

Computer Vision

Naeemullah Khan
naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

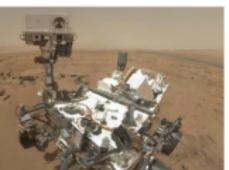
November 19, 2023

Building artificial systems that process, perceive, and reason about visual data

Computer Vision is Everywhere



Left to right:
[Image by Roger H Giesen](#) is licensed under [CC BY 2.0](#)
[Image](#) is CC0 1.0 public domain
[Image](#) is CC0 1.0 public domain
[Image](#) is CC0 1.0 public domain



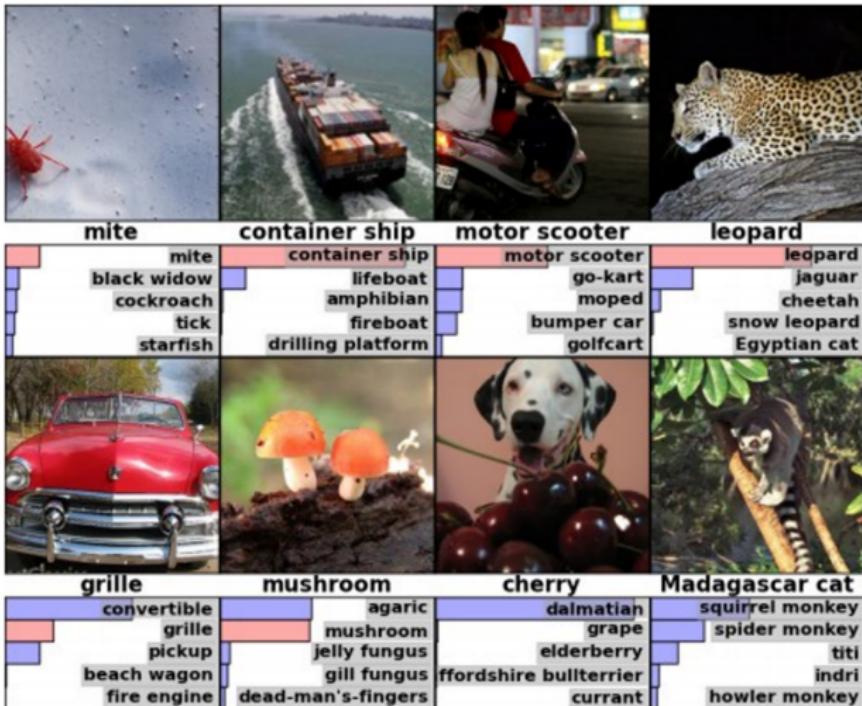
Left to right:
[Image](#) is free to use
[Image](#) is CC0 1.0 public domain
[Image](#) by NASA is licensed under CC BY 2.0
[Image](#) is CC0 1.0 public domain



Bottom row, left to right:
[Image](#) is CC0 1.0 public domain
[Image](#) by Derek Keats is licensed under [CC BY 2.0](#); changes made
[Image](#) is public domain
[Image](#) is licensed under [CC BY 2.0](#); changes made

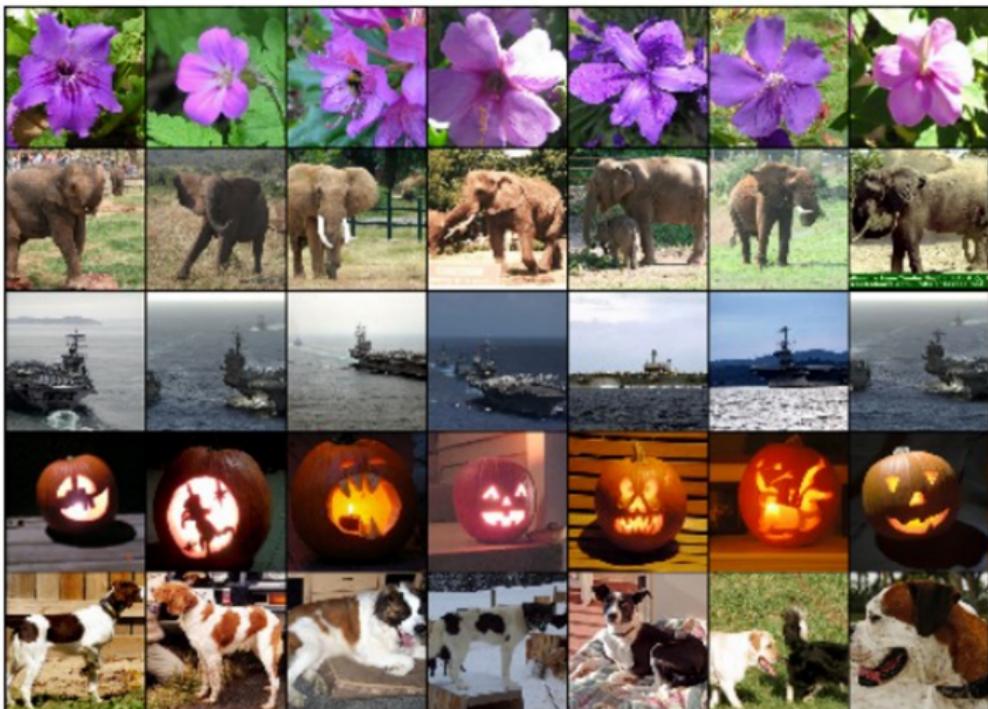
Some Applications

Image Classification



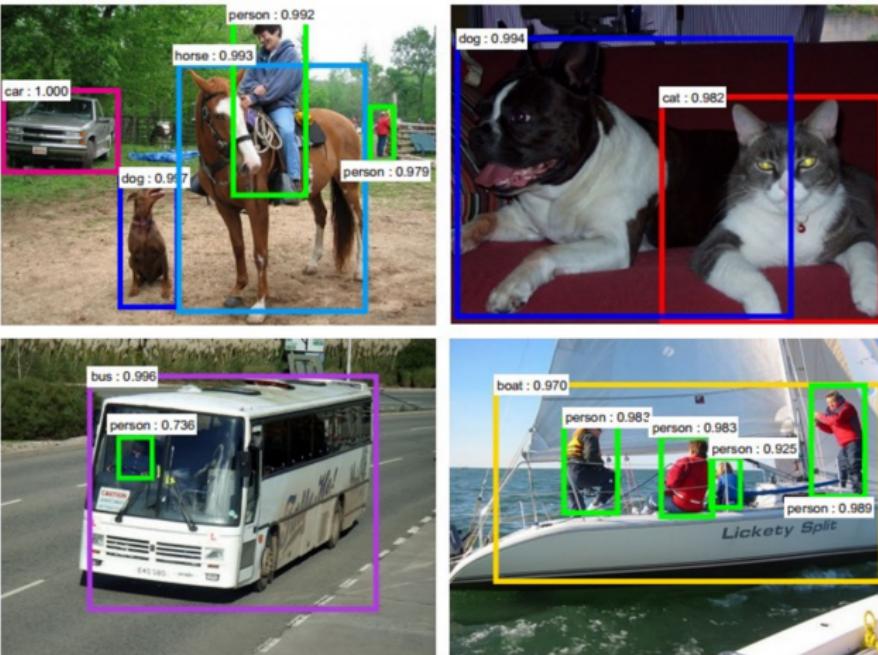
Some Applications (cont.)

Image Retrieval



Some Applications (cont.)

Object Detection



Ren, He, Girshick, and Sun, 2015

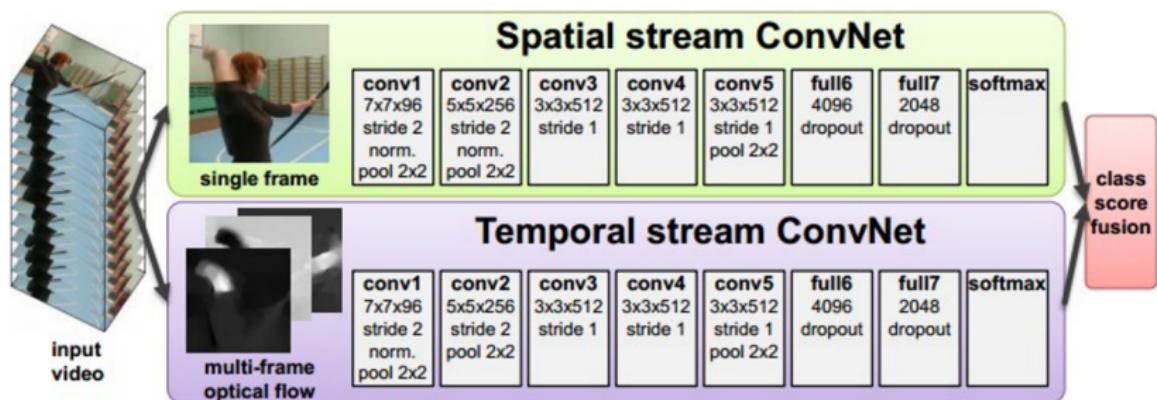
Some Applications (cont.)

