

# Advanced Computer Vision

Naeemullah Khan

[naeemullah.khan@kaust.edu.sa](mailto:naeemullah.khan@kaust.edu.sa)



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

KAUST Academy  
King Abdullah University of Science and Technology

June 15, 2023

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

# Computer Vision is Everywhere



Left to right:  
[Image by Google](#), [in domain](#), [in domain](#)  
[in domain](#), [in domain](#)  
[in domain](#), [in domain](#), [in domain](#)  
[in domain](#), [in domain](#), [in domain](#)



Left to right:  
[Image by Google](#)  
[in domain](#), [in domain](#), [in domain](#)  
[in domain](#), [in domain](#), [in domain](#)  
[in domain](#), [in domain](#), [in domain](#)

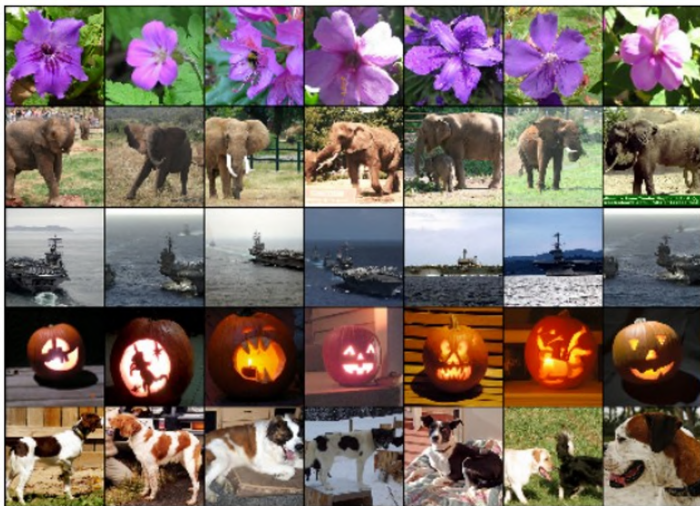


Bottom row, left to right:  
[Image by Google](#), [in domain](#)  
[Image by Google](#), [in domain](#)  
[Image by Google](#), [in domain](#)  
[Image by Google](#), [in domain](#)  
[Image by Google](#), [in domain](#)

