# From Human Vision to Computer Vision: Rolling in the deep with an image

**Submitted By :**

Nimisha Tiwari

Project Associate

Translational Bioinformatics Group

ICGEB, New Delhi

**Email:** nimisha@icgeb.res.in

# Content

- What is an Image, how does a computer perceives it?

- Intro Google Colab.

- Image Preprocessing.

- The algorithmic story of Convolution Neural Network.

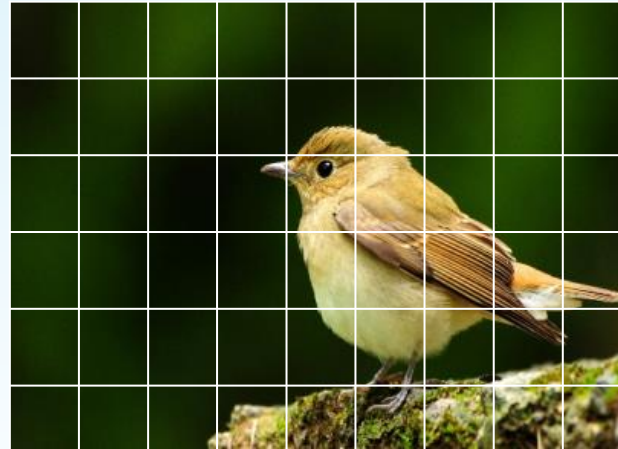- CNN architecture Models : Transfer Learning For Image Classification

**What is an Image ?**

*{What is difference in the vision of a human and a computer}*

**What a human sees**

**What a computer sees**



Color/RGB Image

6x9 pixels

| 28 | 34 | 32 | 30 | 29 | 30 | 31 | 33 | 30 |
|----|----|----|----|----|----|----|----|----|
| 31 | 33 | 32 | 31 | 34 | 31 | 35 | 34 | 31 |
| 32 | 31 | 28 | 29 | 90 | 88 | 79 | 33 | 32 |
| 32 | 30 | 27 | 31 | 99 | 75 | 64 | 48 | 33 |
| 31 | 32 | 29 | 30 | 45 | 68 | 54 | 50 | 36 |
| 30 | 31 | 30 | 52 | 48 | 54 | 55 | 56 | 58 |

- 256 pixels
- 0 – 255
- Black - 0
- White - 255

# Consider learning an image:

➡ Some patterns are much smaller than the whole image

**Can represent a small region with fewer parameters**



"beak" detector

# Convolution Neural Network

- Deep learning explores the possibility of learning features directly from input data, avoiding hand-crafted features.

- A deep net is trained by feeding it input and letting it compute layer-by-layer to generate the final output for comparison with the correct answer.

- After computing the error at the output, this error flows backward through the net by backpropagation.

- At each step backward the model parameters are tuned in a direction that tries to reduce the error.

- helps in model improvement, training is an iterative process that involves multiple passes of the input data until the model converges.

# Convolution Neural Network: CNN

- There are three layers used to build CNN architectures:-

  ➥ Convolutional layer,

  ➥ Pooling layer, and

  ➥ Fully connected layer.

# Convolution

These are the network parameters to be learned

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

6 x 6 image

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

⋮  ⋮

Each filter detects a small pattern (3 x 3).

# Convolution

|   |   |   |
|---|---|---|
| 1 | -1 | -1 |
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

stride=1

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

Dot product →  3    -1

The distance between the applications of filters is called stride.

6 x 6 image

# Convolution

If stride=2

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

6 x 6 image

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

3    -3

Stride hyper parameter is smaller than the filter size the convolution is applied in overlapping windows

# Convolution

stride=1

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

6 x 6 image

| 1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| 3 | -1 | -3 | -1 |
|---|---|---|---|
| -3 | 1 | 0 | -3 |
| -3 | -3 | 0 | 1 |
| 3 | -2 | -2 | -1 |

# Convolution

| | | |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

stride=1

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

6 x 6 image

Repeat this for each filter

| | | | |
|---|---|---|---|
| -1 | -1 | -1 | -1 |
| -1 | | | 1 |
| -1 | -1 | -2 | 1 |
| -1 | 0 | -4 | 3 |