Image Segmentation



Fabaret et al, 2012

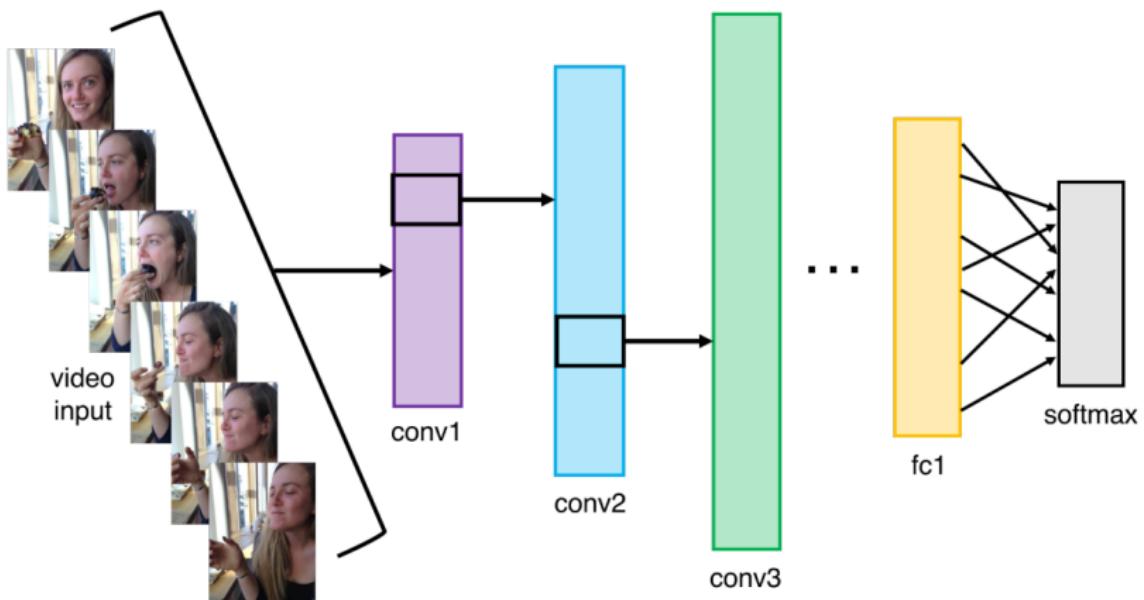
Video Classification



Simonyan et al, 2014

Some Applications (cont.)

Activity Recognition



Some Applications (cont.)

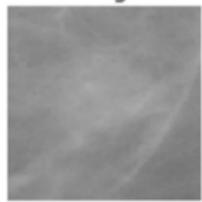
Pose Recognition (Toshev and Szegedy, 2014)



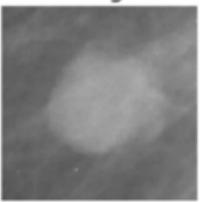
Some Applications (cont.)

Medical Imaging

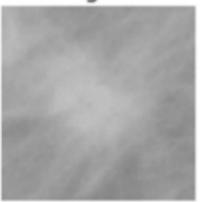
Benign



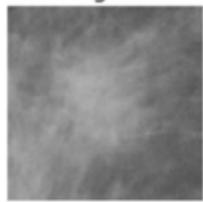
Benign



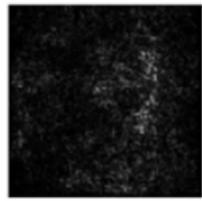
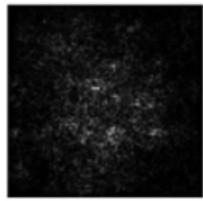
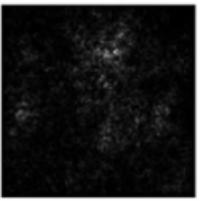
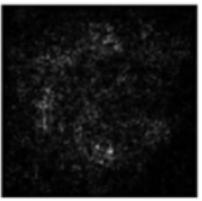
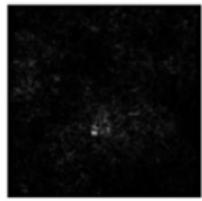
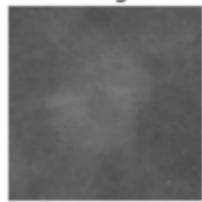
Malignant



Malignant



Benign



Some Applications (cont.)



Image Captioning

Vinyals et al, 2015

Karpathy and Fei-Fei, 2015



A man in a baseball uniform throwing a ball



A woman is holding a cat in her hand



A man riding a wave on top of a surfboard



A cat sitting on a suitcase on the floor



A woman standing on a beach holding a surfboard

All images are CC0 Public domain:

http://yohay.com/coco/teddy_bear_in_grass_1603000/
http://yohay.com/coco/baseball_player_throwing_ball_in_diamond_1423400/
http://yohay.com/coco/cat_sitting_in_trunk_3030000/
http://yohay.com/coco/surfer_riding_a_wave_4900000/
http://yohay.com/coco/woman_standing_on_sandholding_surfboard_2000000/

Some Applications (cont.)

Image Generation



“Teddy bears working on new AI research underwater with 1990s technology”

DALL-E 2

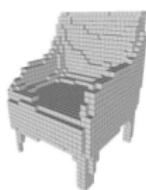
Some Applications (cont.)



Style Transfer

Some Applications (cont.)

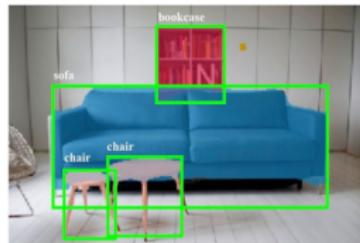
3D Vision



Choy et al., 3D-R2N2: Recurrent Reconstruction Neural Network (2016)



Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)



Gkioxari et al., "Mesh R-CNN", ICCV 2019

How to represent an image?

- ▶ Images are represented as Matrices with elements in [0, 255]
- ▶ Grayscale images have one channel while RGB images have 3 channels



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	34	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	261	237	239	239	228	227	87	71	201
172	105	207	238	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	168	139	76	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	282	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	258	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	79	1	81	47	0	6	217	258	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	206	175	13	95	218

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	34	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	261	237	239	239	228	227	87	71	201
172	105	207	238	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	168	139	76	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	282	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	258	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	79	1	81	47	0	6	217	258	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	206	175	13	95	218

⁰<https://www.v7labs.com/blog/image-recognition-guide>

Fully-Connected Neural Networks

Deep Neural Network

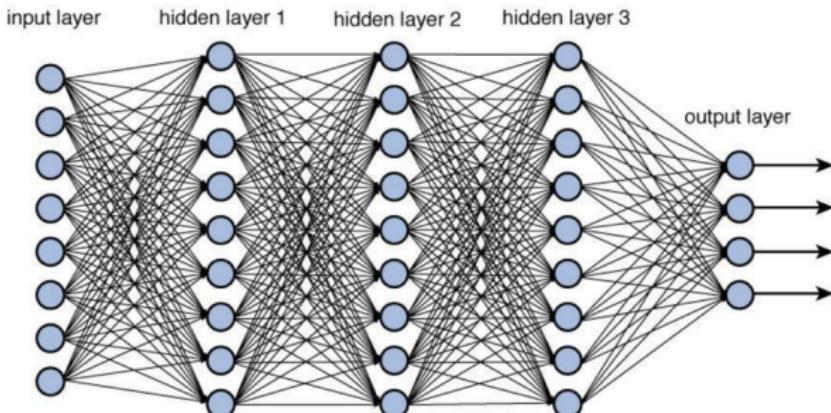


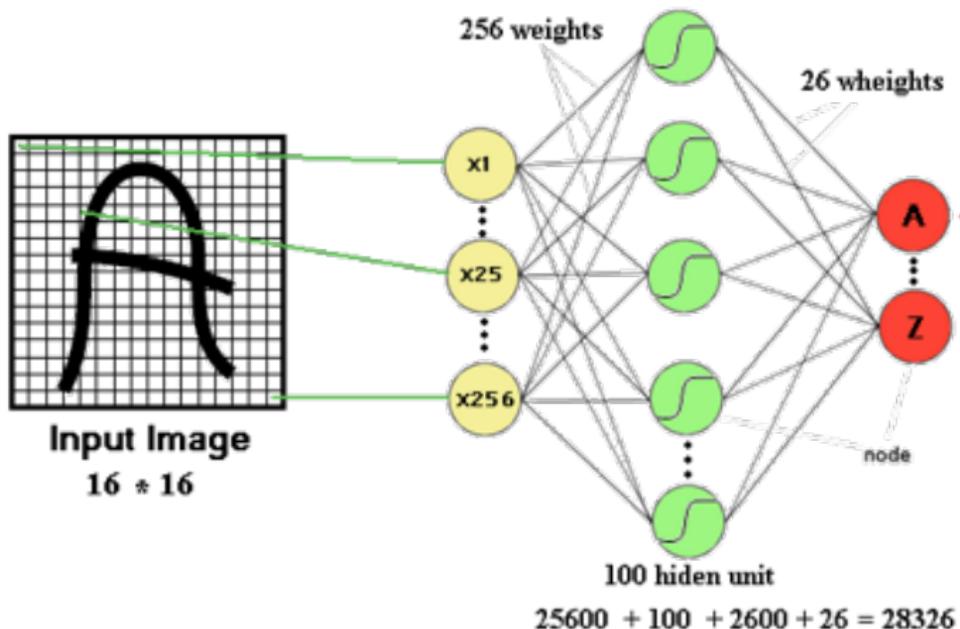
Figure 12.2 Deep network architecture with multiple layers.