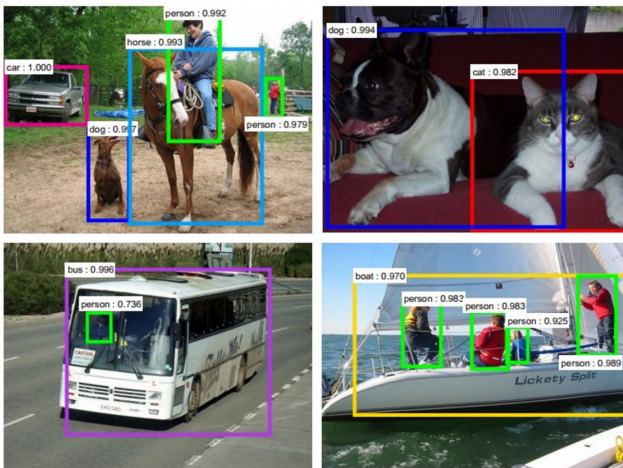
## Image Classification



## Image Retrieval

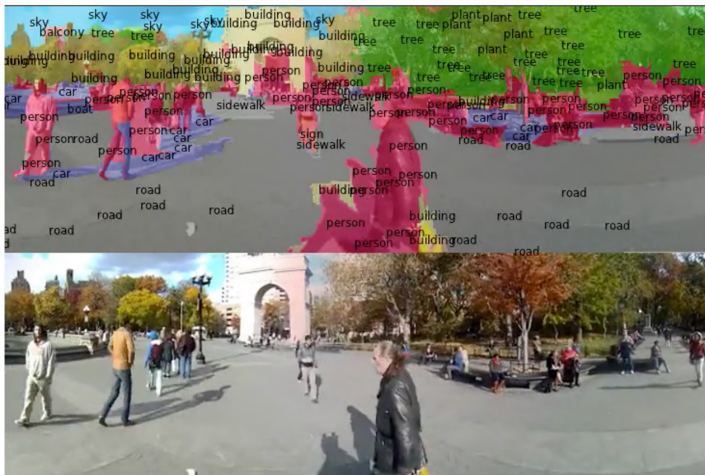


## Object Detection



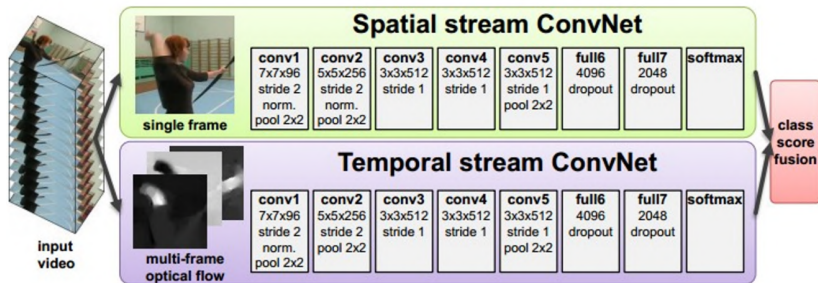
Ren, He, Girshick, and Sun, 2015

## Image Segmentation



Fabaret et al, 2012

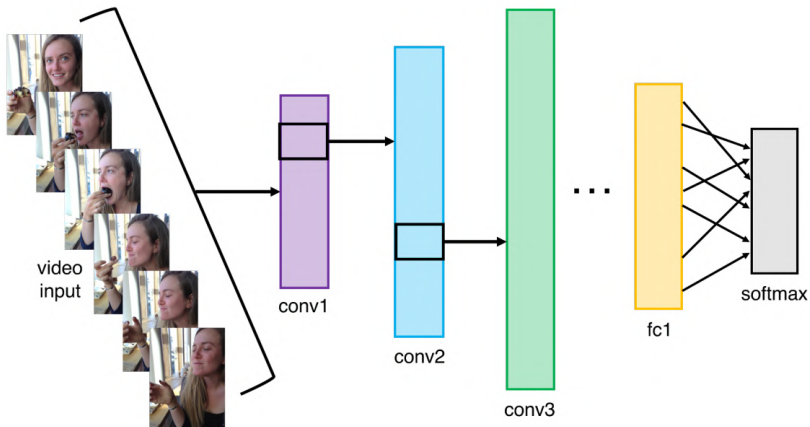
## Video Classification



Simonyan et al, 2014



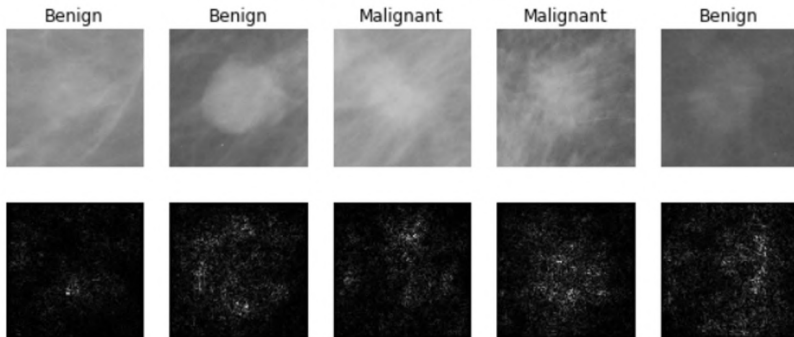
## Activity Recognition



Pose Recognition (Toshev and Szegedy, 2014)



## Medical Imaging



## Image Captioning

Vinyals et al, 2015

Karpathy and Fei-Fei, 2015

© All images are CC0 public domain.  
[https://commons.wikimedia.org/wiki/File:White\\_teddy\\_bear\\_sitting\\_in\\_the\\_grass](https://commons.wikimedia.org/wiki/File:White_teddy_bear_sitting_in_the_grass)  
[https://commons.wikimedia.org/wiki/File:A\\_baseball\\_player\\_in\\_a\\_blue\\_uniform\\_throwing\\_a\\_ball](https://commons.wikimedia.org/wiki/File:A_baseball_player_in_a_blue_uniform_throwing_a_ball)  
[https://commons.wikimedia.org/wiki/File:A\\_woman\\_holding\\_a\\_cat\\_in\\_her\\_hand](https://commons.wikimedia.org/wiki/File:A_woman_holding_a_cat_in_her_hand)  
[https://commons.wikimedia.org/wiki/File:A\\_man\\_riding\\_a\\_wave\\_on\\_top\\_of\\_a\\_surfboard](https://commons.wikimedia.org/wiki/File:A_man_riding_a_wave_on_top_of_a_surfboard)  
[https://commons.wikimedia.org/wiki/File:A\\_cat\\_sitting\\_on\\_a\\_suitcase\\_on\\_the\\_floor](https://commons.wikimedia.org/wiki/File:A_cat_sitting_on_a_suitcase_on_the_floor)  
[https://commons.wikimedia.org/wiki/File:A\\_woman\\_standing\\_on\\_a\\_beach\\_holding\\_a\\_surfboard](https://commons.wikimedia.org/wiki/File:A_woman_standing_on_a_beach_holding_a_surfboard)



*A white teddy bear  
sitting in the grass*



*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*

## Image Generation



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

DALL-E 2

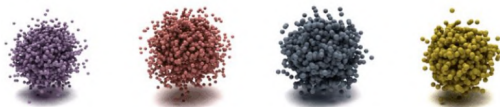


Style Transfer

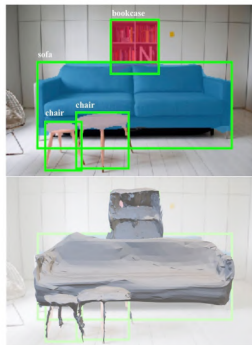
## 3D Vision



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



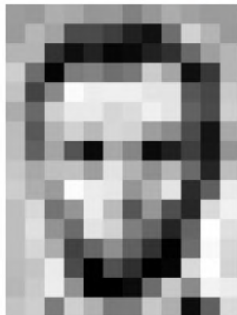
Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)



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	164
180	180	50	14	54	6	10	33	48	106	159	181
206	109	5	134	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	96	74	206
188	88	179	209	185	215	211	158	139	75	30	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	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	168	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	164
180	180	50	14	54	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	96	74	206
188	88	179	209	185	215	211	158	139	75	30	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	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	168	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218



