



Advanced Artificial Intelligence: Computer Vision



May 6, 2023

Course Website



<https://naeemullah-khan.github.io/AI23>



KAUST Academy

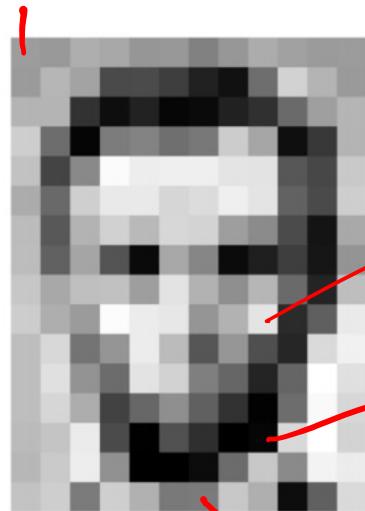
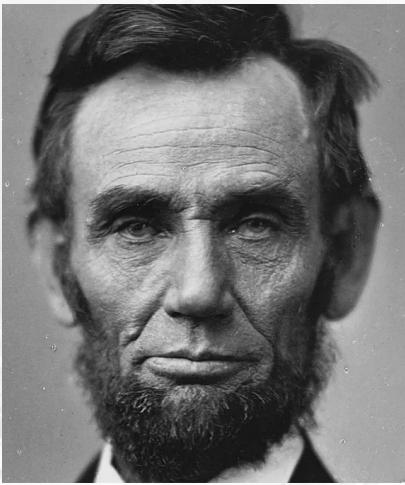


Convolutional Neural Networks

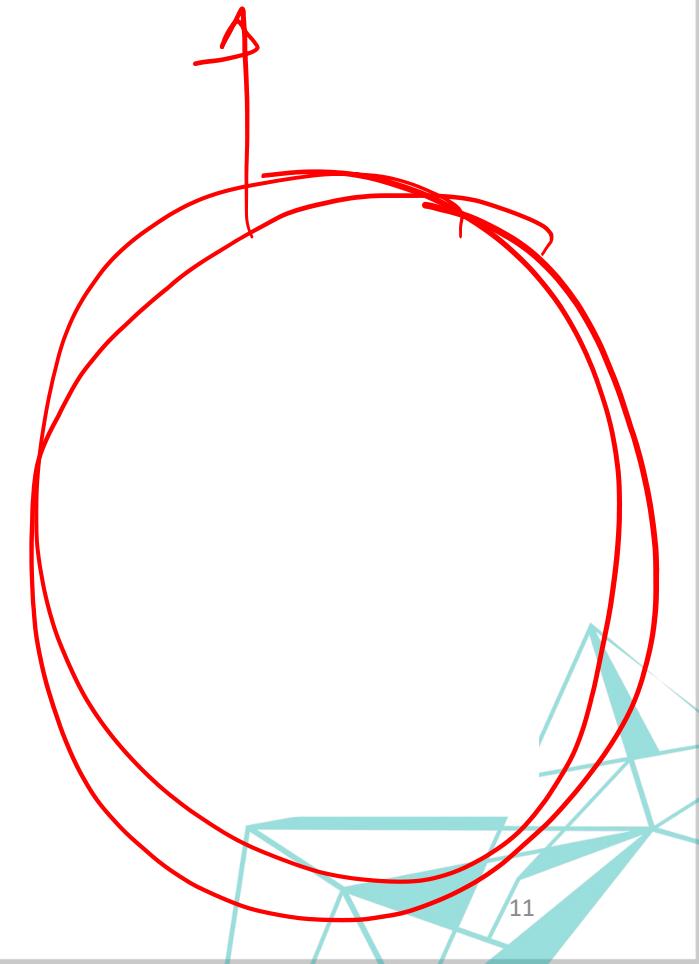
How to represent an image?

- Matrix with elements in $[0, 255]$

$(0, \dots, 255)$



Naeemullah Khan



How to represent an image?

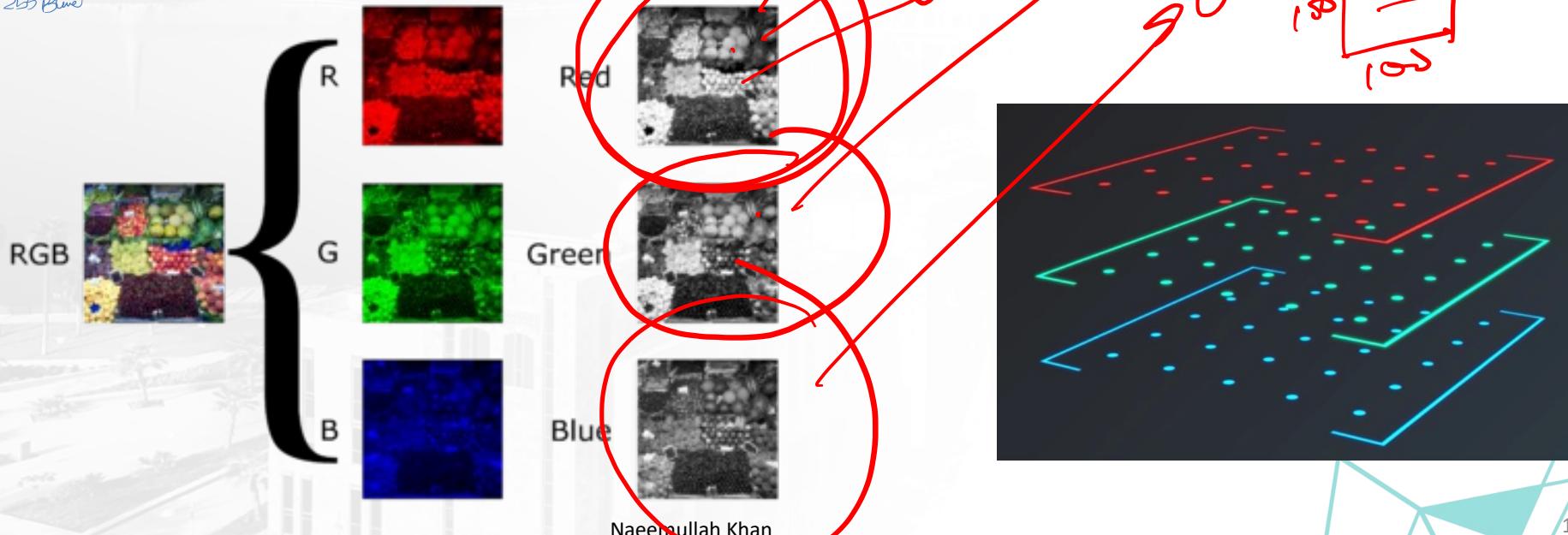
- Matrix with elements in [0,255]
- Tensor with RGB channels, each being a matrix



How to represent an image?

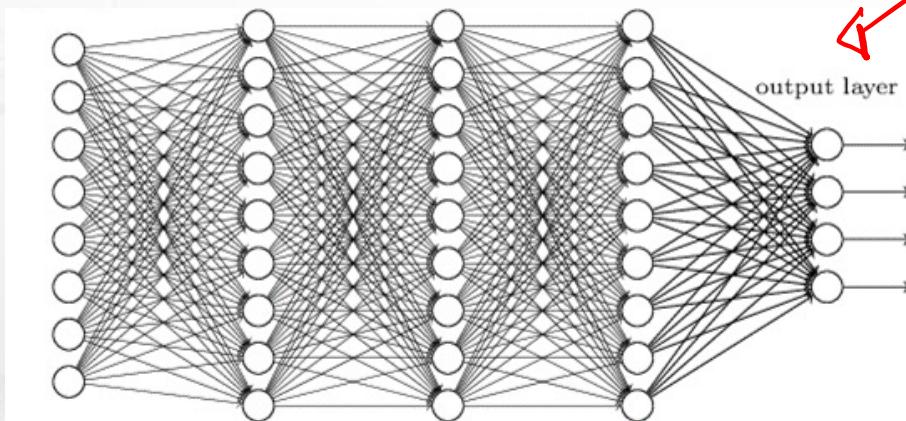
- Grayscale images are matrices
- Color images are Tensors

3 channels (Matrices)
↓
255 Red
255 Green
255 Blue



Convolutional Neural Network (CNN)

Full-connected network



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

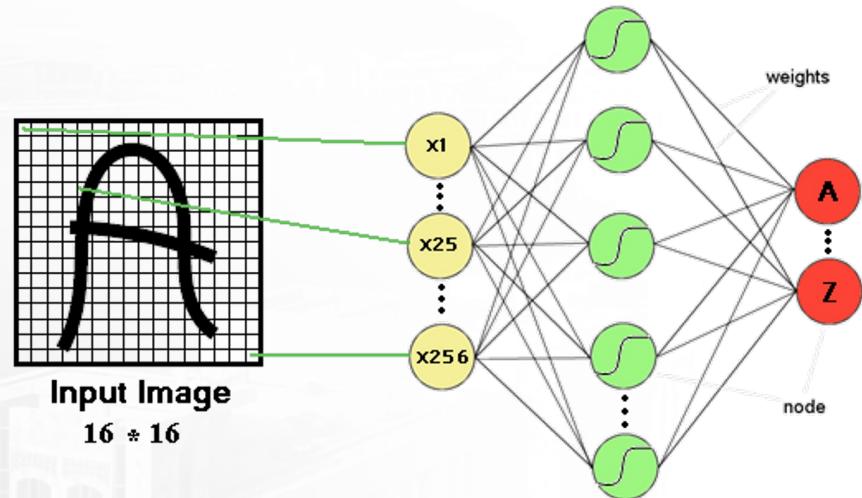
Convolutional network



$$z_{ij} = W * x_{i,j} = \sum_{a=0} \sum_{b=0} W_{ab} x_{(i+a)(j+b)}$$

Multi-layer perceptron and image processing

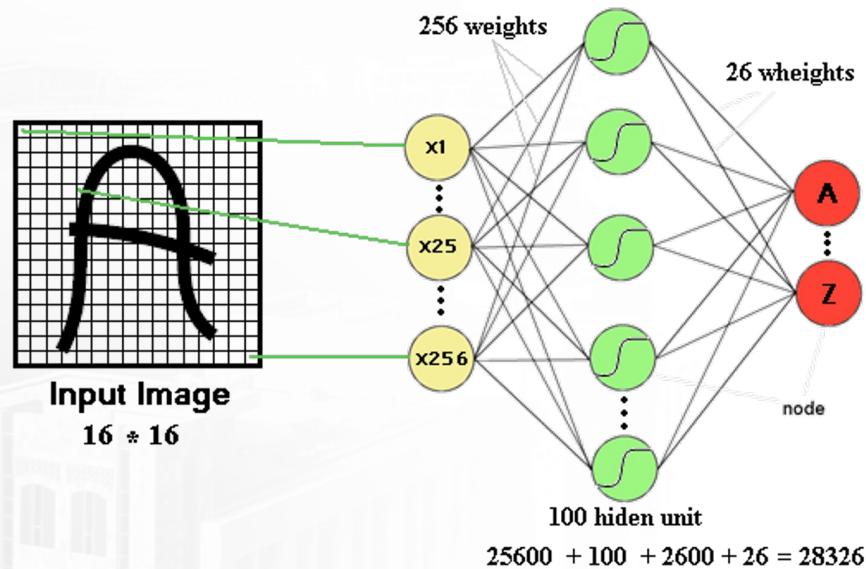
- ⦿ One or more hidden layers
- ⦿ Sigmoid activations functions



