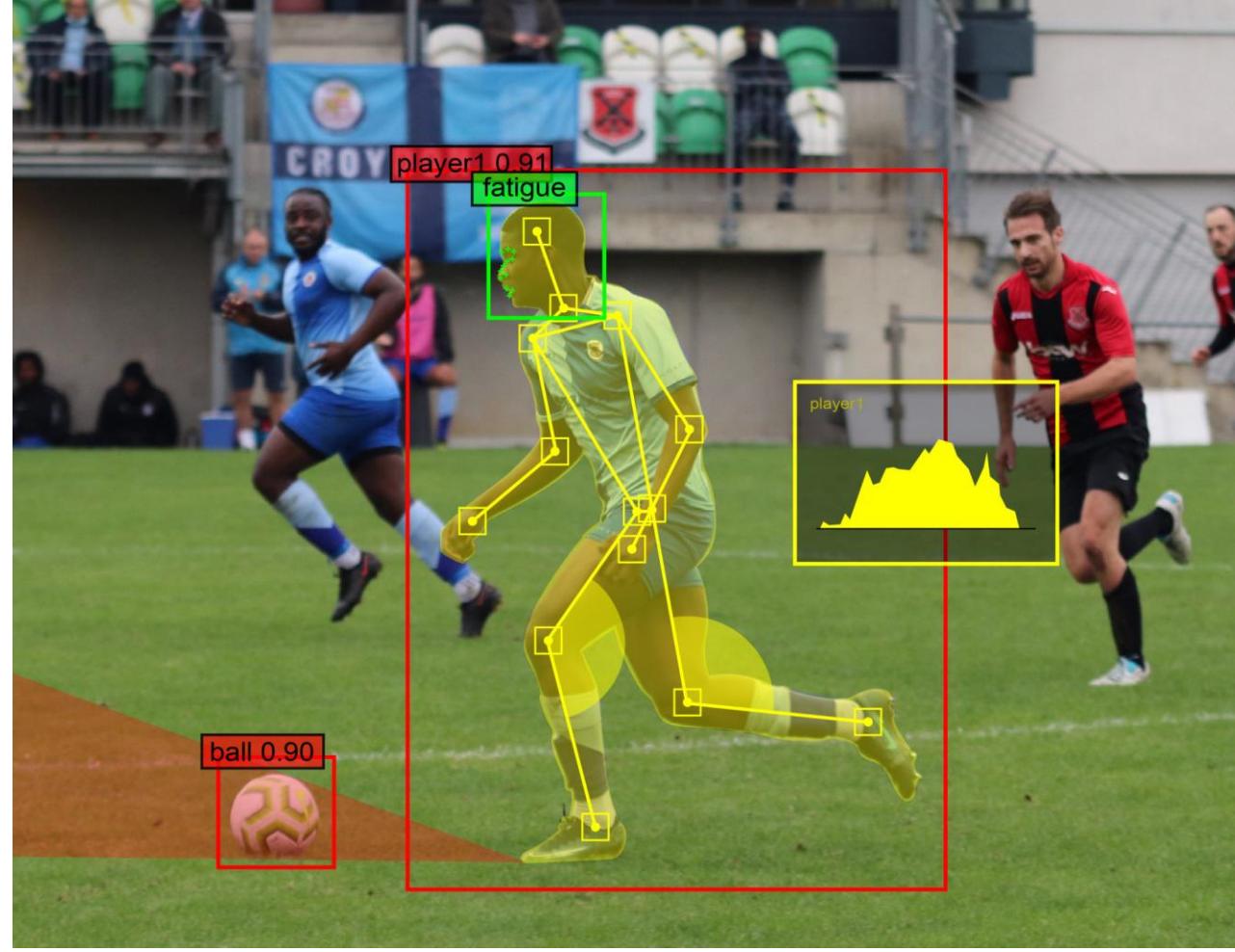






2018.12.28 13:53:17

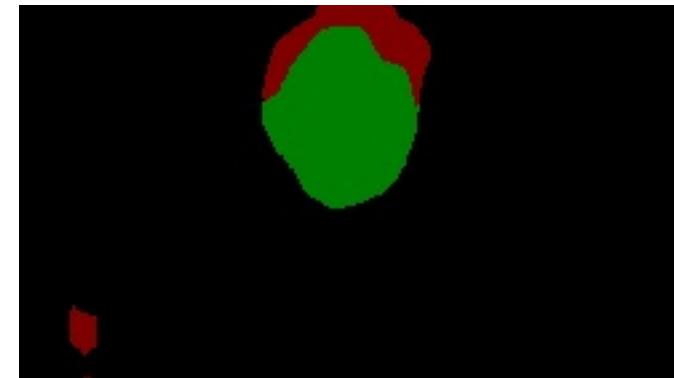