$$z = W_1x_1 + W_2x_2 + \cdots + W_nx_n + b$$

⁰<https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964>

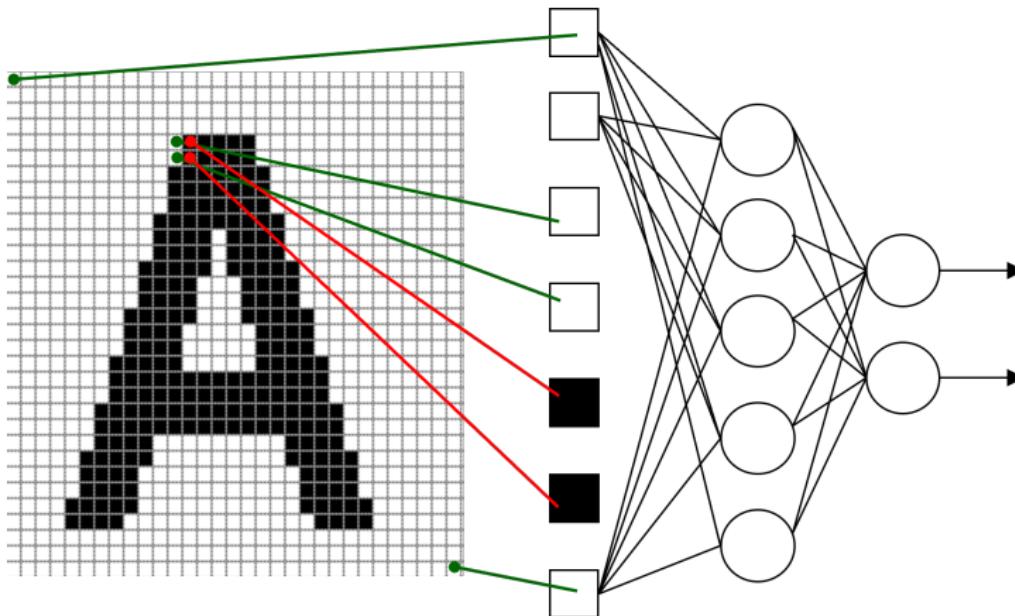
Drawbacks of Fully-Connected Neural Networks

- ▶ The number of trainable parameters becomes extremely large



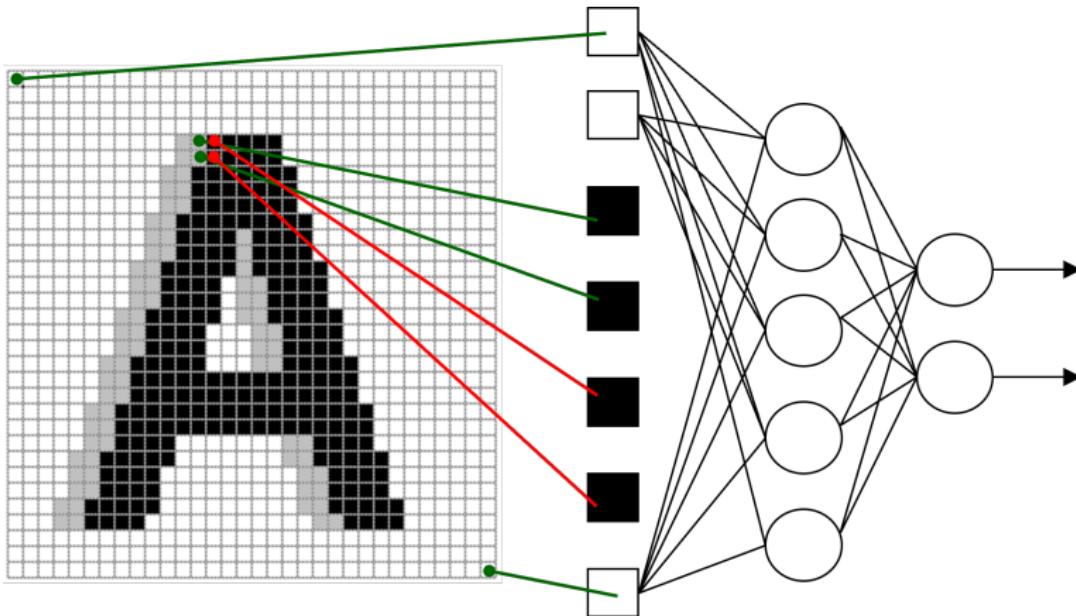
Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ Little or no invariance to shifting, scaling, and other forms of distortion



Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ Little or no invariance to shifting, scaling, and other forms of distortion

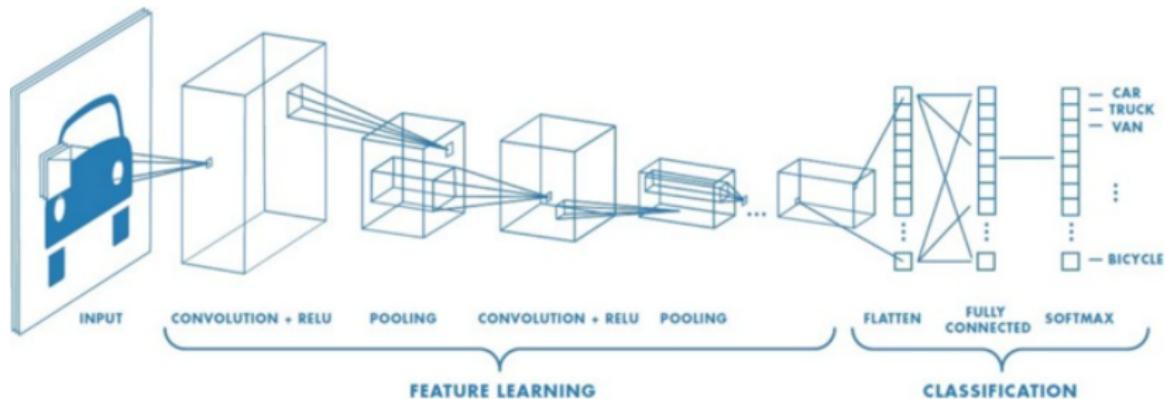


Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ The topology of the input data is completely ignored
- ▶ For a 32×32 image, we have
 - Black and white patterns: $2^{32 \times 32} = 2^{1024}$
 - Grayscale patterns: $256^{32 \times 32} = 256^{1024}$



Convolutional Neural Networks (CNNs)



$$z = W * x_{i,j} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} W_{ab} x_{(i+a)(j+b)}$$

How Convolution Works?

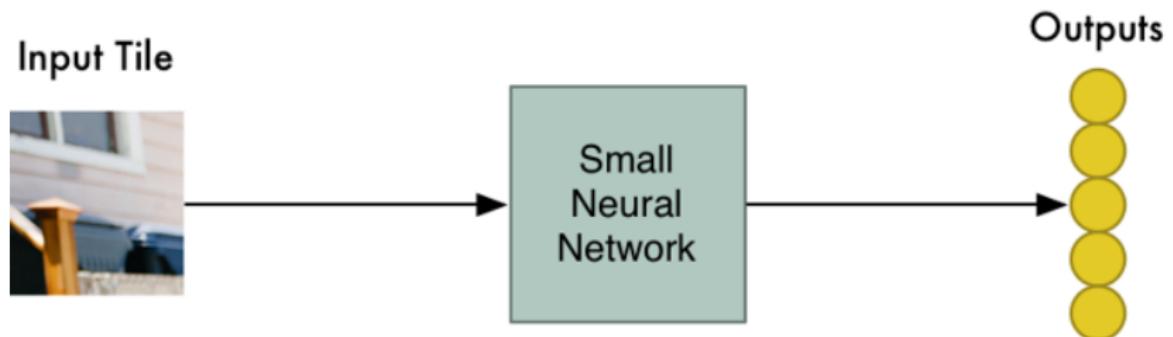


How Convolution Works? (cont.)

