

CHAPTER 1

INTRODUCTION

1.1 Overview

Artificial intelligence (AI) refers to the simulation of human intelligence in machines that are programmed to think like humans and mimic their actions. The term may also be applied to any machine that exhibits traits associated with a human mind such as learning and problem-solving. The ideal characteristic of artificial intelligence is its ability to rationalize and take actions that have the best chance of achieving a specific goal.

A subset of artificial intelligence is machine learning, which refers to the concept that computer programs can automatically learn from and adapt to new data without being assisted by humans. Deep learning techniques enable this automatic learning through the absorption of huge amounts of unstructured data such as text, images, or video.

Artificial intelligence is based on the principle that human intelligence can be defined in a way that a machine can easily mimic it and execute tasks, from the most simple to those that are even more complex. The goals of artificial intelligence include learning, reasoning, and perception. As technology advances, previous benchmarks that defined artificial intelligence become outdated. For example, machines that calculate basic functions or recognize text through optical character recognition are no longer considered to embody artificial intelligence, since this function is now taken for granted as an inherent computer function.

AI is continuously evolving to benefit many different industries. Machines are wired using a cross-disciplinary approach based on mathematics, computer science, linguistics, psychology, and more.

1.2 Problem Statement

Sign language uses lots of gestures so that it looks like movement language which consists of a series of hands and arms motion. For different countries, there are different sign languages and hand gestures. Also, it is noted that some unknown words are translated by simply showing gestures for each alphabet in the word. In addition, sign language also

includes specific gestures to each alphabet in the English dictionary and for each number between 0 and 9. Based on these sign languages are made up of two groups, namely static gesture, and dynamic gesture. The static gesture is used for alphabet and number representation, whereas the dynamic gesture is used for specific concepts. Dynamic also include words, sentences, etc. The static gesture consists of hand gestures, whereas the latter includes the motion of hands, head, or both. Sign language is a visual language and consists of 3 major components, such as finger-spelling, word-level sign vocabulary, and non-manual features. Finger-spelling is used to spell words letter by letter and convey the message whereas the latter is keyword-based. But the design of a sign language translator is quite challenging despite many research efforts during the last few decades. Also, even the same signs have significantly different appearances for different signers and different viewpoints. This work focuses on the creation of a static sign language translator by using a Convolutional Neural Network. We created a lightweight network that can be used with embedded devices/standalone applications/web applications having fewer resources

1.3 Motivation

The various advantages of building a Sign Language Recognition system includes:

- ✓ Sign Language hand gestures to text/speech translation system or dialog systems which are used in specific public domains such as airports, post offices.
- ✓ Sign Language Recognition (SLR) can help to translate the video to text or speech enables inter-communication between normal and deaf people.

1.4 Technologies

Machine Learning

Machine Learning is a subset of Artificial Intelligence that uses statistical learning algorithms to build systems that can automatically learn and improve from experiences without being explicitly programmed. ML algorithms can be broadly classified into three categories:

- Supervised Learning

In supervised learning we have input variables (x) and an output variable (Y) and we use an algorithm to learn the mapping from input to output. In other words, a

supervised learning algorithm takes a known set of input dataset and its known responses to the data (output) to learn the regression/classification model. A learning algorithm then trains a model to generate a prediction for the response to new data or the test datasets.

- **Unsupervised Learning**

Unsupervised Learning is used when we do not have labelled data. Its main focus is to learn more about the data by inferring patterns in the dataset without reference to the known outputs. It is called unsupervised because the algorithms are left on their own to group the unsorted information by finding similarities, differences and patterns in the data. Unsupervised learning is mostly performed as a part of exploratory data analysis. It is most commonly used to find clusters of data and for dimensionality reduction.

- **Reinforcement Learning**

In simple terms, reinforcement learning can be explained as learning by continuously interacting with the environment. It is a type of machine learning algorithm in which an agent learns from an interactive environment in a trial and error way by continuously using feedback from its previous actions and experiences. The reinforcement learning uses rewards and punishments, the agents receive rewards for performing correct actions and penalties for doing it incorrectly.

Deep Learning

Deep learning is a machine learning technique that is inspired by the way a human brain filters information, it is basically learning from examples. It helps a computer model to filter the input data through layers to predict and classify information. Since deep learning processes information in a similar manner as a human brain does, it is mostly used in applications that people generally do. It is the key technology behind driver-less cars, that enables them to recognize a stop sign and to distinguish between a pedestrian and lamp post. Most of the deep learning methods use neural network architectures, so they are often referred to as deep neural networks. Deep Learning is basically mimicking the human brain, it can also be defined as a multi neural network architecture containing a large number of parameters and layers.

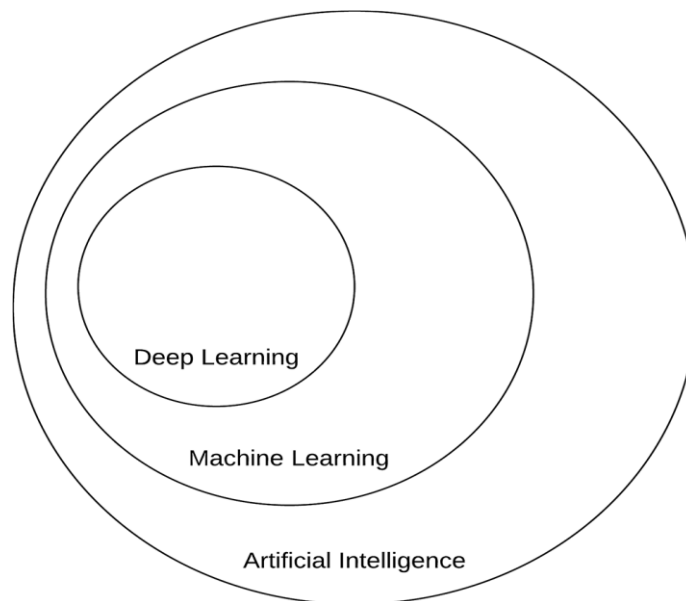


Fig 1.1 Relation between Artificial Intelligence, Machine and Deep Learning

Computer vision

Computer vision is an interdisciplinary field that deals with how computers can be made to gain high-level understanding from digital images or videos. From the perspective of engineering, it seeks to automate tasks that the human visual system can do." Computer vision is concerned with the automatic extraction, analysis and understanding of useful information from a single image or a sequence of images. It involves the development of a theoretical and algorithmic basis to achieve automatic visual understanding”.

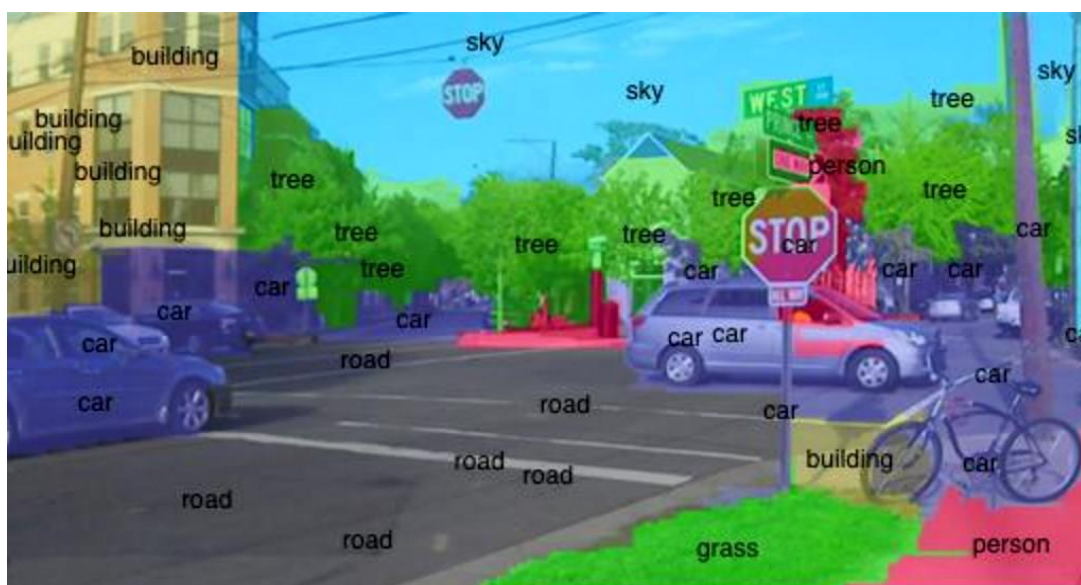


Fig 1.2 Computer Vision – Important manifestation of Artificial Intelligence

As a scientific discipline, computer vision is concerned with the theory behind artificial systems that extract information from images. The image data can take many forms, such as video sequences, views from multiple cameras, or multi-dimensional data from a medical scanner. As a technological discipline, computer vision seeks to apply its theories and models for the construction of computer vision systems.