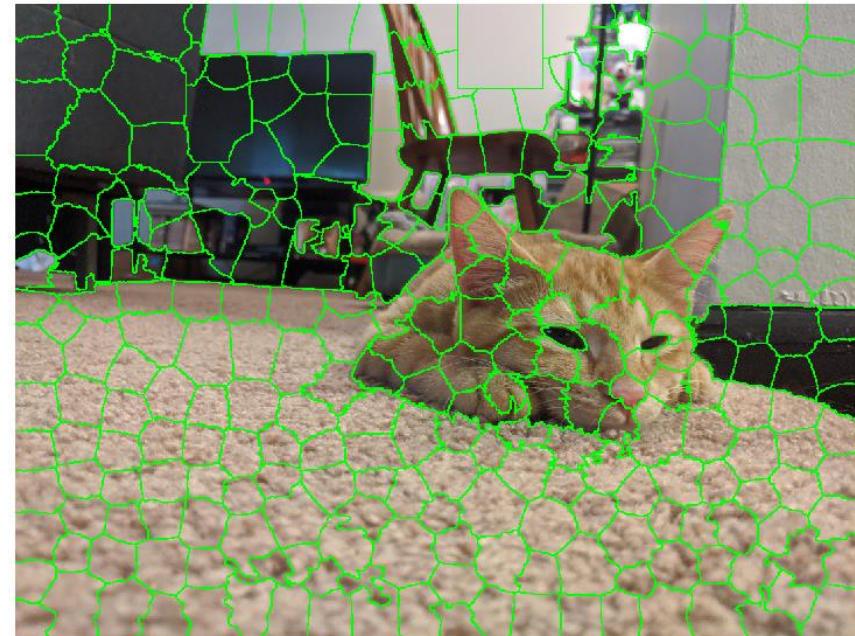
4.5 CNN in segmentation

face segmentation

segments an image into three areas of background, hair