Multi-layer perecptron is another name for Neural Networks

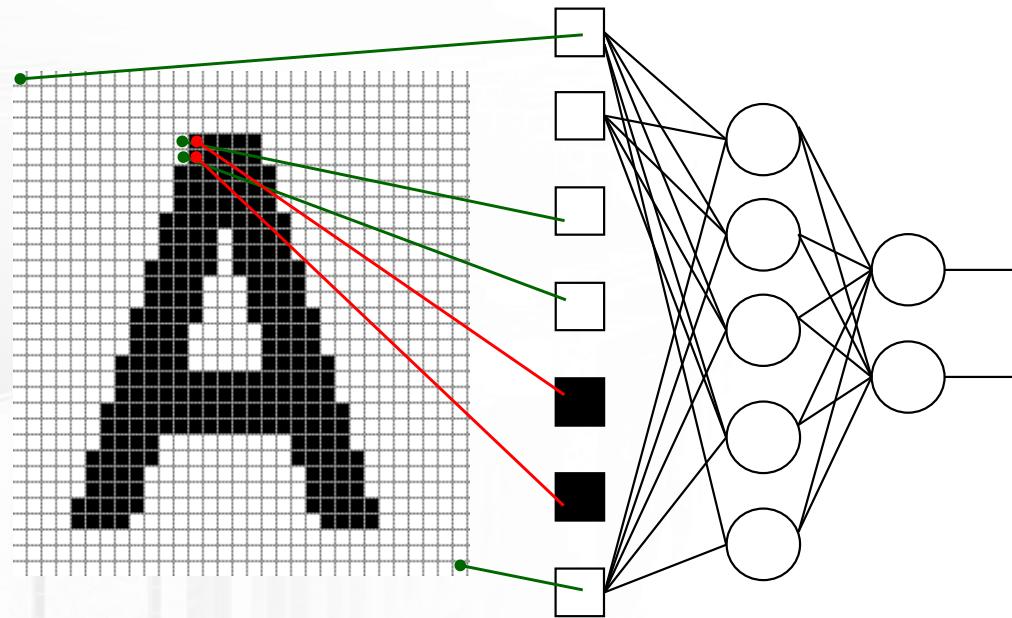
Drawbacks of previous neural networks

- the number of **trainable parameters** becomes extremely large



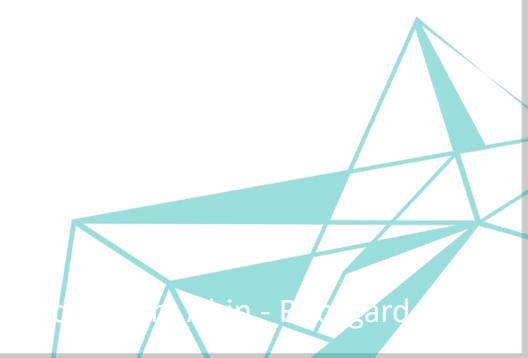
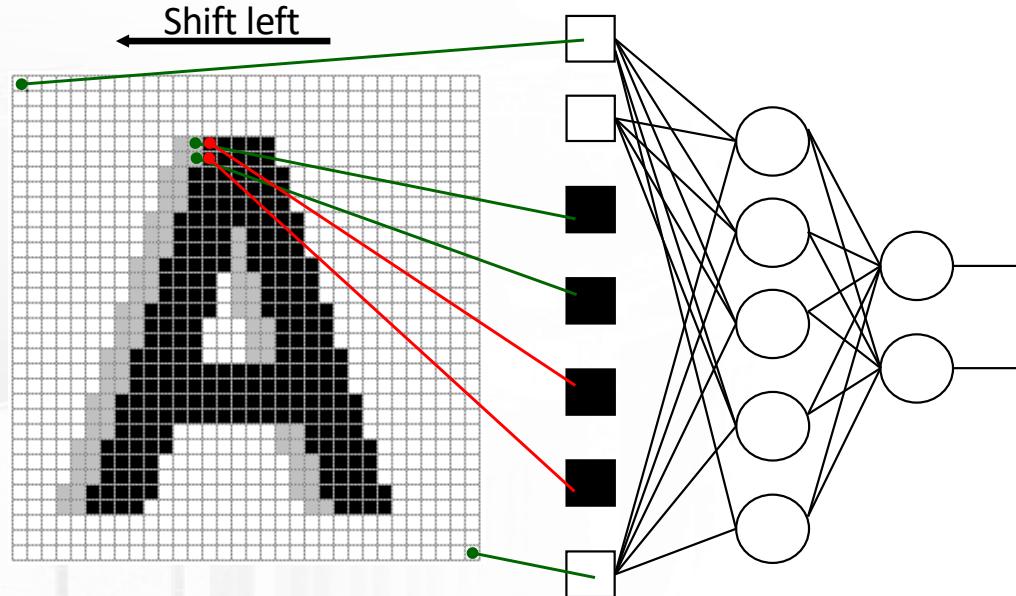
Drawbacks of previous neural networks

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

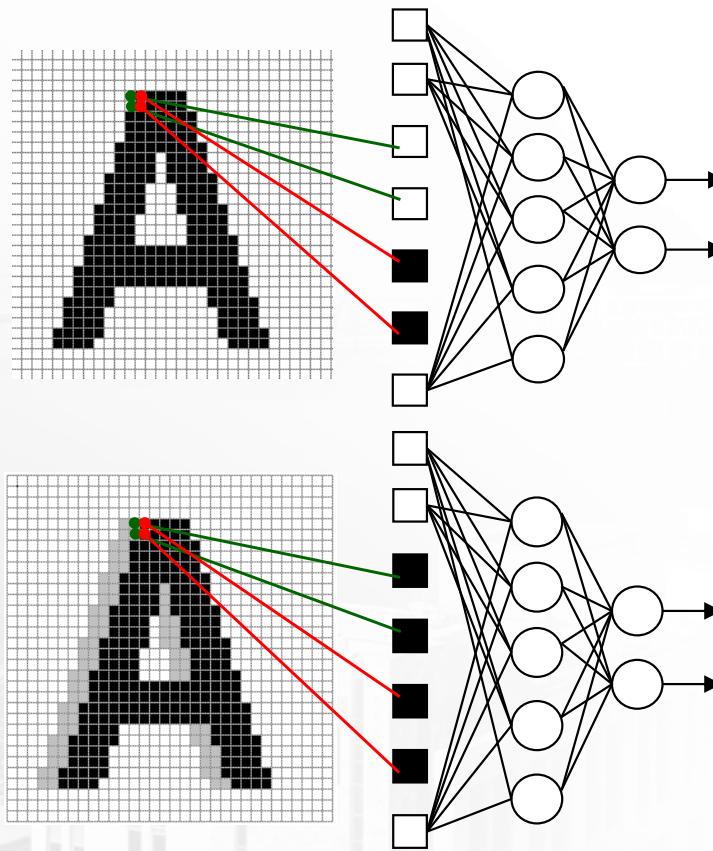


Drawbacks of previous neural networks

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

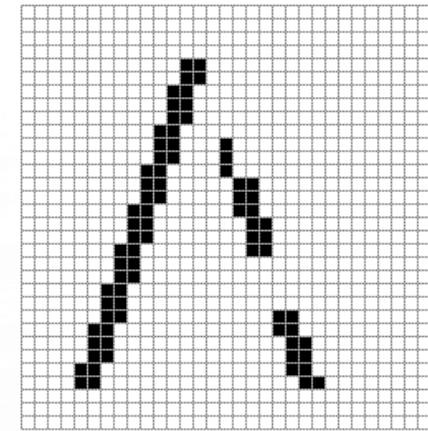


Drawbacks of previous neural networks



~~Topology~~: Where does the data come from

- NN does not care for topology
but we want it to



154 input change
from 2 shift left

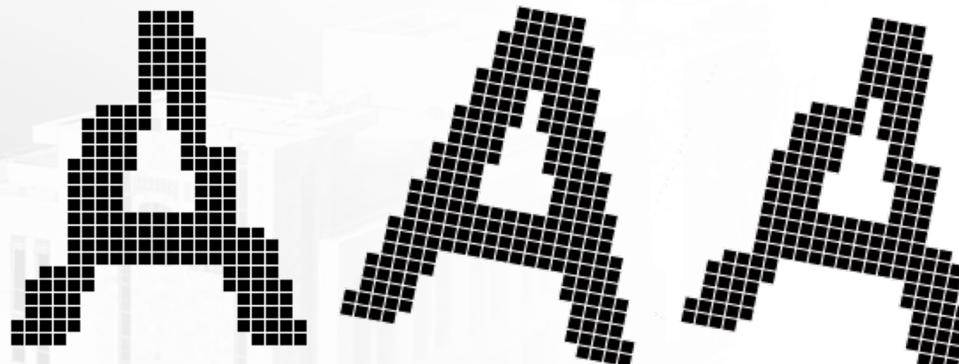
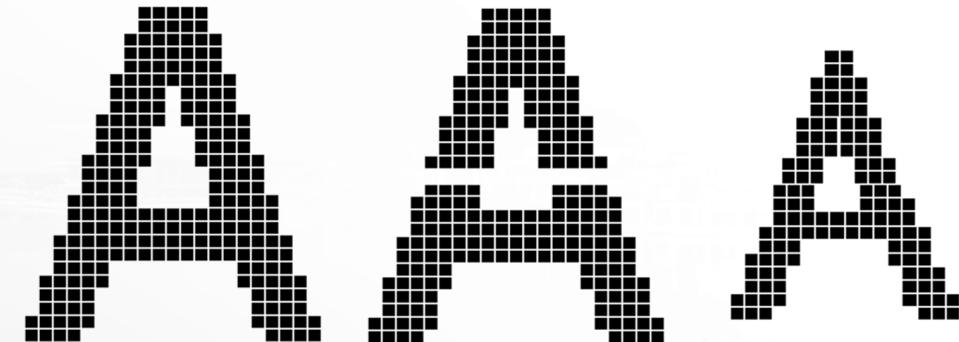
77 : black to white

77 : white to black

if I did this the output
should not change
& NN should know this

Drawbacks of previous neural networks

- scaling, and other forms of distortion



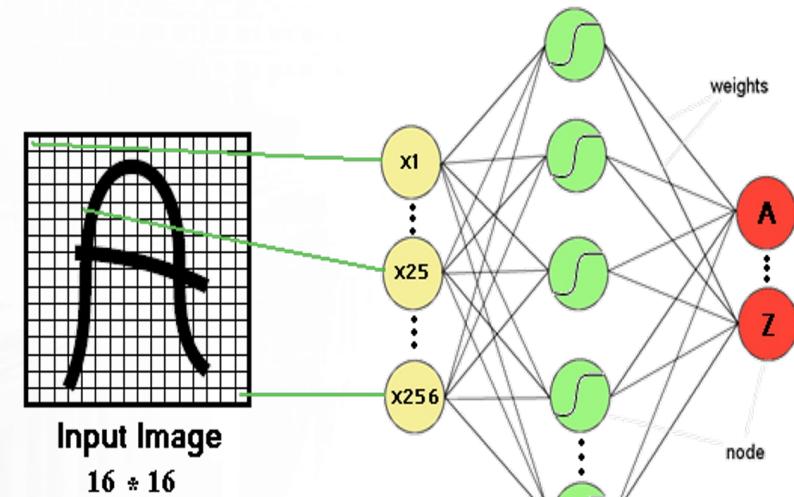
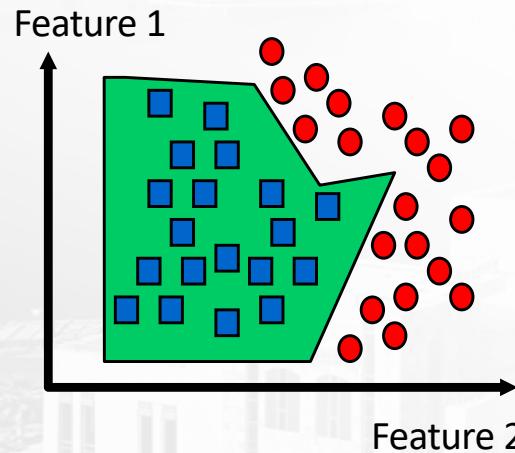
AAA
AAA



Drawbacks of previous neural networks

- the topology of the input data is completely ignored
- work with raw data.

If we look at data pixel wise we would never recognise a thing. It's the pattern that it counts



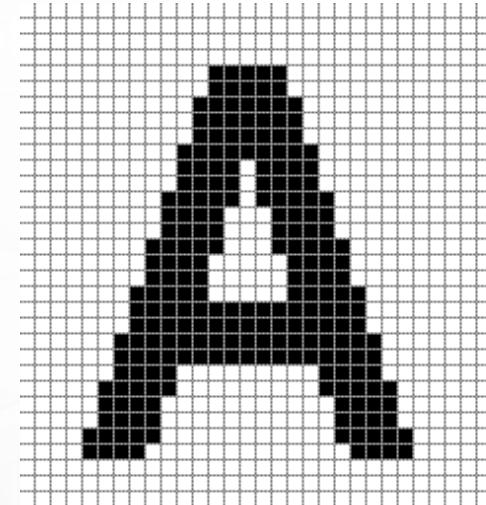
Drawbacks of previous neural networks

Black and white patterns:

