



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 8 Issue: IX Month of publication: September 2020

DOI: <https://doi.org/10.22214/ijraset.2020.31693>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Comparative Analysis of different Convolutional Neural Network Algorithm for Image Classification

H. Agrawal¹, M. Kalantri², A. Bansal³

^{1,2}Student, Department of Computer Science and Information Technology, Symbiosis University of Applied Sciences, Indore, Madhya Pradesh, India

³Professor, Department of Computer Science and Information Technology, Symbiosis University of Applied Sciences, Indore, Madhya Pradesh, India

Abstract: In today's fast and furious world where new techniques and models are being developed every day for the sake of increasing the efficiency and performance of the neural network, this research conducts an in-depth study about some of the most popular and important Convolutional neural network models. This paper contains a detailed study and data-rich analysis of the 5 most popular Convolutional Neural networks (CNNs) for Image Detection and Identification. These 5 CNNs are LeNet, AlexNet, VGGNet16, ResNet50, and GoogLeNet. The performance of these neural networks is evaluated and benchmarked using well known and most commonly used Cifar10 and Cifar100 datasets. This research aims to make the process of understanding different neural networks and working with them easy. With rigorous training on the high end and graphics enabled machine for several months continuously the data and information gathered have been compiled in this research paper with all the obligatory information required to comprehend Convolutional Neural Networks.

Keywords: Neural Network; Deep Learning; Perceptron; Convolutional Neural Network; Object detection; Object classification

I. INTRODUCTION

In today's tech-savvy and automated world, where most of the work and functioning of mankind has become dependent on technology and its efficaciousness, Artificial Intelligence is picking up the heat to furthermore augment these technologies in order to achieve a world where the requirement of human interaction is negligible. Artificial Intelligence is "the theory and development of computer systems able to perform tasks normally requiring human intelligence, such as visual perception, speech recognition, decision-making, and translation between languages". One of the most prominent tasks in which Artificial Intelligence is being used nowadays is of "Image/Object detection and Recognition" and to perform these tasks there is a subdomain of Artificial Intelligence that comes into play which Deep learning and Deep Neural networks [DNN]. A DNN is an artificial neural network having multiple hidden layers. It propagates through the network and finds the correct mathematical manipulation to turn the input into the output. This research paper presents an in-depth study of a subdomain of DNN known as Convolutional Neural Networks [CNN]. CNN has been presenting an operative class of models for better understanding of contents present in an image, therefore resulting in better image recognition, segmentation, detection, and retrieval. CNN's are efficiently and effectively used in many patterns and image recognition applications, for example, gesture recognition, face recognition, object classification, and generating scene descriptions. Due to advances and development in learning algorithms for deep network construction and moderately to the open-source large labelled data set available for experimentation purpose, the successful integration of all the stated applications is possible. CNN has well known trained networks that use datasets like CIFAR 10, 100, MNIST etc. which are available in open-source networks and increases its efficacy of classification after getting trained over millions of images contained in the datasets of CIFAR-10 and CIFAR-100. The datasets used are composed of millions of tiny images. Therefore, they can simplify well and accurate and hence successfully categorize the classes' out-of-sample examples. When comparisons are made on a large data set such as CIFAR-10, 100 etc., The neural network classification and prediction accuracy and error rates are all most comparable to that of humans This work aims at analyzing the capability of convolutional neural networks to categories the scene in videos on the basis of identified objects. A variety of image categories are included in CIFAR- 100 and CIFAR 10 datasets for training the CNN. The test datasets are videos of different categories and subjects. The contradiction branches out because of the feature extraction capabilities of different CNN. Talking about different CNN's, in this research 5 types of networks are used. These networks used for our study are constructed using existing neural networks namely LeNET, AlexNet, VGGNET-16, ResNET50 and GoogLeNET and each of these networks have different layers, therefore their performance varies considerably. The paper starts with the Literature survey and giving insights about the concepts used and progresses further by representing and explaining the results obtained and concluded by discussing learning outcomes.

II. LITERATURE SURVEY

The different architectures of Convolutional Neural Network are attained by using deep learning methodologies of which, the perceptron is the basic unit.

A. Perceptron

The single artificial neuron, a perceptron is a simulation of biological neurons. The observation in Fig. 1 shows the analogy between Biological and Artificial neurons:

- 1) The small circles denoted by (1) are the data inputs $x(i)$
- 2) Arrows denoted by (2) are the synapse, used to send inputs multiplied with weights, $w(i)$. These weights can be assumed as the measure of 'importance' of a particular input data, $w(i) x(i)$.
- 3) Now, for better representation of data, bias is added W to the sum of $w(i) x(i)$ ($0 \leq i \leq n$). Bias helps in shifting of a graph left, right for the better fit (representation) of data points. These operations are performed in (3), the actual cell of a perceptron.
- 4) The cell, denoted by (3) will produce output as a single value. In case of more than one output requirement, we apply a function to the output of the cell, (3). This function is used to produce distinct output according to the given input. These functions are known as 'activation functions', denoted by (4) in Fig. 1.
- 5) After multiplying weights, adding bias and applying activation functions to the input given, the neuron produces output. Output, denoted by (5), is given on the grounds of a number of distinct classes.

