

PREMIER UNIVERSITY

Department of Computer Science & Engineering



Final Year Thesis Report

On

**Human Action Recognition in Library Environment Using
CNN-RNN Model**

SUBMITTED BY

Name: Mahmudur Rashid Mehraj

ID: 1903610201794

Name: Raj Dev Ullash

ID: 1903610201804

Name: Bonnhi Shikha Parna

ID: 1903610201820

In partial fulfillment for the degree of
Bachelor of Science in Computer Science & Engineering
under the Supervision of

Anik Sen

Assistant Professor

Department of Computer Science & Engineering

Premier University, Chattogram

September, 2023

Author's Declaration of Originality

We hereby declare that the thesis work entitled “**Human Action Recognition in Library Environment Using CNN-RNN Model**” submitted to the Premier University, is a record of an original work done by us under the guidance of Mr. Anik Sen, Lecturer, Department of Computer Science & Engineering, Premier University, Chattogram and this work is submitted for fulfillment of the degree of Bachelor Science in Computer Science & Engineering. We can assure that the result of this thesis has not been submitted to any other university.

Mahmudur Rashid Mehraj
ID: 1903610201794

Raj Dev Ullash
ID: 1903610201804

Bonnhi Shikha Parna
ID: 1903610201820

CERTIFICATION

The thesis entitled “**Human Action Recognition in Library Environment Using CNN-RNN Model**” submitted by, Mahmudur Rashid Mehraj, ID: 1903610201794, Raj Dev Ullash, ID: 1903610201804, Bonnhi Shikha Parna, ID: 1903610201820 has been accepted as satisfactory in partial fulfillment of the degree of Bachelor Science in Computer Science & Engineering (CSE) to be awarded by Premier University, Chattogram.

Anik Sen
Assistant Professor
Department of Computer Science & Engineering
Premier University, Chattogram

Dr. Shahid Md. Asif Iqbal
Chairman and Associate Professor
Department of Computer Science & Engineering
Premier University, Chattogram

*It is my genuine gratefulness and warmest regard that
I dedicate this work to my beloved
Father and Mother*

Table of Contents

List of Figures	vi
List of Tables	vii
1 Introduction	1
1.1 Introduction	1
1.2 Motivations	2
1.3 Application	2
1.4 Challenges	3
1.5 Summary	3
2 Literature Review	4
2.1 Introduction	4
2.2 Background Study	4
2.2.1 Neural Network	4
2.2.2 Convolutional Neural Network	5
2.2.3 Long Short-Term Memory (LSTM)	6
2.2.4 Transfer Learning (TL)	6
2.3 Related Review for Human Action Recognition	6
2.4 Summary	9
3 Methodology/Proposed Method/Workflow	10
3.1 Introduction	10
3.2 Dataset Collection	11
3.3 Overall Working Procedure	12
3.3.1 VGG16 Model	12
3.3.2 DenseNet121	14
3.3.3 InceptionV3 Model	15
3.3.4 A Hybrid Architecture for Human Action Recognition	15
3.4 Summary	16
4 Experimental Results and Discussion	17
4.1 Introduction	17
4.2 System Configuration	18
4.3 Confusion Matrix	18
4.4 Classification Report	18

4.4.1	Precision:	19
4.4.2	Recall:	19
4.4.3	F1 score:	19
4.5	Dataset Description	20
4.6	Parameter Calculation of Our Best Model VGG16	21
4.6.1	Results and Experiments VGG16 of Approach	24
4.6.2	Results and Experiments InceptionV3 of Approach	26
5	Conclusion	29
5.1	Limitation	30
5.2	Future Work	30
	Bibliography	31

List of Figures

1.1	Student activity monitoring	3
3.1	Workflow of the proposed human action recognition	11
3.2	Overall Process for Human Action Recognition	16
4.1	Reading	20
4.2	Writing	20
4.3	Using Smartphone	20
4.4	Training Accuracy vs Test Accuracy of The Best VGG16 Model	24
4.5	Confusion Matrix of The Best VGG16 Model + Bidirectional LSTM Model	25
4.6	Training Accuracy vs Val Accuracy of The Best InceptionV3 Model . . .	27
4.7	Confusion Matrix of The Best InceptionV3 + Bidirectional LSTM Model	27