$$2^{32 \times 32} = 2^{1024}$$

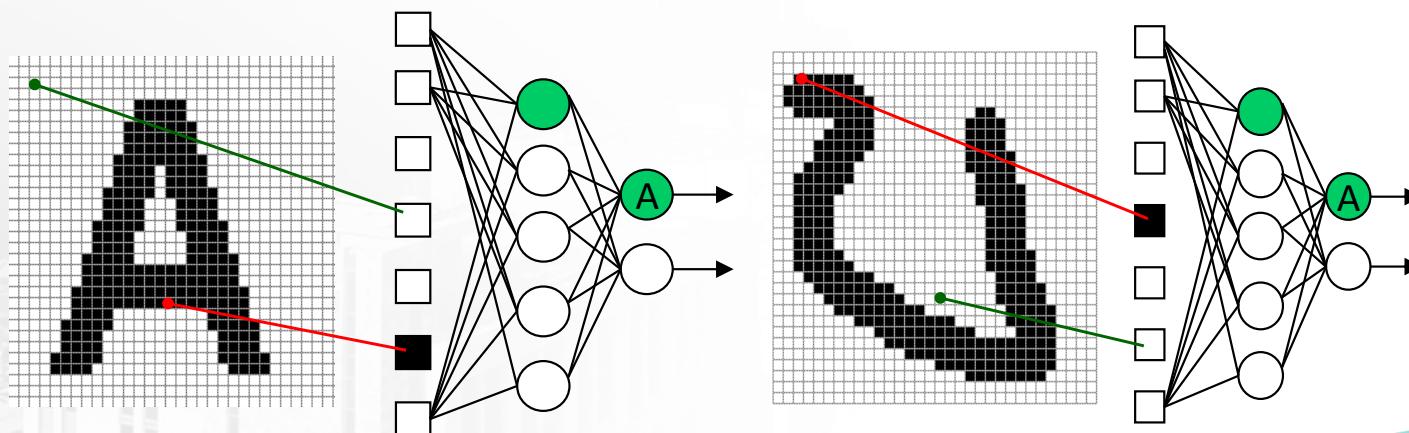
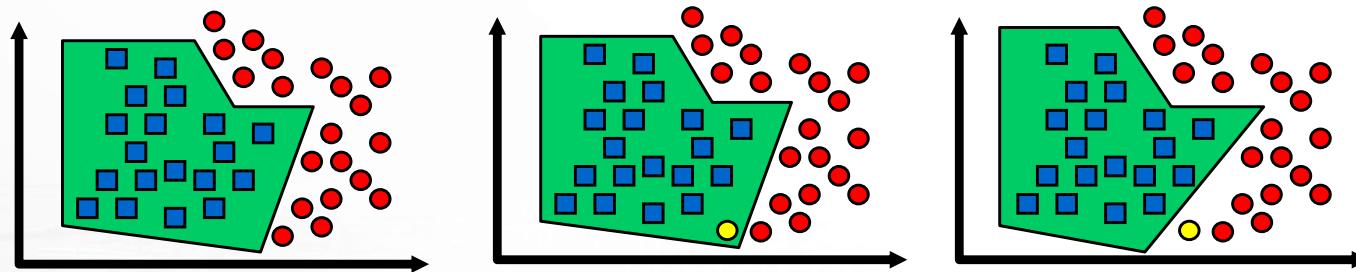
$$256^{32 \times 32} = 256^{1024}$$

Gray scale patterns :



32 * 32 input image

Drawbacks of previous neural networks



Improvement

⌚ Fully connected network of sufficient size can produce outputs that are invariant with respect to such variations.

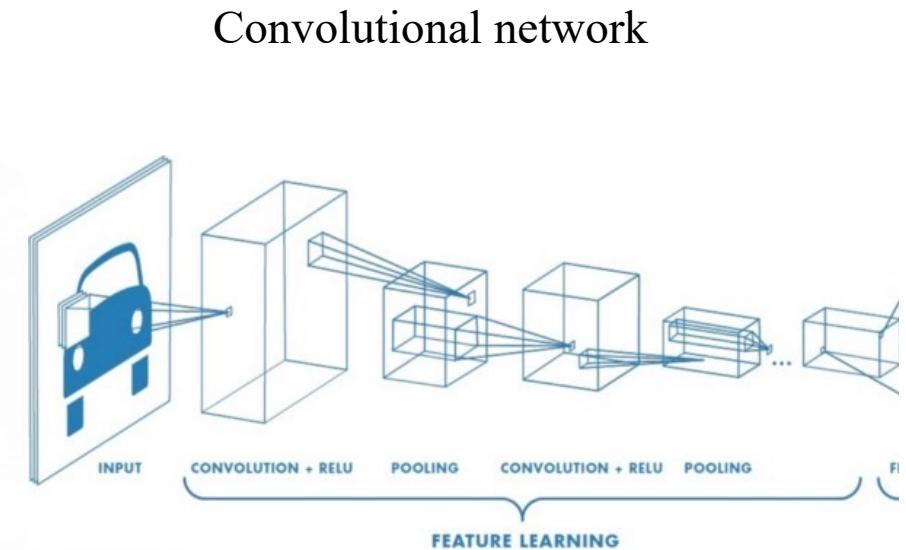
⌚ **Training time**

⌚ **Network size**

⌚ **Free parameters**

Convolutional Neural Network (CNN)

- Local connectivity
- Reuse parameters



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

Convolutional Neural Network

- How Convolution works
- Filters / Padding / Stride
- Activations
- Pooling Layers

How Convolution works

- Convolution on an image



Naeemullah Khan

How Convolution works

- Break the image in overlapping tiles

MaxPooling → yes or No (have you found the pattern?)



Cats/Dogs

If this type of pattern it
Means its a cat

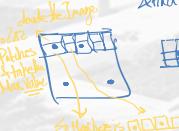
(+) If this pattern
its a dog

Dot Product is the Simplest Way to find
Similarity

If the Value is high it Means that this Part of the
Image is Similar to the Patch

If the Value is low then this Part of the
Image is Not Similar to the Patch

So I will try to take the Dot Product
Define Place with High Value

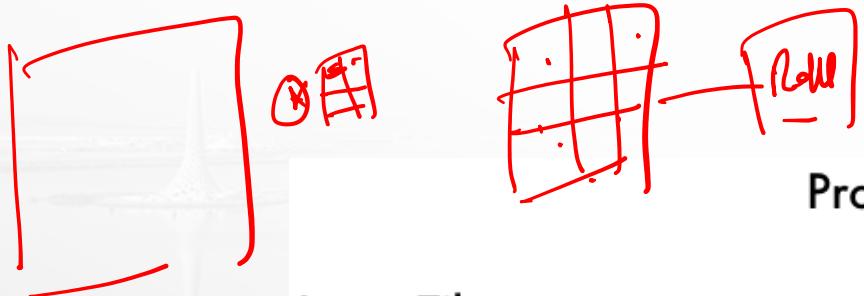


Naeemullah Khan

How Convolution works

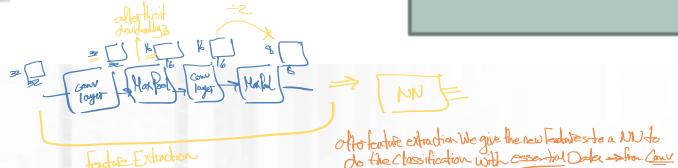
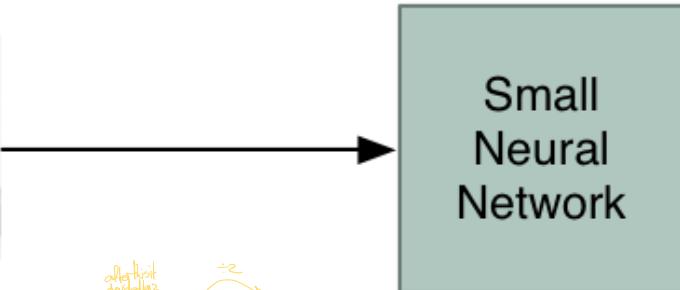
$$\begin{matrix} x^1 & x^2 \\ x^3 & x^4 \end{matrix} \cdot \begin{matrix} w^1 & w^2 \\ w^3 & w^4 \end{matrix} = \begin{matrix} w^1x^1 + w^3x^3 & w^2x^2 + w^4x^4 \end{matrix}$$

- Feed each image tile into small NN



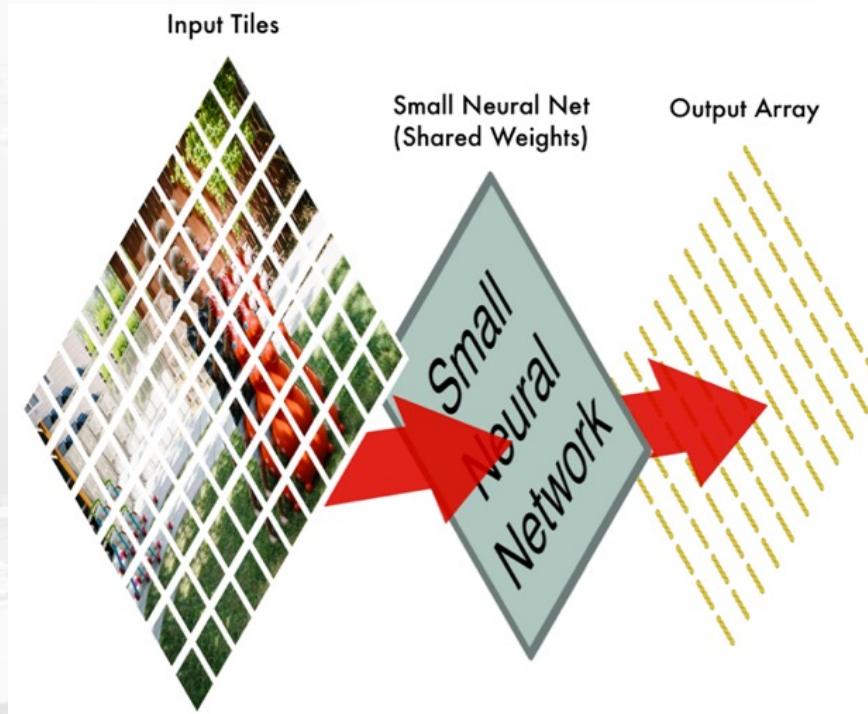
Processing a single tile

Input Tile



How Convolution works

- Store results from each tile into a new array

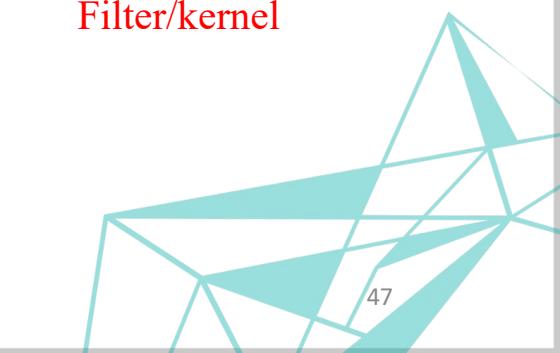


Naeemullah Khan

$$z_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{m-1} W_{ab} x_{(i+a)(j+b)}$$

Filter/kernel

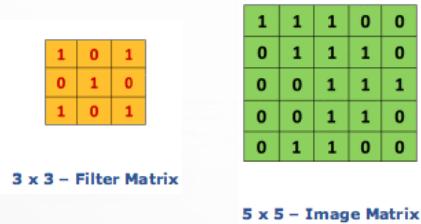
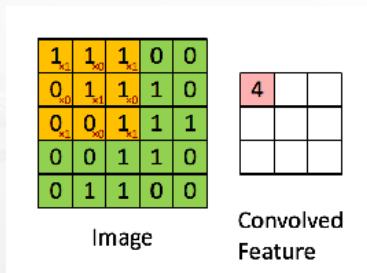
1



CNN – Filters

$$z_{ij} = W * x_{i,j}$$

- Filters and convolution



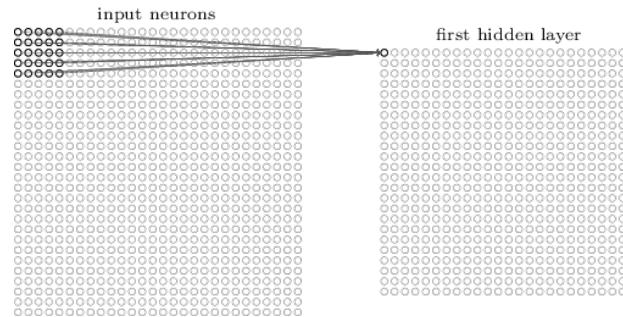
Basically How Conv changes the Image:

If filter is $f \times f$ you will lose $\frac{f-1}{2}$ from the boundary

$$\text{filter } 3 \times 3 \Rightarrow \frac{3-1}{2} = 1, \quad 5 \times 5 \Rightarrow \frac{5-1}{2} = 2$$

here we will lose 2 row, column from the boundary of the Image

You can add Padding if you want to keep the output size
your image



highlighting the more activated part of the image

Smoothing

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

Multiple filters

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

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

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

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

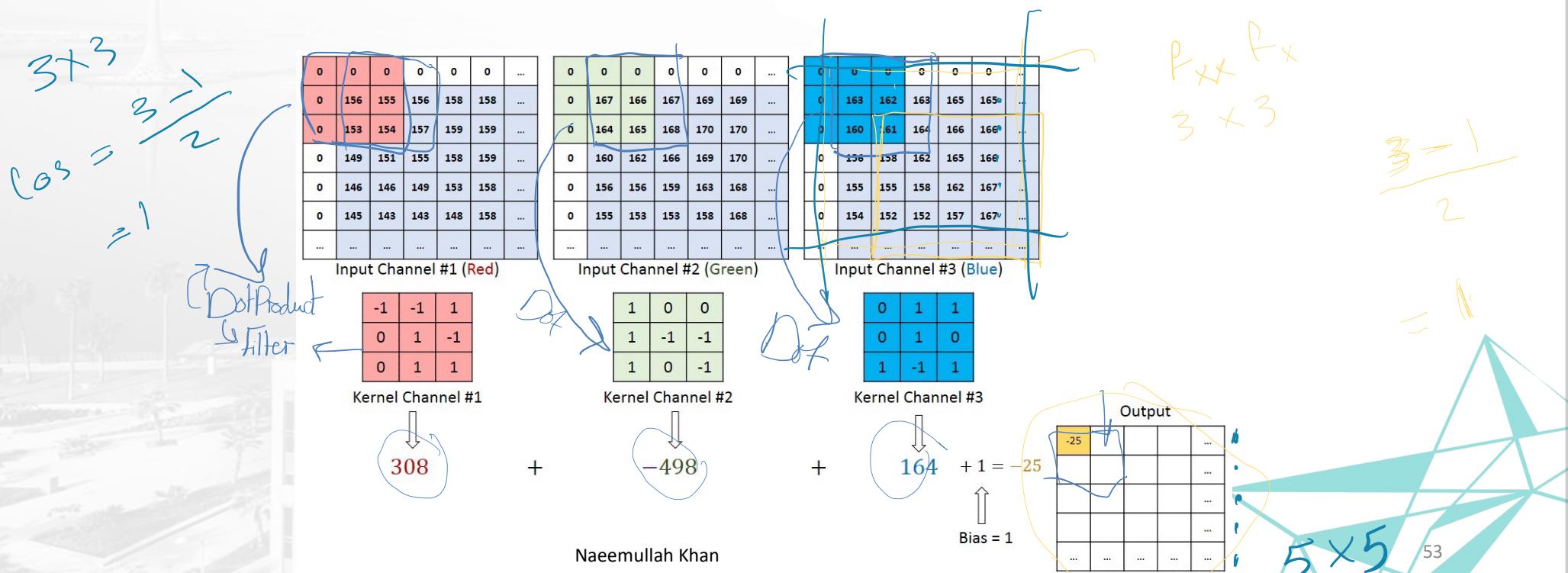
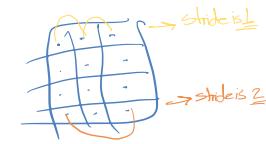


CNN – Filters

$f_1 \times f_2 \}$ size of the filter

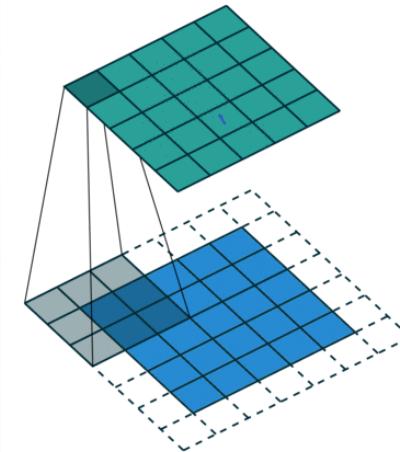
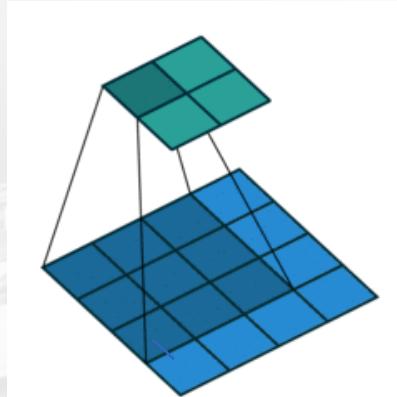
- Multiple channels
 - Ex: RGB color

$$\text{Loss} = \frac{f_1 - 1}{2} =$$



CNN - Padding

- What happen to the borders?
 - Reduce size of the output
 - Pad the picture with zeros (zero-padding)

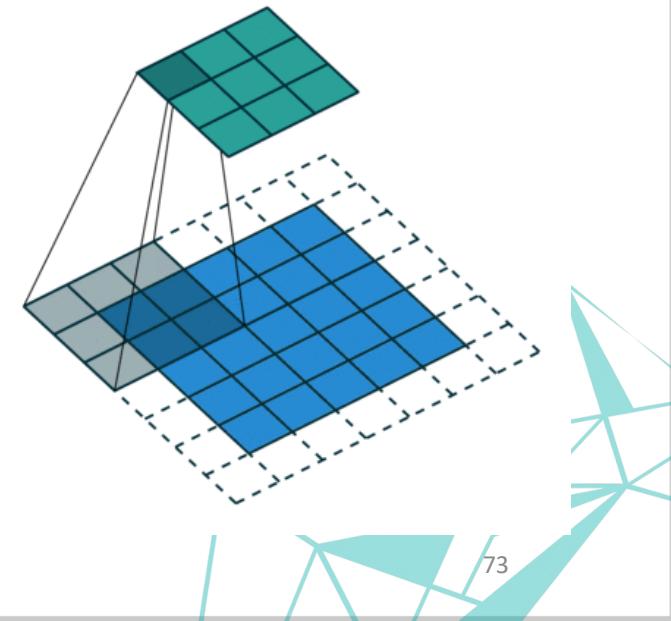


Naeemullah Khan

More on: http://deeplearning.net/software/theano/tutorial/conv_arithmeti.html

CNN – Strides

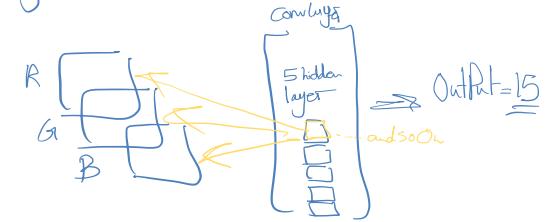
- If not need to keep all output
 - Skip some pixels
 - Stride: number of skipped pixels
 - Output with smaller resolution



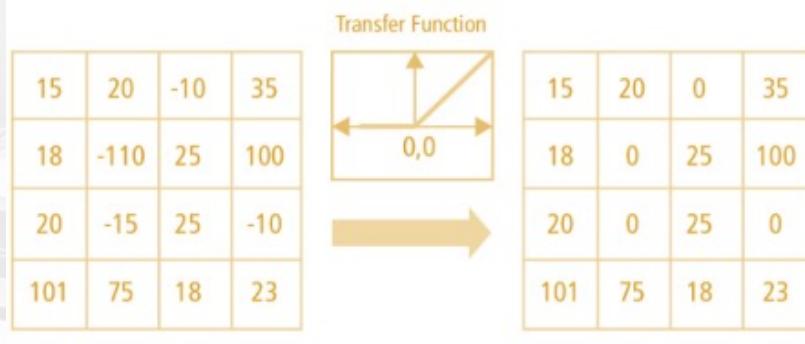
CNN – Activations

- Any previous activation fits
- Breaks linearity
- Ex: Rectified Linear Unit (ReLU)

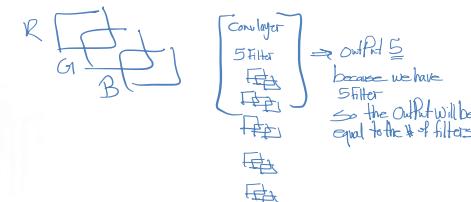
Image has 3 channels



$$\sigma(x) = \max(0, x)$$

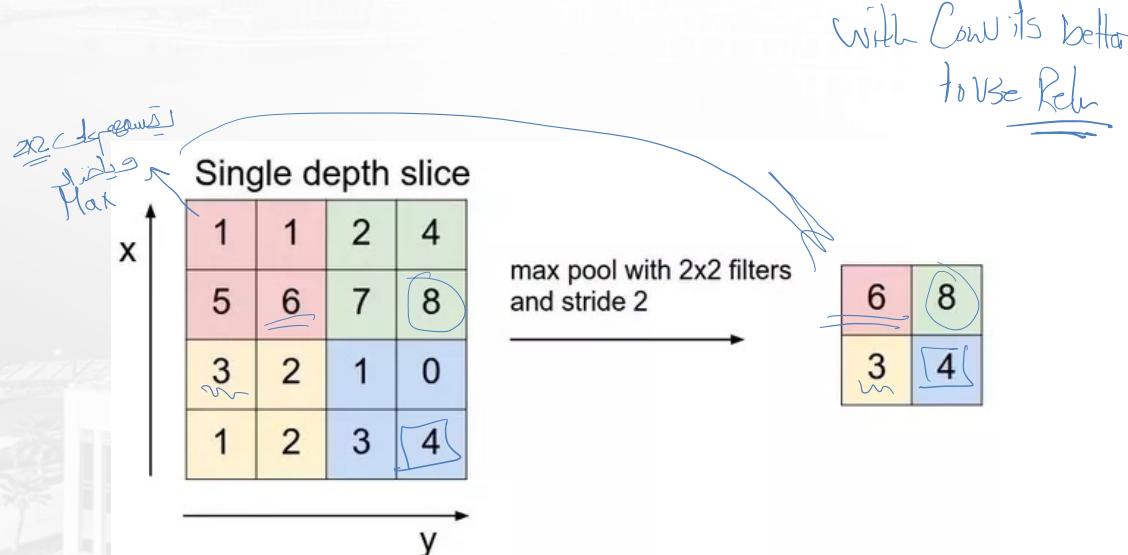


Naeemullah Khan



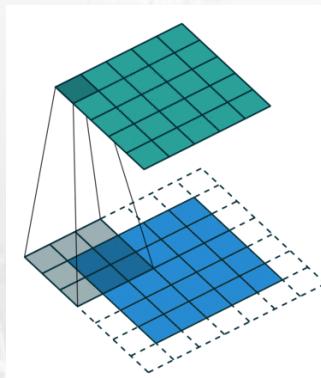
CNN – Pooling Layer

- Max / Mean / Sum Pooling
- Reduce the resolution



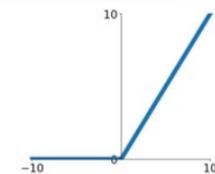
A convolution layer

Convolution



Activation

ReLU
 $\max(0, x)$



Pooling

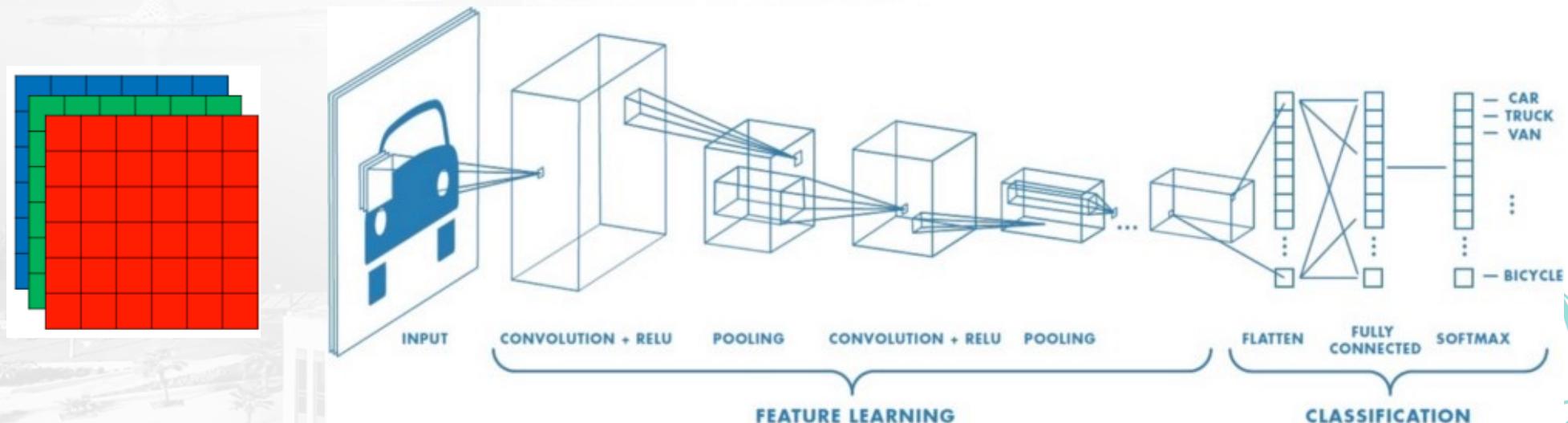
Single depth slice			
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters
and stride 2

6	8
3	4

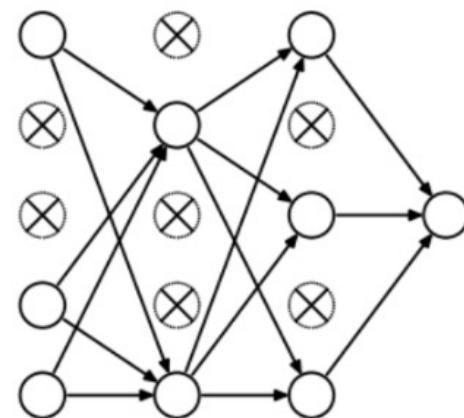
Deep Convolutional network

- Multiple convolutional layers



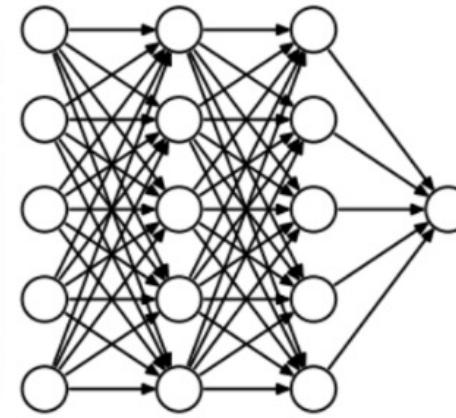
Deep Convolutional network

- Add redundancy to the network
- A way of regularization



Training

Naeemullah Khan



Testing

CNN – Most notable network

- LeNet-5 (1998)

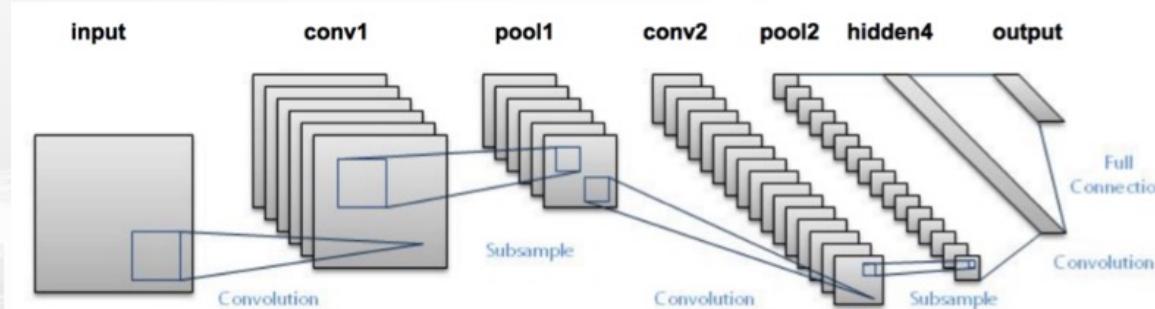
- 2 convolutional layers, 2 pooling layers, 2 fully connected layers
- Used for handwritten numbers on checks

[PDF] LeNet - Yann LeCun

yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf ▾

by Y LeCun - 1998 - Cited by 17209 - Related articles

OF THE IEEE, NOVEMBER 1998. 1 ... Real-life document recognition systems are composed Convolutional Neural Network called LeNet - 5 described in.