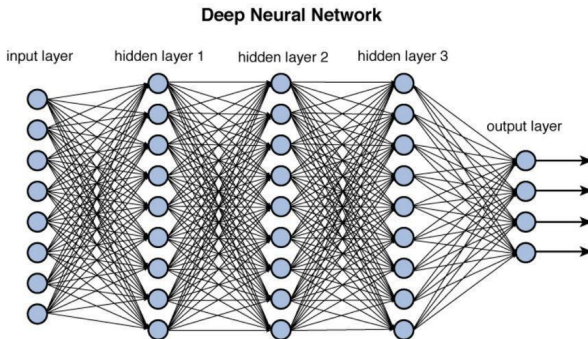
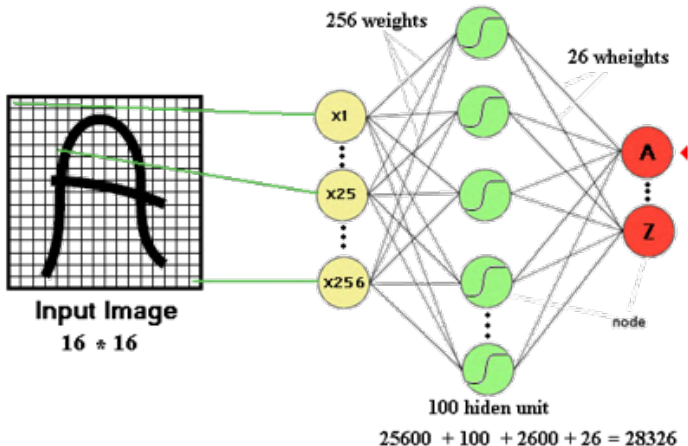


Figure 12.2 Deep network architecture with multiple layers.

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

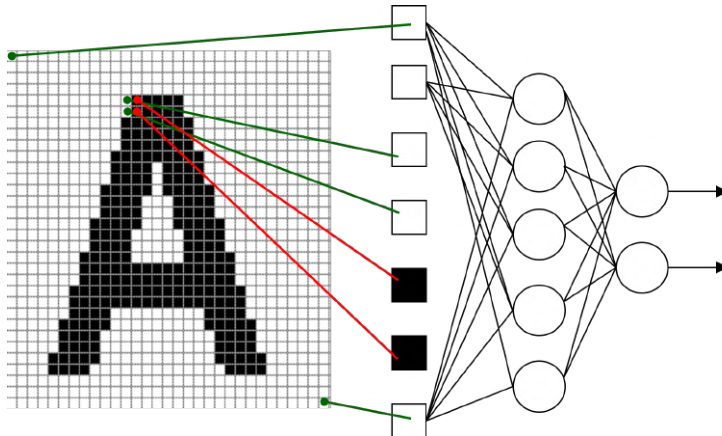
# Drawbacks of Fully-Connected Neural Networks

- The number of trainable parameters becomes extremely large



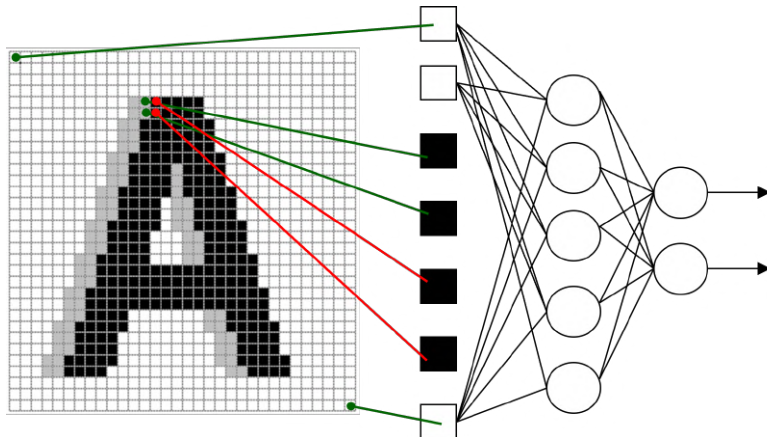
# Drawbacks of Fully-Connected Neural Networks (cont.)

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



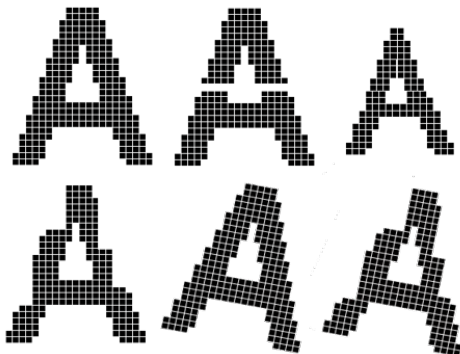
# Drawbacks of Fully-Connected Neural Networks (cont.)

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

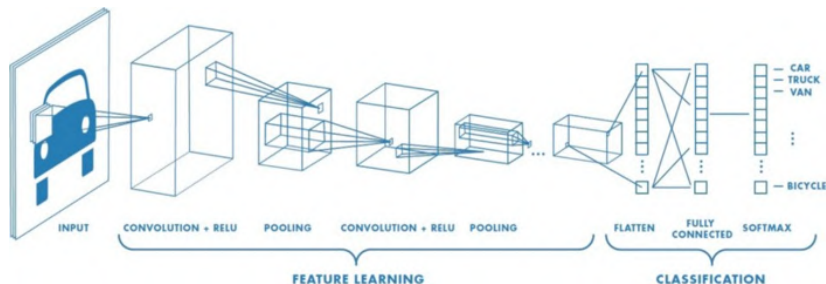


# Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ The topology of the input data is completely ignored
- ▶ For a  $32 \times 32$  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?