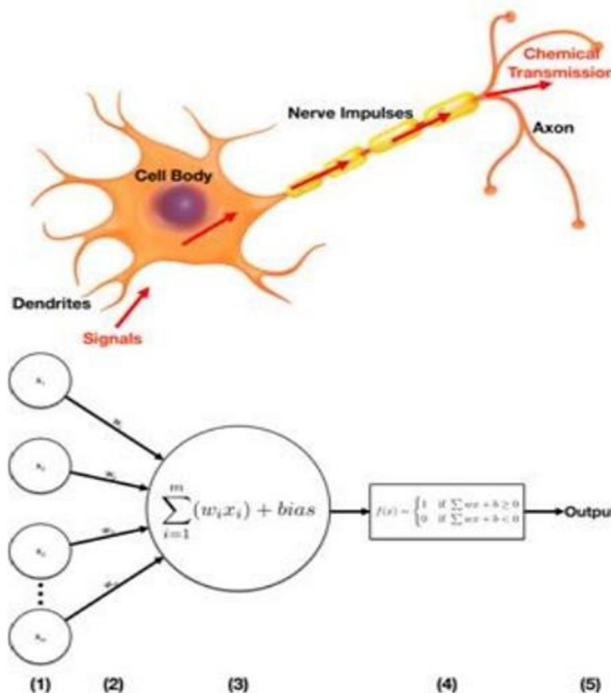


Fig. 1. Perception and Biological Neuron

B. Deep Learning

Deep learning is an artificial intelligence function that imitates the workings of the human brain in processing data and creating patterns for use in decision making. Deep learning is a subset of machine learning in artificial intelligence (AI) that has networks capable of learning unsupervised from data that is unstructured or unlabeled. Also known as deep neural learning or deep neural network. Deep learning has evolved hand-in-hand with the digital era, which has brought about an explosion of data in all forms and from every region of the world. This data, known simply as big data, is drawn from sources like social media, internet search engines, e-commerce platforms, and online cinemas, among others. This enormous amount of data is readily accessible and can be shared through fintech applications like cloud computing. However, the data, which normally is unstructured, is so vast that it could take decades for humans to comprehend it and extract relevant information. Companies realize the incredible potential that can result from unravelling this wealth of information and are increasingly adapting to AI systems for automated support.

Various use cases of Deep Learning mainly include:

- 1) Face recognition and identification.
- 2) Handwritten text recognition
- 3) Speech recognition
- 4) Language translation
- 5) Building game strategies
- 6) Control and drive self-driving cars.
- 7) Natural Language Processing

Not only these, but there are several other areas where Deep Learning models are challenging human capabilities and conventional methods.

C. Neural Networks

A neural network is a network or circuit of neurons or in a modern sense, an artificial neural network, composed of artificial neurons or nodes. Thus, a neural network is either a biological neural network, made up of real biological neurons, or an artificial neural network, for solving artificial intelligence (AI) problems. The connections of the biological neuron are modelled as weights. All inputs are modified by weight and summed. Finally, an activation function controls the amplitude of the output. For example, an acceptable range of output is usually between 0 and 1, or it could be -1 and 1. These artificial networks may be used for predictive modelling, adaptive control and applications where they can be trained via a dataset. Self-learning resulting from experience can occur within networks, which can derive conclusions from a complex and seemingly unrelated set of information.

1) Types of Neural Network

- a) Feedforward Neural Network – Artificial Neuron
- b) Radial basis function Neural Network
- c) Multilayer Perceptron
- d) Recurrent Neural Network (RNN) – Long Short-Term Memory
- e) Convolutional Neural Network
- f) Modular Neural Network

D. Convolutional Neural Networks

A convolutional neural network (CNN, or ConvNet) is a class of deep neural networks, most commonly applied to analyzing visual imagery. They have applications in image and video recognition, recommender systems, image classification, medical image analysis, natural language processing, and financial time series. Traditional neural networks are not great at image processing and must be fed images in reduced-resolution pieces. CNN has their “neurons” arranged more like the area responsible for processing visual stimuli in humans and other animals. The layers of neurons are arranged in a way that covers the entire visual field avoiding the piecemeal image processing problem of traditional neural networks.

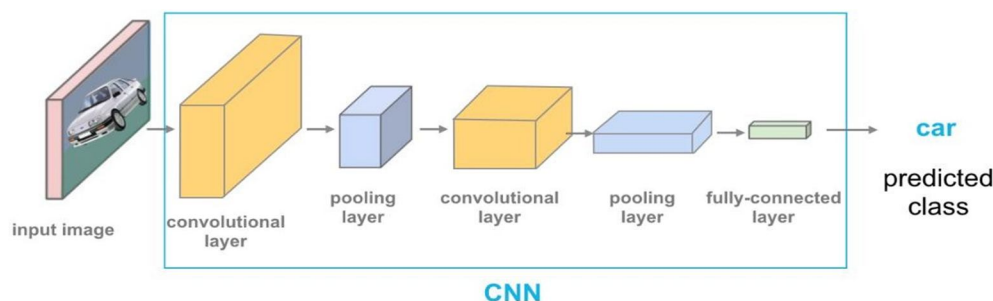


Fig. 2. The Architecture of a typical CNN

CNN uses a system like a multilayer perceptron designed for reduced processing requirements. A simple ConvNet is a sequence of layers, and every layer of a ConvNet transforms one volume of activations to another through a differentiable function. The sequence of layers of a CNN consist of an input layer, an output layer and a hidden layer that includes multiple convolutional layers, pooling layers, fully connected layers and normalization layers.

A convolutional layer contains a set of filters whose parameters need to be learned. The height and weight of the filters are smaller than those of the input volume. Each filter is convolved with the input volume to compute an activation map made of neurons. The output volume of the convolutional layer is obtained by stacking the activation maps of all filters along the depth dimension.

A pooling layer is usually incorporated between two successive convolutional layers. The pooling layer reduces the number of parameters and computation by down-sampling the representation. The pooling function can be max or average. Max pooling is commonly used as it works better. Fully connected layers in a neural network are those layers where all the inputs from one layer are connected to every activation unit of the next layer.