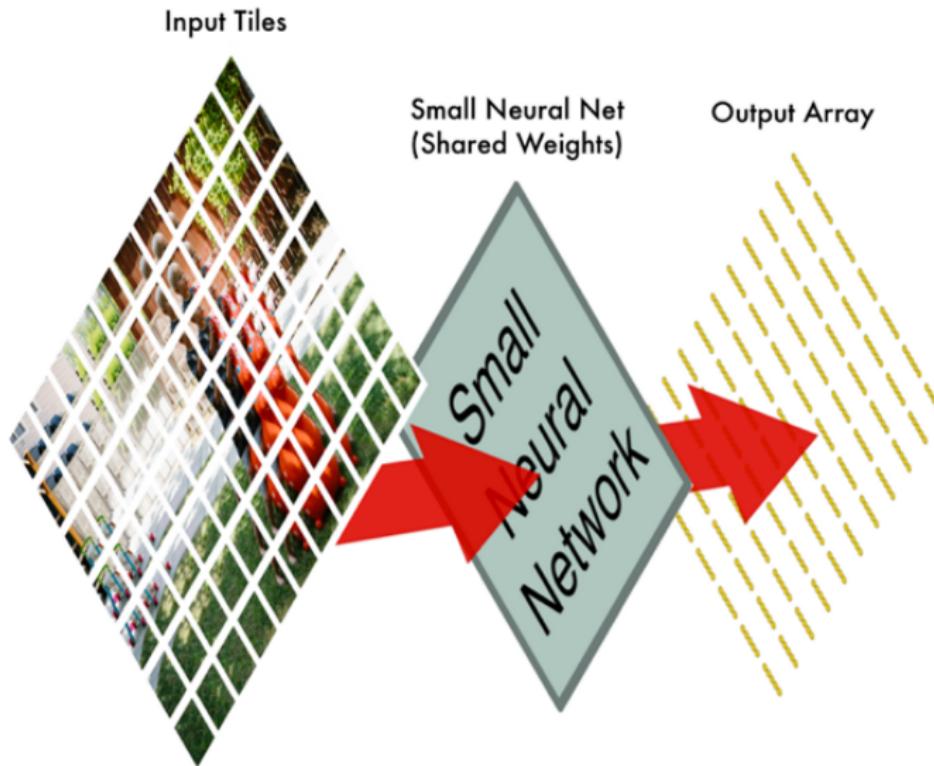


How Convolution Works? (cont.)

Processing a single tile

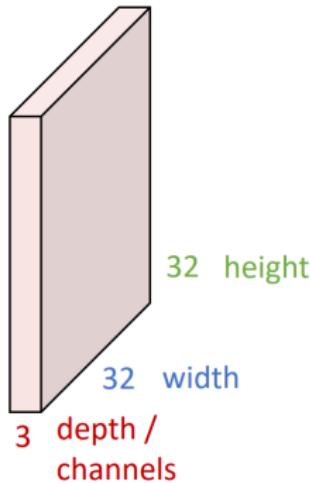


How Convolution Works? (cont.)



How Convolution Works? (cont.)

3x32x32 image

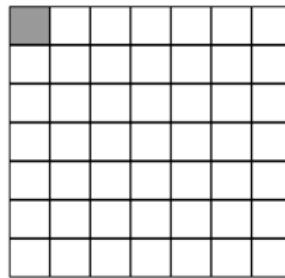
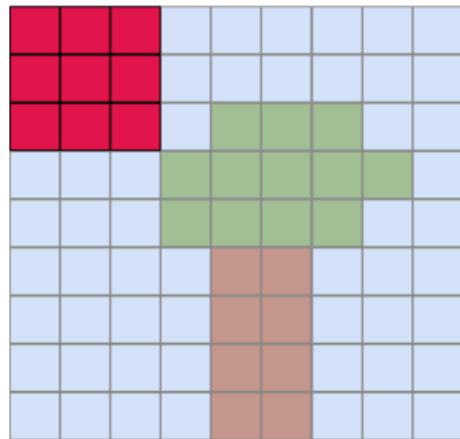


3x5x5 filter



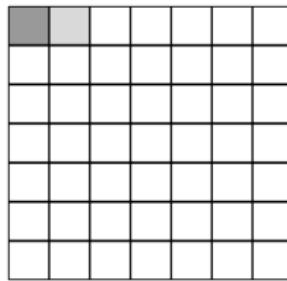
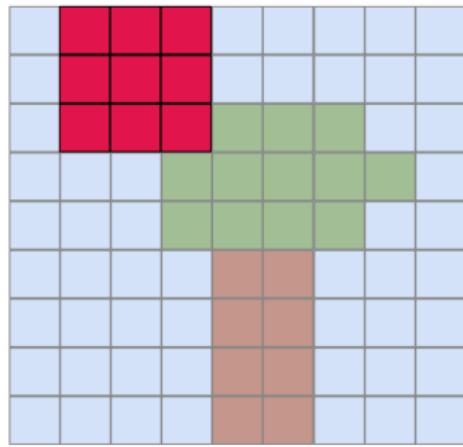
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

How Convolution Works? (cont.)



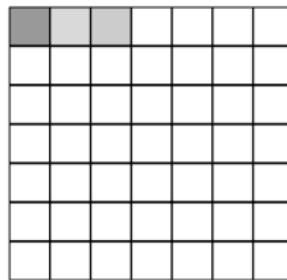
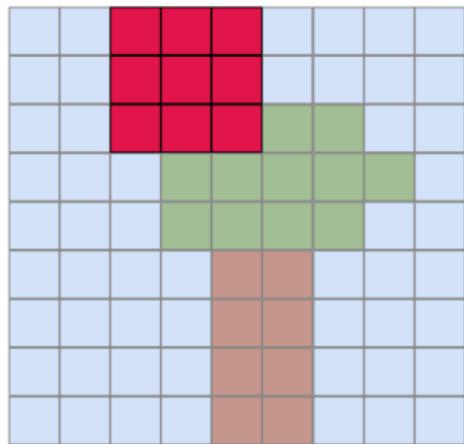
The **kernel** slides across the image and produces an output value at each position

How Convolution Works? (cont.)



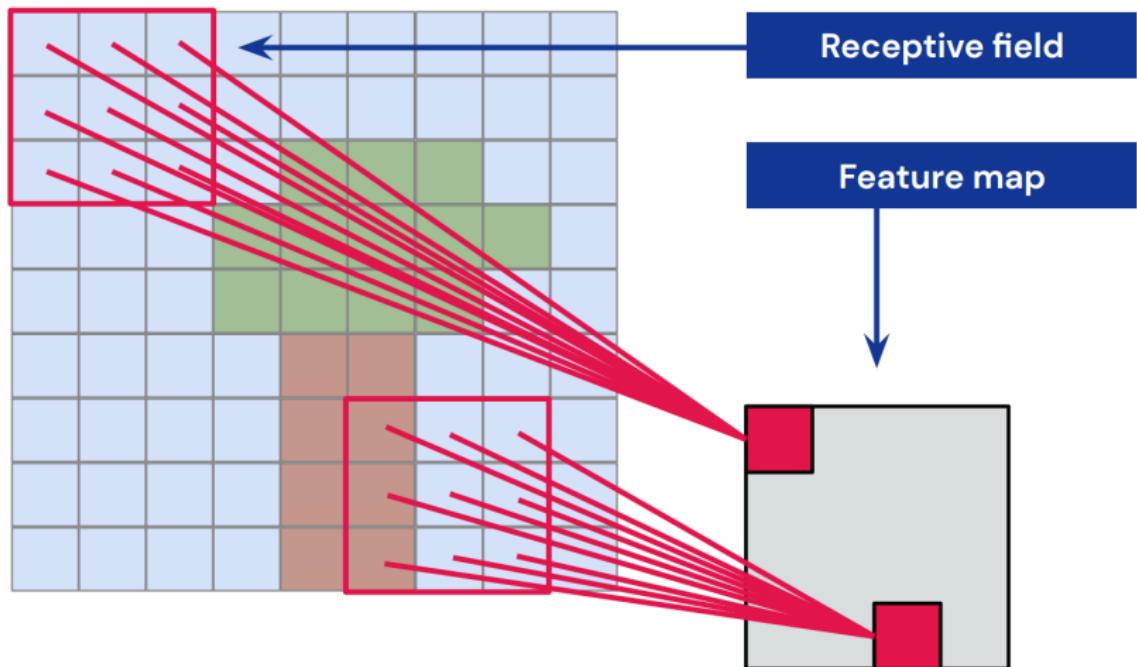
The **kernel** slides across the image and produces an output value at each position

How Convolution Works? (cont.)