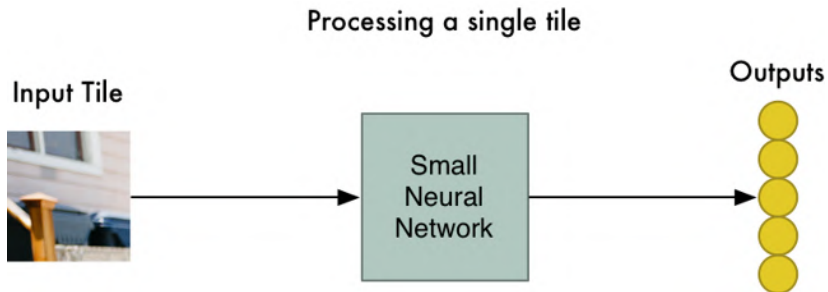


# How Convolution Works? (cont.)

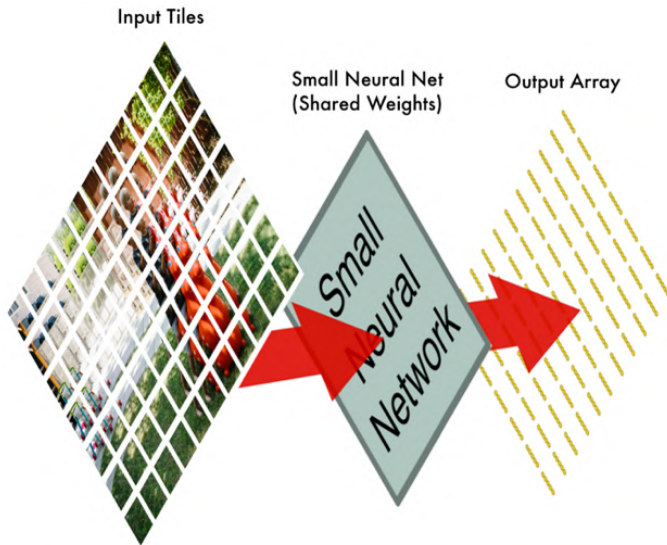




# How Convolution Works? (cont.)

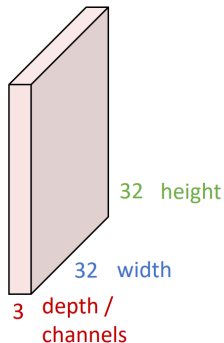


# How Convolution Works? (cont.)



# How Convolution Works? (cont.)

3x32x32 image

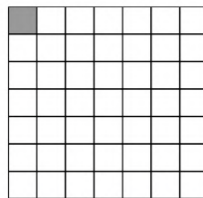
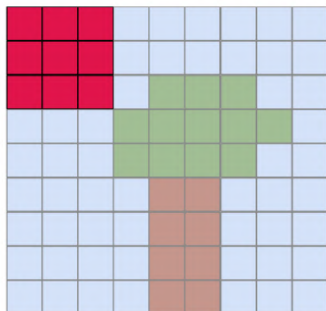


3x5x5 filter



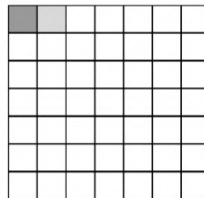
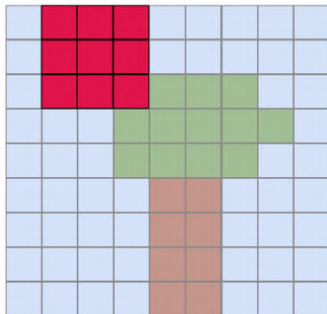
**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# How Convolution Works? (cont.)



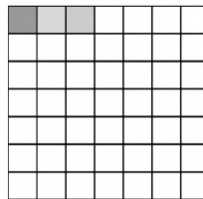
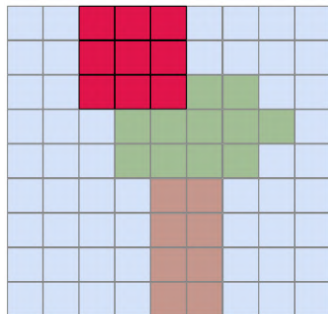
The **kernel** slides across the image and produces an output value at each position

# How Convolution Works? (cont.)



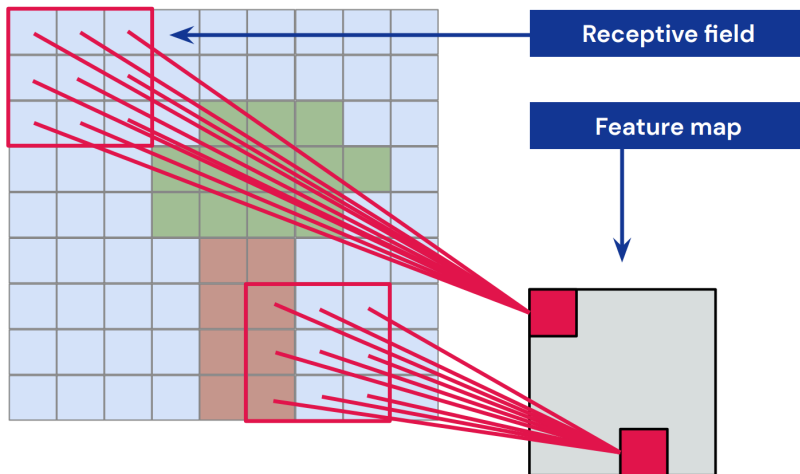
The **kernel** slides across the image and produces an output value at each position