Normalization is an approach which is applied during the preparation of data in order to change the values of numeric columns in a dataset to use a common scale when the features in the data have different ranges. In this article, we will discuss the various normalization methods which can be used in deep learning models. There are many types of normalization layers like Batch normalization and layer normalization.

Neurons are the units which take input, perform some functions and pass on the output to another neuron. The functions used in the neuron are called as activation functions. There are several activations functions like Step function, Sigmoid function, Tanh and Relu Activation function. The removal of limitations and increase in efficiency for image processing results in a system that is far more effective and simpler to train.

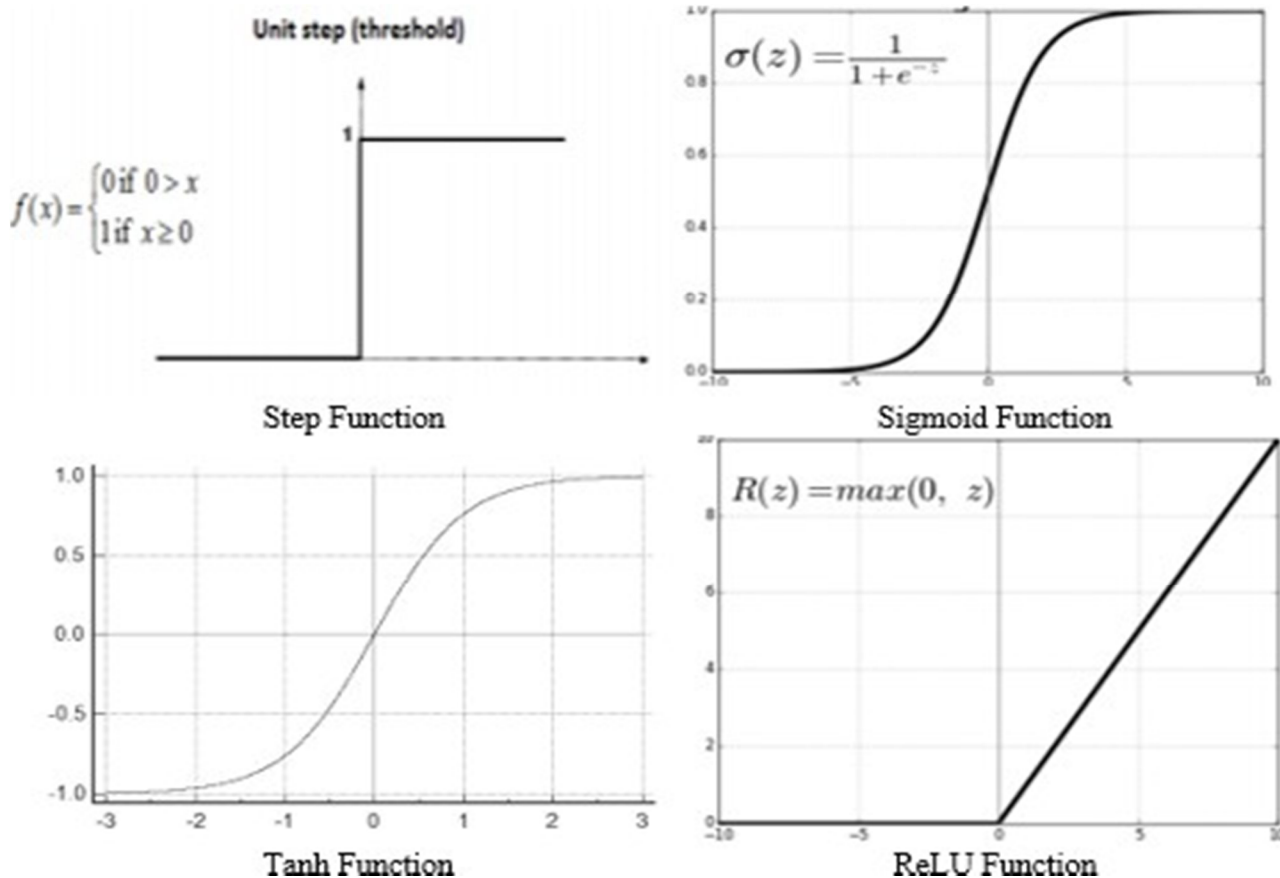


Fig. 3. Activation functions commonly used in CNN

III.ARCHITECTURE

The different architectures studied are: LeNet -5, AlexNet, VGGNet16, ResNet50 and GoogLeNet.

A. LeNet – 5

LeNet-5, a pioneering 7-level convolutional network by LeCun et al in 1998, that classifies digits, was applied by several banks to recognize hand-written numbers on checks (cheques) digitized in 32x32 pixel greyscale input images. The ability to process higher-resolution images requires larger and more convolutional layers, so this technique is constrained by the availability of computing resources.

The LeNet-5 architecture consists of two sets of convolutional and average pooling layers, followed by a flattening convolutional layer, then two fully-connected layers and finally a softmax classifier. The input for LeNet-5 is a 32×32 grayscale image which is passed through the LeNet-5 model. The output is softmax output layer \hat{y} with 10 possible values corresponding to the digits from 0 to 9. The total number of parameters is around 60,000.