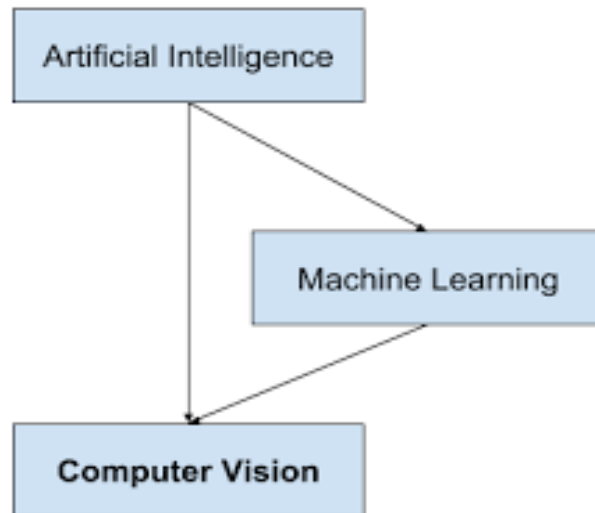


Fig 1.2 Relation between AI, Machine Learning and Computer Vision

1.5 Applications

The applications for artificial intelligence are endless. The technology can be applied to many different sectors and industries. AI is being tested and used in the healthcare industry for dosing drugs and different treatment in patients, and for surgical procedures in the operating room. Other examples of machines with artificial intelligence include computers that play chess and self-driving cars. Each of these machines must weigh the consequences of any action they take, as each action will impact the end result. In chess, the end result is winning the game. For self-driving cars, the computer system must account for all external data and compute it to act in a way that prevents a collision.

Artificial intelligence also has applications in the financial industry, where it is used to detect and flag activity in banking and finance such as unusual debit card usage and large account deposits—all of which help a bank's fraud department. Applications for AI are also being used to help streamline and make trading easier. This is done by making supply, demand, and pricing of securities easier to estimate.

CHAPTER 2

LITERATURE SURVEY

2.1 INTRODUCTION

A literature survey in a project report is that section which shows the various analyses and research made in the field of your interest and the results already published taking into account the various parameters of the project and the extent of the project. A Literature survey refers to getting the content from the books which are related to the topic or a given project. It should be referred from some research paper that is related to the topic.

Any material which is related to the project from the internet which is valuable for the student and has helped the student to enhance the report status as well as the calculation, analysis and tabulation majorly reflect in the survey. So, in this way, one can select the literature survey. It is necessary to emphasize that it is the most important part in the project report. It is the most important part of the report as it gives the students a direction in the area of their research. It helps the students to set a goal for analysis - thus giving them their problem statement.

When one writes a literature review in respect of their project, they have to write the researches made by various analysts - their methodology (which is their abstract) and the conclusions they have arrived at. One should also give an account of how this research has influenced their thesis.

Literature surveys are needed for:

- To see what has and has not been investigated.
- To identify data sources that other researchers have used.
- To learn how others have defined and measured key concepts.
- To develop alternative research projects.
- To put one's perspective into work.
- To contribute to the field by moving research forward.
- To provide evidence that may be used to support your own findings.

2.2 Literature Survey

There are many systems developed in colleges and industries to keep track of attendance. But there are performance and stability problems.

[1] Kang, Byeongkeun, Subarna Tripathi, and Truong Q. Nguyen. "Real- time sign language fingerspelling recognition using convolutional neural networks from depth map." arXiv preprint arXiv: 1509.03001 (2015).

This work focuses on static fingerspelling in American Sign Language. A method for implementing a sign language to text/voice conversion system without using handheld gloves and sensors, by capturing the gesture continuously and converting them to voice. In this method, only a few images were captured for recognition. The design of a communication aid for the physically challenged.

[2] Suganya, R., and T. Meeradevi. "Design of a communication aid for physically challenged." In Electronics and Communication Systems (ICECS), 2015 2nd International Conference on, pp. 818-822. IEEE, 2015.

The system was developed under the MATLAB environment. It consists of mainly two phases via training phase and the testing phase. In the training phase, the author used feed-forward neural networks. The problem here is MATLAB is not that efficient and also integrating the concurrent attributes as a whole is difficult.

[3] Sruthi Upendran, Thamizharasi. A," American Sign Language Interpreter System for Deaf and Dumb Individuals", 2014 International Conference on Control, Instrumentation, Communication and Computa.

The discussed procedures could recognize 20 out of 24 static ASL alphabets. The alphabets A, M, N and S couldn't be recognized due to occlusion problem. They have used only a limited number of images.

CHAPTER 3

SYSTEM AND SOFTWARE REQUIREMENTS AND SPECIFICATIONS

The program works on Desktop PC and is executed using a Django Rest Frameworks interface which interacts with a PostgreSQL database running on localhost.

3.1 FUNCTIONAL REQUIREMENTS

A description of the facility or feature required. Functional requirements deal with what the system should do or provide for users. They include description of the required functions, outlines of associated reports or online queries, and details of data to be held in the system.

3.1.1 Interface Requirements:

- ✓ The system shall provide user to upload image via frontend.
- ✓ The system shall give the prediction output in the web page.
- ✓ The system shall provide option to convert the predicted text to speech and all the details will be stored in the database.
- ✓ The system shall provide option to make prediction on both web application and standalone application.
- ✓ The system shall provide option to capture live hand gesture via camera and detect the corresponding English alphabets and then store in the database.

3.2 NON-FUNCTIONAL REQUIREMENTS:

Non-functional requirements define the overall qualities or attributes of the resulting system.

3.2.1 Usability

Usability is the ease with which a user can learn to operate the Online Shopping Website system and get results.

3.2.2 Security

Security requirements are included in a system to ensure

- ✓ All inventory and supplier information are well secured
- ✓ SQL injection is prevented

3.2.3 Reliability

Reliability is the ability of a system to perform its required functions under stated conditions for a specific period of time. Constraints on the run-time behaviour of the system can be considered under two separate headings:

- ✓ Availability: is the system available for service when requested by end-users.
- ✓ Failure rate: how often does the system fail to deliver the service as expected by end- users.

3.2.4. Efficiency

The comparison of what is actually produced or performed with what can be achieved with the same consumption of Clouds (money, time, labor, etc.). It is an important Factor in Determination of Productivity.

3.3 SOFTWARE REQUIREMENTS

Programming language	:	Python
Operating system	:	ANY OS
Application required	:	Web Application
Coding language	:	Core Python OOPS

3.4 HARDWARE REQUIREMENTS

CPU	:	Pentium IV 2.4 GHz
Memory (Primary)	:	512 MB, 1 GB or above
Hard Disk	:	40 GB, 80GB, 160GB or above
Monitor	:	15 VGA color

CHAPTER 4

SYSTEM DESIGN AND METHODOLOGY

4.1 Methodology

Image recognition is done using TensorFlow

Techniques used:

- ✓ Convolution Neural Network
- ✓ OpenCV
- ✓ TensorFlow
- ✓ Django Rest Frameworks