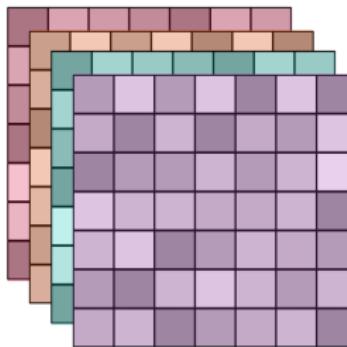
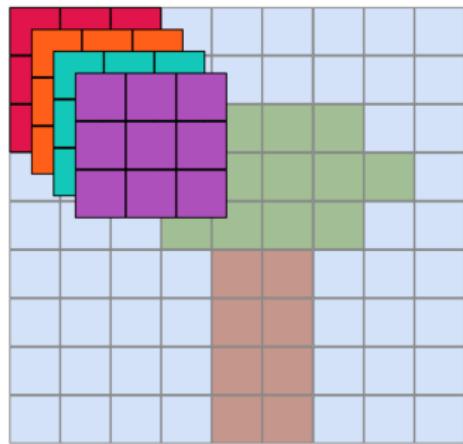


The **kernel** slides across the image and produces an output value at each position

How Convolution Works? (cont.)



How Convolution Works? (cont.)

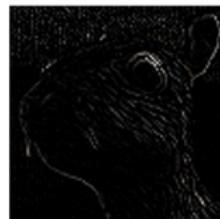
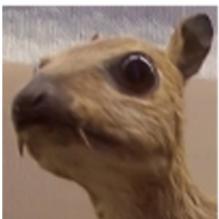


We convolve multiple kernels and obtain multiple feature maps or **channels**

How Convolution Works? (cont.)

$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$$

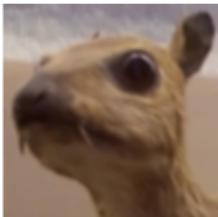
$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



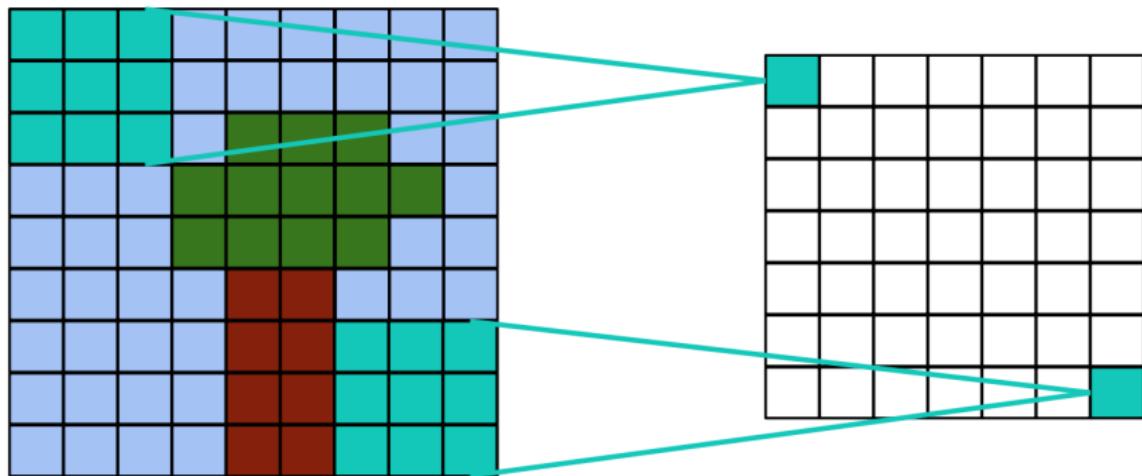
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

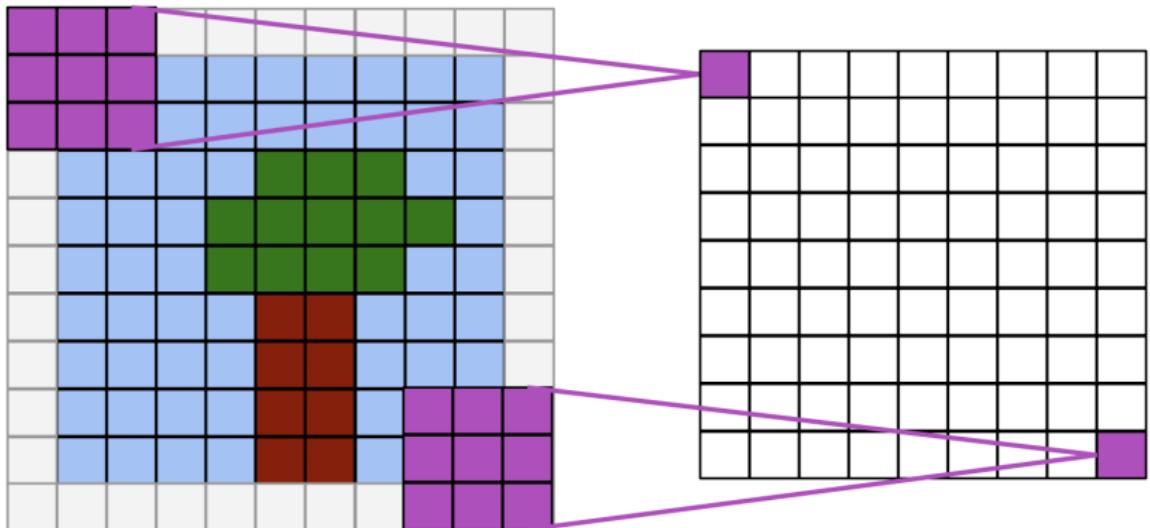


- ▶ Applying Convolution as such reduces the size of the borders.
- ▶ Sometimes this is not desirable.
- ▶ We can pad the border with zeros.



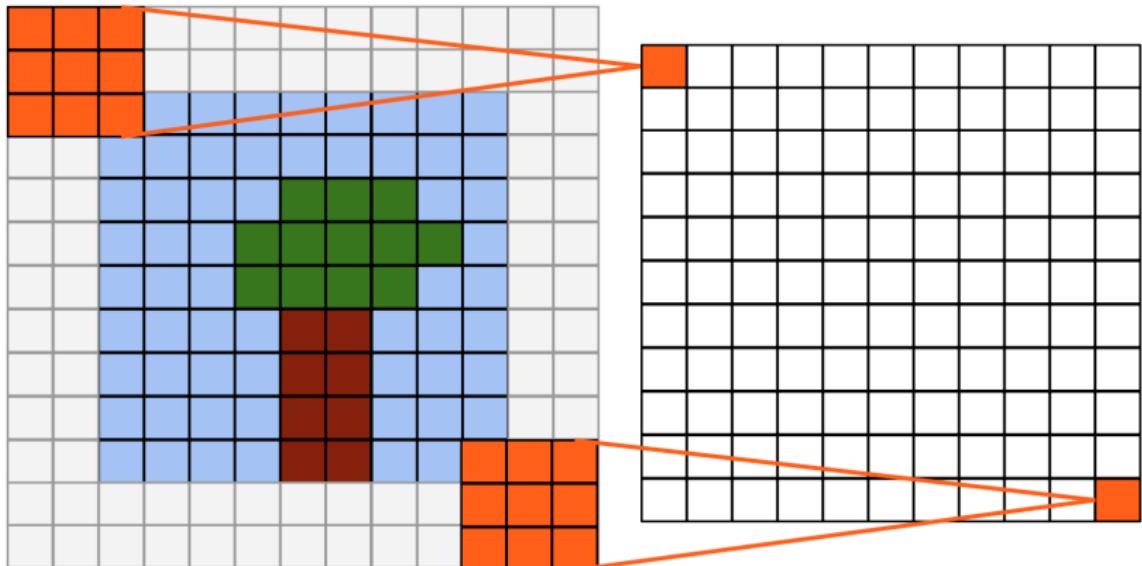
Padding (cont.)