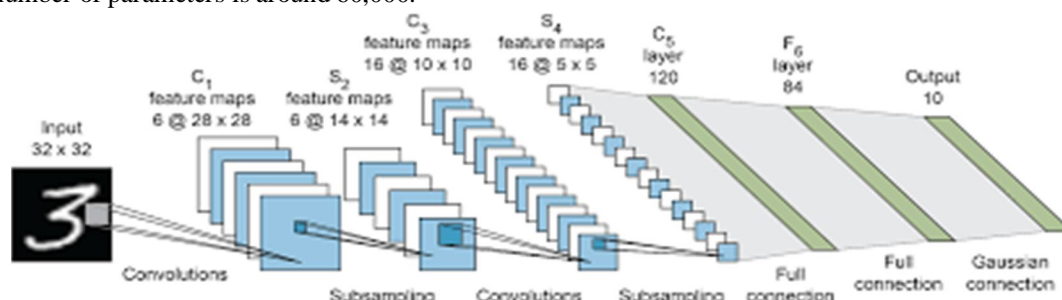


Fig. 4. LeNet – 5 architecture

B. AlexNet

AlexNet developed by Krizhevsky et al in 2012 and had a large impact on the field of machine learning, specifically in the application of deep learning to machine vision. It famously won the 2012 ImageNet LSVRC-2012 competition by a large margin (15.3% VS 26.2% (second place) error rates). The network was very similar to LeNet but was much deeper and had around 60 million parameters. The input was an RGB image of size 227×227 and the output a softmax layer predicting values for each class. The key points in the success of AlexNet are that it used Relu activation function is used instead of Tanh to add non-linearity which accelerates the speed by 6 times at the same accuracy. Dropout was used instead of regularization to deal with overfitting. However, the training time is doubled with a dropout rate of 0.5. It also overlapped pooling to reduce the size of the network.

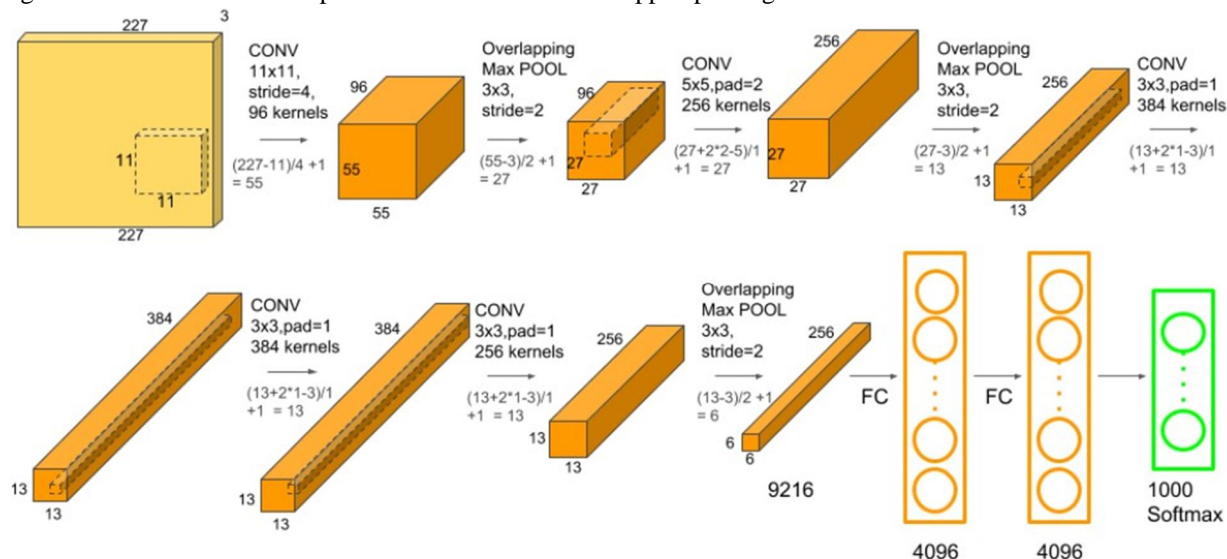


Fig. 5. AlexNet architecture

C. VGGNet16

VGGNet was developed by Simonyan et al and was the runner up of 2014 Imagenet challenge. It consists of 16 convolutional layers and because of the simplicity of its uniform architecture, it appeals to a new-comer as a simpler form of a deep convolutional neural network. This network is one of the most used choices for feature extraction from images (taking images and converting them to a smaller dimensional array that contains important information regarding the image). The weight configuration of the VGGNet-16 is publicly available and has been used in many other applications and challenges as a baseline feature extractor. VGGNet consists of 138 million parameters, which can be a bit challenging to handle. The input image was an RGB image of 224×224 pixels and the total parameters are around 138 million.

VGGNet has 2 simple rules of thumb to be followed:

- 1) Each Convolutional layer has configuration — kernel size = 3×3 , stride = 1×1 , padding = same. The only thing that differs is the number of filters.
- 2) Each Max Pooling layer has configuration — windows size = 2×2 and stride = 2×2 . Thus, we half the size of the image at every Pooling layer.