and face



superpixels



Labeled Faces in the Wild

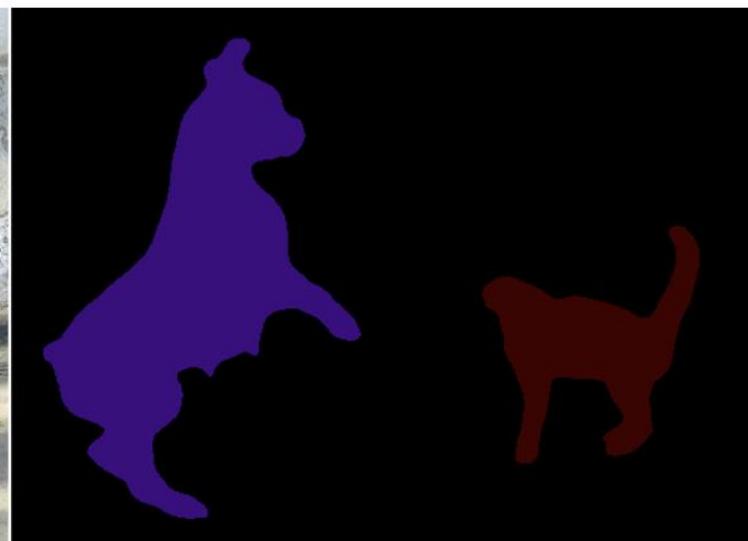


superpixels