List of Tables

3.1	Distribution of Samples of Our Dataset	12
3.2	Pre-trained VGG16 Architecture for 224×224 Input Shape	13
3.3	Summary of VGG16 + Bidirectional LSTM Model	13
3.4	Summary of DenseNet121 + Bidirectional LSTM Model	14
3.5	Summary of InceptionV3 + Bidirectional LSTM Model	15
4.1	System Configuration	18
4.2	Distribution of Samples of Our Dataset	20
4.3	Paramter Table Summary of VGG16	21
4.4	Experiments of the VGG16 Approach	24
4.5	Calculations of Confusion Matrix from	25
4.6	Experiments of the InceptionV3 Approach	26
4.7	Calculations of Confusion Matrix from	28

Abstract

Video-based Human Action recognition (HAR) has been popular in recent years for its real-time activity monitoring. In this paper, we present a new deep learning-based human activity recognition technique. For educational institutions that have libraries, understanding how students engage with study materials can provide valuable insights into learning behaviors. It is challenging to automatically learn video actions in a library environment due to its fine-grain nature. As a result, we proposed a hybrid CNN-RNN model as a novel classifier. Human performance accompanies human body shapes and their relative motions. In this thesis, a novel dataset of 1301 videos has been developed. We have collected and preprocessed video data which includes fine-grained actions of three different categories; Reading, Writing, and Using Smartphone. Our collected dataset of videos where each video clip contains a human performing a specific action from two views. To address the challenge by collecting and preprocessing a large amount of fine-grained video data, hybrid CNN-RNN is proposed here. The integration of CNN and RNN architectures brings together the strengths of library action data analysis. Furthermore, to leverage transfer learning, we built multiple pre-trained models with ImageNet weights including VGG16, InceptionV3, and DenseNet121 then conducted hyperparameter tuning. Our proposed architecture VGG16 achieves the highest test accuracy of 84% on the dataset.

Keywords: Human Action recognition (HAR), Library Actions Dataset, Fine-grained video classification, Human actions dataset

CHAPTER 1

INTRODUCTION

1.1 Introduction

Human action recognition is a vital field related to computer vision that enables researchers to identify actions performed by humans in different environments. Though image classification has been implemented to assign labels or categories video-based classification is popular now. Video-based human action recognition is the task of identifying and categorizing human actions or activities performed in videos. It involves analyzing the temporal sequence of video frames to determine what action is being performed by the human subject. This task is important for various applications, including video surveillance, sports analysis, human-computer interaction, and more. In this thesis, we implement how human action recognition is used in library environments and classify it. We collect a dataset of videos where each video clip contains a human performing a specific action. To address the challenge of collecting and preprocessing a large amount of fine-grained video data, a hybrid CNN-RNN is proposed here. We split the dataset into training and testing sets. To improve the work capability for a model we again split the training set to get another validation test set. CNN to extract spatial features and RNN to sequential patterns from the frames. Finally, we use the softmax activation function to generate a set of probabilities for each action in the video data. Using transfer learning techniques, we can provide an effective classification in three categories: Reading, Writing, and Using Smartphones.

1.2 Motivations

We look into human activity recognition in library settings, focusing on behaviors like reading, writing, and smartphone use. However, monitoring and understanding the intricate actions that take place within library environments remains a challenge due to the absence of human behavior data. It is a difficult problem perspective of computer vision and deep learning. We are motivated by other's work to emphasize the importance of real-time processing for action recognition. There is also need for timely decision-making and the increasing demand for applications that require real-time action recognition, such as surveillance and security.

To automatically learn video actions in the library environment, our goal is to identify the reasoning behind these actions and recognize the actions, which will contribute to a more thorough knowledge of human behavior in such situations. No published work regarding human action recognition in a library environment. To address this gap, the proposed approach of developing a hybrid Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) model holds great promise. This hybrid model offers the advantage of learning complex patterns in both spatial and temporal dimensions, enabling it to effectively classify fine-grained actions unique to libraries. The challenge in human activity recognition is to efficiently recognize various actions in complex situations, provide a high accuracy recognition rate, and to simplify implementation in real-time application while using less computing power motivated us.