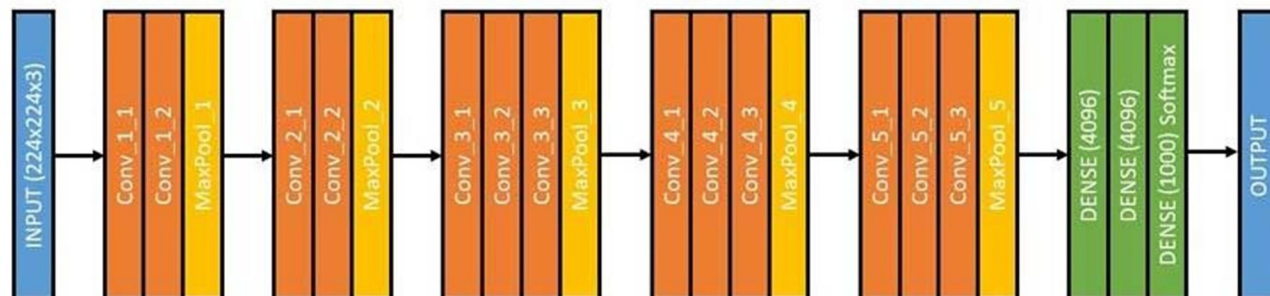


Fig. 6. VggNet-16 architecture

D. ResNet50

ResNet introduced by Kaiming et al in 2015 imagenet competition brought about a top-5 error rate of 3.57%, which is lower than the human error on top-5. The “50” refers to the number of layers it has. The network introduced a novel approach called — “skip connections”. Ideally, deeper models should perform better than shallower ones but that was not the case. Deep networks often suffer from a problem of vanishing gradients, i.e.: as the model backpropagates, the gradient gets smaller and smaller. These tiny gradients can make learning intractable.

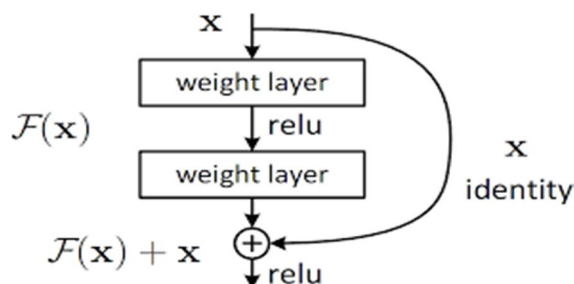


Fig.7. Identity Block

The skip connection in the diagram above is labelled “identity.” It allows the network to learn the identity function, which allows it to pass the input through the block without passing through the other weight layers. This allows you to stack additional layers and build a deeper network, offsetting the vanishing gradient by allowing your network to skip through layers of it feels they are less relevant in training.

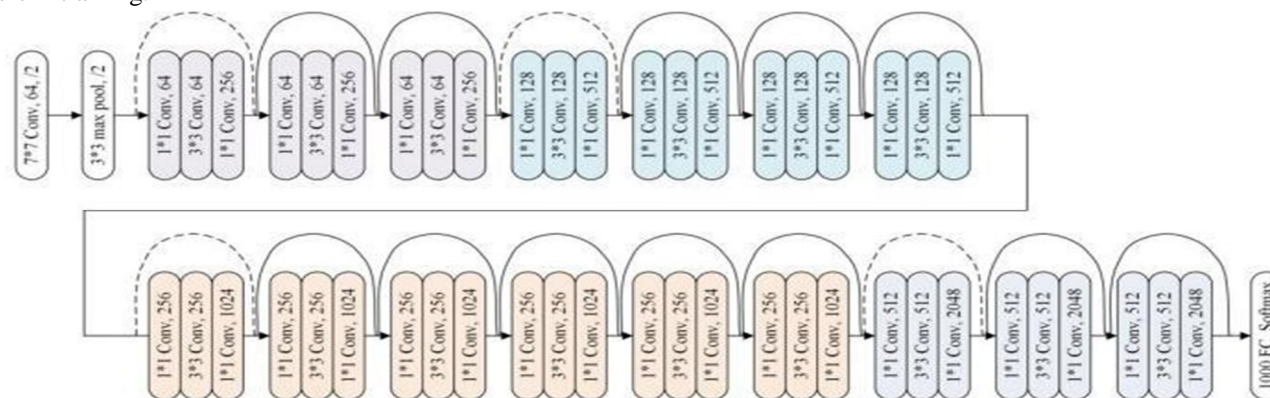


Fig. 8. ResNet-50 architecture

E. GoogLeNet

GoogLeNet or Inception developed by Szegedy et al was the winner of the 2014 ImageNet competition. It used an inception module, a novel concept, with smaller convolutions that allowed the reduction of the number of parameters to a mere 4million.

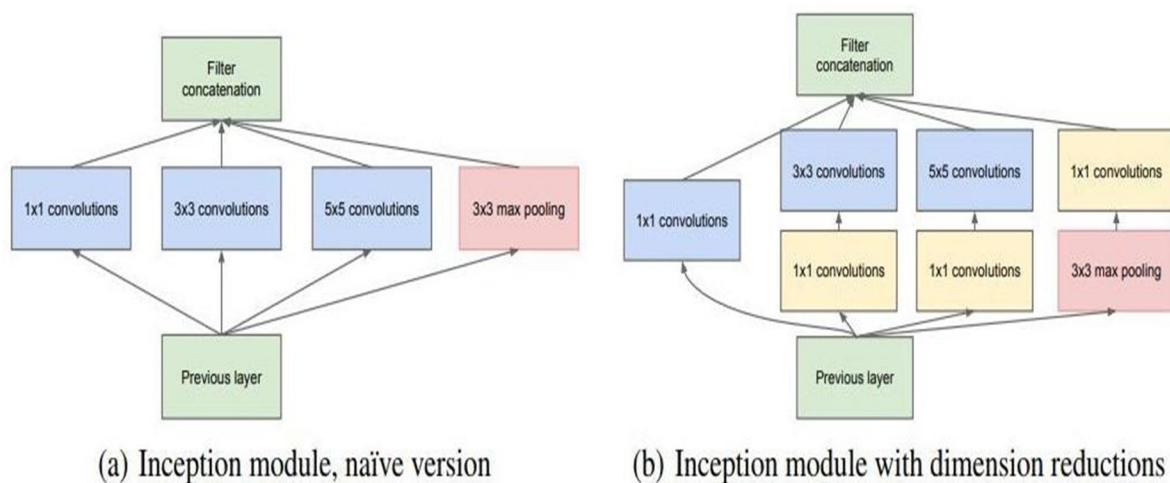


Fig. 9. Inception modules

There are 22 layers in total. It is already a very deep model compared with previous AlexNet, Auxiliary classifiers are also used in this model that is connected to the intermediate layers which are expected to encourage discrimination in the lower stages in the classifier, increase the gradient signal that gets propagated back, and provide additional regularization. These auxiliary classifiers are used during training only.

