# Computer Vision

## Naeemullah Khan
naeemullah.khan@kaust.edu.sa



جامعة الملك عبدالله
للعلوم والتقنية
King Abdullah University of
Science and Technology

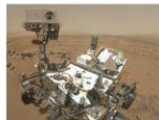KAUST Academy
King Abdullah University of Science and Technology

November 19, 2023

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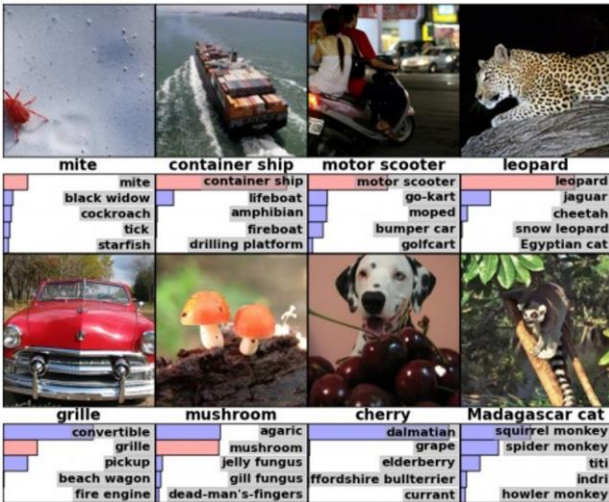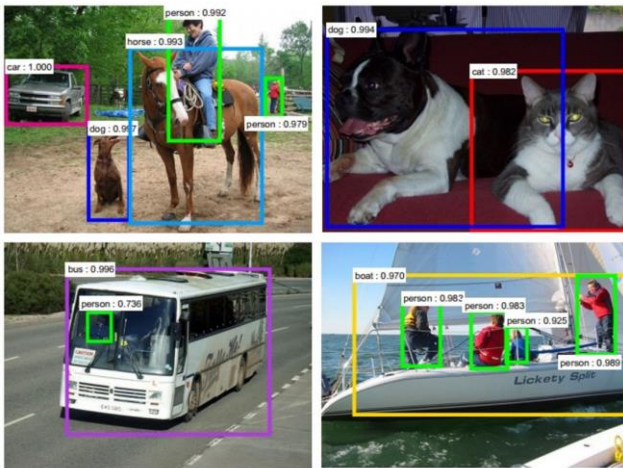
# Image Classification

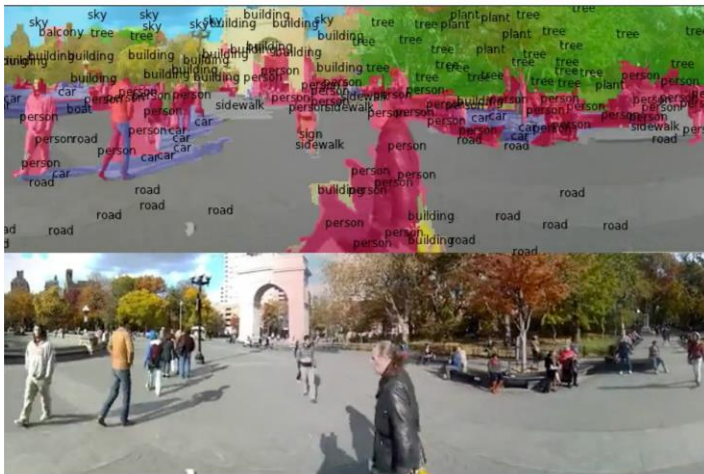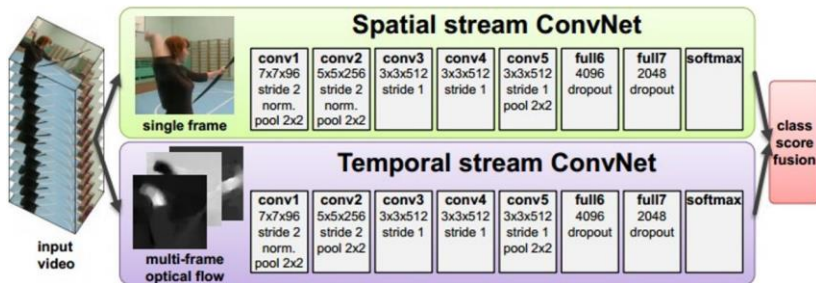Image Retrieval