**Feature Map**

Two 4 x 4 images
Forming 2 x 4 x 4 matrix

# Color image: RGB 3 channels

# Convolution v.s. Fully Connected



image

convolution

Fully-connected

$x_1$

$x_2$

$x_{36}$

Filter 1

6 x 6 image

**Fewer Parameters**

1
2
3
4
⋮
7
8
9
10
⋮
13
14
15
16
⋮

3

Only connect to 9 inputs, not fully connected

Filter 1

6 x 6 image

Fewer parameters

Even fewer parameters

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

-1

Shared weights

# Why Pooling

➡ Subsampling pixels will not change the object

bird



Subsampling

bird

We can subsample the pixels to make image smaller

➡ fewer parameters to characterize the image

# Max Pooling

| | | |
|---|---|---|
| 1 | -1 | -1 |
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| | | |
|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 2

| | |
|---|---|
| 3 | -1 |
| -3 | 1 |

| | |
|---|---|
| -3 | -1 |
| 0 | -3 |

| | |
|---|---|
| -3 | -3 |
| 3 | -2 |

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

| | |
|---|---|
| -1 | -1 |
| -1 | -1 |

| | |
|---|---|
| -1 | -1 |
| -2 | 1 |

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

| | |
|---|---|
| -2 | 1 |
| -4 | 3 |

# Max Pooling

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

6 x 6 image

**Conv**

**Max Pooling**

New image but smaller

-1   1

0   3

2 x 2 image

Each filter is a channel

# The whole CNN

# Flattening



Flattened

**Fully Connected Feedforward network**

# CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D tensor)*

input

```
model2.add( Convolution2D( 25,3,3,
            input_shape=(28,28,1)) )
```

There are 25 3x3 filters.

Input_shape = ( 28 , 28 , 1)

28 x 28 pixels          1: black/white, 3: RGB

```
model2.add(MaxPooling2D((2,2)))
```

| 3 | -1 |
|---|----|
| -3 | 1 |

→ 3

**Convolution**

**Max Pooling**

**Convolution**

**Max Pooling**

# CNN in Keras

Input

1 x 28 x 28

Convolution

```
model2.add( Convolution2D( 25,3,3,
          input_shape=(28,28,1)) )
```

How many parameters for each filter?

**25x9 = 225**

25 x 26 x 26

Max Pooling

```
model2.add(MaxPooling2D((2,2)))
```

25 x 13 x 13

Convolution

```
model2.add(Convolution2D(50,3,3))
```

50 x 11 x 11

Max Pooling

```
model2.add(MaxPooling2D((2,2)))
```

50 x 5 x 5

# CNN Architectures

## LeNet-5
[LeCun et al., 1998]



Input
Image Maps
Convolutions
Subsampling
Fully Connected
Output

- Conv filters were 5x5, applied at stride 1

- Subsampling (Pooling) layers were 2x2 applied at stride 2

- i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

# AlexNet

[Krizhevsky et al. 2012]