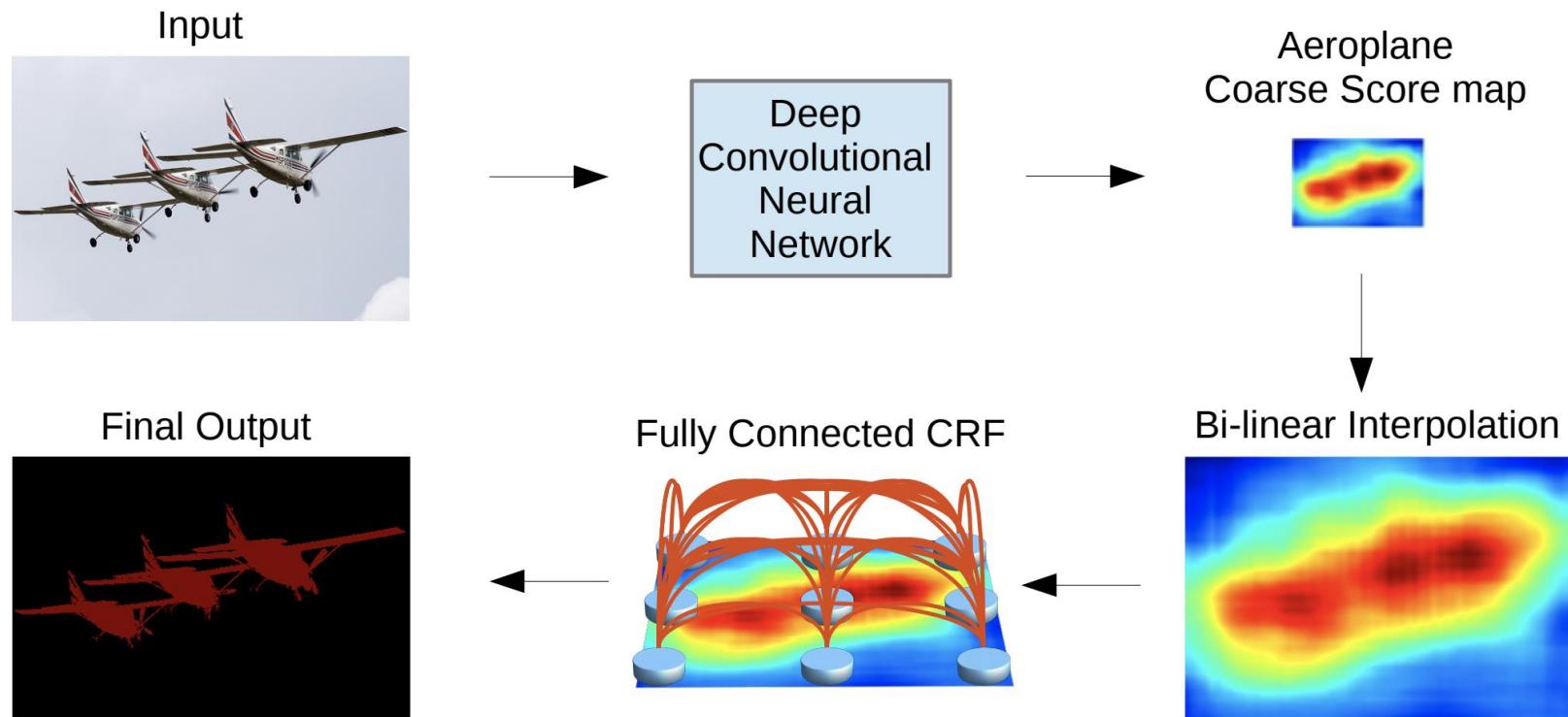
groundtruth

semantic segmentation

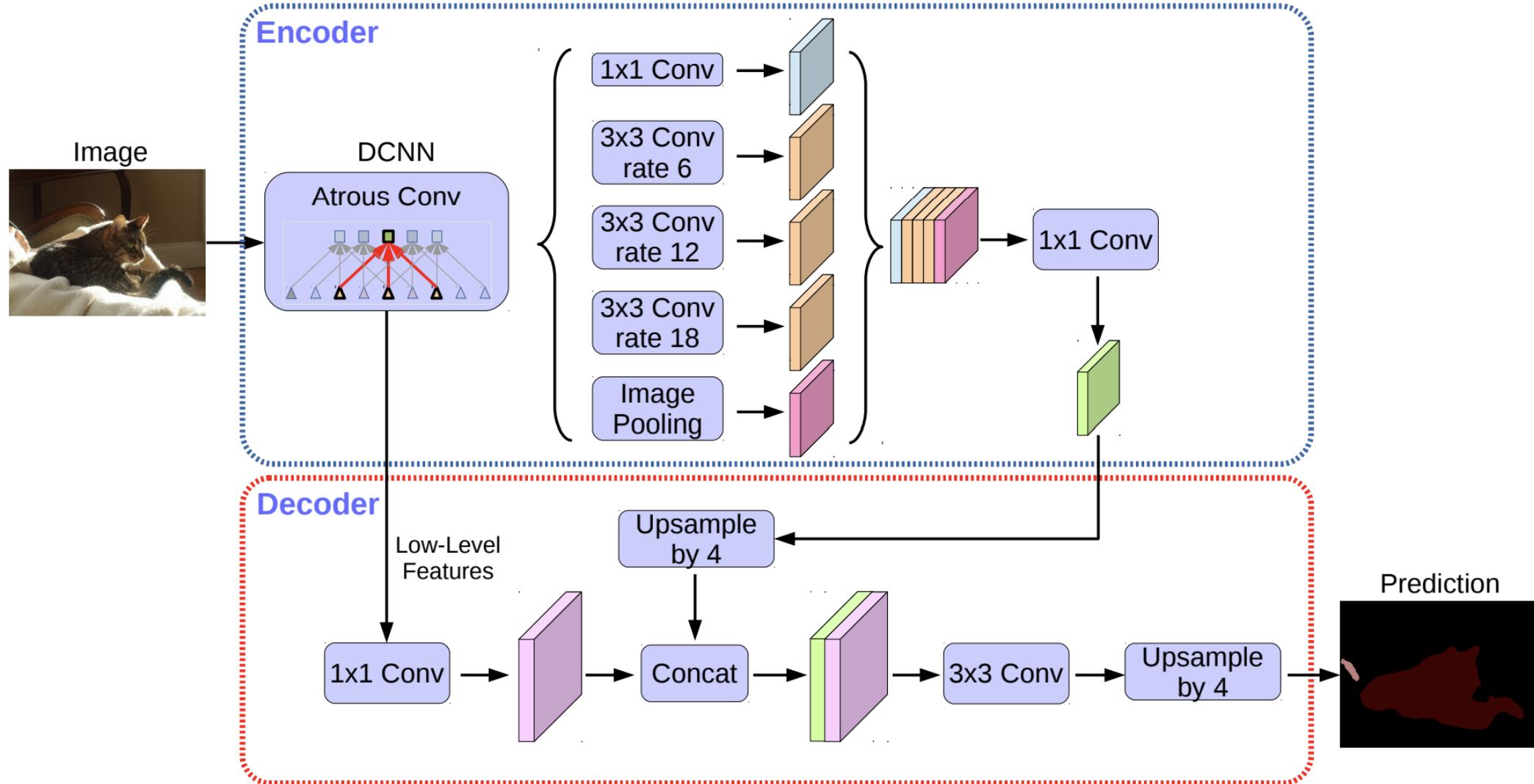