## Object Detection



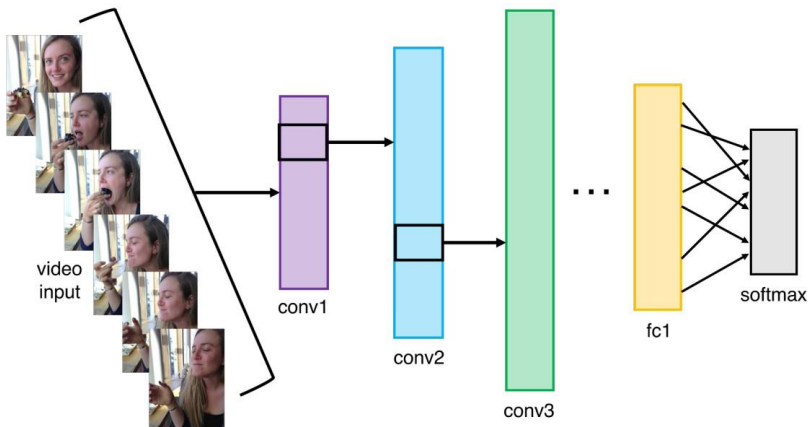Ren, He, Girshick, and Sun, 2015

## 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

*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)

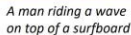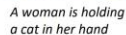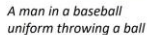Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)

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

- Images are represented as Matrices with elements in $[0, 255]$

- Grayscale images have one channel while RGB images have 3 channels

**Deep Neural Network**

Figure 12.2 Deep network architecture with multiple layers.

$$z = W_1 x_1 + W_2 x_2 + \cdots + W_n x_n + b$$

► The number of trainable parameters becomes extremely large



256 weights

26 wheights

Input Image
16 * 16

x1

x25

x256

A

Z

node

100 hiden unit

25600 + 100 + 2600 + 26 = 28326

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

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

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

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

Processing a single tile

Input Tile



Small
Neural
Network

Outputs

3x32x32 image

3x5x5 filter

32 height

32 width

3 depth / channels

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

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

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

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

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

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



$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \qquad \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} \qquad \frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

| 10 | 10 | 10 | 0 | 0 | 0 |
|----|----|----|---|---|---|
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |

$*$

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

| 0 | 30 | 30 | 0 |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

$*$

vertical edges

horizontal edges

- 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: output size = input size + kernel size - 1

► Kernel slides along the image with a step > 1

► Kernel slides along the image with a step > 1

# Pooling

► Compute mean or max over small windows to reduce resolution

4 x 4

6 x 6 x 3

*

3 x 3 x 3

=

4 x 4

*

3 x 3 x 3

=

4 x 4

6 x 6 x 3

*

3 x 3 x 3

*

3 x 3 x 3

## Single depth slice



max pool with 2x2 filters
and stride 2

- No learnable parameters
- Introduces spatial invariance

# Dilated Convolution

► Kernel is spread out, step > 1 between kernel elements

# Activation

► 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): $\sigma(x) = max(0, x)$



Transfer Function

| 15 | 20 | -10 | 35 |
|----|----|-----|-----|
| 18 | -110 | 25 | 100 |
| 20 | -15 | 25 | -10 |
| 101 | 75 | 18 | 23 |

0,0

| 15 | 20 | 0 | 35 |
|----|----|---|-----|
| 18 | 0 | 25 | 100 |
| 20 | 0 | 25 | 0 |
| 101 | 75 | 18 | 23 |

ReLU Layer

| 10 | 10 | 10 | 0 | 0 | 0 |
|----|----|----|---|---|---|
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |

*

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

| 0 | 30 | 30 | 0 |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

**Sparsity of connections:** In each layer, each output value depends only on a small number of inputs.

## Convolution Layers

## Pooling Layers

224x224x64 pool 112x112x64

224 downsampling 112
224 112

## Fully-Connected Layers

x h s

## Activation Function

## Normalization

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

If you have 10 filters that are 3 x 3 x 3 in one layer of a neural network, how many parameters does that layer have?

$$\text{Floor}\left(\frac{(W-F+2P)}{S} + 1\right)$$

# Most Notable CNNs

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

32×32 ×1    5 × 5, s = 1    28×28×6    avg pool, f = 2, s = 2    14×14×6    5 × 5, s = 1    10×10×16    avg pool, f = 2, s = 2    5×5×16    120    84    $\hat{y}$

[LeCun et al., 1998. Gradient-based learning applied to document recognition]

- 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

[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

- Stack of three 3x3 conv (stride 1) layers has same effective receptive field as one 7x7 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

# InceptionNet

- ► 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

The problem of computational cost



$28 \times 28 \times 192$  →  CONV $5 \times 5$, same, $32$  →  $28 \times 28 \times 32$

Using 1x1 convolution



$28 \times 28 \times 192$

CONV
$1 \times 1$,
16,
$1 \times 1 \times 192$

$28 \times 28 \times 16$

CONV
$5 \times 5$,
32,
$5 \times 5 \times 16$

$28 \times 28 \times 32$

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



Inception module

$1 \times 1$

$3 \times 3$

$5 \times 5$

MAX-POOL

$28 \times 28 \times 192$

28

28

32

32

128

64

[Szegedy et al. 2014. Going deeper with convolutions]