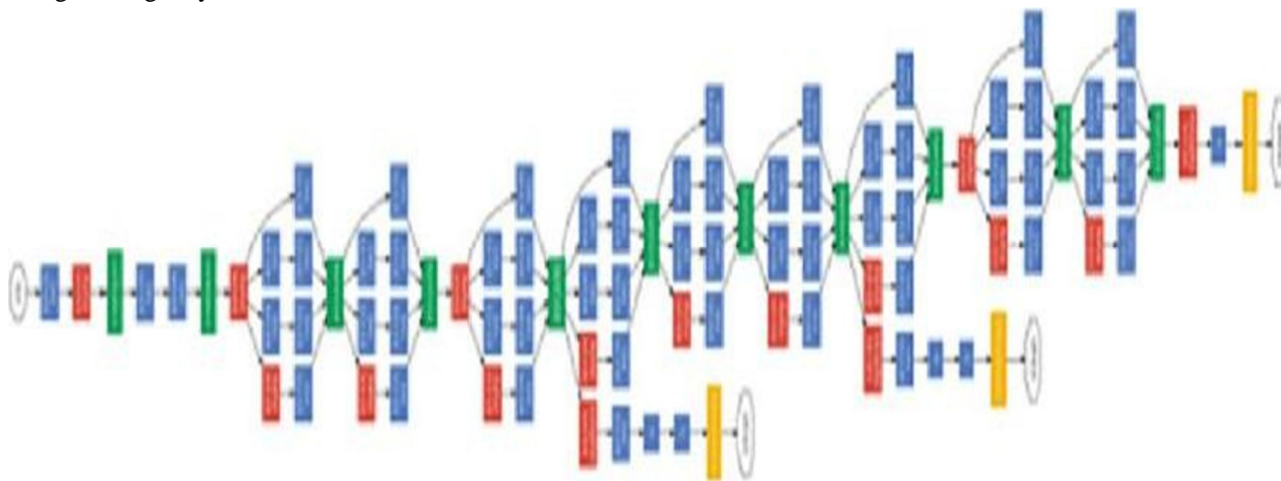


Fig. 10. GoogLeNet architecture

IV.DATASET

Image dataset of CIFAR- 100 which has numerous super-classes of general object images and several subclass categories of each superclass. CIFAR-100 has 100 classes of images with each class having 600 images each. These 600 images are divided into 500 training images and 100 testing images for each class, therefore, making a total of 60,000 different images. These 100 classes are clubbed together into 20 super-classes. The second dataset used was CIFAR-10. The dataset is a collection of images that are commonly used to train machine learning and computer vision algorithms. It is one of the most widely used datasets for machine learning research. The CIFAR-10 dataset contains 60,000 32x32 color images in 10 different classes. The 10 different classes represent airplanes, cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks. There are 6,000 images of each class. CIFAR-10 is a set of images that can be used to teach a computer how to recognize objects. Since the images in CIFAR-10 are low-resolution (32x32), this dataset allows to quickly try different algorithms to see what works. CIFAR-10 is a labelled subset of the 80 million tiny images dataset.

V. PROCESS AND METHODOLOGY

The flow of the research includes the followings steps:

- A. The first step includes Data. Data is the soul of every Machine Learning, Deep Learning project. For every project, data is divided into two parts. One, training set, another, test set.
- B. Both the training and test set aren't decided by the user but selected randomly. This is done to avoid the overfitting or underfitting of the model. This is the most important step which decides the working of the whole project and model.
- C. We use only the training set further to build the model. This is to avoid overfitting of the model and test the model further on another set.
- D. After this, we write a deep learning algorithm. This algorithm depends on the type of data, type of output, number of classes and various other factors.
- E. As soon as we develop an algorithm, we evaluate it using a test set. This is a measure of accuracy, input data, parameters and memory used. After this, rigorous testing of the model starts.
- F. After training and testing the CNN on both the datasets for multiple epochs and with repetitive calibrating and hyper tuning, the final test set is executed.
- G. The results are recorded and comparative analysis of all the 5 models is done.

VI. RESULTS

The performance analysis of CNN's is done by testing each of the networks on CIFAR-100 and CIFAR-10 datasets.