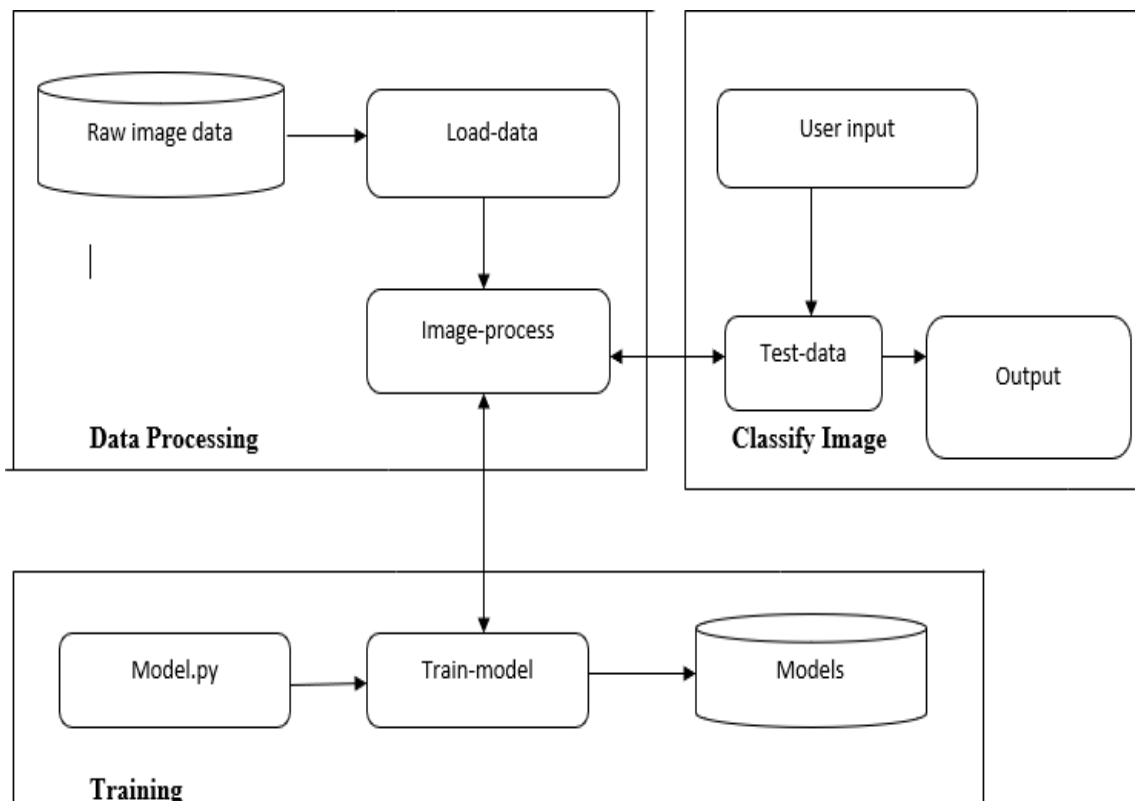


Fig 4.1: Methodology of Proposed Model

The methodology has 3 major steps as shown in the Fig 4.1:

1. Data preprocessing
2. Training
3. Image classification

Data preprocessing is the first step of the process. Here the raw data is obtained, loaded and then it is processed. Training is the second step of the process. Here, the

preprocessed images are given to the YOLO model and the model is trained on the images. Once the model is trained, the input image is given through the camera, the input image is processed and the processed image is trained to model. The model classifies the image and the output is obtained. The model is deployed using Django Rest Frameworks.

Steps involved are:

- ✓ Initializing the model using Sequential class.
- ✓ Adding the Input layer that is images and applying the Convolution 2D method.
- ✓ Adding the Hidden layer and applying the activation function.
- ✓ Max pooling method is applied then.
- ✓ Flattening the image using the Flatten method.
- ✓ Adding the output layer and applying a suitable activation function. This is called full connection method.
- ✓ Compiling and training the object detection model
- ✓ Predicting the model and evaluating the accuracy.

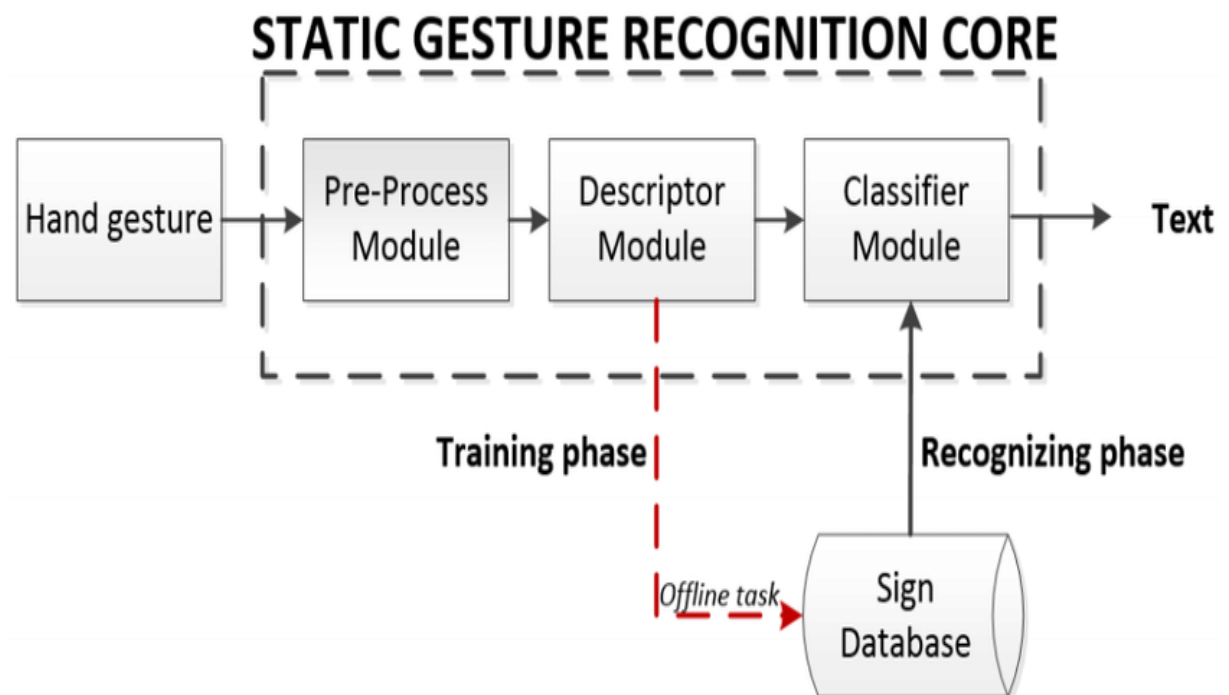


Fig. 4.2 Flow chart of the Proposed Model

Web Developing steps:

- ✓ The frontend is done using Html, CSS, and JavaScript.
- ✓ Backend is done using Django
- ✓ The database used PostgreSQL

Finally, to compound all these challenges, there is the issue of signer independence. While larger data sets are starting to appear, few allow true tests of signer independence over long continuous sequences. Maybe this is one of the most urgent problems in SLR that of creating data sets which are not only realistic, but also well annotated to facilitate machine learning. Despite these problems recent uses of SLR include translation to spoken language, or to another sign language when combined with avatar technology. SLR is also set to be used as an annotation aid, to automate annotation of sign video for linguistic research, currently a time-consuming and expensive task. Fig 4.2 shows the flow of the model proposed.

CHAPTER 5

IMPLEMENTATION

5.1 Dataset

We have used multiple datasets and trained multiple models to achieve good accuracy.

5.1.1 ASL Alphabet

The data is a collection of images of the alphabet from the American Sign Language, separated into 29 folders that represent the various classes.

The training dataset consists of 87000 images which are 200x200 pixels. There are 29 classes of which 26 are English alphabets A-Z and the rest 3 classes are SPACE, DELETE, and, NOTHING. These 3 classes are very important and helpful in real-time applications.

5.1.2 Sign Language Gesture Images Dataset

The dataset consists of 37 different hand sign gestures which include A-Z alphabet gestures, 0-9 number gestures, and also a gesture for space which means how the deaf (hard hearing) and dumb people represent space between two letters or two words while communicating.

Each gesture has 1500 images which are 50x50 pixels, so altogether there are 37 gestures which mean there 55,500 images for all gestures. Convolutional Neural Network (CNN) is well suited for this dataset for model training purposes and gesture prediction.

5.2 Data Pre-processing

An image is nothing more than a 2-dimensional array of numbers or pixels which are ranging from 0 to 255. Typically, 0 means black, and 255 means white. Image is defined by mathematical function $f(x, y)$ where 'x' represents horizontal and 'y'

the pixel value at that point of an image.

Image Pre-processing is the use of algorithms to perform operations on images. It is important to Pre-process the images before sending the images for model training. For example, all the images should have the same size of 200x200 pixels. If not, the model cannot be trained.