semantic segmentation

- DeepLab



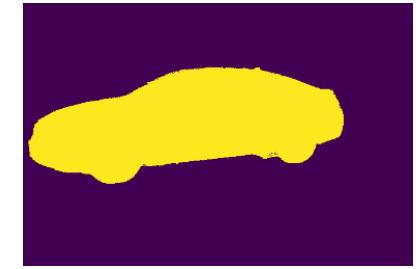
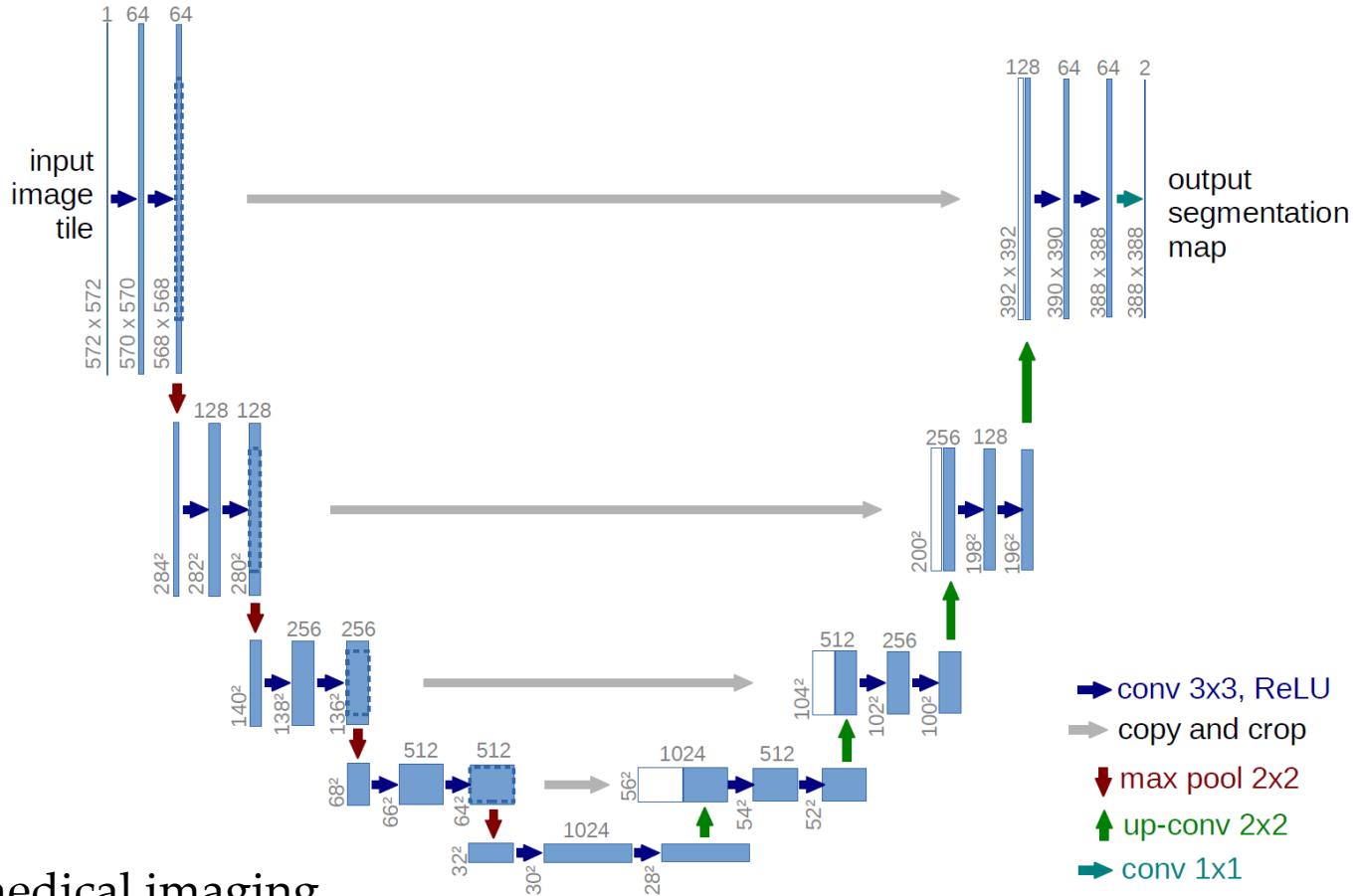
DeepLabv3