CNN – Most notable network

- AlexNet (2012)

- First big improvement in image classification
- CNN / MaxPool / Dropout / ReLU
- 5 convolutional layers, followed by max-pooling layers; the last three are fully connected layers

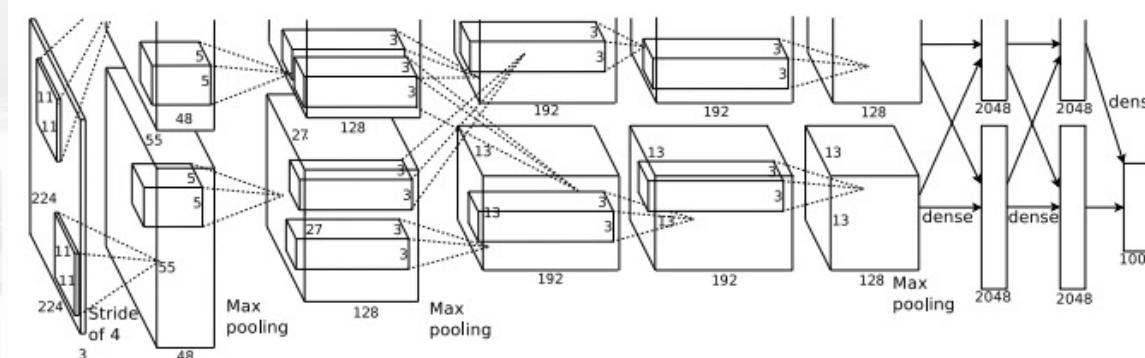
[PDF]

[ImageNet Classification with Deep Convolutional ... - NIPS Proceedings](#)

<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neu...>

by A Krizhevsky - 2012 - Cited by 36889 - Related articles

The specific contributions of this paper are as follows: we trained one of the largest convolutional neural networks to date on the subsets of ImageNet used in the ...



Naeemullah Khan

CNN – Most notable network

- VGGNet-16 (2014)
 - 13 convolutional layers, 3 fully-connected layers

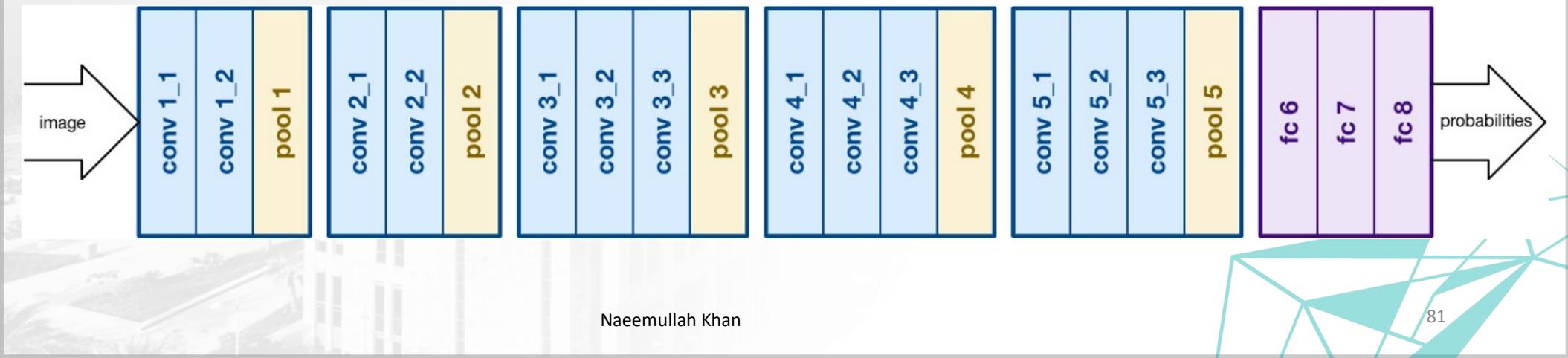
[PDF] Very Deep Convolutional Networks for Large-Scale Image Recognition

<https://arxiv.org/pdf/1409.1556.pdf>

by K. Simonyan - 2014 - Cited by 20568 - Related articles

Apr 10, 2015 - CVJ 10 Apr 2015. Published as a conference paper at ICLR 2015 ... Google DeepMind.

http://www.robots.ox.ac.uk/~vgg/research/very_deep/. 1 ...



CNN – Most notable network

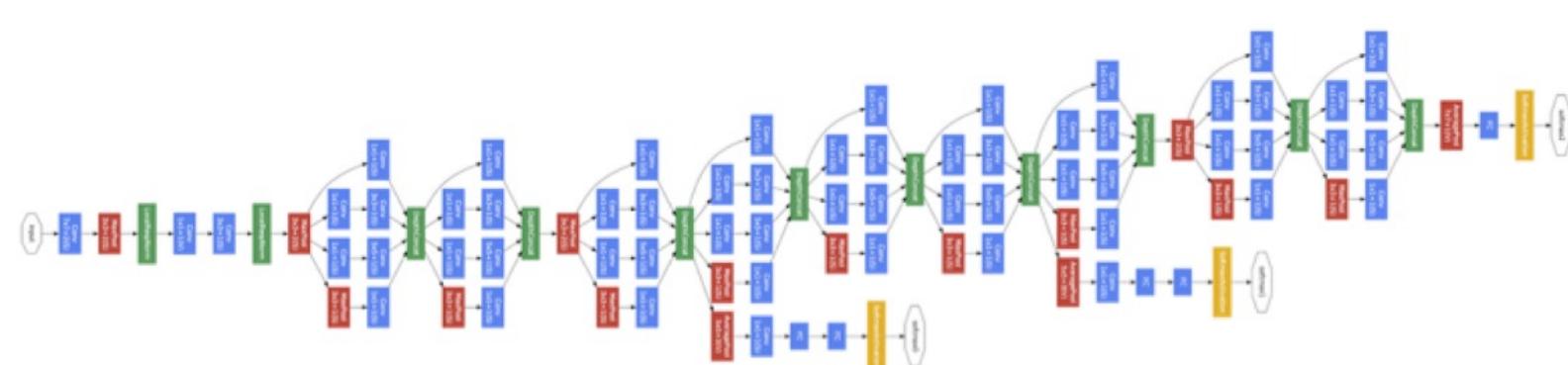
- GoogleNet/Inception (2014)
 - Going Deep: 22 layers
 - Very small convolutions
 - Drastically reduce the number of parameters
 - From 138M (VGG) to 4M parameters

[Going Deeper with Convolutions](#)

<https://arxiv.org/> > cs

by C Szegedy - 2014 - Cited by 12157 - Related articles

Sep 17, 2014 - Abstract: We propose a deep convolutional neural network architecture codenamed "Inception", which was responsible for setting the new state ...



Naeemullah Khan

CNN – Most notable network

- ResNet (2015)
 - Residual layer
 - Up to 152 layers (50-101-152)
 - Requires huge resources

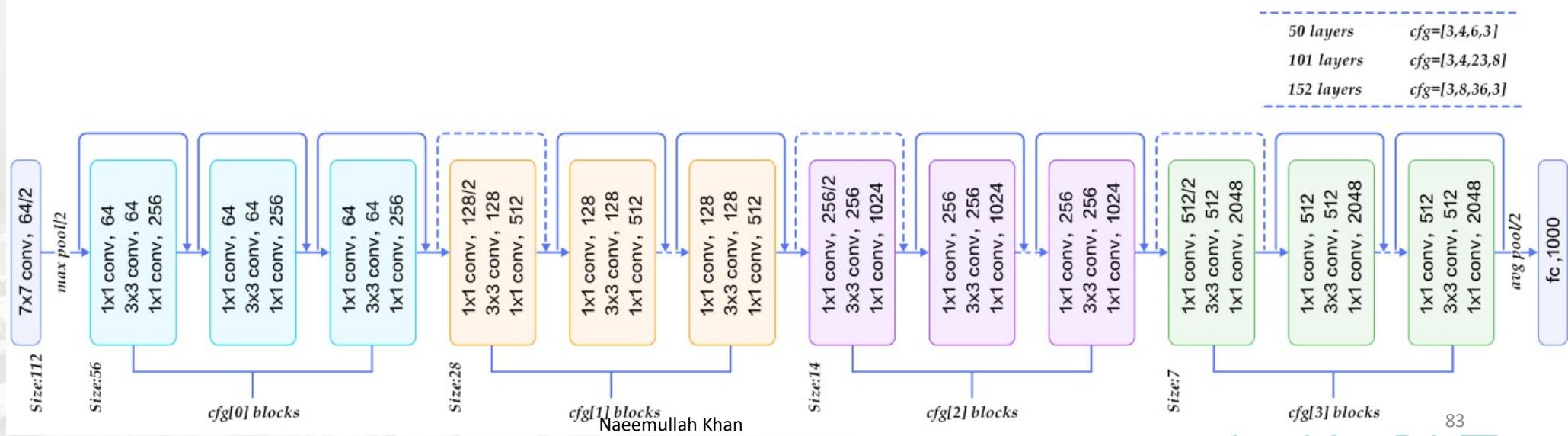
[ResNet - arXiv](#)

www.arxiv.org/abs/1512.03385

by K He - 2015 - Cited by 19344 - Related articles

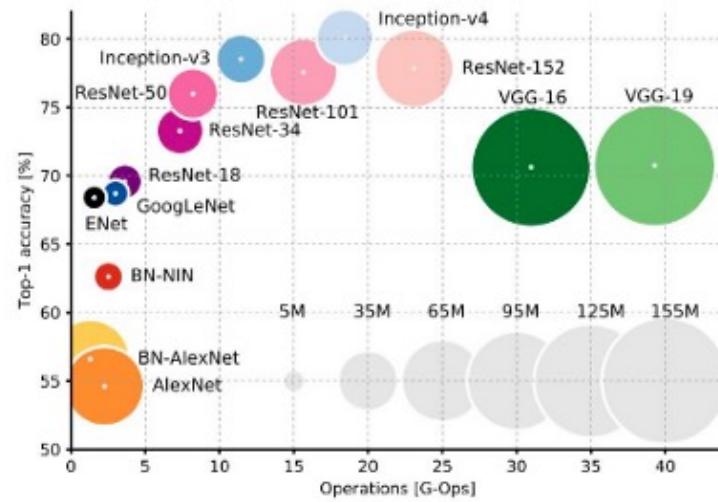
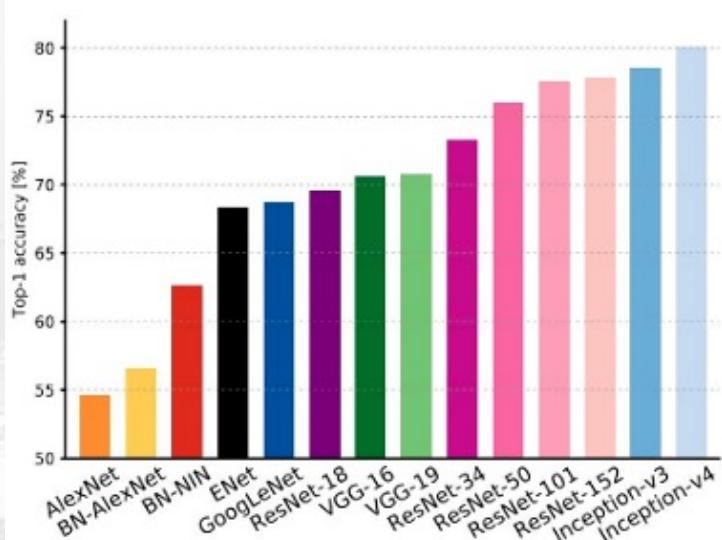
No information is available for this page.