- ▶ Same Convolution: Output is the same size as input



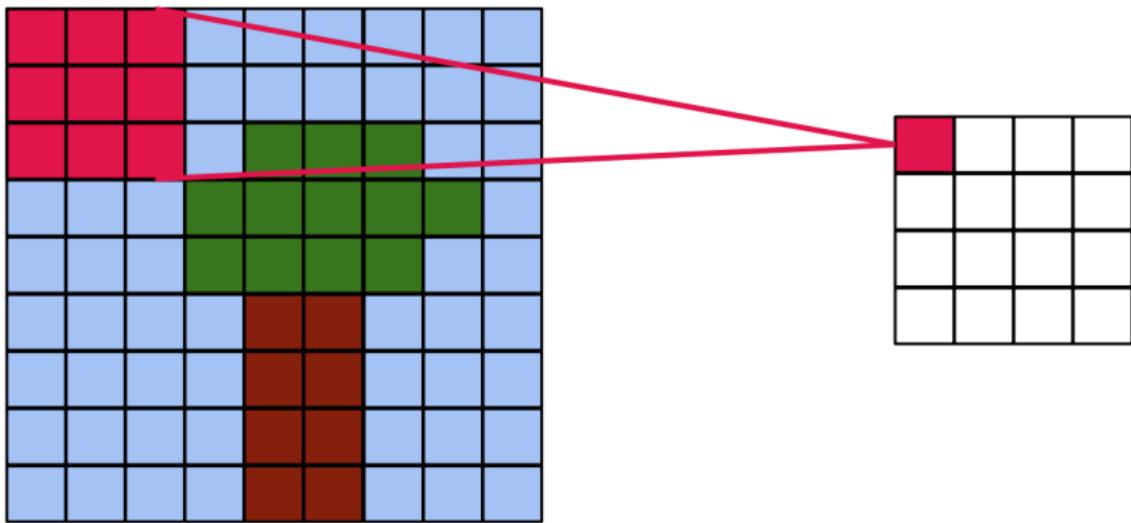
Padding (cont.)

- ▶ Full Convolution: $\text{output size} = \text{input size} + \text{kernel size} - 1$



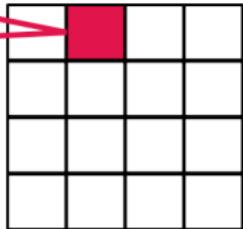
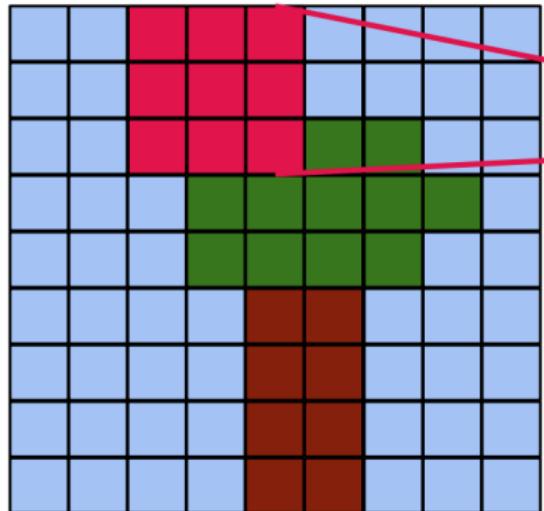
Strided Convolution

- ▶ Kernel slides along the image with a step > 1



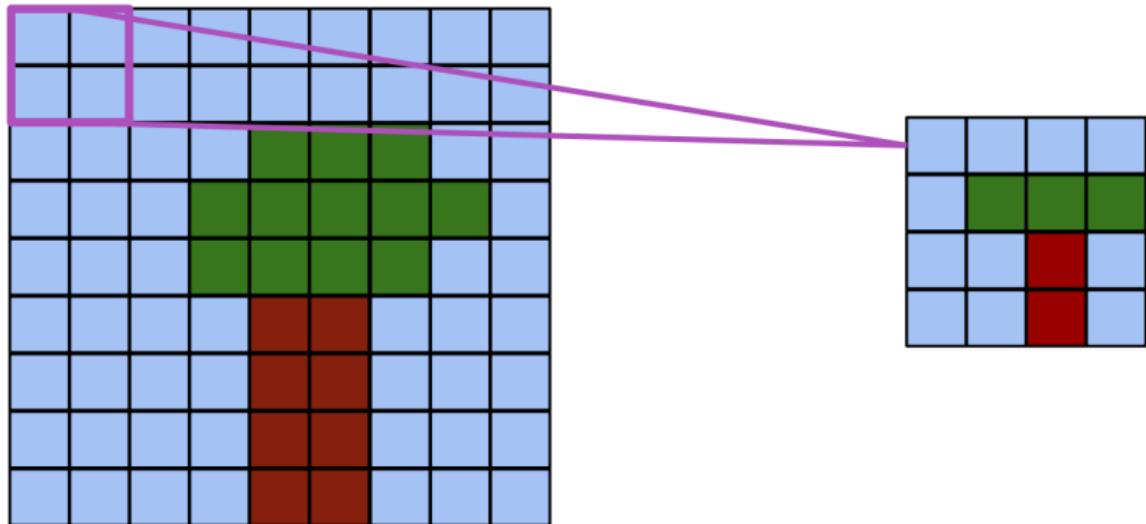
Strided Convolution (cont.)

- ▶ Kernel slides along the image with a step > 1

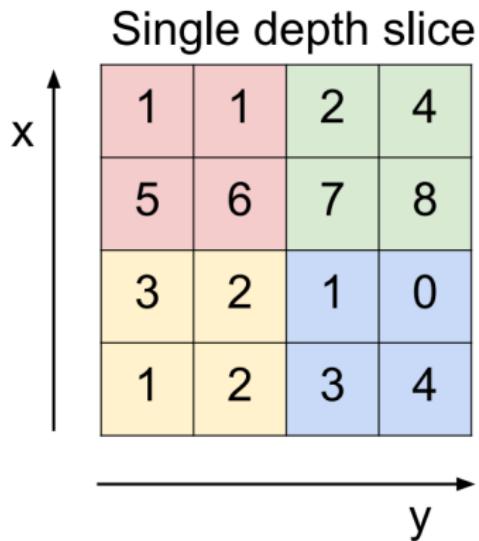


Pooling

- ▶ Compute mean or max over small windows to reduce resolution



Pooling (cont.)



max pool with 2x2 filters
and stride 2

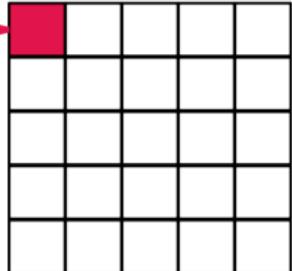
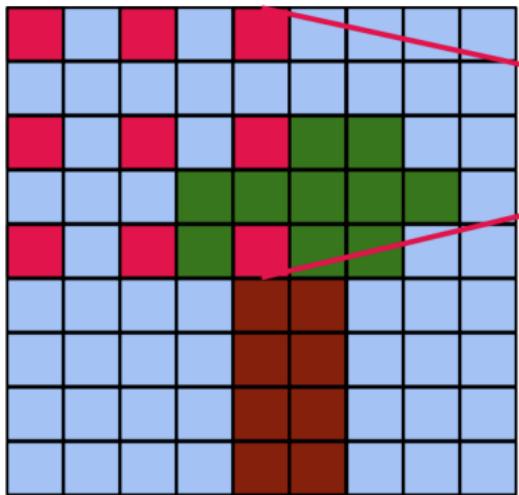


6	8
3	4

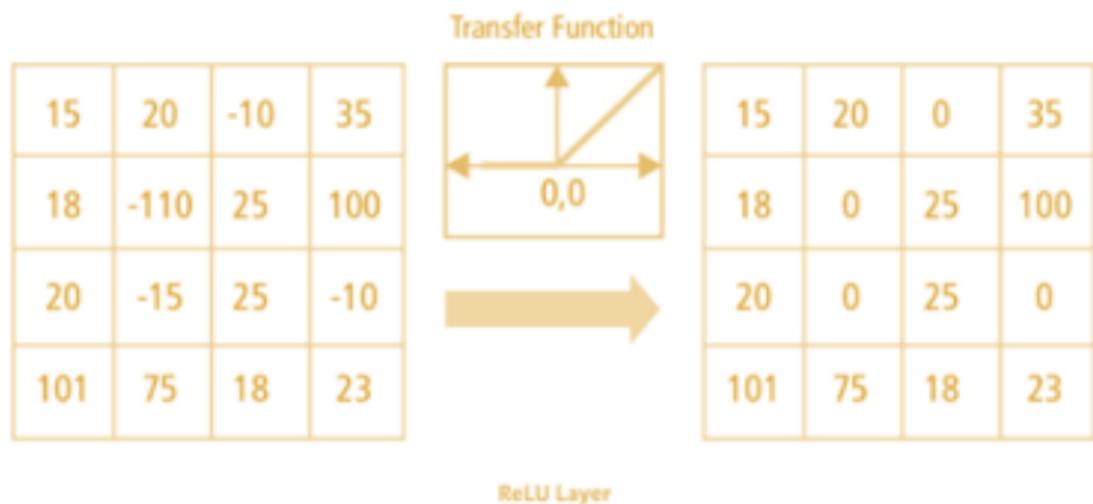
- No learnable parameters
- Introduces spatial invariance