Fig. 5.1. Sample Image without Pre-processing

The steps we have taken for image Pre-processing are:

- ✓ Read Images.
- ✓ Resize or reshape all the images to the same
- ✓ Remove noise.
- ✓ All the image pixels array are converted to 0 to 255 by dividing the image array by 255.



Fig. 5.2. Pre-Processed Image

5.3 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are one of the most popular architectures of deep learning which simulate biological nervous system like Artificial Neural Networks (ANNs). AlexNet, GoogleNet, Squeeze Net and ResNet are the most common architectures of CNN. In comparison with ANNs, CNNs take the advantage of local connections instead of fully connections in all layers except the last layer.

CNNs structure contains series of most common layers as shown in Fig 4.3:

- ✓ Convolution Layer
- ✓ Pooling Layer
- ✓ Fully connected layer

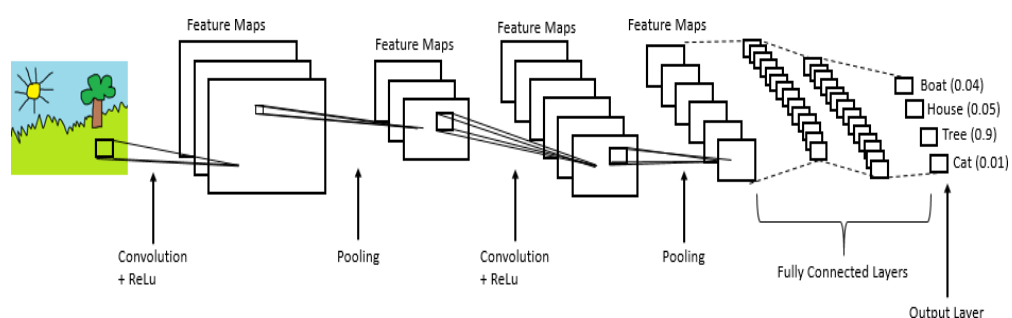


Fig. 5.3. CNN Layers

5.3.1 Convolution Layer

The first layer is convolutional layer. In this layer, each region that contains feature maps is connected to the feature maps of a local region in the previous layer by calculating weights that are known as kernels (filter banks). Sum of all local weights goes through a non-linearity function, e.g. Relu.

Input image, feature detector and feature maps are the three essential elements that enter into the convolution operation. The basic idea of the convolution is Input image x Feature detector = Feature Map. The main advantage of using feature maps is reducing the size of the input image, making it easier to read. A convolution matrix to adjust an image. It can be used to blur, edge detect and sharpen the image. The Eqn 1 represents the convolution function

$$(f*g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau \quad (\text{Eqn 1})$$

The two main down sides of convolution are shrinking of outputs and the loss of information, especially in the corners of the image. In order to overcome this, padding is done. In simple words, it can be described as an additional layer that can be added to the border of the image. This helps in increasing the accuracy.

$$n - f + 1 = 6$$

$$n = 6 - 1 + f$$

$$n = 5 + 3 = 8$$

The value of n obtained here is 8. If a padding p=1 is applied then, we get an image of size 8x8 is obtained. The following is the result obtained on multiplying the image with the filter.

$$\text{Result} = n + 2p - f + 1$$

$$\text{Result} = 6 + 2 - 3 + 1 = 6$$

The resultant matrix, thus obtained is 6x6 and there is no loss of data. The original size of the image can be retained if padding is applied. Any number of convolution and padding can be applied in order to retain the size of the image.

5.3.2 Pooling Layer

Pooling layer is the second common layer of CNNs which is used to decrease the spatial dimension of the output of convolutional layer without any change in depth. The advantage of using pooling layer is that by decreasing computational operations it prevents the overfitting in training process. Pooling layer contains min, max and average operations as shown in fig 4.4. However max pooling layer has achieved reasonable results in most fields.

The main objective of max pooling is to reduce dimensionality, down-sample an input representation and allowing the assumptions made about the features contained in the sub-regions. It can be described as a sample-based discretization process. In max pooling, only the high pixel values are being taken and hence, we are clearly able to recognize the face in the image.

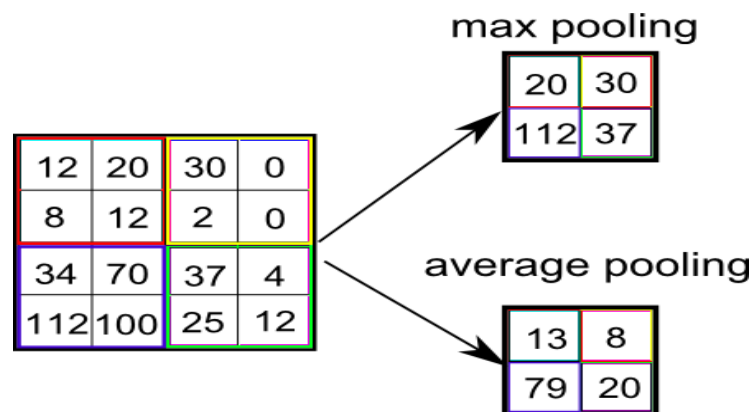


Fig. 5.4. Types of Pooling

5.3.3 Fully connected layer

The third common layer is fully connected layer. Before this, flattening must be done as shown in fig 4.5. Flattening can be defined as the process of conversion all the resultant two-dimensional arrays into a single long continuous linear vector.

And it is connected to the final classification model, which is called a fully-connected layer as shown in fig 4.6. In other words, all the pixel data is put in one line and make connections with the final layer.

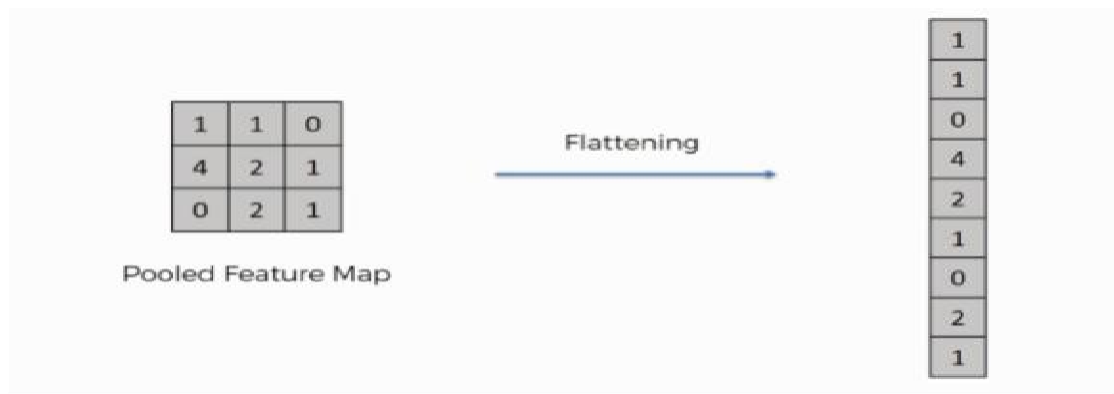


Fig. 5.5. Flattening

In the fully connected layer shown in fig 4.6 ,each neuron is not only connected to the all neurons of the previous layer but also calculated scores of dataset's classes are given in this layer. Moreover, generally in the last most convolutional layers softmax function is utilized to calculate the probable distribution through the labels of the class.