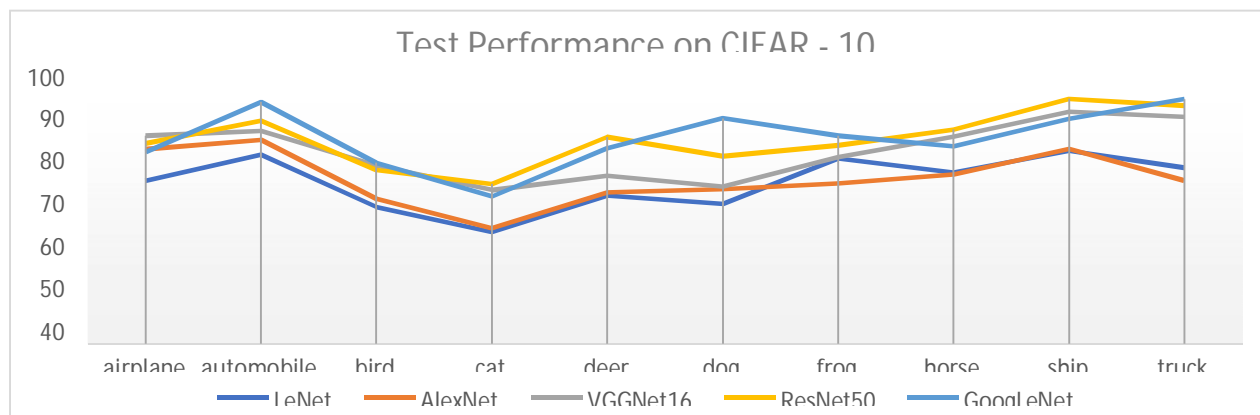
A. For CIFAR-10

From table 1 it can be concluded that objects present in the "automobile", "ship" and "Truck" category are classified efficiently by all the Networks. Whereas, the "cat" category was the most poorly predicted. Moreover, "Horse" and "Airplane" were predicted with accuracy well above 60%. Accuracy of "Bird", "deer" and "Dog" categories varied from one network to the other.

CLASS	LeNet	AlexNet	VGGNet16	ResNet50	GoogLeNet
airplane	60.3	72	77	74.1	70.9
automobile	69.9	75.4	78.8	82.6	89.4
bird	50.7	53.6	65.9	64.2	66.9
cat	41.3	42.8	56.9	59.2	54.5
deer	54.9	56.1	62.2	76.5	72.3
dog	51.8	57.1	58.2	69.4	83.4
frog	68.6	59.4	69	73.4	77
horse	63.4	62.6	76.6	79.3	73.1
ship	71.4	72.1	85.9	90.6	83.3
truck	65.1	60.4	84	88.1	90.5

Table 1. Performance of CNN's on CIFAR-10

GoogLeNet and ResNet50 performed very well on almost all the categories with VGGNet16 lacking behind a bit. Whereas, AlexNet and LeNet because of having a lesser number of layers and low training optimization lacked behind the other three by a very visible margin.

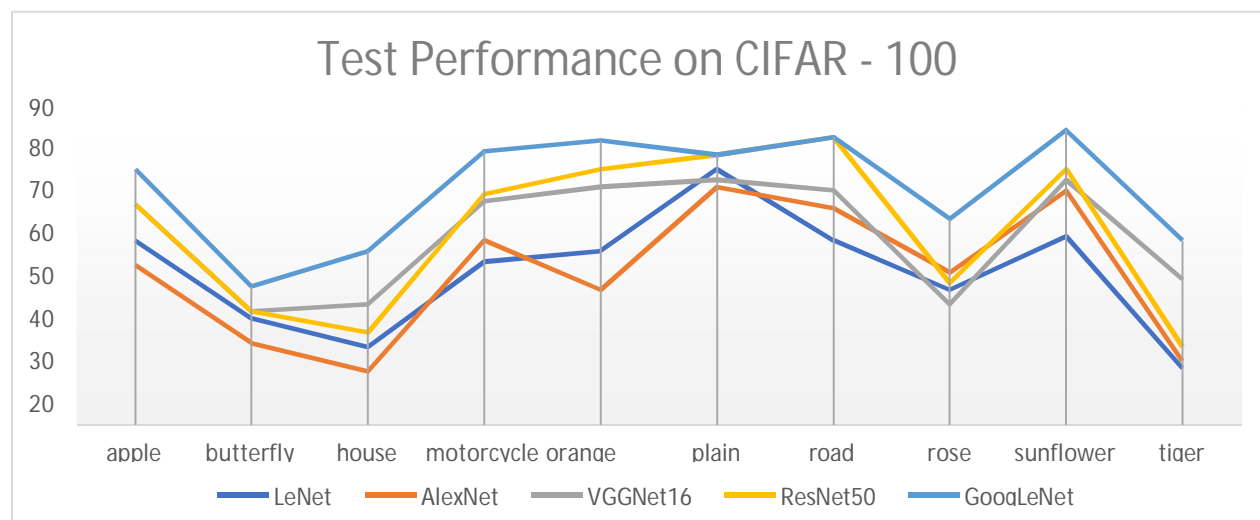


B. For CIFAR-100

From table 2 it can be concluded that objects present in the “sunflower”, “road” and “orange” category are classified efficiently by all the Networks. Whereas, the “House” category was the most poorly predicted. Moreover, “Motorcycle”, “Apple” and “Plain” were predicted with accuracy well above 60%. Accuracy of “Rose”, “Tiger” and “Butterfly” categories varied from one network to the other. GoogLeNet and ResNet50 performed very well on almost all the categories with VGGNet16 lacking behind a bit. Whereas, AlexNet and LeNet because of having a lesser number of layers and low training optimization lacked behind the other three by a very visible margin.

CLASS	LeNet	AlexNet	VGGNet16	ResNet50	GoogLeNet
apple	52	45	62	62	72
butterfly	30	23	32	32	39
house	22	15	34	26	49
motorcycle	46	52	63	65	77
orange	49	38	67	72	80
plain	72	67	69	76	76
road	52	61	66	81	81
rose	38	43	34	40	58
sunflower	53	66	69	72	83
tiger	16	18	41	22	52

Table 2. Performance of CNN's on some categories of CIFAR-100



C. For Live Video Feeds

The real-time analysis of convolutional neural networks shows that the performance of CNN's on images vary substantially in live testing results. In live testing, CNNs get confused between few objects. The overall accuracy for LeNet, AlexNet, VGGNet, ResNet and GoogLeNet models were 14.56%, 18.32%, 31.88%, 39.35% and 56.19% respectively. The results proved that the accuracy of GoogLeNet and ResNet was pretty good compared to other models.