**Architecture:**
CONV1
MAX POOL1
NORM1
CONV2
MAX POOL2
NORM2
CONV3
CONV4
CONV5
Max POOL3
FC6
FC7
FC8

## What is ImageNet DataSet

- It is a large dataset of annotated photographs.
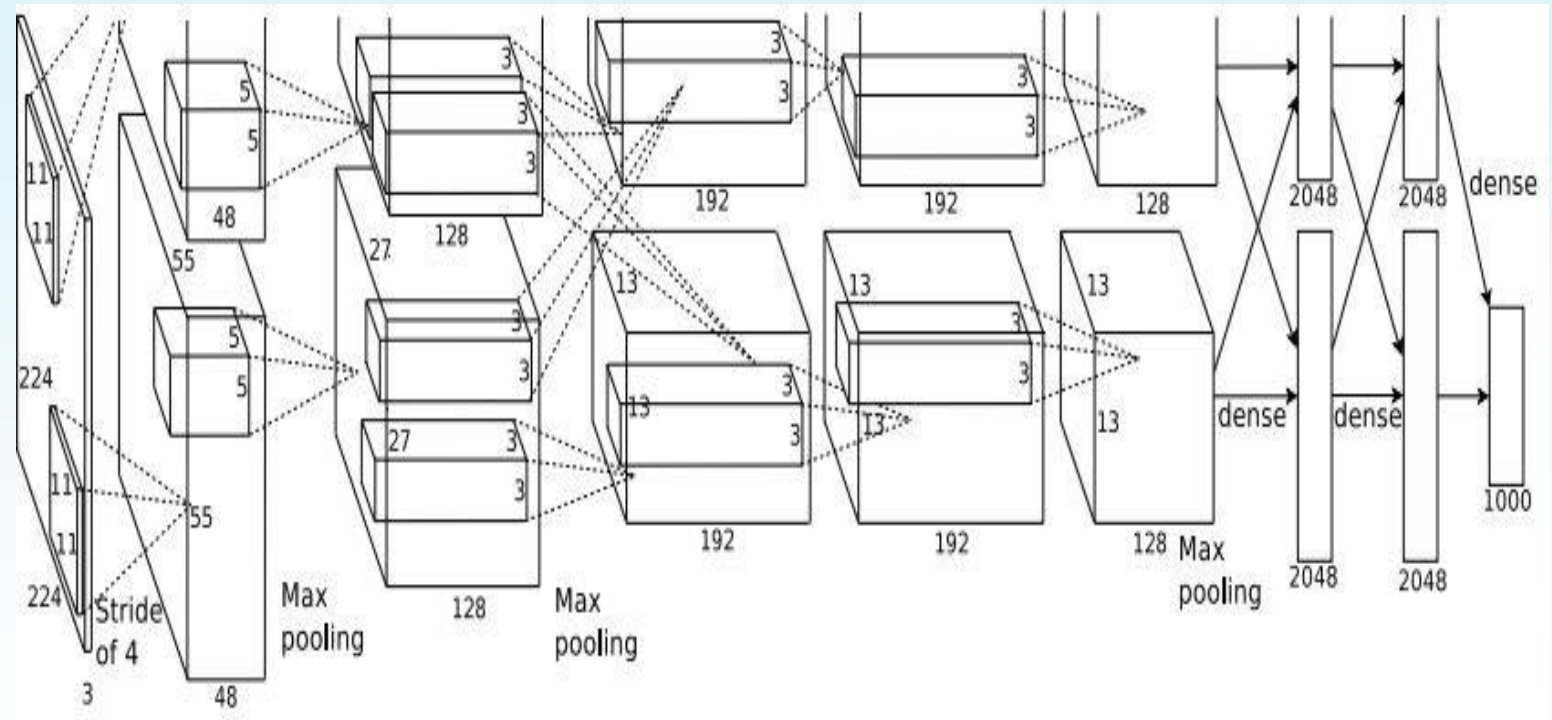- This dataset consists of about 14 million images, more than 21000 groups or classes and more than 1 million images that have bounding box annotation.

## ImageNet large scale visual recognition Challenge (ILSVRC)

- ImageNet large scale visual recognition Challenge for short ILSVRC. The goal of this challenge is to train a model that can correctly classify an image into a class out of 1000 separate object categories.