Dilated Convolution

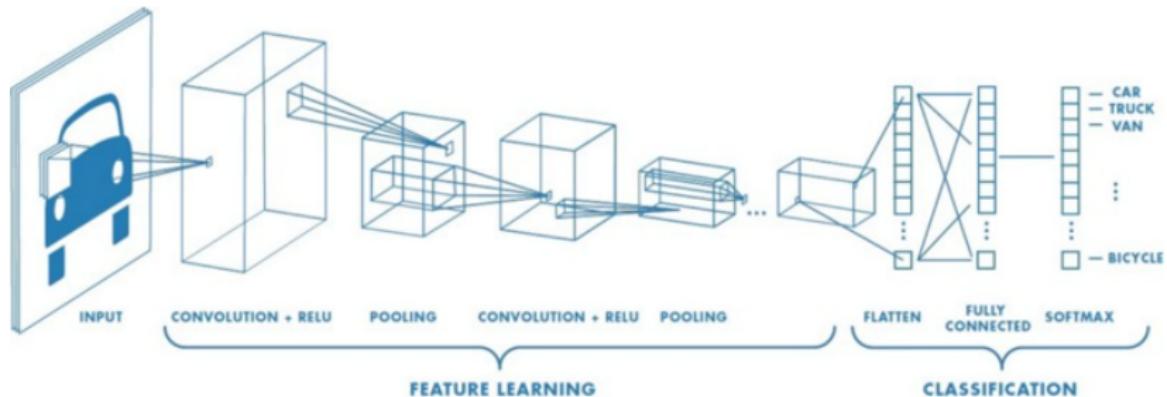
- ▶ Kernel is spread out, step > 1 between kernel elements



- ▶ Just like Fully-Connected Neural Networks, we can apply an activation over convolutional layer outputs
- ▶ It helps break linearity
- ▶ For example, Rectified Linear Unit (ReLU): $\sigma(x) = \max(0, x)$

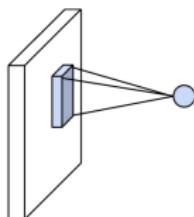


Convolutional Neural Networks

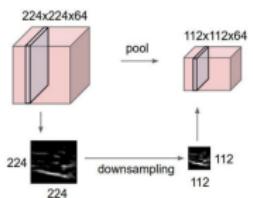


Components of a CNN

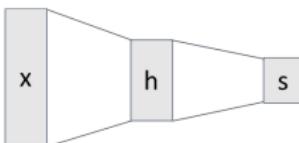
Convolution Layers



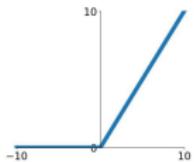
Pooling Layers



Fully-Connected Layers



Activation Function



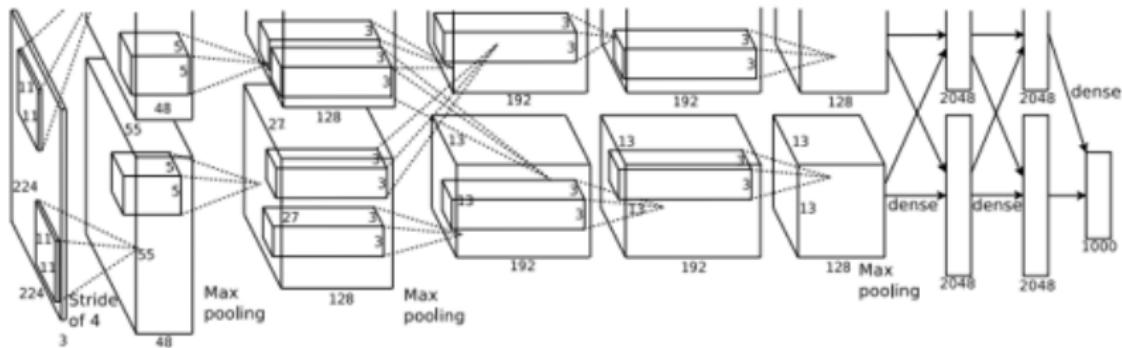
Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

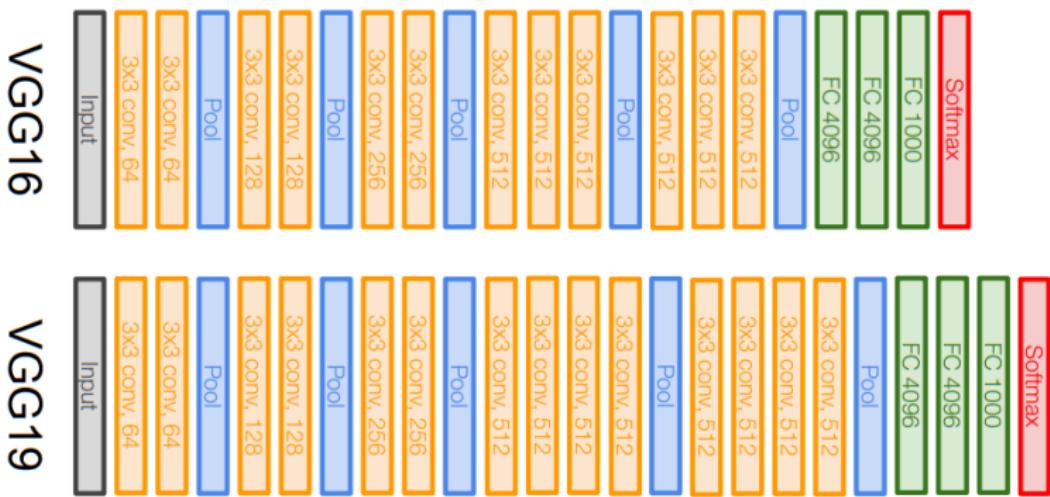
Most Notable CNNs

- ▶ AlexNet [*Krizhevsky et al. 2012*]
- ▶ VGGNet [*Simonyan and Zisserman, 2014*]
- ▶ InceptionNet (GoogLeNet) [*Szegedy et al., 2014*]
- ▶ ResNet [*He et al., 2015*]

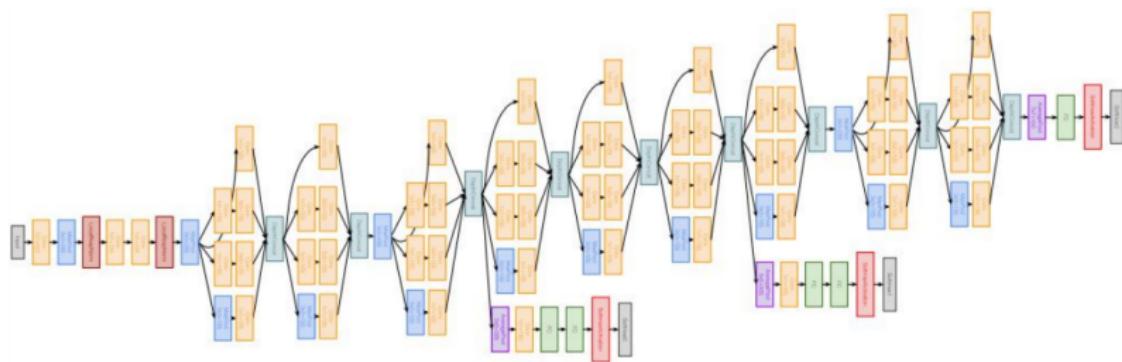
- ▶ First big improvement in image classification
- ▶ Made use of CNN, pooling, dropout, ReLU and training on GPUs.
- ▶ 5 convolutional layers, followed by max-pooling layers; with three fully connected layers at the end



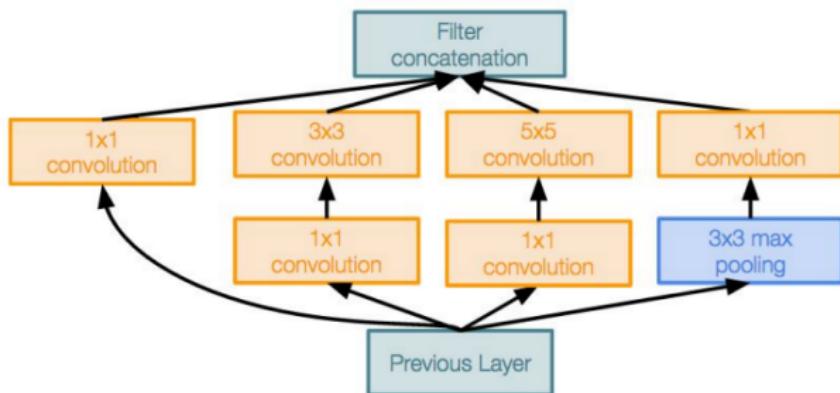
- ▶ Stack of three 3×3 conv (stride 1) layers has same effective receptive field as one 7×7 conv layer
- ▶ But deeper, more non-linearities and lesser parameters
- ▶ 13 or 16 conv layers with 3 fully-connected layers. Most params in the fully connected layer