CATEGORY	LeNet	AlexNet	VGGNet16	ResNet50	GoogLeNet
airplane	15	18	20	39	52
automobile	19	17	21	34	59
bird	10	15	16	36	43
cat	12	20	24	41	63
deer	16	21	26	43	60
dog	9	19	19	51	72
frog	18	11	27	26	44
horse	13	14	29	37	58
ship	21	13	23	40	65
truck	11	18	21	33	46

Table 3. Performance of CNN's trained on CIFAR-10 on Live Video Feeds

CATEGORY	LeNet	AlexNet	VGGNet16	ResNet50	GoogLeNet
apple	25	21	29	41	45
butterfly	16	24	22	34	38
house	15	19	25	55	65
motorcycle	22	23	34	36	42
orange	18	19	27	29	36
plain	12	15	19	34	40
road	16	19	18	43	62
rose	19	18	31	30	39
sunflower	23	22	33	38	51
tiger	13	17	20	31	68

Table 4. Performance of CNN's trained on CIFAR-100 on Live Video Feeds for the same categories as in Table 2.

VII. CONCLUSION

The Research done evaluated and analyzed the accuracy of prediction of 5 different types of Convolutional Neural Networks on the most popular datasets for training and testing namely CIFAR-10 and CIFAR-100. The main focus of this research was to find out the accuracy of different CNN models on the same datasets and to classify the efficaciousness and consistency of predictions done by the networks. The result obtained presents a thorough analysis of all the 5 Convolutional neural networks and both the datasets. The results show that neural networks can classify objects accurately given that the image quality is optimum and the training set is large enough. When the training set is small or the details visible in the image is degraded the neural network struggles to identify the object and the accuracy drops sharply. Moreover, with the increase in the layers of the neural network, the accuracy gets augmented. It was registered that GoogLeNet and ResNet50 were able to detect and classify objects on broad-spectrum whereas LeNet and AlexNet struggled to distinguish between objects of the same classes. The conclusion can be derived from this is that the number of layers in a neural network is directly proportional to the accuracy of its prediction. Also, for better classification and prediction a neural network required the image dataset which is provided to be of correct size format and optimum quality so that the matrix operation can be performed over it and prediction classes can be formed. The hardware requirement for training these neural networks rigorously is high and when these hardware requirements are met then this prediction can do wonders in day to day task and can increase the productivity and efficacy of human beings.

REFERENCES

- [1] Gradient-Based Learning Applied to Document Recognition [online] Available at: <<http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>>
- [2] ImageNet Classification with Deep Convolutional Neural Networks [online] Available at: <https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks>
- [3] Very Deep Convolutional Networks for Large-Scale Image Recognition [online] Available at: <https://arxiv.org/abs/1409.1556>
- [4] Deep Residual Learning for Image Recognition [online] Available at: <https://arxiv.org/abs/1512.03385>
- [5] Going deeper with convolutions [online] Available at: <https://arxiv.org/abs/1409.4842>
- [6] Rethinking the Inception Architecture for Computer Vision [online] Available at: <https://arxiv.org/abs/1512.00567>
- [7] An Analysis of Convolutional Neural Networks for Image Classification [online] Available at: <<https://www.sciencedirect.com/science/article/pii/S1877050918309335>>
- [8] The CIFAR-10 and CIFAR-100 dataset: Available at: <<https://www.cs.toronto.edu/~kriz/cifar.html>>
- [9] Wikipedia. Convolutional Neural Network. [online] Available at: https://en.wikipedia.org/wiki/Convolutional_neural_network/
- [10] CS231n Convolutional Neural Networks for Visual Recognition [online] Available at: <https://cs231n.github.io/convolutional-networks/>
- [11] Learn Open CV: Understanding AlexNet [online] Available at: <<https://www.learnopencv.com/understanding-alexnet/>>
- [12] PyImageSearch: LeNet Convolutional Neural Network [online] Available at: <<https://www.pyimagesearch.com/2016/08/01/lenet-convolutional-neural-network-in-python/>>
- [13] NeuroHive: VGG16-Convolutional network for Classification and Detection [online] Available at: <https://neurohive.io/en/popular-networks/vgg16/>
- [14] Techwasti: CNN using Gluon [Online] Available at: <https://www.techwasti.com/cnn-convolutional-neural-network-using-gluon-88f4f7ccdb8/>
- [15] Investopedia: Deep Learning [online] Available at: <https://www.investopedia.com/terms/d/deep-learning.asp>
- [16] ScienceDirect: An Analysis Of Convolutional Neural Networks For Image Classification [online] Available at: <https://www.sciencedirect.com/science/article/pii/S1877050918309335>
- [17] Wikipedia: Neural Network [online] Available at: https://en.m.wikipedia.org/wiki/Neural_network
- [18] Wikipedia: Convolutional Neural Network [online] Available at: https://en.wikipedia.org/wiki/Convolutional_neural_network
- [19] Search Enterprise AI: convolutional neural network [online] Available at: <https://searchenterpriseai.techtarget.com/definition/convolutional-neural-network>
- [20] Wikipedia: CIFAR-10 [online] Available at: <https://en.wikipedia.org/wiki/CIFAR-10>
- [21] ScienceDirect: Convolutional Layer [online] Available at: <https://www.sciencedirect.com/topics/engineering/convolutional-layer>
- [22] AnalyticsDimag: Understanding normalization methods in deep learning [online] Available at: <<https://analyticsindiamag.com/understanding-normalization-methods-in-deep-learning/>>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)