1.3 Application

We can apply this research for real-time student activity monitoring. It offers an exciting opportunity to reshape education. By embracing emerging technologies and data-driven insights, educational institutions can create an environment where each student's journey is tracked, analyzed, and enhanced in real-time. This approach empowers educators to provide personalized support, fosters student success, and contributes to the advancement of educational practices through ongoing research and innovation. Real-time monitoring is visualized below Figure 1.1.



Figure 1.1. Student activity monitoring

1.4 Challenges

It was challenging to collect and preprocess a large amount of video data. As fine-grained action, there were huge similarities among the actions.

1.5 Summary

In this chapter, we have established the context and rationale for investigating human action recognition in library environments. It outlines the research objectives and scope, highlighting the potential benefits of the study's findings for enhancing library design, services, and user experiences. The results of our experiments reveal that our technique is efficient and delivers a flawless detection rate. This chapter serves as a foundational guide for the reader, setting the stage for the comprehensive analysis presented in the following chapters. The development of a hybrid CNN-RNN model to classify fine-grained video actions in library environments addresses the need for a data-driven approach to understanding activities in these spaces. By leveraging the power of deep learning and video data collection, this research not only has the potential to benefit libraries but also contributes to the broader field of computer vision and human activity recognition.

CHAPTER 2

LITERATURE REVIEW

2.1 Introduction

In recent years, deep learning models have emerged as a critical pillar within the machine learning community. Gradually, they have transformed into extensively employed computational tools, particularly in the specialized realm of machine learning. A prominent benefit of deep learning is its capacity to effectively process massive volumes of data. We implement the CNN model with the help of LSTM as part of RNN. Based on these models, this review serves as an essential precursor to our study, Human Action Recognition in a library environment, as it offers a comprehensive assessment of the existing research and its implications. By examining the extensive body of literature available, we aim to identify gaps, contradictions, and opportunities for further investigation. Moreover, we seek to highlight emerging trends, theoretical frameworks, and methodological approaches that have emerged since earlier studies.

2.2 Background Study

2.2.1 Neural Network

A neural network is a computational model inspired by the structure and functioning of the human brain. It consists of interconnected artificial neurons (also called nodes or units) organized in layers to process and learn from data. Neural networks are the foundation of deep learning, a subset of machine learning that focuses on training models with multiple

layers to automatically learn features from data and make predictions.

2.2.2 Convolutional Neural Network

A Convolutional Neural Network (CNN) is a specialized type of artificial neural network designed to process and analyze visual data, particularly images and videos. It's inspired by the human visual system's ability to recognize patterns and features in visual information. CNNs are highly effective for tasks such as image recognition, object detection, facial recognition, and more. Convolutional Neural Network (CNN) model for processing video frames is a concept commonly employed in video analysis tasks to capture temporal information and reduce the computational complexity of training and inference. This approach leverages the fact that neighboring frames in a video often share similar patterns and features.

Key components we use a CNN include:

Convolutional Layers These layers perform convolution operations on the input data using learnable filters or kernels. The filters slide over the input data, extracting features by detecting patterns, edges, textures, and other important information.

Activation Functions Non-linear activation functions, such as Rectified Linear Activation (ReLU), introduce non-linearity into the network. ReLU replaces negative values with zeros and keeps positive values unchanged. We use Softmax activation function in the last output layer convert into probability distributions.

Fully Connected Layers After multiple layers of convolution and pooling, the extracted features are flattened and fed into one or more fully connected layers. These layers make final predictions or classifications based on the extracted features. We use 3 dense layer.

CNNs excel at capturing hierarchical features in visual data. Lower layers detect simple patterns like edges and corners, while deeper layers recognize more complex shapes and objects. This hierarchical feature extraction enables CNNs to automatically learn meaningful representations from raw visual data. CNNs have revolutionized various fields, including computer vision, where they have achieved state-of-the-art performance on tasks like image classification, object detection, image segmentation, and even artistic style transfer. They've also been applied in other domains, such as natural language processing and bioinformatics, by adapting their architecture to suit the specific characteristics of the data.