- ▶ Going Deep: 22 layers
- ▶ Only 5 million parameters! (12x less than AlexNet and 27x less than VGGNet)
- ▶ Introduced efficient "Inception module"
- ▶ Introduced "bottleneck" layers that use 1x1 convolutions to reduce feature channel size and computational complexity

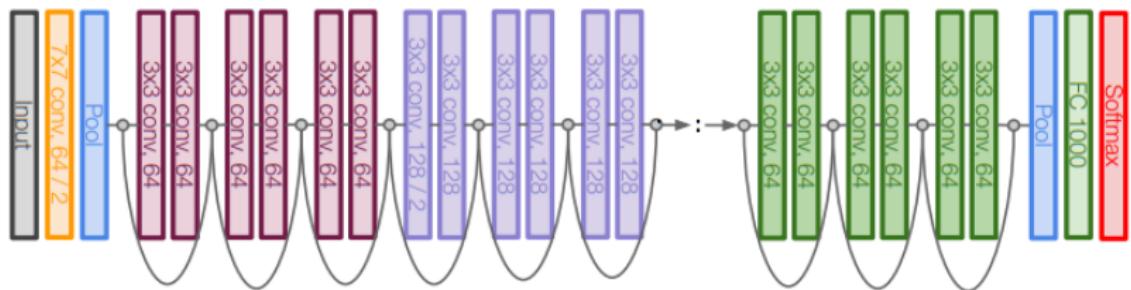


- ▶ **Inception module:** design a good local network topology (network within a network) and then stack these modules on top of each other



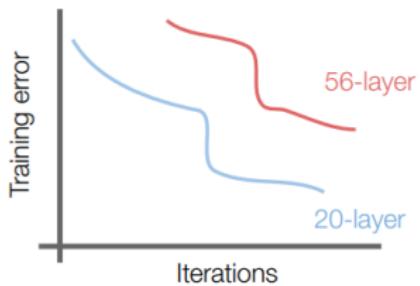
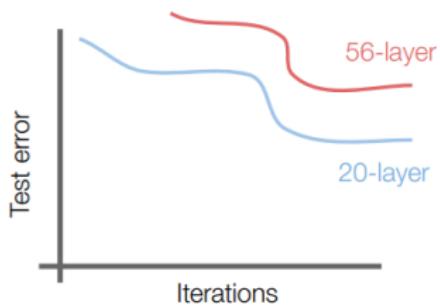
Inception module

- ▶ Very deep networks using residual connections
- ▶ 152-layer model for ImageNet
- ▶ Stacked Residual Blocks

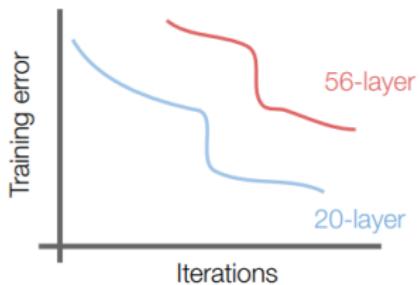
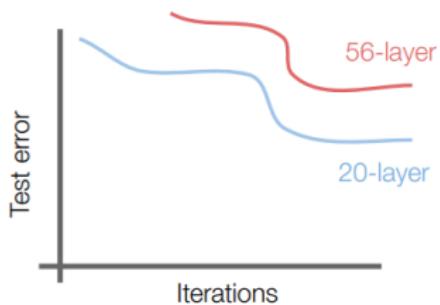


- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

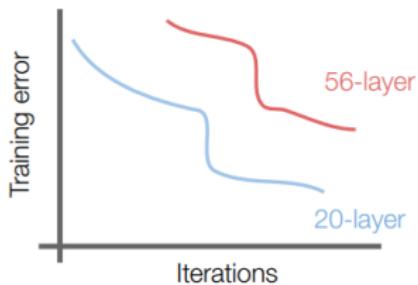
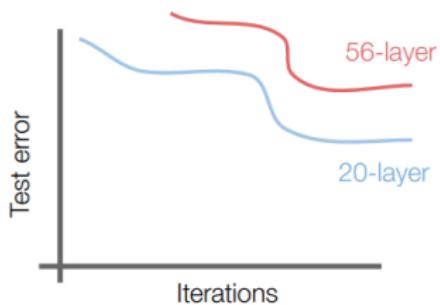


- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



- ▶ 56-layer model performs worse on both test and training error

- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



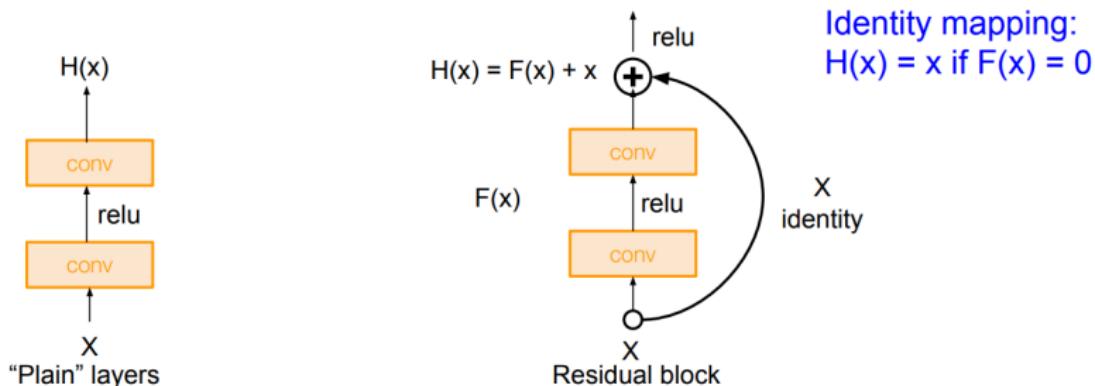
- ▶ 56-layer model performs worse on both test and training error
- ▶ The deeper model performs worse, but it's not caused by overfitting!

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize
- ▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

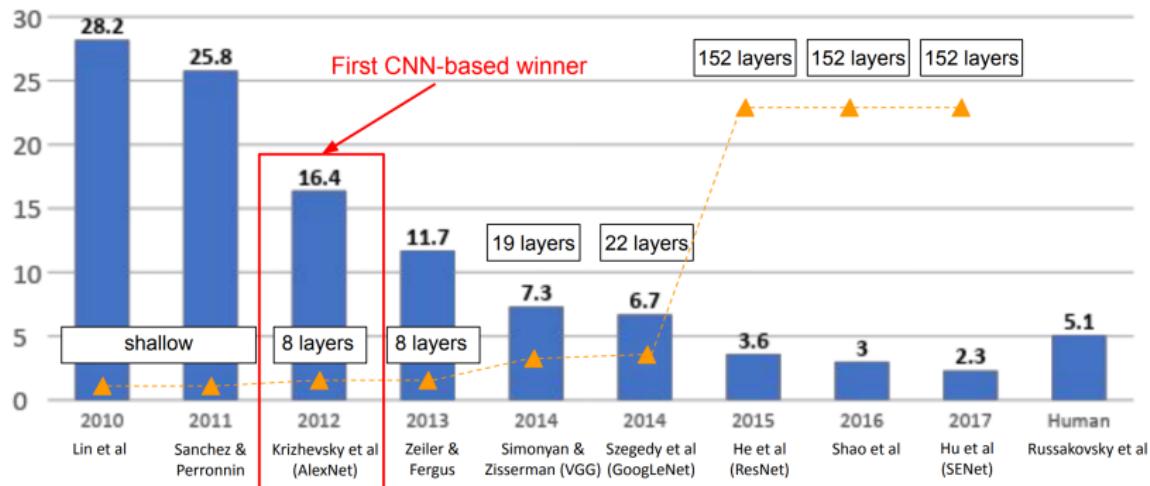
- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize
- ▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



- ▶ The most extensive data for Image Classification
- ▶ 3 RGB channels from 0 to 255
- ▶ 14,197,122 images
- ▶ 1000 classes



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: Deep Learning for Computer Vision
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV Deep Learning for Computer Vision: Fundamentals and Applications
- ▶ Justin Johnson, UMich EECS 498.008/598.008: Deep Learning for Computer Vision
- ▶ Sander Dieleman, Deepmind: Deep Learning Lecture Series 2020