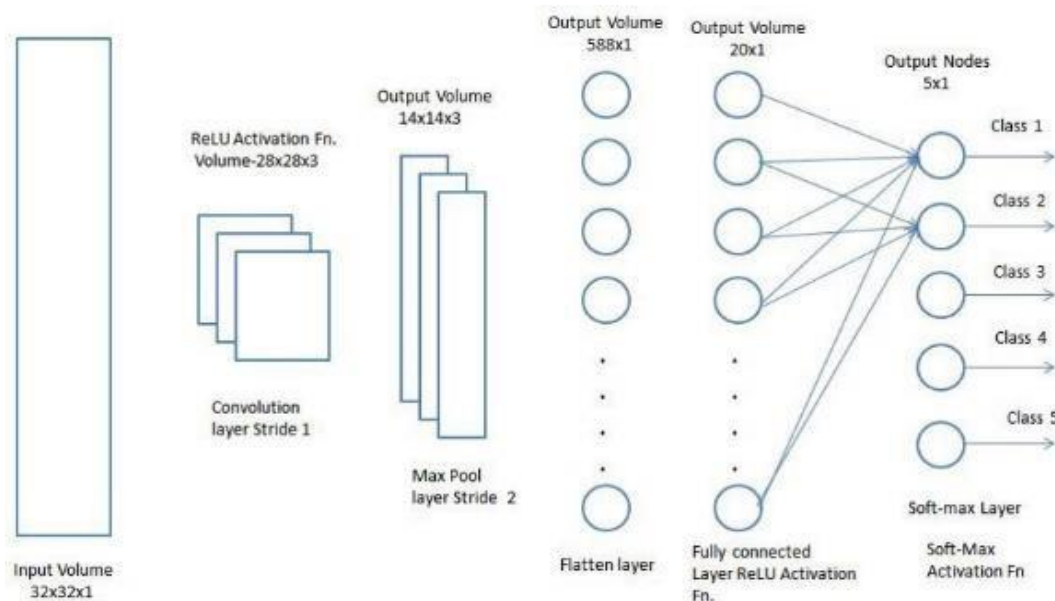


Fig 5.6. Fully connected layer

5.4 Convolution Neural Network (CNN) Architectures

5.4.1 LeNet-5 Architecture

LeNet-5 is one of the simplest architectures. It has 2 convolutional and 3 fully-connected layers. The LeNet-5 architecture consists of two pairs of convolutional and average pooling

layers, followed by a flattening convolutional layer, then two fully connected layers, and finally a SoftMax classifier.