# ResNet

- 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

# ResNet

- **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

# ResNet

- **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 $F(x) = 0$

- Low computational cost at deployment
- Useful for mobile and embedded vision applications
- Key idea: Normal vs. depthwise-separable convolutions

$6 \times 6 \times 3$   *   $3 \times 3 \times 3$   =   $4 \times 4 \times 5$

Computational cost   =   #filter params   x   # filter positions   x   # of filters

Normal Convolution

Depthwise Separable Convolution



Depthwise    Pointwise

# Depthwise convolution
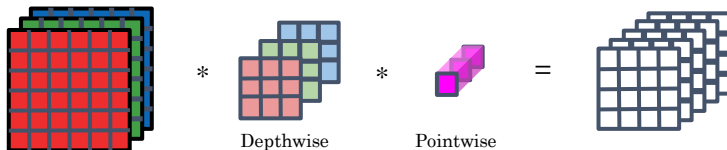


6 x 6 x 3     *     3 x 3     =     4 x 4 x 3

Computational cost     =     #filter params     x     # filter positions     x     # of filters

4 x 4 x 3        *        1 x 1 x 3        =        4 x 4 x 5

Computational cost    =    #filter params    x    # filter positions    x    # of filters

Depthwise Convolution

Pointwise Convolution

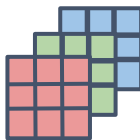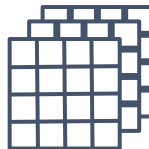# MobileNet



MobileNet v1

MobileNet v2

Residual Connection

Expansion    Depthwise    Projection

[Sandler et al. 2019, MobileNetV2: Inverted Residuals and Linear Bottlenecks]

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

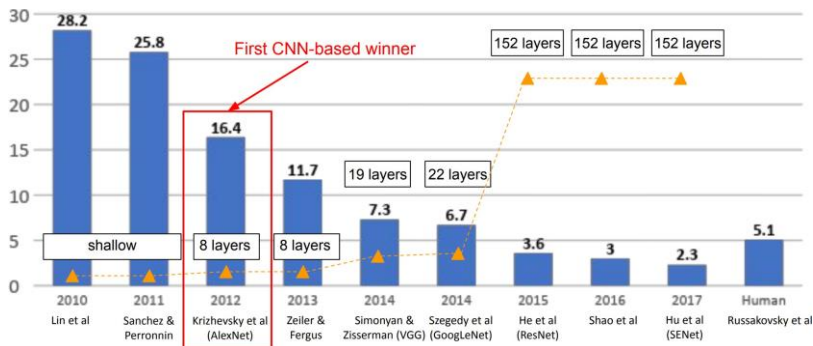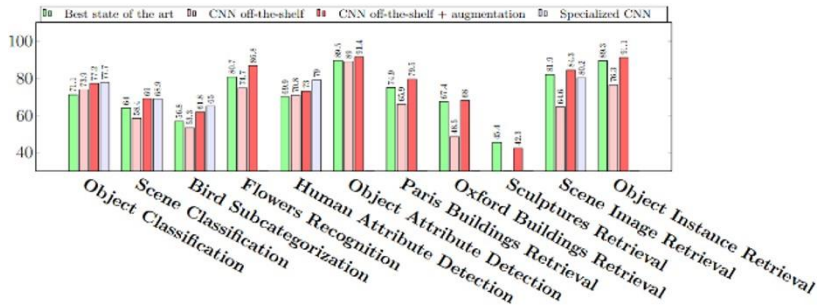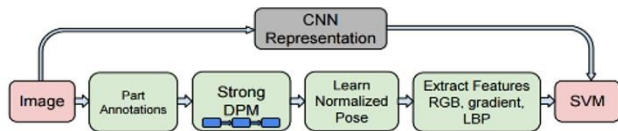- Improvement of learning in a **new** task through the transfer of knowledge from a **related** task that has already been learned.

- We will look at one strategy of transfer learning called Fine-Tuning

# When to fine-tune your model?

- ► New dataset is small with distribution similar to original dataset.
  - Keep the feature extraction part fixed and fine-tune the classifier part of the network
- ► New dataset is large with similar distribution to the original dataset
  - Fine tune both the feature extractor and the classifier part of the network
- ► New dataset is small but different distribution from the original dataset
  - Use SVM classifier on the features extracted from the feature extractor part of the Network
- ► New dataset is large and different distribution from the original dataset
  - Fine tune both the feature extractor and the classifier part of the network

0 Razavian et al. 2014

Figure 2: Learning and Transferring Mid-Level Image Representations using Convolutional Neural Networks

[0]Oquab et al. CVPR 2014

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