Learn why



CNN – Most notable network

- Totally more than 100 K citations



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Naeemullah Khan

Outline

- From machine learning to deep learning
- Deep neural network
- Train a deep neural network
- Convolutional neural networks
- Application on computer vision

Application on Computer Vision

Classification



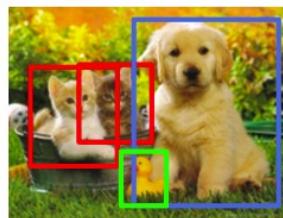
CAT

Classification + Localization



CAT

Object Detection



CAT, DOG, DUCK

Instance Segmentation



CAT, DOG, DUCK

CLASSIFICATION



ELLEN

LOCALIZATION



ELLEN

SINGLE OBJECT

DETECTION



ELLEN, JULIA, PETER, JENNIFER, BRADLEY,
BRAD, MERYL, KEVIN, LUPITA, CHANNING

SEGMENTATION

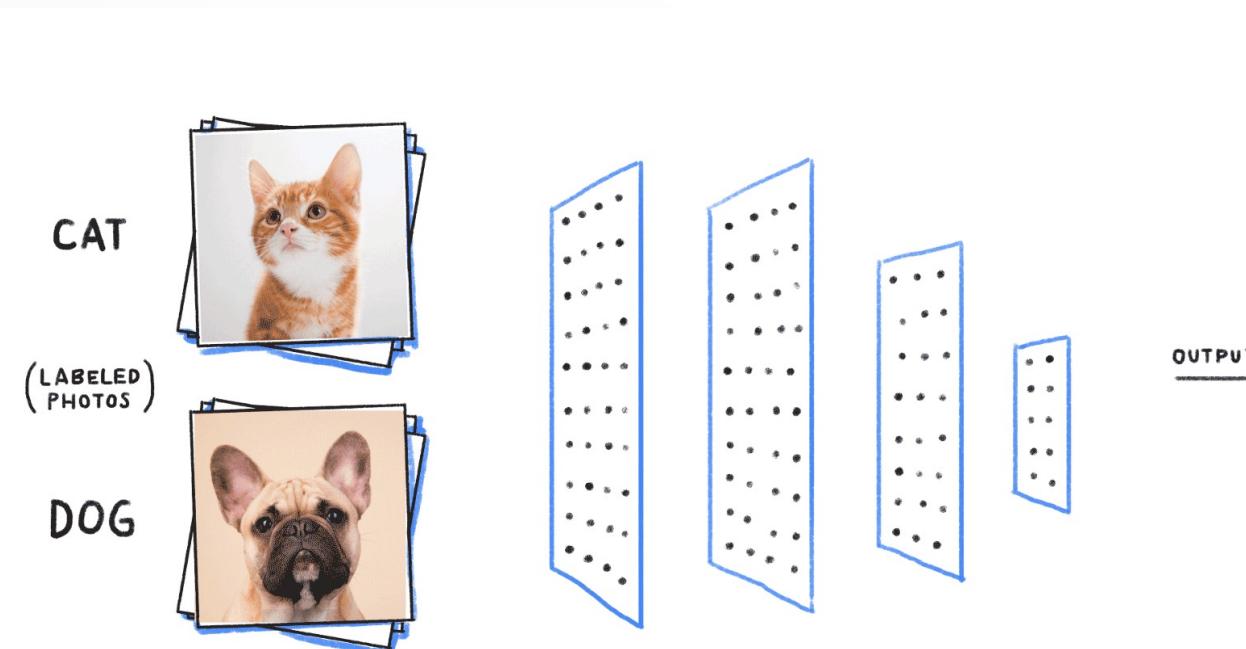


ELLEN, JULIA, PETER, JENNIFER, BRADLEY,
BRAD, MERYL, KEVIN, LUPITA, CHANNING

MULTIPLE OBJECTS

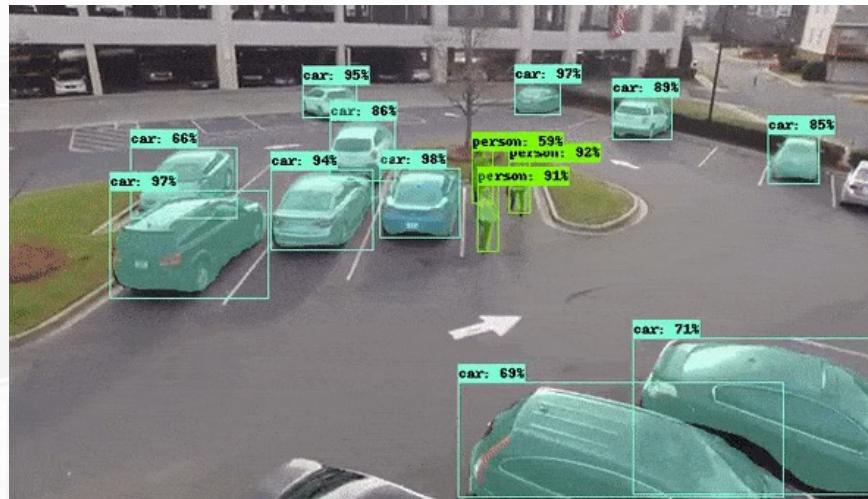
Application on Computer Vision

Classification



Application on Computer Vision

Object detection



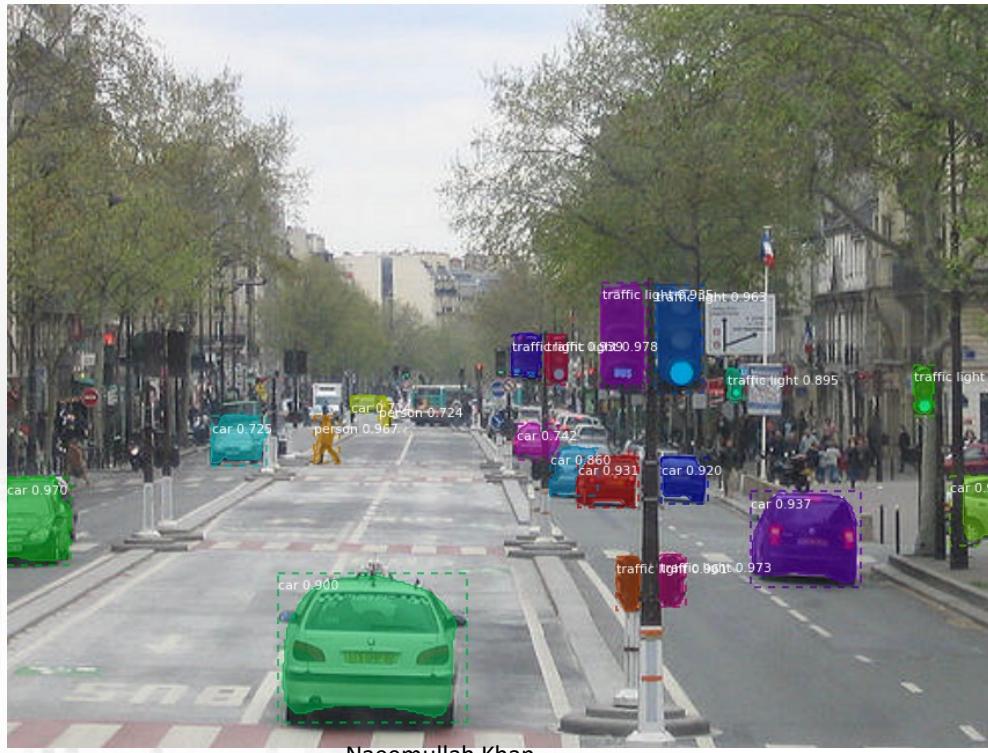
Application on Computer Vision

Semantic segmentation



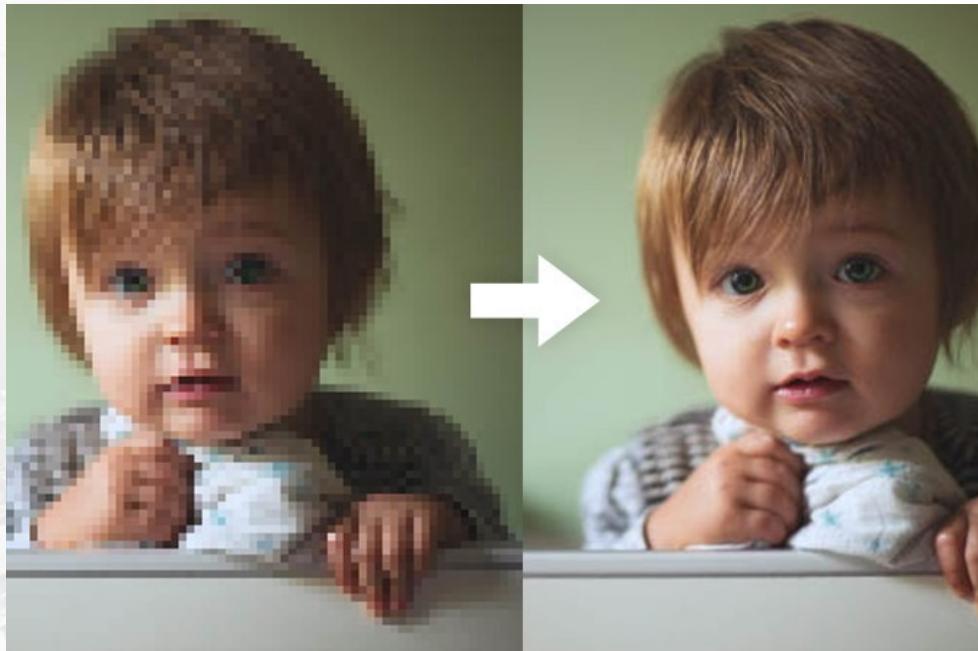
Naeemullah Khan

Instance segmentation



Application on Computer Vision

Image super-resolution



Naeemullah Khan

Dataset and Challenges

Classification

- MNIST

- Digit Classification (handwritten)
- Grayscale intensity from 0 to 255
- $28 \times 28 = 768$ pixels
- 10 classes
- 10,000 images



CAT

000000000000000000
111111111111111111
222222222222222222
333333333333333333
444444444444444444
555555555555555555
666666666666666666
777777777777777777
888888888888888888
999999999999999999

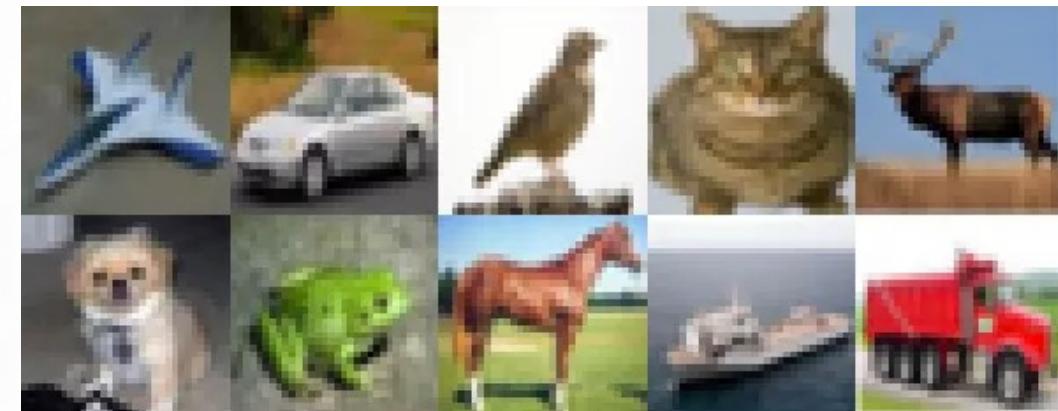
Dataset and Challenges

Classification

- CIFAR-10
 - Object Classification
 - 3 RGB channels from 0 to 255
 - $32 \times 32 = 1024$ pixels
 - 10 classes
 - 60,000 images



CAT



Dataset and Challenges

Classification

- ImageNet
 - Object Classification
 - 3 RGB channels from 0 to 255
 - Any resolution
 - 1000 classes
 - 14,197,122 images