2.2.3 Long Short-Term Memory (LSTM)

Long Short-Term Memory (LSTM) is a type of recurrent neural network (RNN) architecture designed to handle sequences of data. It addresses the vanishing gradient problem that can occur in traditional RNNs and enables the network to learn and capture long-range dependencies within sequences. LSTMs are particularly effective for tasks involving sequential data, such as time series forecasting, natural language processing, speech recognition, and more. The LSTM architecture effectively addresses the issue of vanishing gradients by allowing information to flow through the memory cells without significant degradation over time. This makes LSTMs capable of capturing and retaining long-term dependencies in sequences, which is crucial for tasks where context over extended periods is important.

2.2.4 Transfer Learning (TL)

A unusual phenomena is shared by many deep neural networks trained on images: early layers of a deep learning model attempt to learn low level information, such as identifying edges, colors, fluctuations in intensities, etc. No matter what sort of picture we are processing, whether for lion or automobile detection, such characteristics don't seem to be relevant to a given dataset or job. In both situations, we must find these inconspicuous qualities. Regardless of the precise cost function or picture dataset, all these traits exist. . It utilizes the knowledge a model has comprehended from a task with many available labeled training data in a new task that does not have much data. Rather than starting the learning process from scratch, we start with patterns comprehended from unraveling an affiliated task. Transfer learning is not a machine learning technique, but can be seen as a “design methodology” within the field, for example, active learning.

2.3 Related Review for Human Action Recognition

Author [1] targets recognition of human actions in videos. This problem is complicated due to the high complexity of human actions such as appearance variation, motion pattern variation, occlusions, etc. The author claims that the ability to process a video in real-time will be a key factor in future action recognition applications. They have also evaluated their methods on publicly available datasets; CAD60, and CAD120. Their experiments show that the methods proposed in this thesis improve state-of-the-art results. Author [2] works with the automatic recognition of actions in videos which is a challenging task To obtain good classification rates, it is necessary to work with spatial information (for example, shapes found in a single frame of the video) and temporal information (for example, move-

ment patterns found throughout the frames in the video). In this thesis, new methods are proposed for the automatic recognition of human actions based on deep learning features extracted from videos. There is a high demand for the development of new techniques for automatic pattern recognition in videos, for example for the automatic recognition of human actions, this demand is motivated by the advances in the technologies of production, storage, transmission, and sharing of videos, such advances triggered the production of a huge volume of videos that need to be automatically processed to be useful. Author [3], focuses on automating the recognition of human actions in videos, a challenging task due to factors like diverse human motion, occlusion, and changing environments. It reviews and categorizes existing techniques into handcrafted and deep-learning approaches. The proposed framework integrates novel algorithms for both types. A handcrafted method addresses "viewpoint variations" using silhouette-based features, achieving top-tier results without compromising efficiency. Two deep learning methods are introduced: one employs transfer learning from pre-trained models, showcasing better accuracy than handcrafted methods on limited data. The second utilizes unsupervised Deep Belief Networks for feature extraction, achieving high recognition accuracy on complex videos. The thesis presents a comprehensive exploration of automatic human action recognition, offering contributions in both handcrafted and deep learning domains. Discussions and future research directions round up the work, emphasizing its relevance and potential impact. Author [4], the rapid proliferation of surveillance technology has led to the emergence of diverse digital monitoring devices across various contexts. The installation of these devices has resulted in a wealth of surveillance videos encompassing different human actions and behaviors, including pedestrian and vehicle tracking, object recognition, and human behavior analysis. These videos hold the potential to aid in distinguishing normal and abnormal behaviors, thus enhancing public safety. However, the traditional approach to analyzing these videos relies heavily on manual efforts, posing challenges such as high labor costs and energy consumption. As working hours increase, human energy levels decline, underscoring the need for more efficient and automated methods. Author [5], human motion analysis is gaining attention due to diverse applications like body part segmentation, joint tracking, and 3D structure recovery for athletics and medical diagnosis. Automated activity monitoring in security zones is vital for law enforcement. Interpreting videos aids digital libraries, user interfaces, and video conferencing. Action recognition research focuses on robust techniques for complex scenes. Author [6], tackles human action recognition in real-world video data, such as movies and online videos. Accurate action recognition has diverse applications, from surveillance to medical diagnosis. Challenges arise due to varying factors like appearance, backgrounds, lighting, and viewpoints. Local space-time features are successful but lack structure. The thesis enhances Bag-of-Features video representations by introducing more discriminative features and supervision, aiming to enhance action recognition in challenging videos. Various local