deeplab



UNet

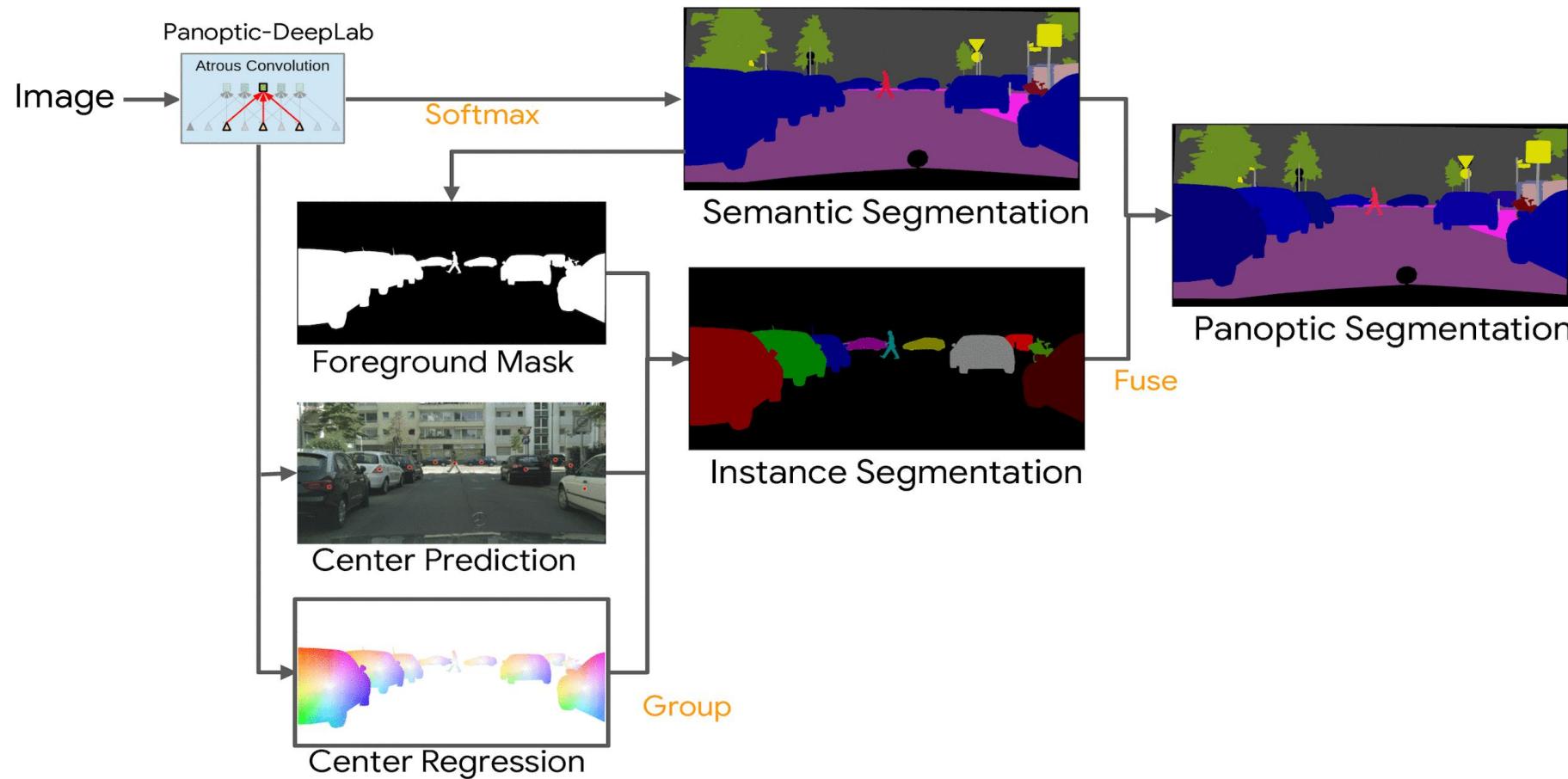


widely used for medical imaging

instance segmentation and panoptic segmentation



Panoptic-DeepLab



Thank you!

Reference

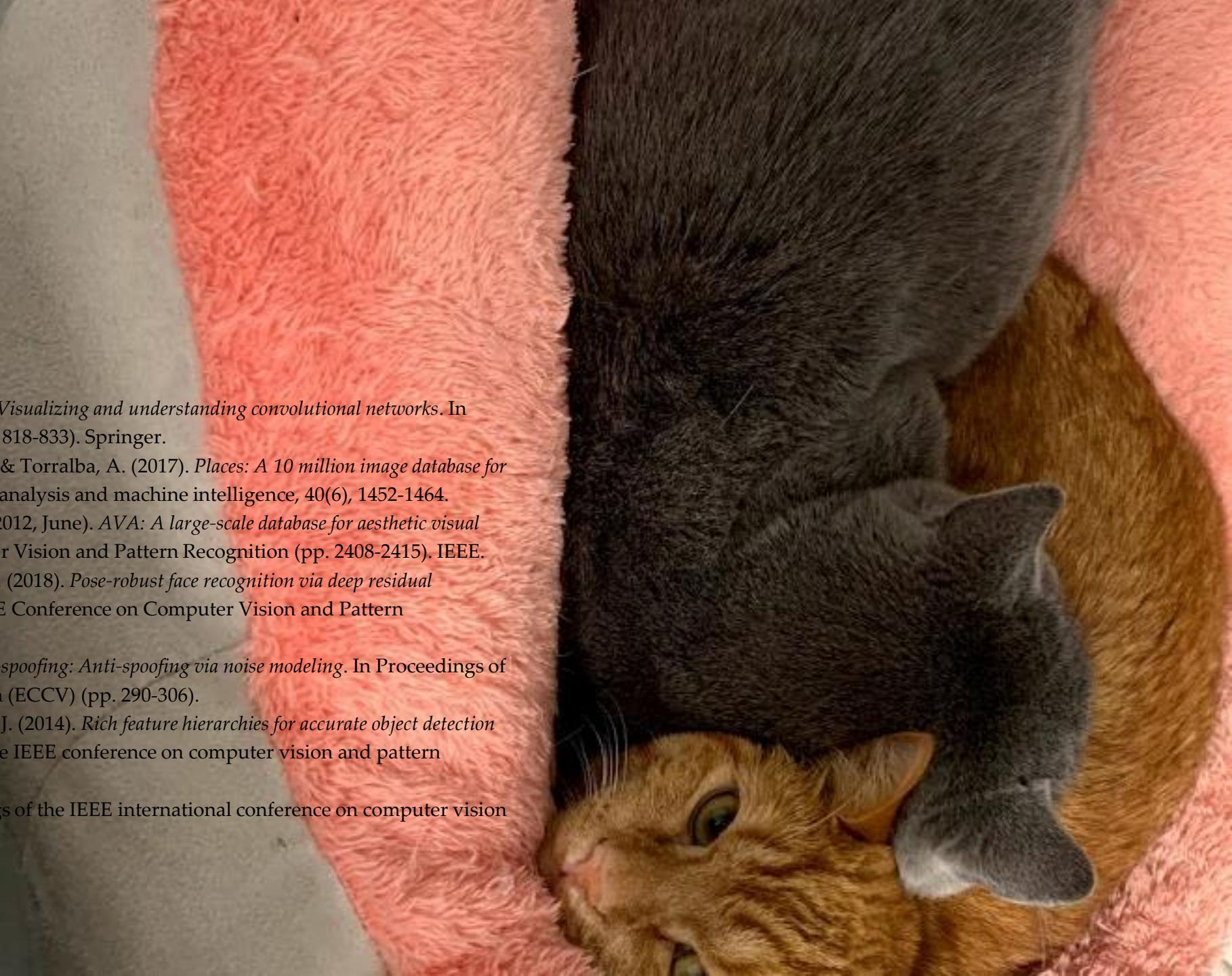
- Ch 5, 7,13, *Deep Learning for Coders with fastai and PyTorch*.
- Ch 9, *Deep Learning*.
- Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In *European conference on computer vision* (pp. 818-833). Springer.
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. In *ICLR 2015*.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... & Rabinovich, A. (2015). Going deeper with convolutions. In *CVPR* (pp. 1-9).
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *CVPR* (pp. 770-778).