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

- AlexNet
- GoogLeNet (e.g. InceptionN),
-  VGGNet (e.g. VGG16 or VGG19),
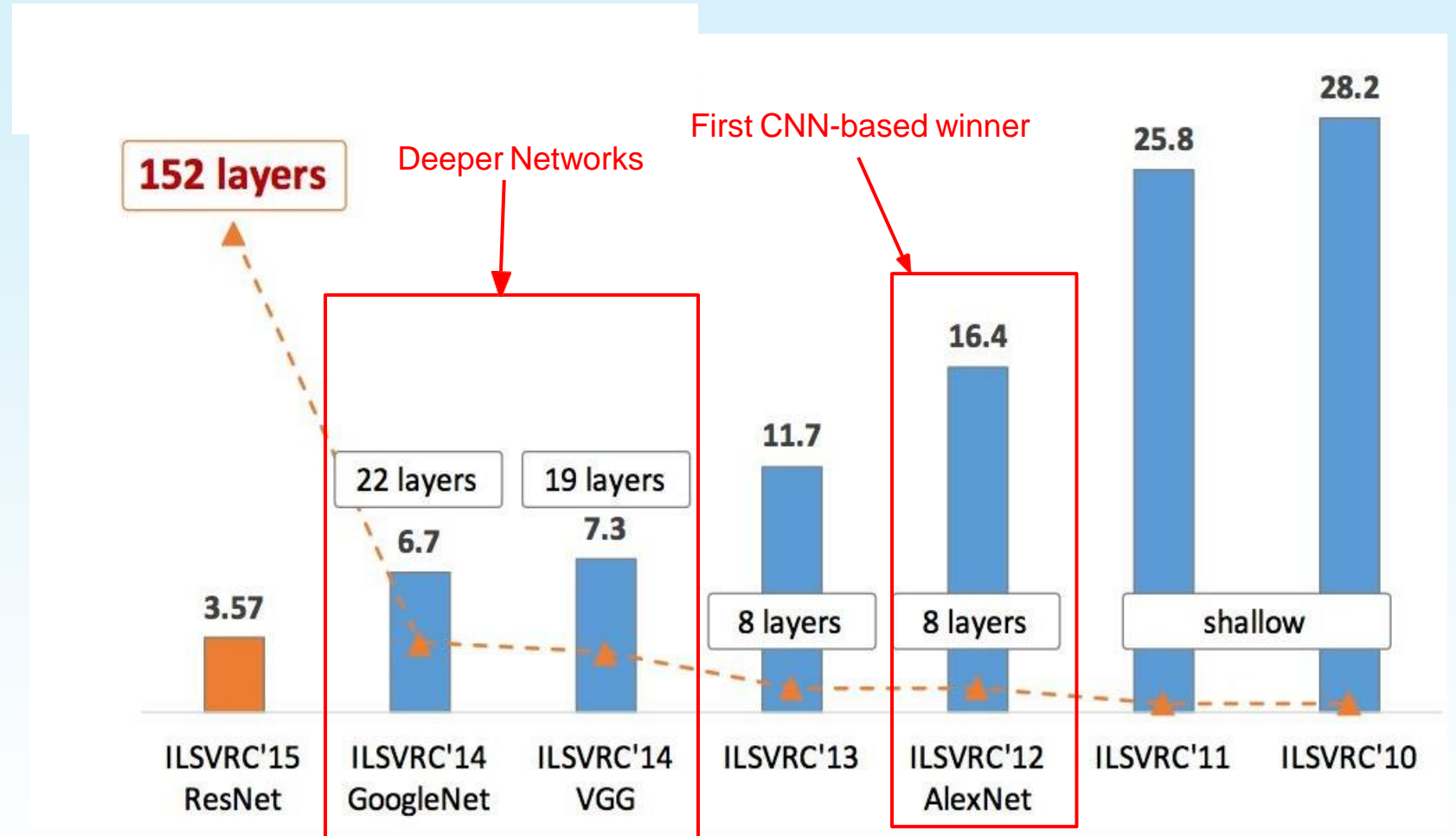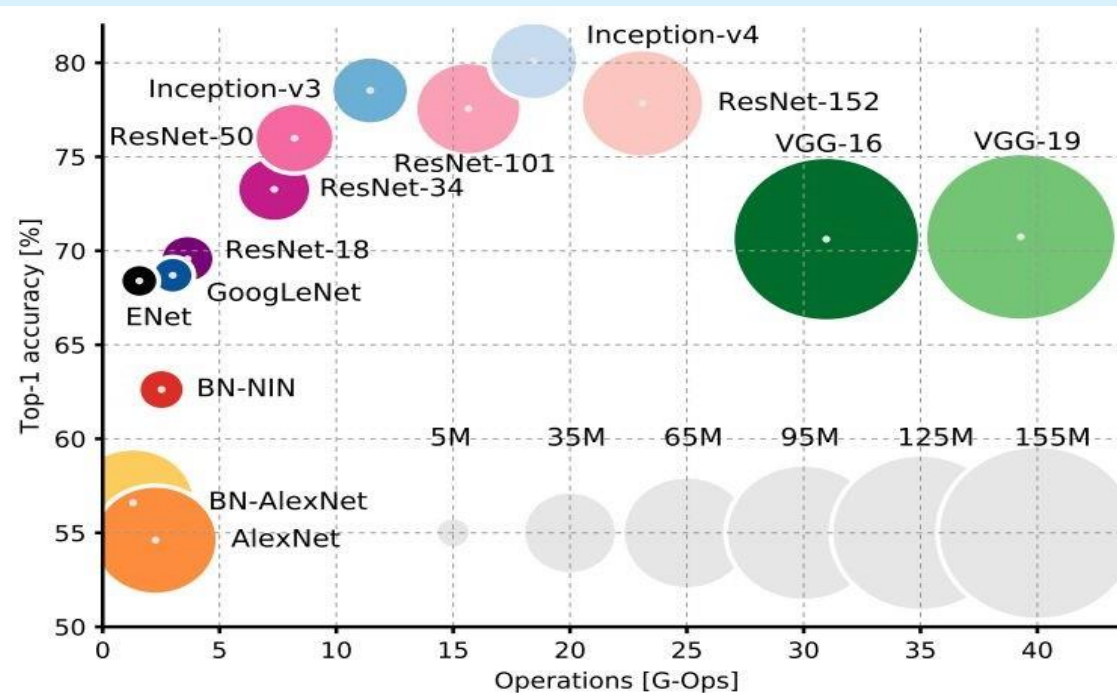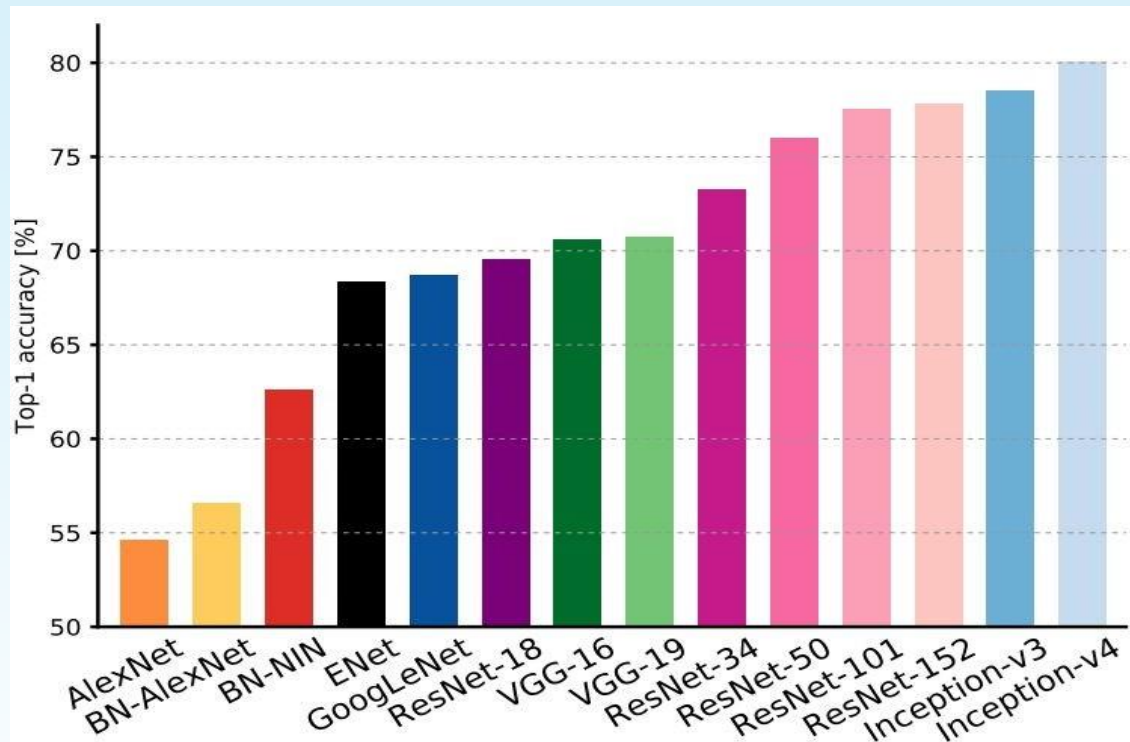- Residual Network (e.g. ResNetN)



Fig: ILSVRC winners (with training error % vs Years)

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

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

VGG: Highest memory, most operations.
GoogLeNet: Most efficient.
AlexNet: Smaller compute, still memory heavy, lower accuracy.
ResNet: Moderate efficiency depending on model, highest accuracy

# Transfer Learning Model

- Transfer learning is a technique whereby a neural network model is first trained on a problem similar to the problem that is being solved.

- One or more layers from the trained model are then used in a new model hence, transfer learning is a method of reusing a pre-trained model knowledge for another task.

- There are perhaps a dozen or more top-performing models for image recognition that can be used. AlexNet, ResNet, VGG, and Inception etc. are some of the CNN based transfer learning model.

34

Thank You