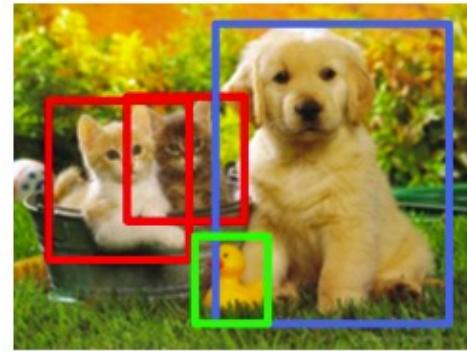
CAT



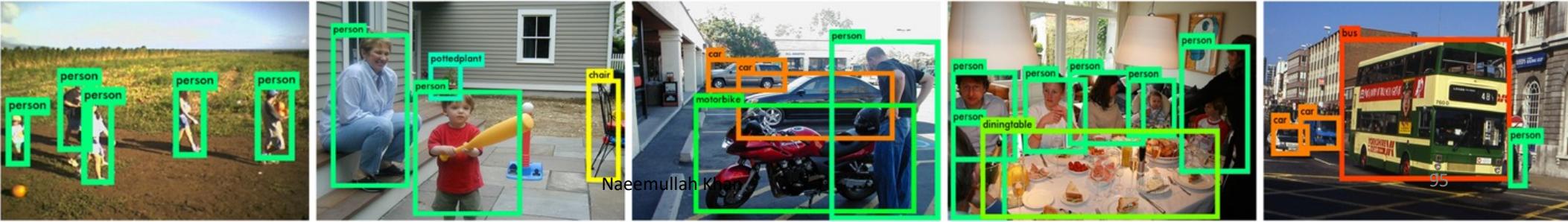
Dataset and Challenges

- Pascal VOC (2007)
 - Object Detection
 - 9,963 images
 - 24,640 annotated objects

Object Detection



CAT, DOG, DUCK



Dataset and Challenges

- Pascal VOC (2012)
 - Object Segmentation
 - 11,530 images
 - 27,450 objects
 - 6,929 segmentations



CAT, DOG, DUCK

Instance
Segmentation



Image Classification – Applications

- Image/Photo search (Google Lens)

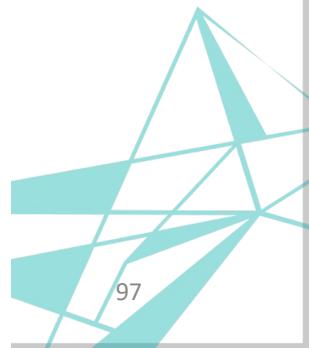
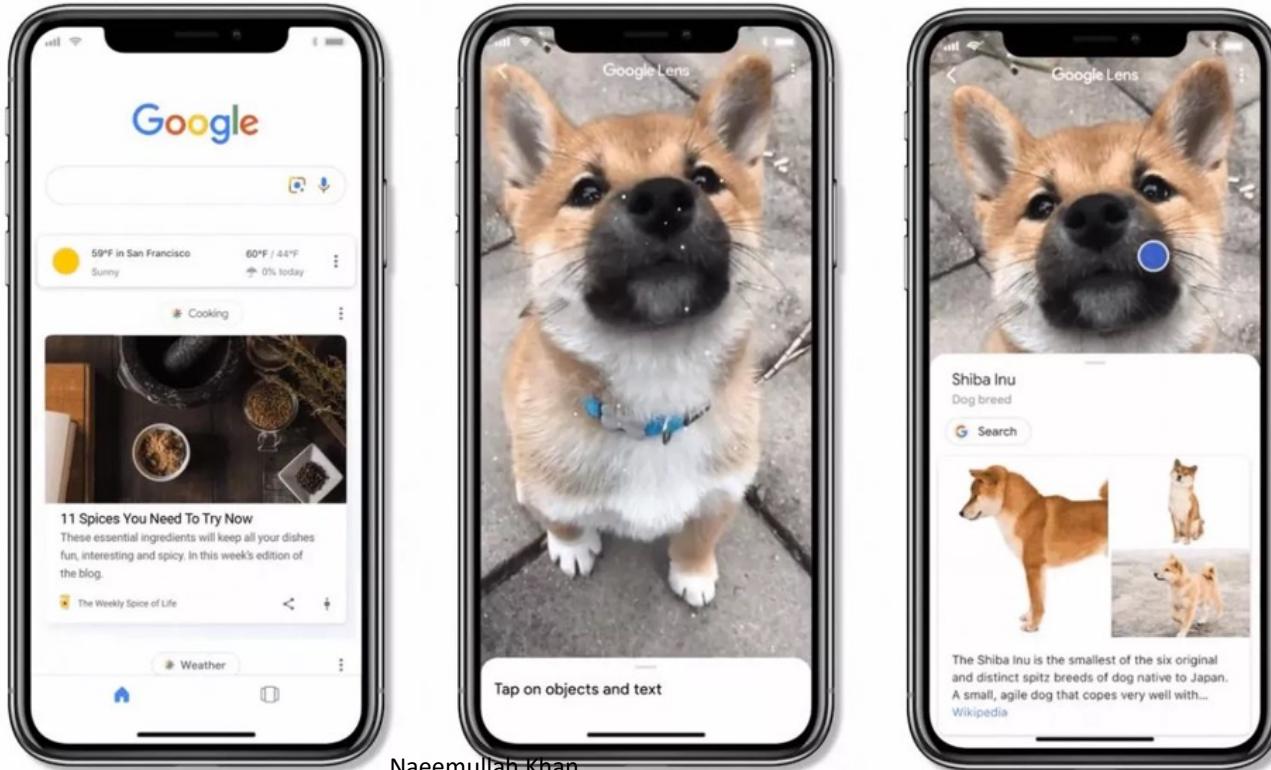




Image Classification – Applications

- Agriculture



Maize



Common wheat



Sugar beet



Scentless Mayweed



Chickweed



Shepherd's Purse



Cleavers



Charlock



Fat Hen



Cranesbill

Naeemullah Khan Black-grass



Loose Silky-bent

PlantNet Plant Identification



Download on the
App Store

ANDROID APP ON
Google play

Object Localization



- Detection as a classification
 - Sliding window

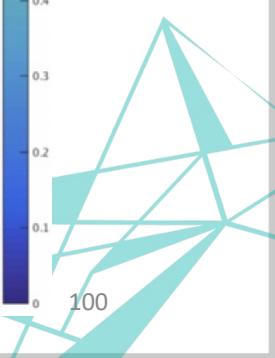
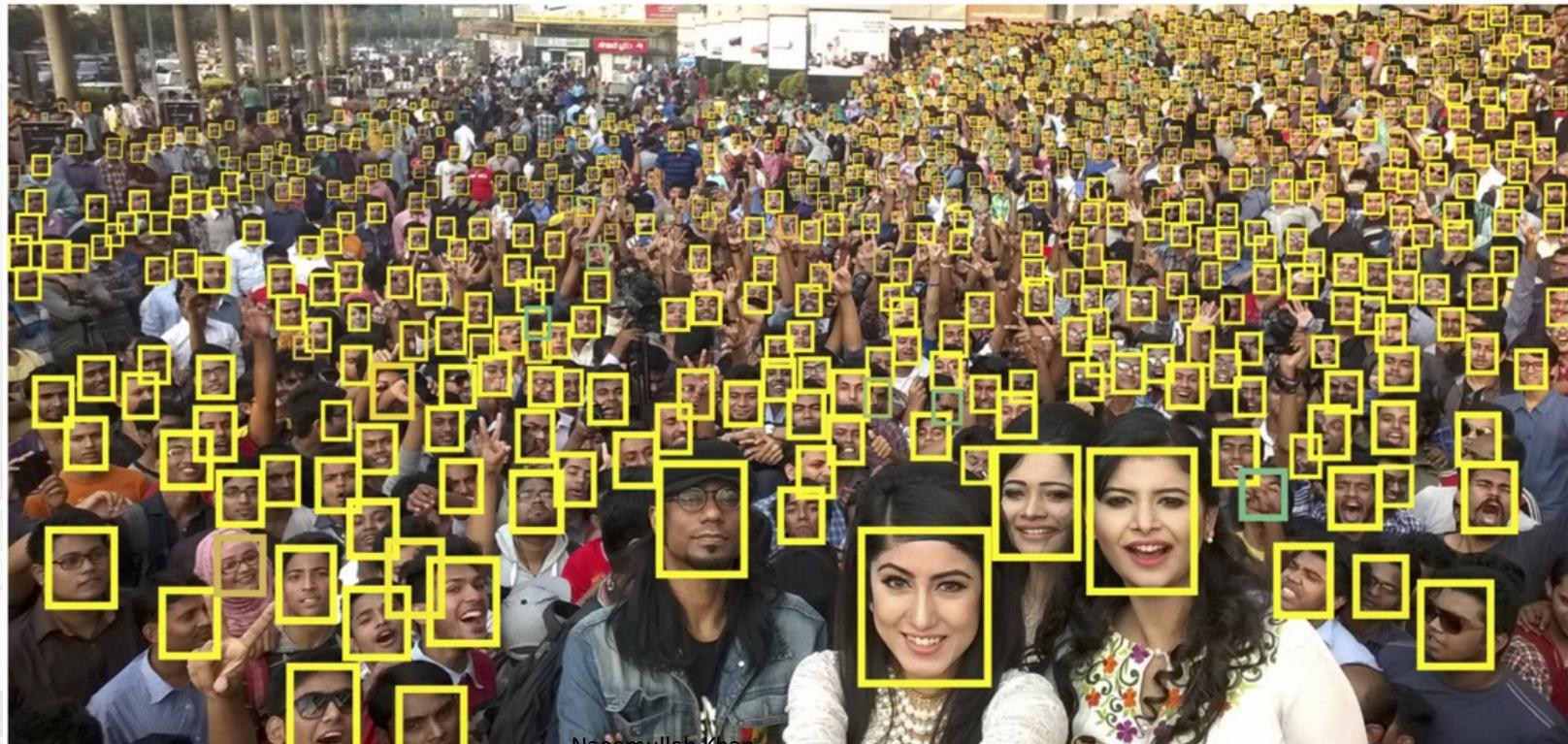


Naeemullah Khan

Object Localization - Application



- Finding tiny faces in the wild



Object Detection – Application



- Drone Application
 - Parking Lot Vehicle Detection Using DL



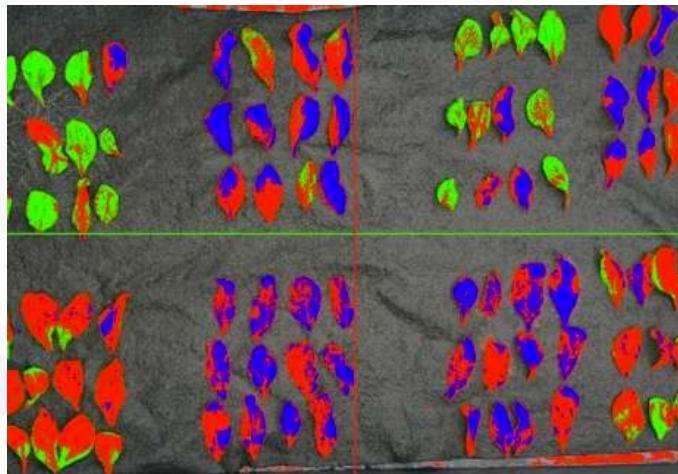
Naeemullah Khan

More on: <https://medium.com/geoai/parking-lot-vehicle-detection-using-deep-learning-49597917bc4a>