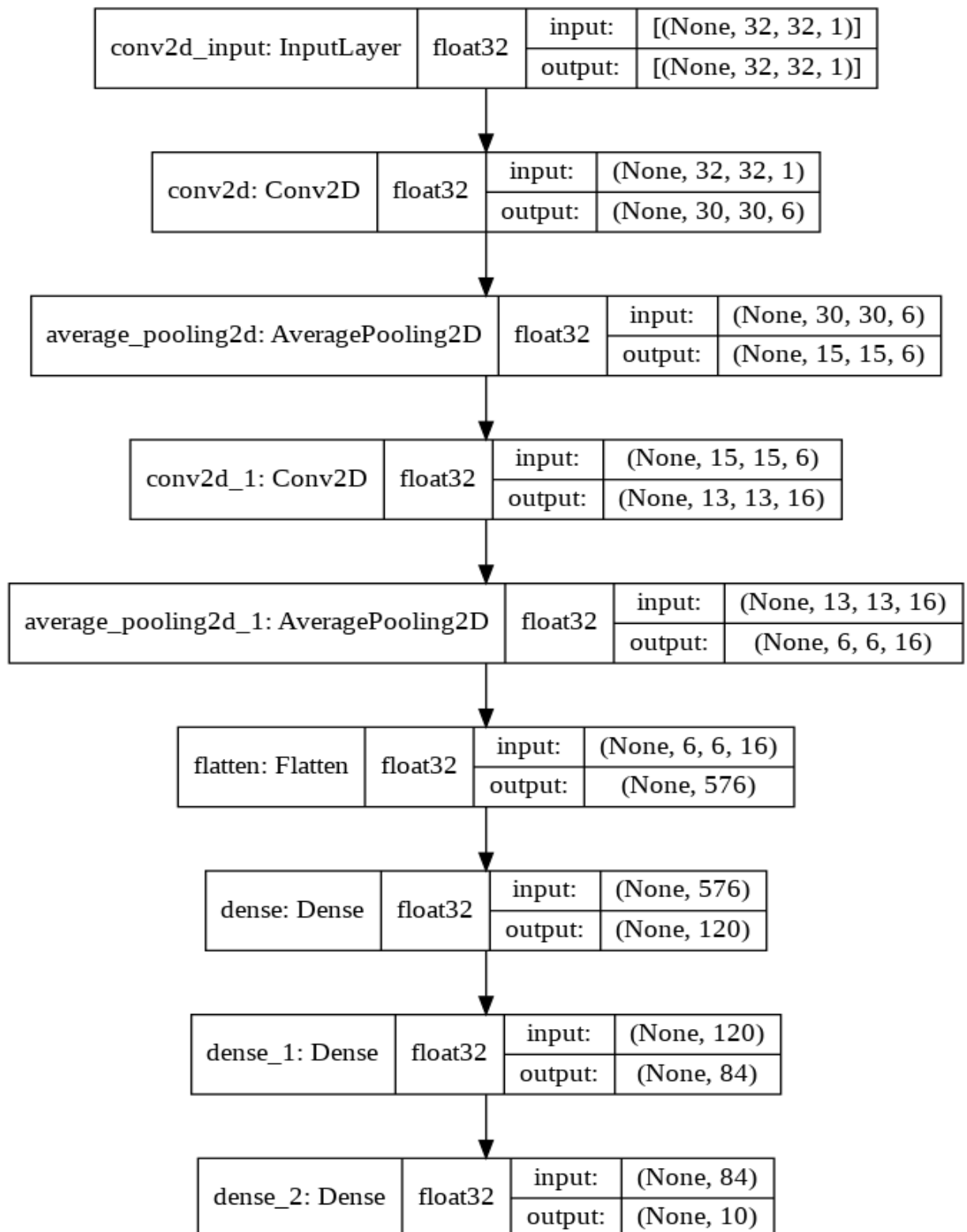


Fig 5.7. LeNet Architecture

5.4.2 MobileNet V2 Architecture

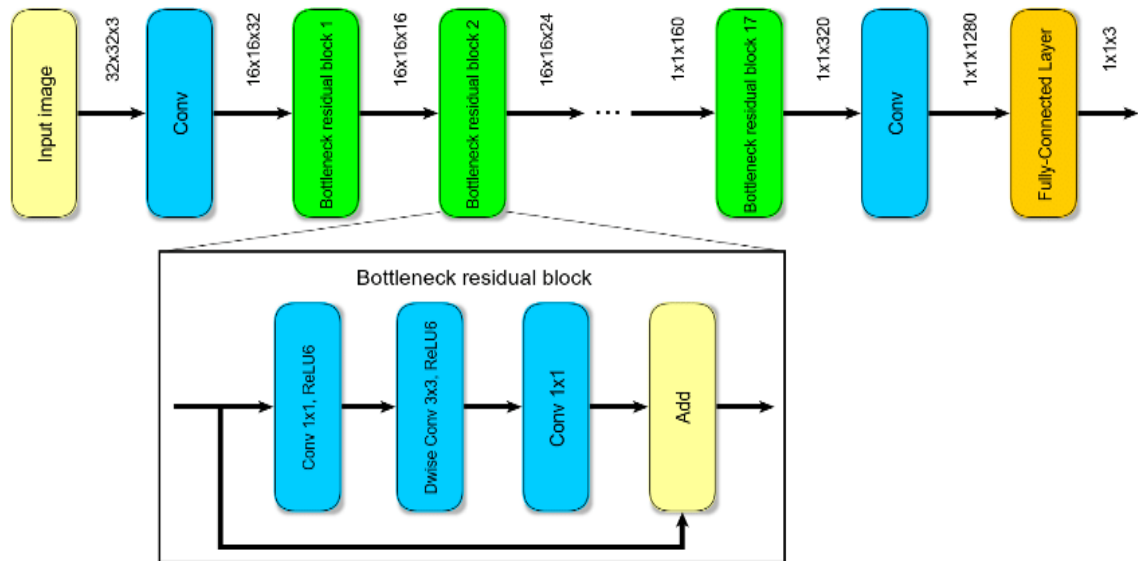


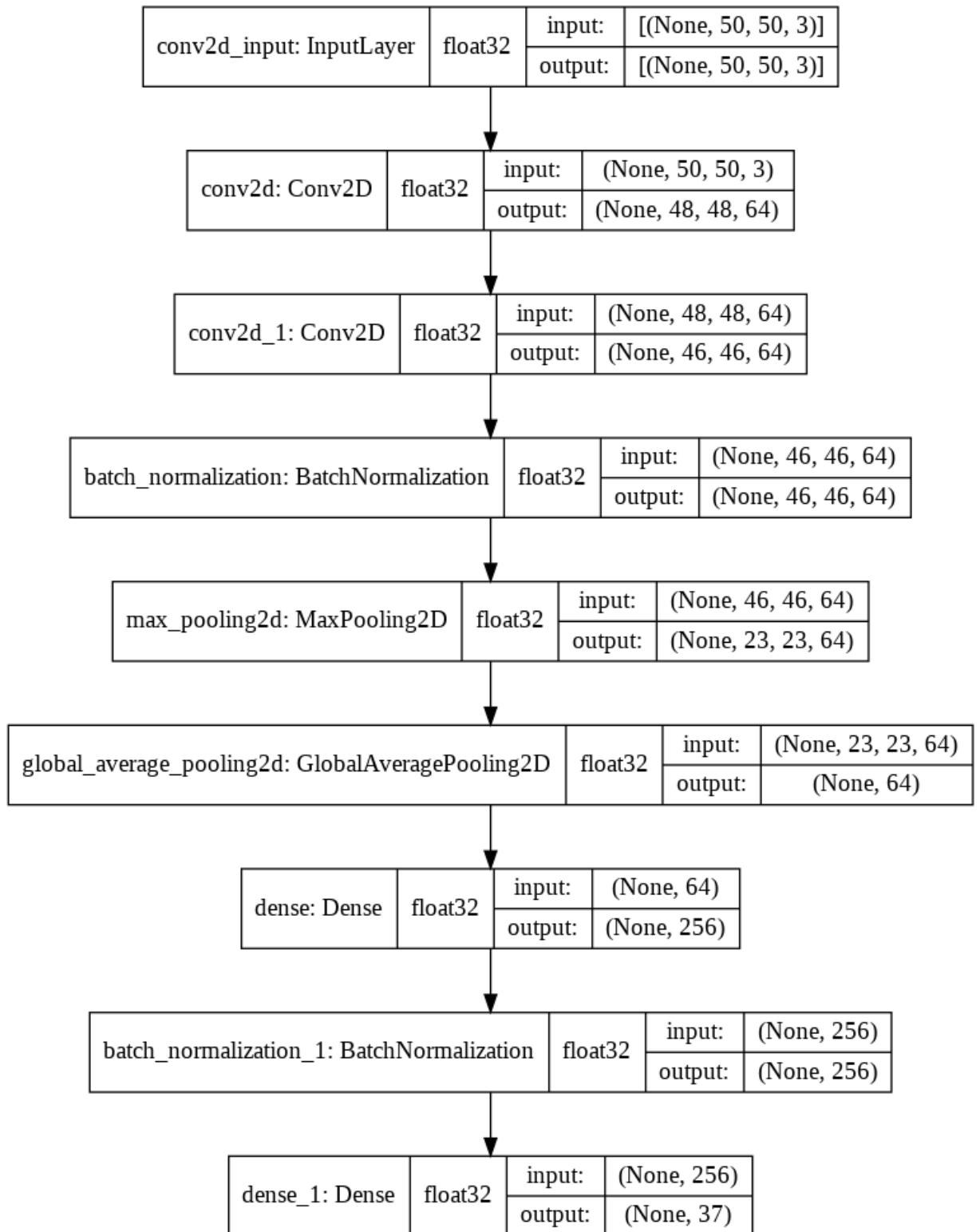
Fig. 5.8. MobileNet V2 Implementation

In the previous version MobileNetV1, Depthwise Separable Convolution is introduced which dramatically reduce the complexity cost and model size of the network, which is suitable to Mobile devices, or any devices with low computational power. In MobileNetV2, a better module is introduced with inverted residual structure. Non-linearities in narrow layers are removed this time. With MobileNetV2 as backbone for feature extraction, state-of-the-art performances are also achieved for object detection and semantic segmentation.

MobileNetV2 is a convolutional neural network architecture that performs well on mobile devices. The architecture of MobileNetV2 contains the fully convolution layer with 32 filters, followed by 19 residual bottleneck layers. This network is lightweight and efficient.

5.4.3 Own Architecture

In our own architecture, we have implemented 3 convolution layers followed by batch normalization and max pooling, followed by global average pooling with dense layer and batch normalization, and a final dense layer for classification.

**Fig. 5.9. Own Architecture Implementation**

5.5 Proposed Model

We have trained 3 models with 2 different datasets to perform well on unseen datasets. We have trained LeNet-5, MobineNetV2 and, our own architectures. We have not taken the best model out of 3 we have taken all 3 models and made a final model that will perform an ensemble of these 3 models.

5.5.1 Neural Network Ensemble Horizontal Voting

In Machine Learning we have an ensemble technique where we train multiple sub-models and average them. Random Forest algorithm is an example where it uses multiple Decision tree algorithms. Similarly, we can perform ensemble for Neural Networks [4] as well. There are a lot of ensemble techniques for Neural Network like Stacked generalization [8], Ensemble learning via negative correlation [9] and, Probabilistic Modelling with Neural Networks [10] [11].

We have implemented the Horizontal Voting Ensemble method to improve the performance of neural networks. Horizontal voting is an ensemble technique for neural networks where we train several sub-models and make predictions using these sub-models. For the final predictions, we make predictions from all the sub-models and see which class has got maximum votes.

The final prediction will be the class that has the maximum votes. For this, we have used 3 models that are an odd number of sub-models to avoid an even number of votes for two classes in worst cases.

For MobileNetV2 model we have taken Adam optimizer with learning rate = 0.001, eta_1=0.9, beta_2=0.999, and epsilon=1e-07 For all the models while training we have used ReduceLROnPlateau callback with factor = 0.2, patience = 2, min_lr = 0.001. By using this horizontal ensemble technique, we have achieved 99.8% accuracy.

We have deployed all the models in Django Web Frameworks and built a simple frontend to accept the image from users and send the response. We have also built an API that will pop up the live camera and detects the hand gestures and then converts them to the corresponding English alphabets.