Thank you!

Reference

- Zeiler, M. D., & Fergus, R. (2014, September). *Visualizing and understanding convolutional networks*. In European conference on computer vision (pp. 818-833). Springer.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., & Torralba, A. (2017). *Places: A 10 million image database for scene recognition*. IEEE transactions on pattern analysis and machine intelligence, 40(6), 1452-1464.
- Murray, N., Marchesotti, L., & Perronnin, F. (2012, June). *AVA: A large-scale database for aesthetic visual analysis*. In 2012 IEEE Conference on Computer Vision and Pattern Recognition (pp. 2408-2415). IEEE.
- Cao, K., Rong, Y., Li, C., Tang, X., & Loy, C. C. (2018). *Pose-robust face recognition via deep residual equivariant mapping*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 5187-5196).
- Jourabloo, A., Liu, Y., & Liu, X. (2018). *Face de-spoofing: Anti-spoofing via noise modeling*. In Proceedings of the European Conference on Computer Vision (ECCV) (pp. 290-306).
- Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). *Rich feature hierarchies for accurate object detection and semantic segmentation*. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 580-587).
- Girshick, R. (2015). *Fast R-CNN*. In Proceedings of the IEEE international conference on computer vision (pp. 1440-1448).



Thank you!

Reference

- Ren, S., He, K., Girshick, R., & Sun, J. (2015). *Faster R-CNN: Towards real-time object detection with region proposal networks*. Advances in neural information processing systems, 28, 91-99.
- Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). *You only look once: Unified, real-time object detection*. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 779-788).
- Redmon, J., & Farhadi, A. (2018). *Yolov3: An incremental improvement*. arXiv preprint arXiv:1804.02767.
- Kae, A., Sohn, K., Lee, H., & Learned-Miller, E. (2013). *Augmenting CRFs with Boltzmann machine shape priors for image labeling*. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 2019-2026).
- Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2017). *Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs*. IEEE transactions on pattern analysis and machine intelligence, 40(4), 834-848.
- Ma, L., Jia, X., Sun, Q., Schiele, B., Tuytelaars, T., & Van Gool, L. (2017). *Pose guided person image generation*. In Proceedings of the 31st International Conference on Neural Information Processing Systems (pp. 405-415).
- Zhu, J. Y., Park, T., Isola, P., & Efros, A. A. (2017). *Unpaired image-to-image translation using cycle-consistent adversarial networks*. In Proceedings of the IEEE international conference on computer vision (pp. 2223-2232).
- Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., ... & Shi, W. (2017). *Photo-realistic single image super-resolution using a generative adversarial network*. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 4681-4690).