Object Detection – Application



- Agriculture:
 - Early plant disease detection
-



Naeemullah Khan

More on: <https://phys.org/news/2017-04-drones-early-disease-crops.html#jCp>

Image Segmentation



- SegNet

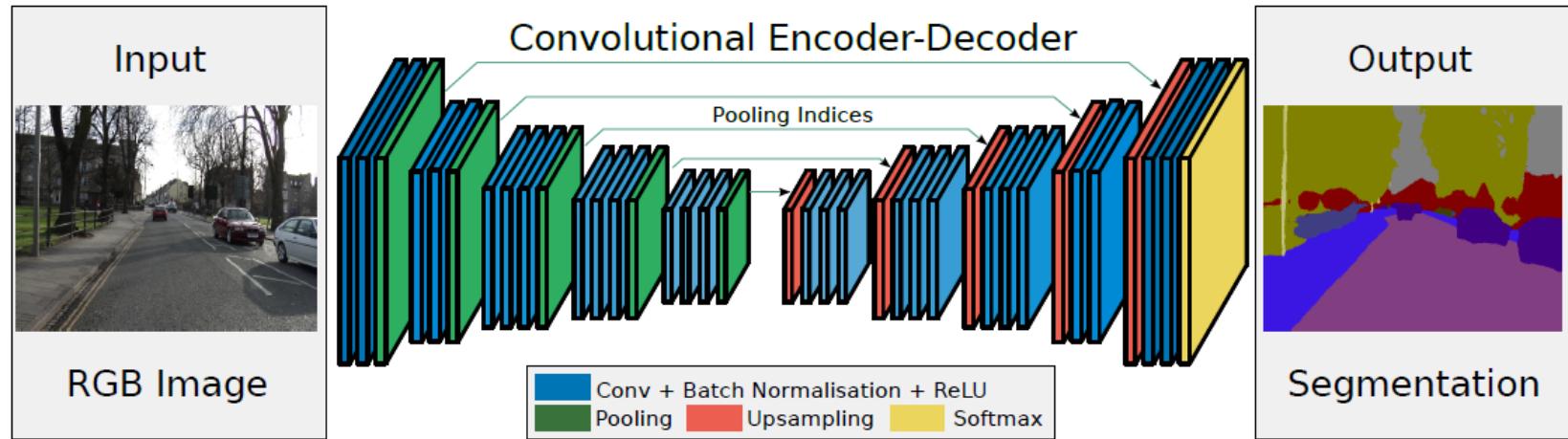
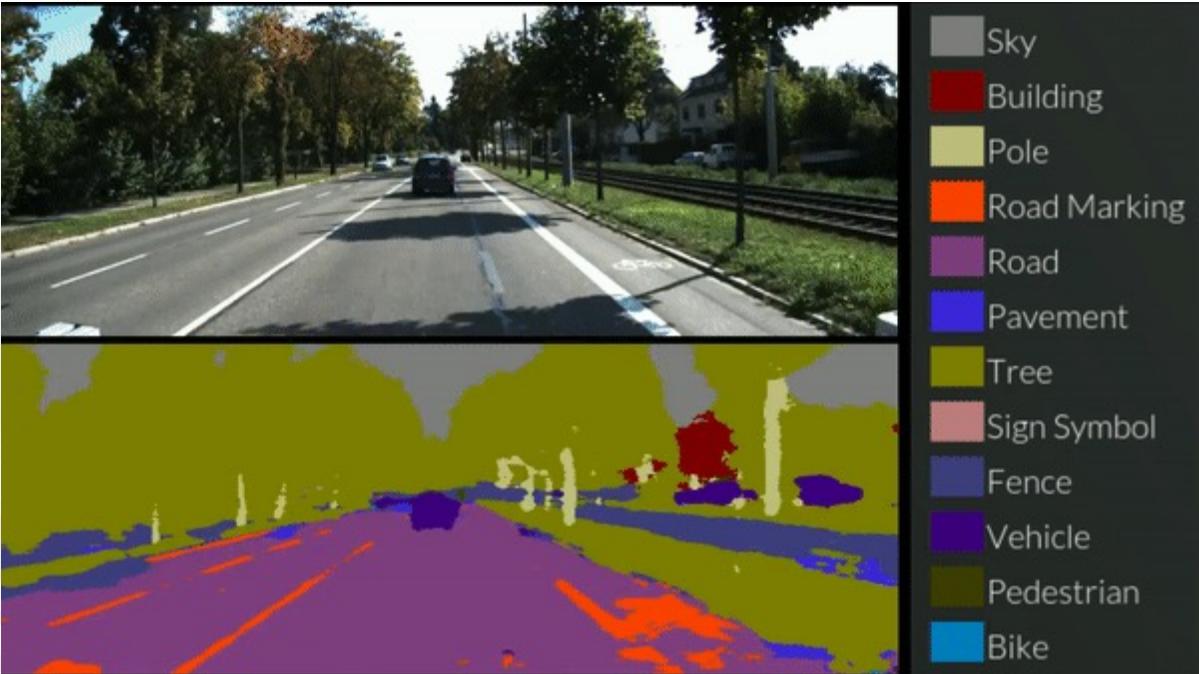




Image Segmentation – Application

- SegNet

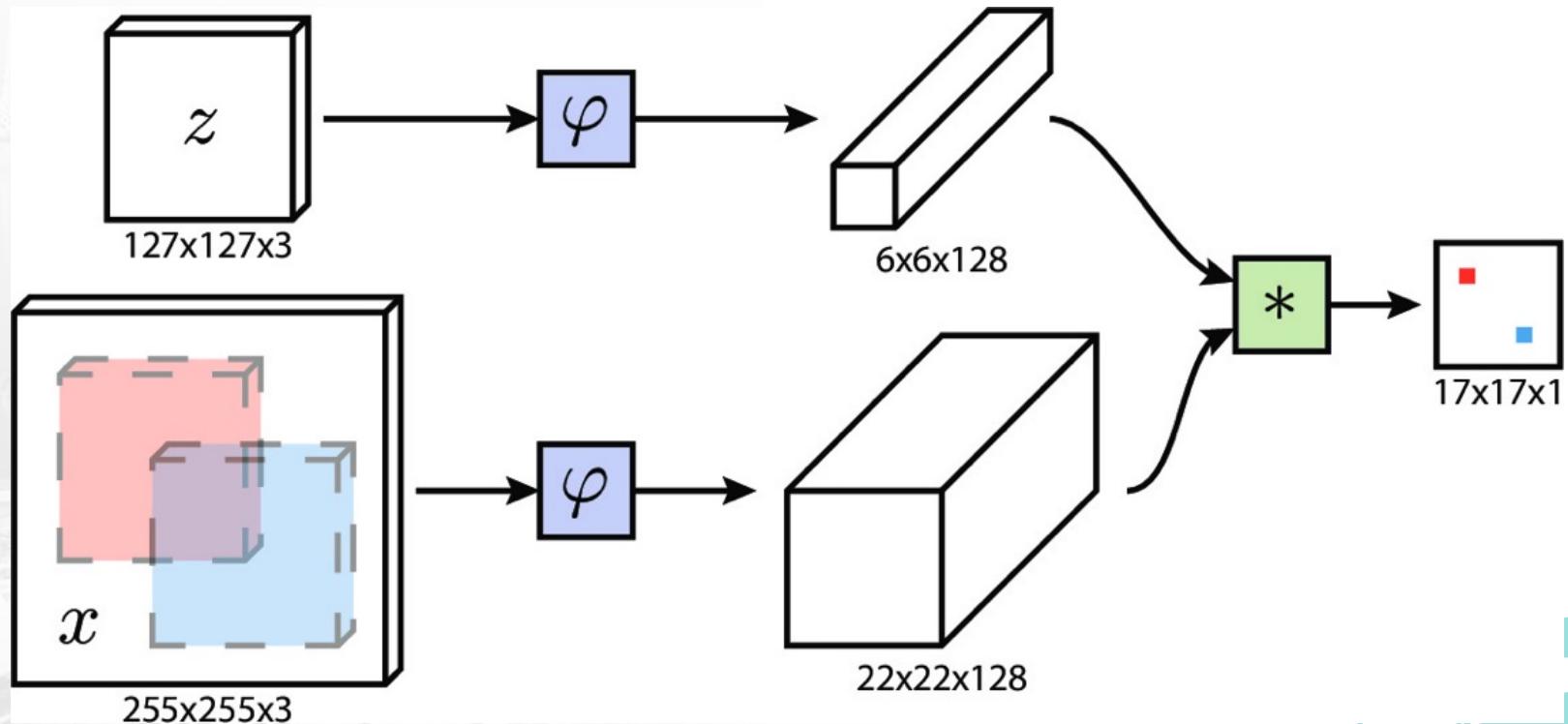


Naeemullah Khan

Object Tracking



- Tracking



Object Tracking – Application



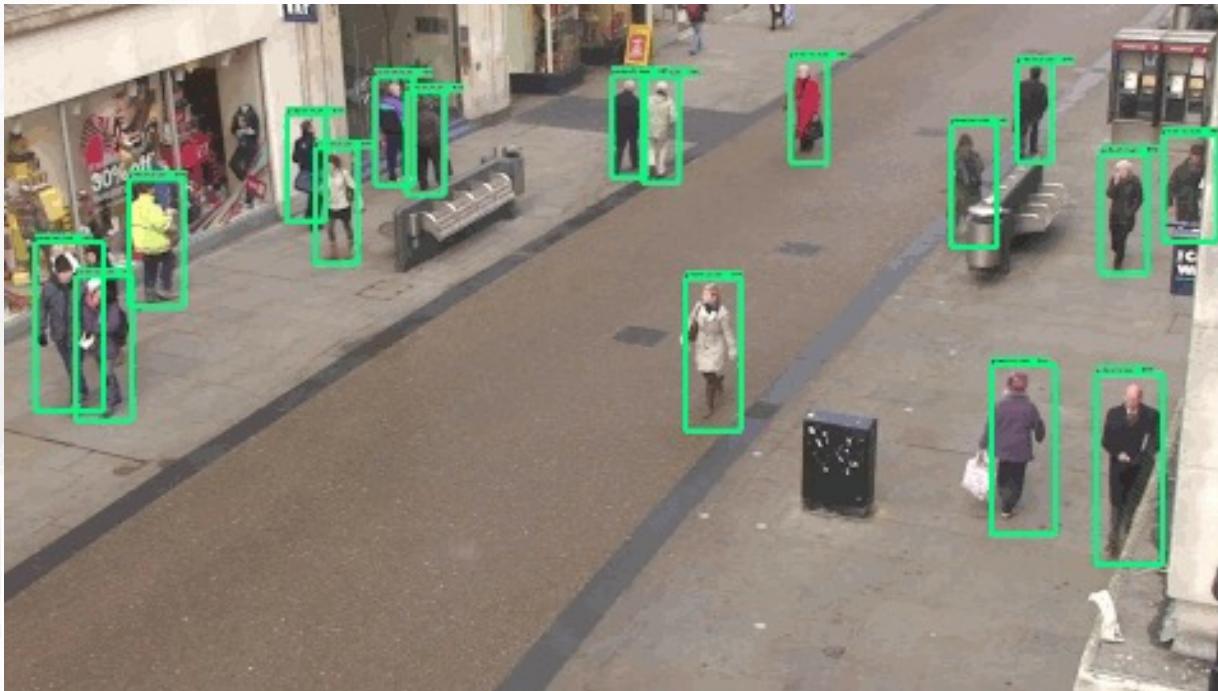
- Tracking



Object Tracking – Application



- People Tracking



Naeemullah Khan

Object Tracking – Application



- UAV tracking



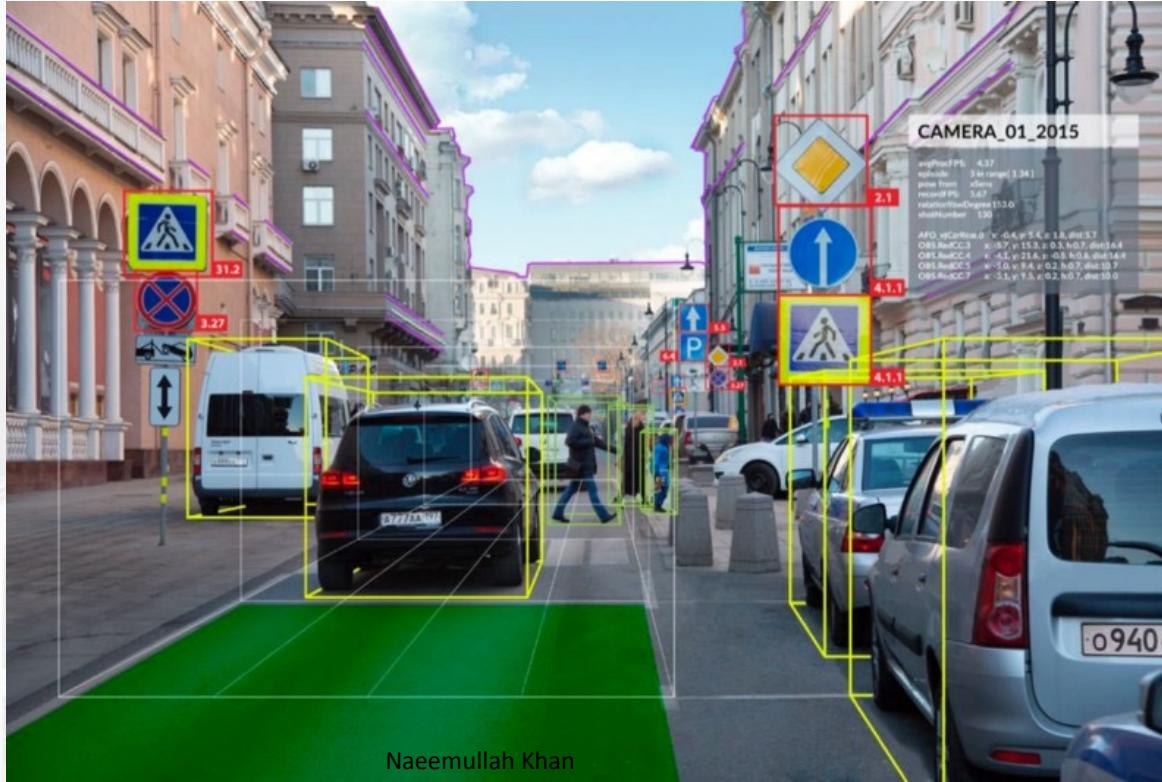
Naeemullah Khan

freetogifmaker.me

112

Computer Vision – Application

- Putting everything together:



ML and AI tools Usage Guide

- Supervised Learning
- Unstructured (Tabular) data Neural Networks
- Grid Data CNNs
- Time-series data (Transformers, RNNs, LSTMs)
- Graph Data (Molecule structure, Social Networks) Graph Neural Networks

Practice Guide

- Use available models architectures
- Data Augmentation
- Fine-Tuning
- Data Collection better than Architecture search

Momentum :