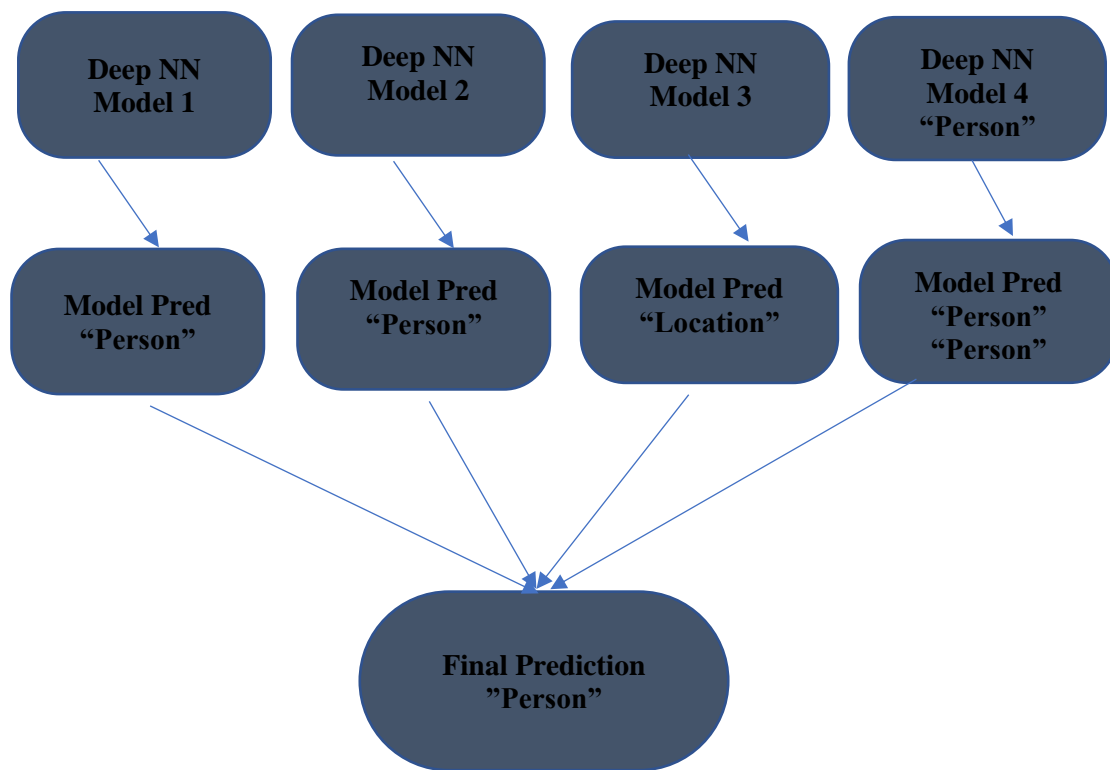


Fig. 5.10. Horizontal Voting

RESULTS

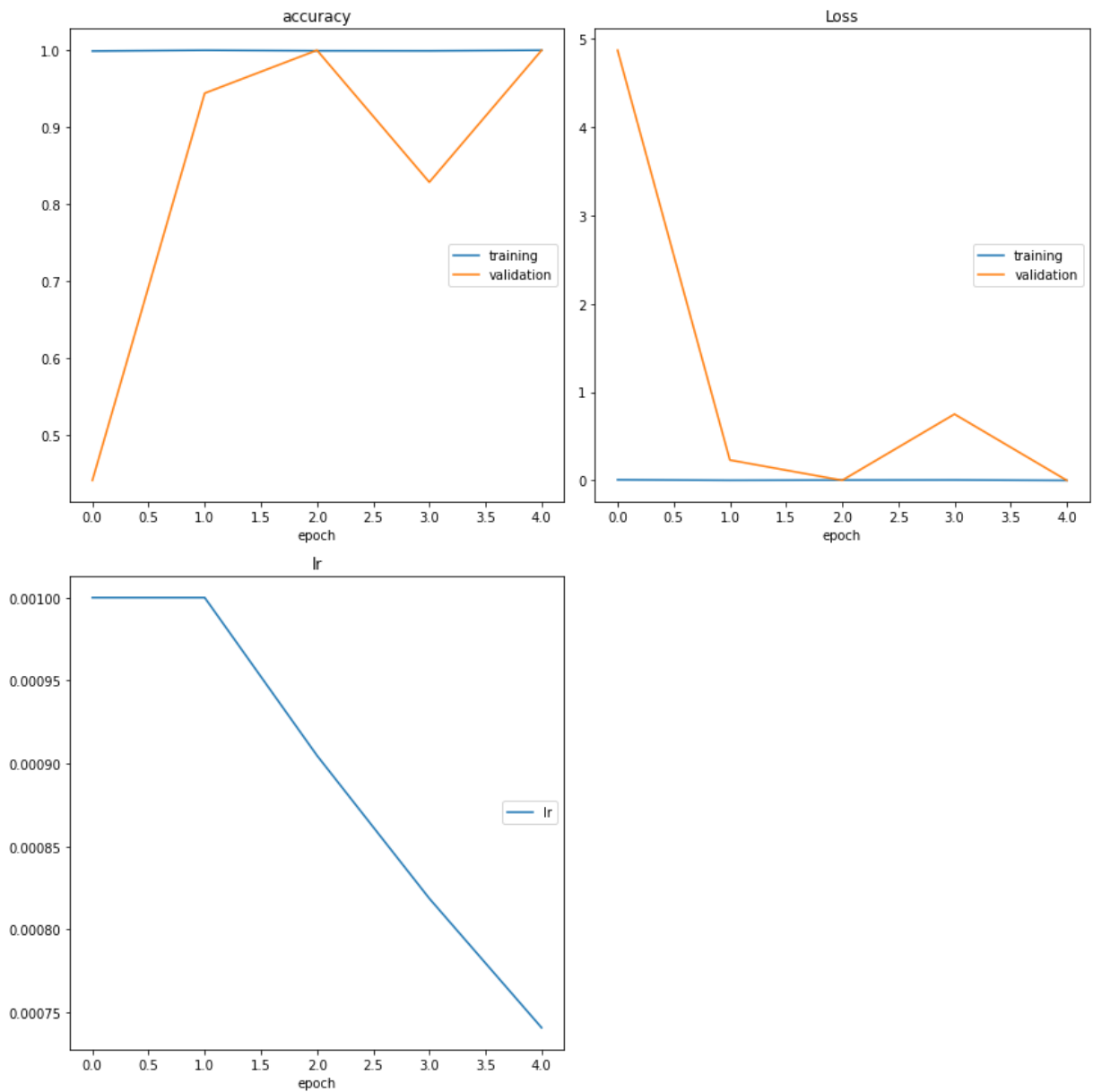


Fig 6.1 Trained Model Plots

Model Analysis at Epoch

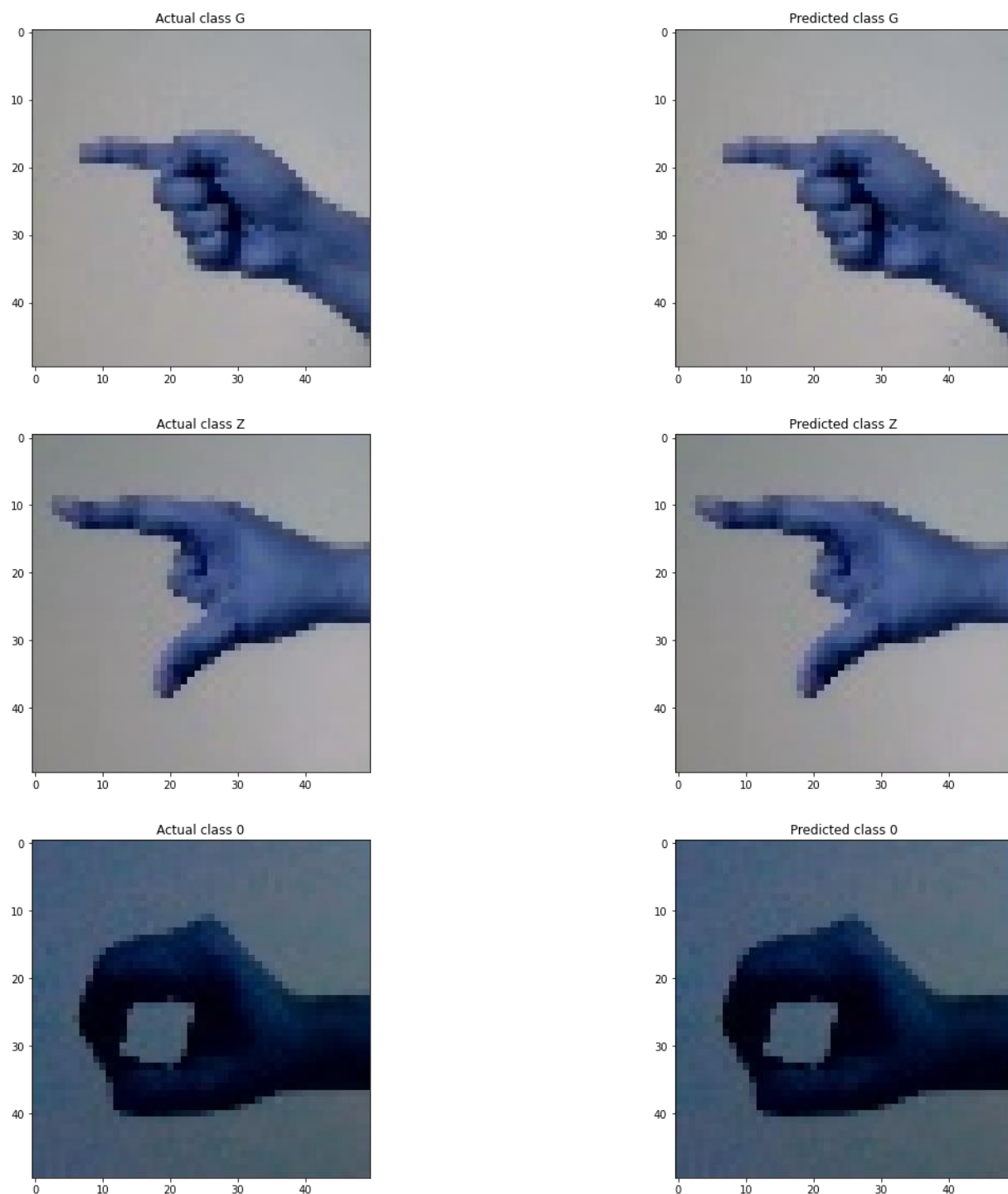


Fig 6.2 Model Prediction

Models	Accuracy
MobileNetV2	98.9%
LeNet-5	97%
Own Model	98%
Ensemble	99.8%

Table. 1. Performance Results

All the models performed well on the test cases. After applying the horizontal voting ensemble technique for these 3 models, we have achieved almost 100% accuracy.