# How Convolution Works? (cont.)

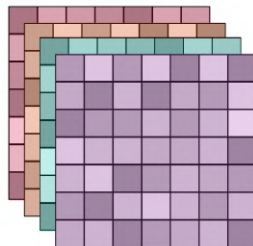
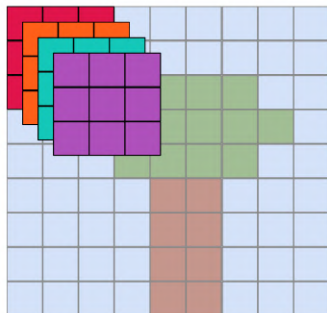


The **kernel** slides across the image and produces an output value at each position

# How Convolution Works? (cont.)



# How Convolution Works? (cont.)

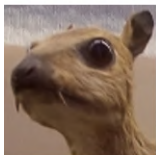


We convolve multiple kernels and obtain multiple feature maps or **channels**

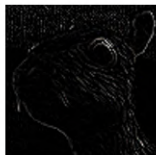


# How Convolution Works? (cont.)

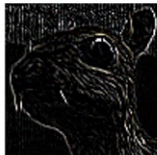
$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$$



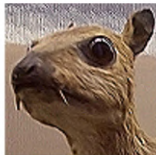
$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



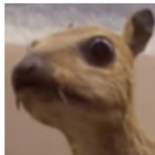
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



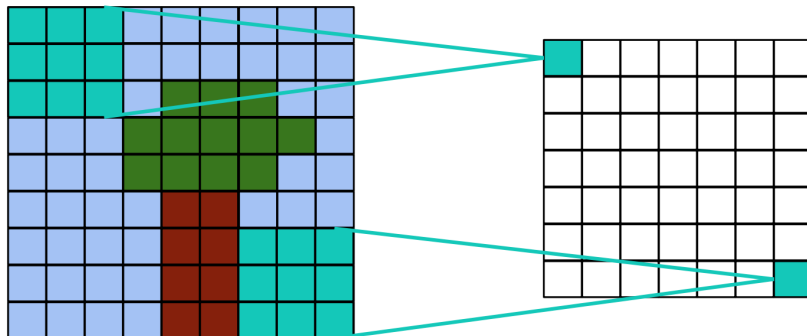
$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$



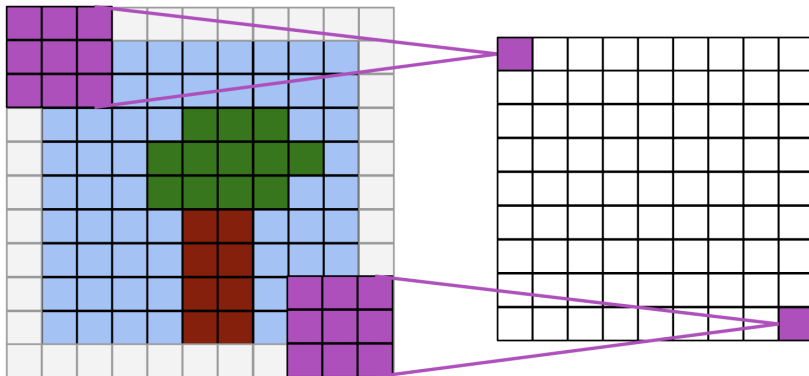
$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



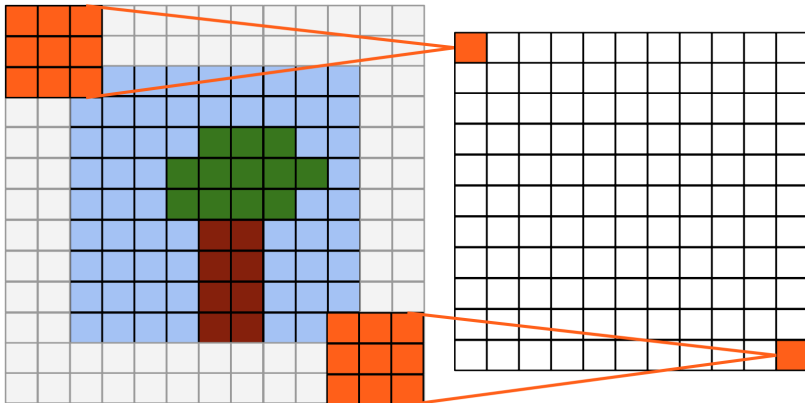
- ▶ Applying Convolution as such reduces the size of the borders.
- ▶ Sometimes this is not desirable.
- ▶ We can pad the border with zeros.



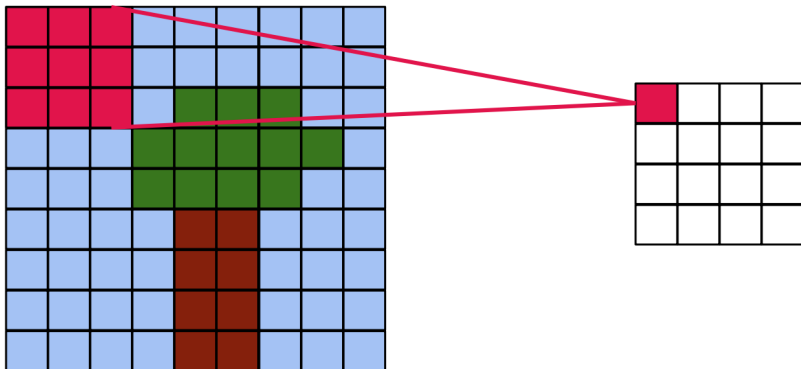
- Same Convolution: Output is the same size as input



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

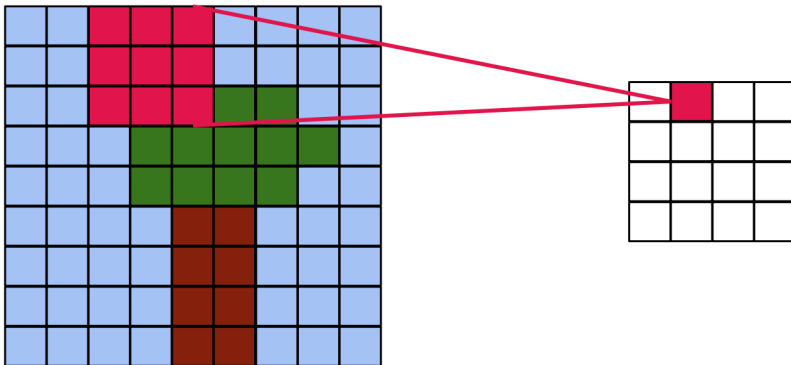


- Kernel slides along the image with a step  $> 1$

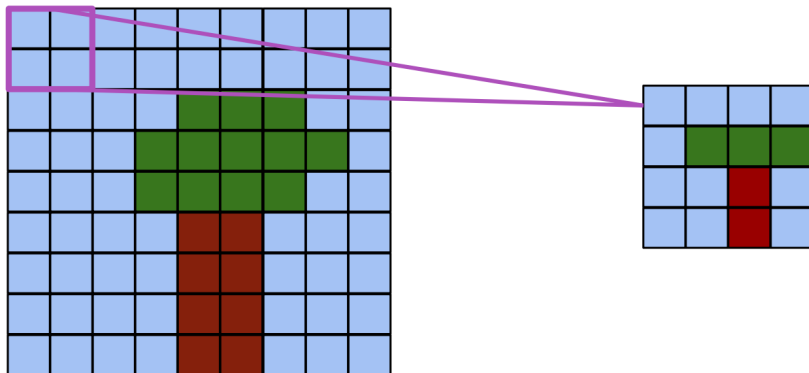


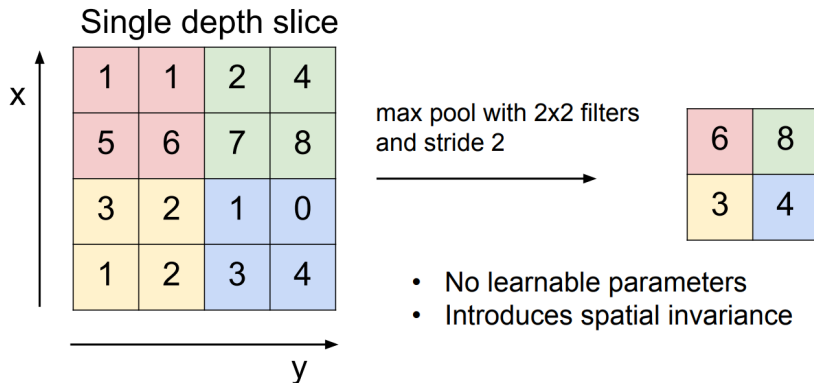
# Strided Convolution (cont.)

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



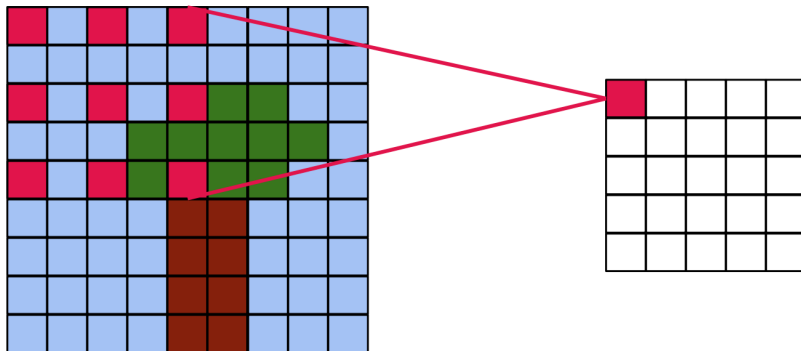
- Compute mean or max over small windows to reduce resolution



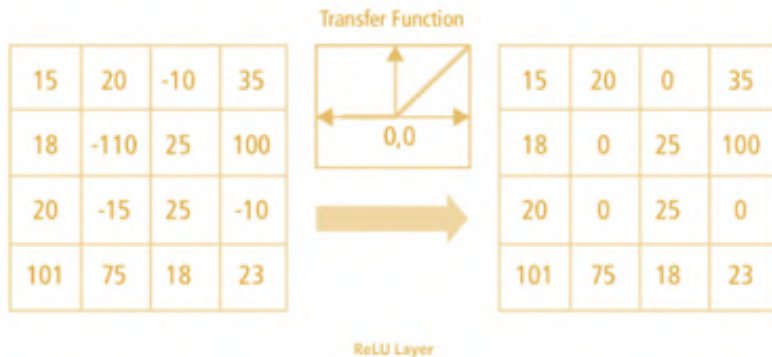




- ▶ Kernel is spread out, step  $> 1$  between kernel elements



- ▶ 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):  $\sigma(x) = \max(0, x)$



- ▶ Consider a single layer  $y = Wx$
- ▶ The following could lead to tough optimization
  - Inputs  $x$  are not centered around zero (need large bias)
  - Inputs  $x$  have different scaling per element (entries in  $W$  will need to vary a lot)

- ▶ Consider a single layer  $y = Wx$
- ▶ The following could lead to tough optimization
  - Inputs  $x$  are not centered around zero (need large bias)
  - Inputs  $x$  have different scaling per element (entries in  $W$  will need to vary a lot)
- ▶ **Idea:** Force inputs to be "nicely scaled" at each layer!

- Consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

- Consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

- **Problem:** What if zero-mean, unit variance is too hard of a constraint?

**Input:**  $x : N \times D$

**Learnable scale and shift parameters:**

$$\gamma, \beta : D$$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$ , will recover the  
identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

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

Normalized x,  
Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is N x D

Estimates depend on minibatch;  
can't do this at test-time!

**Input:**  $x : N \times D$

**Learnable scale and  
shift parameters:**

$$\gamma, \beta : D$$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$ , will recover the  
identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel var, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} \quad \text{Normalized x, Shape is } N \times D$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is } N \times D$$



**Input:**  $x : N \times D$

$\mu_j =$  (Running) average of values seen during training

Per-channel mean, shape is D

**Learnable scale and shift parameters:**

$\sigma_j^2 =$  (Running) average of values seen during training

Per-channel var, shape is D

$\gamma, \beta : D$

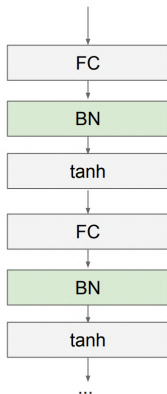
During testing batchnorm becomes a linear operator!  
Can be fused with the previous fully-connected or conv layer

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

Normalized x, Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output, Shape is N x D



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Batch Normalization for  
**fully-connected** networks

$$\mathbf{x}: \mathbf{N} \times \mathbf{D}$$

Normalize



$$\boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} \times \mathbf{D}$$

$$\boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{D}$$

$$\mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$$

Batch Normalization for  
**convolutional** networks  
(Spatial Batchnorm, BatchNorm2D)

$$\mathbf{x}: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W}$$

Normalize



$$\boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$$

$$\boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$$

$$\mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$$

## ► Advantages:

- Makes deep networks much easier to train!
- Improves gradient flow
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!

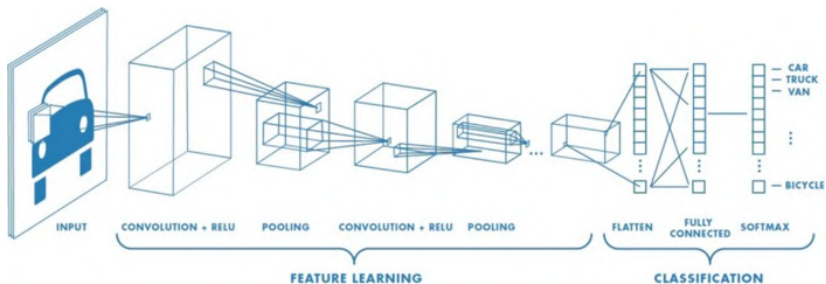
## ► Advantages:

- Makes deep networks much easier to train!
- Improves gradient flow
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!

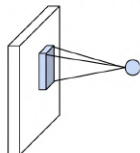
## ► Disadvantages:

- Behaves differently during training and testing: this is a very common source of bugs!

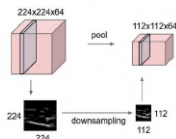
# Convolutional Neural Networks



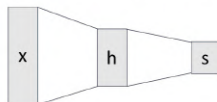
## Convolution Layers



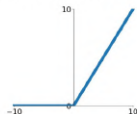
## Pooling Layers



## Fully-Connected Layers



## Activation Function



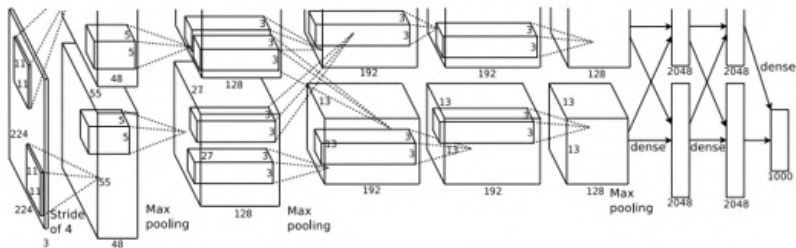
## Normalization

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

- ▶ 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



- The diagram illustrates the VGG16 and VGG19 architectures, showing the sequence of layers and their dimensions. The layers are color-coded: red for Softmax, green for FC layers, blue for Pool layers, orange for 3x3 conv layers, and grey for Input layers.

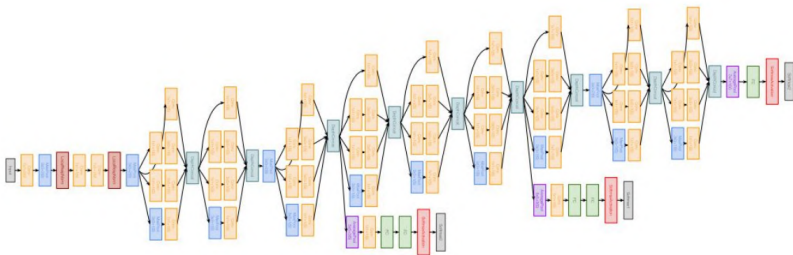
**VGG16 Architecture (16 layers):**

  - Input (grey)
  - 3x3 conv, 64 (orange)
  - 3x3 conv, 64 (orange)
  - Pool (blue)
  - 3x3 conv, 128 (orange)
  - 3x3 conv, 128 (orange)
  - Pool (blue)
  - 3x3 conv, 256 (orange)
  - 3x3 conv, 256 (orange)
  - Pool (blue)
  - 3x3 conv, 512 (orange)
  - 3x3 conv, 512 (orange)
  - Pool (blue)
  - 3x3 conv, 512 (orange)
  - 3x3 conv, 512 (orange)
  - FC 4096 (green)
  - FC 1000 (green)
  - Softmax (red)