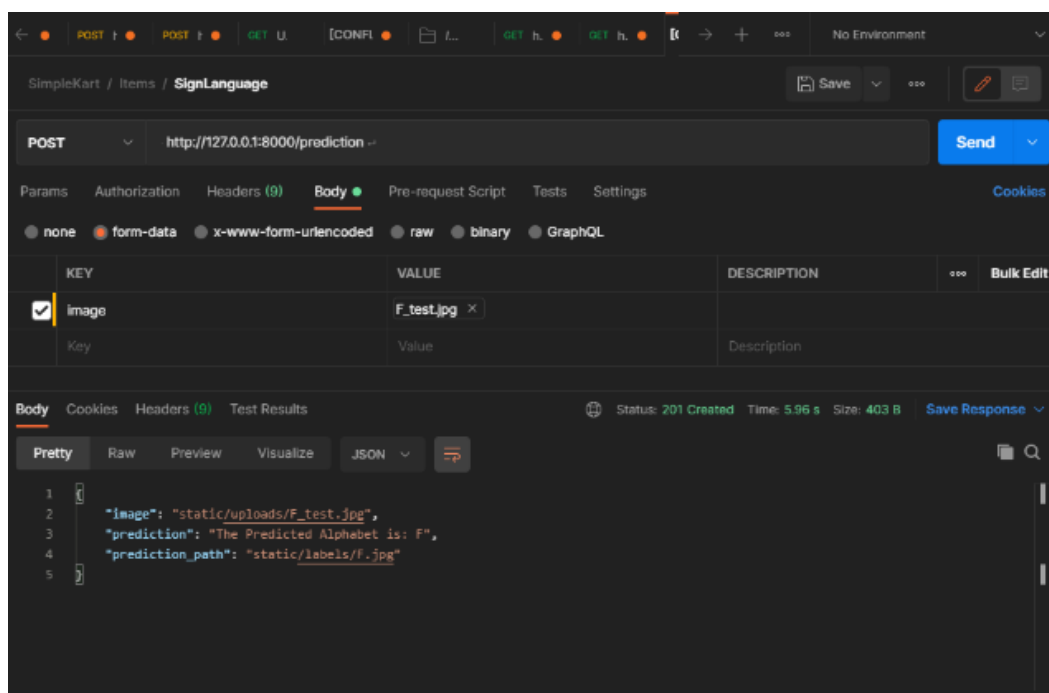


Fig 6.3 JSON Response from the Application

We have used Django Rest Frameworks as a backend for our project. This is the sample JSON response that is sent to the frontend when a user inputs an Image.



Fig 6.4 Prediction

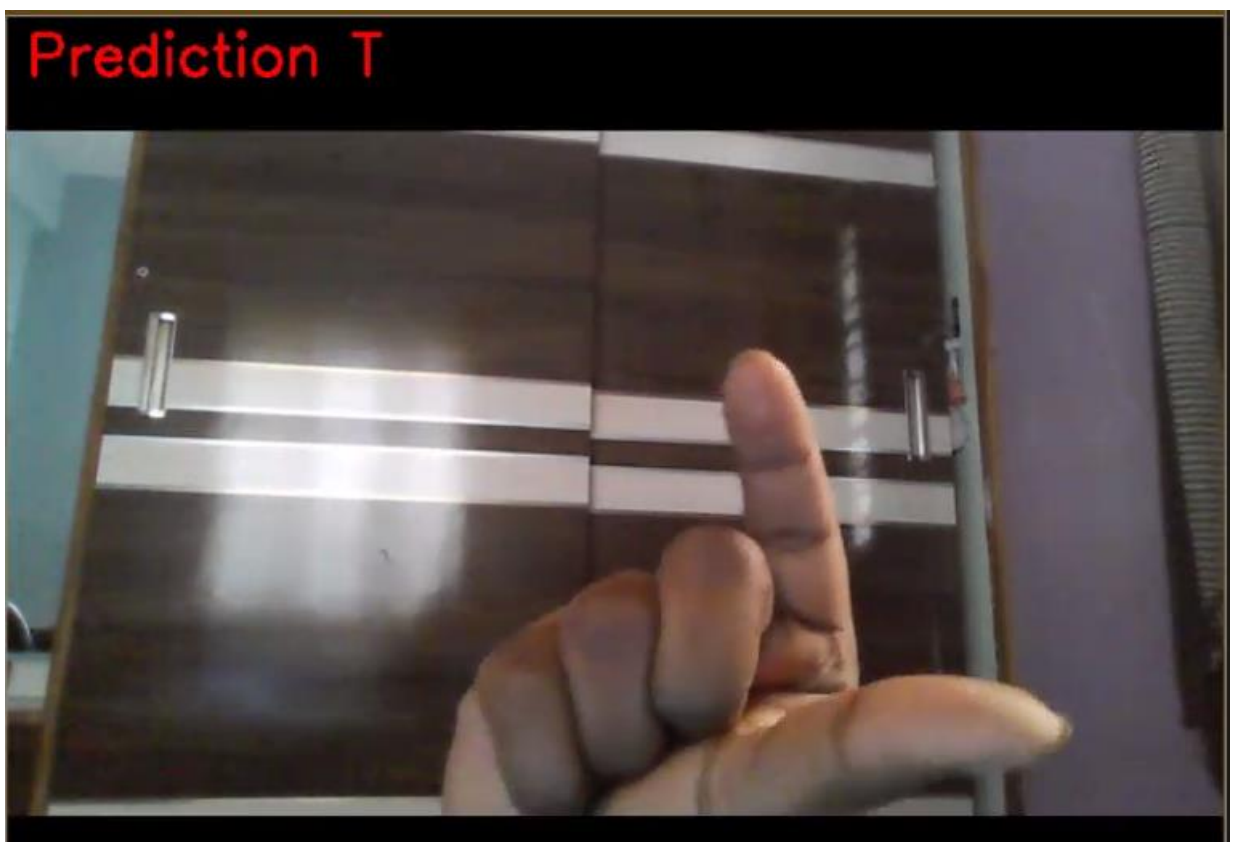


Fig 6.5 Prediction



Fig 6.5 Prediction



Fig 6.6 Prediction

CHAPTER 6

CONCLUSION

Communications between deaf-mute and a normal person have always been a challenging task. The goal of this project is to reduce the barrier of communication by contributing to the field of automatic sign language recognition. Through this work, a CNN classifier is constructed which is capable of recognizing static sign language gestures. A basic GUI application is created to test our classifier in this system. The application allows users to select the static sign gestures as input and it will speak out the words or sentences corresponding to the gesture. We have trained our model for 26 symbols which include alphabets. We were able to achieve an accuracy of 99% for our CNN classifier.

Dept. of CSE, SVIT

REFERENCES

- [1] TensorFlow Documentation
- [2] Django Rest Frameworks Documentation
- [3] CNN Research paper
- [4] L. K. Hansen and P. Salamon, "Neural network ensembles," in IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 12, no. 10, pp. 993-1001, Oct. 1990, doi: 10.1109/34.58871.
- [5] David H. Wolpert, Stacked generalization, Neural Networks, Volume 5, Issue 2, 1992, Pages 241-259, ISSN 0893-6080, [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1).
- [6] Y. Liu, X. Yao, Ensemble learning via negative correlation, Neural Networks, Volume 12, Issue 10, 1999, Pages 1399-1404, ISSN 0893-6080, [https://doi.org/10.1016/S0893-6080\(99\)00073-8](https://doi.org/10.1016/S0893-6080(99)00073-8).
- [7] MacKay D.J.C. (1995) Developments in Probabilistic Modelling with Neural Networks — Ensemble Learning. In: Kappen B., Gielen S. (eds) Neural Networks: Artificial Intelligence and Industrial Applications. Springer, London. https://doi.org/10.1007/978-1-4471-3087-1_37
- [8] Polikar R. (2012) Ensemble Learning. In: Zhang C., Ma Y. (eds) Ensemble Machine Learning. Springer, Boston, MA. https://doi.org/10.1007/978-1-4419-9326-7_1