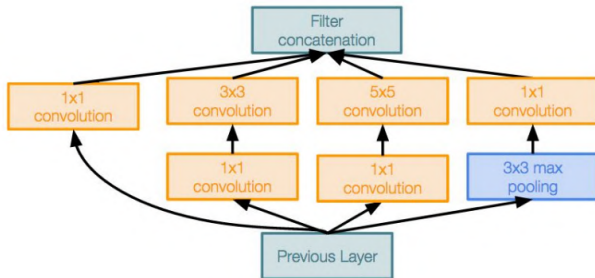
**VGG19 Architecture (19 layers):**

  - Input (grey)
  - 3x3 conv, 64 (orange)
  - 3x3 conv, 64 (orange)
  - Pool (blue)
  - 3x3 conv, 128 (orange)
  - 3x3 conv, 128 (orange)
  - Pool (blue)
  - 3x3 conv, 256 (orange)
  - 3x3 conv, 256 (orange)
  - Pool (blue)
  - 3x3 conv, 512 (orange)
  - 3x3 conv, 512 (orange)
  - Pool (blue)
  - 3x3 conv, 512 (orange)
  - 3x3 conv, 512 (orange)
  - FC 4096 (green)
  - FC 1000 (green)
  - Softmax (red)

- ▶ 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  $1 \times 1$  convolutions to reduce feature channel size and computational complexity

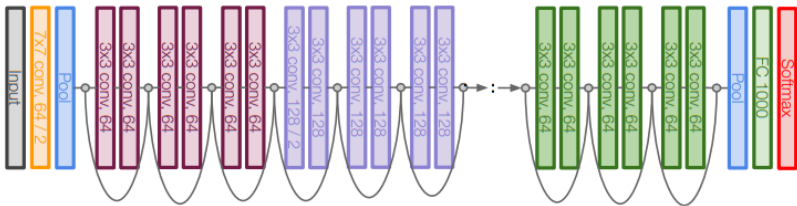


- **Inception module:** design a good local network topology (network within a network) and then stack these modules on top of each other



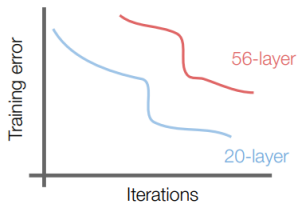
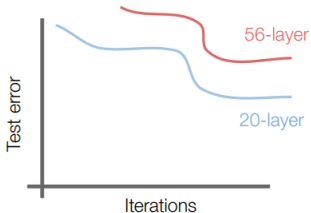
Inception module

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

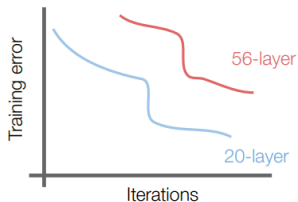
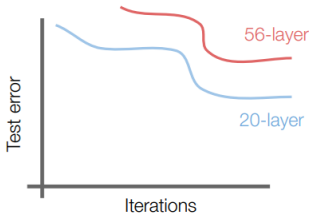


- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

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



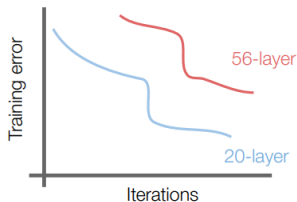
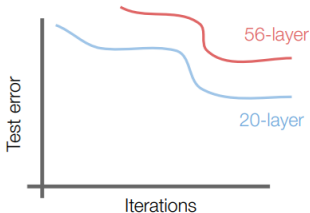
- What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



- 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?



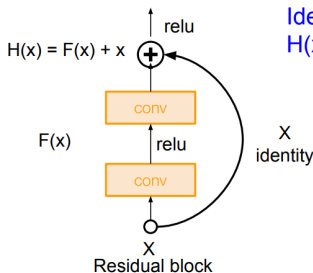
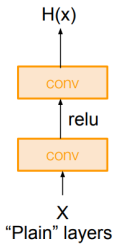
- ▶ 56-layer model performs worse on both test and training error
- ▶ The deeper model performs worse, but it's not caused by overfitting!

- **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

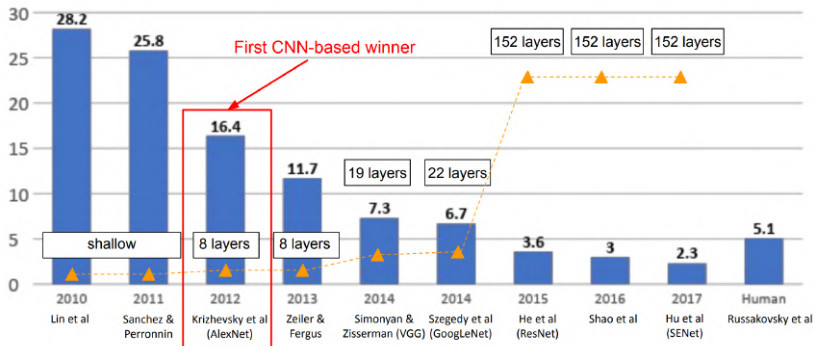


Identity mapping:  
 $H(x) = x$  if  $F(x) = 0$

- ▶ 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](#)