space-time feature detectors and descriptors are evaluated, leading to the integration of non-local region-level information for better recognition. An attribute-based approach encompassing objects and human poses is explored, providing complementary high-level information. Additionally, the thesis introduces Actions, and motion-pattern-based body part detectors trained on synthetic data. Actions show promise for improving action recognition in realistic scenarios. This thesis offers novel supervised video representation methods to advance human action recognition in complex settings. Author [7], focuses on automating the recognition of human actions in videos, a complex task due to diverse variations in appearance, motion, viewpoints, and environmental factors. It reviews, evaluates, and proposes a local feature-based action recognition framework, which is applied across the thesis, along with introducing a new dataset of elder patients' actions. Novel local spatio-temporal descriptors are introduced: one based on a covariance matrix to model linear relations, and another based on Brownian covariance to account for all possible relations between low-level features. The thesis extends beyond local feature techniques with three higher-level representations. The first employs relatively dense trajectories, using object-centric motion trajectory features for spatial awareness. The second encodes appearance and motion relations as pairwise features, utilizing geometric information to characterize their arrangement. The third capture statistics of co-occurring visual words within multi-scale neighborhoods, offering contextual features that encompass density, pairwise relations, and spatio-temporal order among features. These advancements collectively enhance the recognition of human actions in videos. M. Abdellaoui, A. Douik, A. et al, [8] suggested a novel DBN-based HAR approach that aims to improve the accuracy of human activity categorization. As a preliminary step, divide the video segments from the human activity dataset into frames. The result is then transformed into binary frames, which are subsequently treated to some morphological filtering processes to increase their quality. Then, by transforming the new frames into a binary vector, create an input matrix including the training data, the testing data, and their labels. This matrix represents the input data for our DBN architecture. In the last phase, the DBN classifier will be trained using the training data matrix, and the classification outcome will be extracted. Klaser et al. [9] proposed a video sequence local descriptor. Histograms of oriented 3D spatio-temporal gradients serve as the foundation for the proposed descriptor. To compute 3D gradients at any size, the authors first created a memory-efficient technique based on integral films. The researchers then presented a generic 3D orientation quantization based on regular polyhedrons and they evaluated all descriptor settings thoroughly and optimized them for HAR. Their descriptor was applied to several human action datasets, including (KTH 91.4%), Weizmann (84.3%), and Hollywood (24.7%).

2.4 Summary

We introduce methods based on deep +learning features extracted from videos. In this section implementation of all the methods has been discussed. Integrating handcrafted and deep learning methods showcases the benefits of combining different techniques to achieve higher recognition accuracy. To classify and recognize human action recognition, we get a good structure of the model and how to implement it.

CHAPTER 3

METHODOLOGY/PROPOSED METHOD/WORKFLOW

3.1 Introduction

One of the popular algorithms VGG16 model has demonstrated the highest 84% test accuracy than other models on our dataset. The model represents the confluence of several concepts created by numerous scholars over time. VGG16 (Visual Geometry Group 16) is a convolutional neural network architecture. This pre-trained model can capture a wide range of features.

In this thesis, a novel dataset was created, for all experiment 224×224 RGB frames extracted from each video. A CNN architecture was applied to extract spatial features from the video frames. Transfer learning from the pre-trained VGG16 model was then utilized. The extracted features were fed into a Recurrent Neural Network (RNN) to capture the temporal changes of frames. Finally, a softmax layer was introduced to predict the class of each video.

The workflow was divided into several steps, including pre-processing of frames, training and exploitation where training also placed with training and validation set , feature extraction using CNN, design of a hybrid CNN-RNN network, and evaluation of the whole testing model. The entire workflow is illustrated in Figure 3.1.

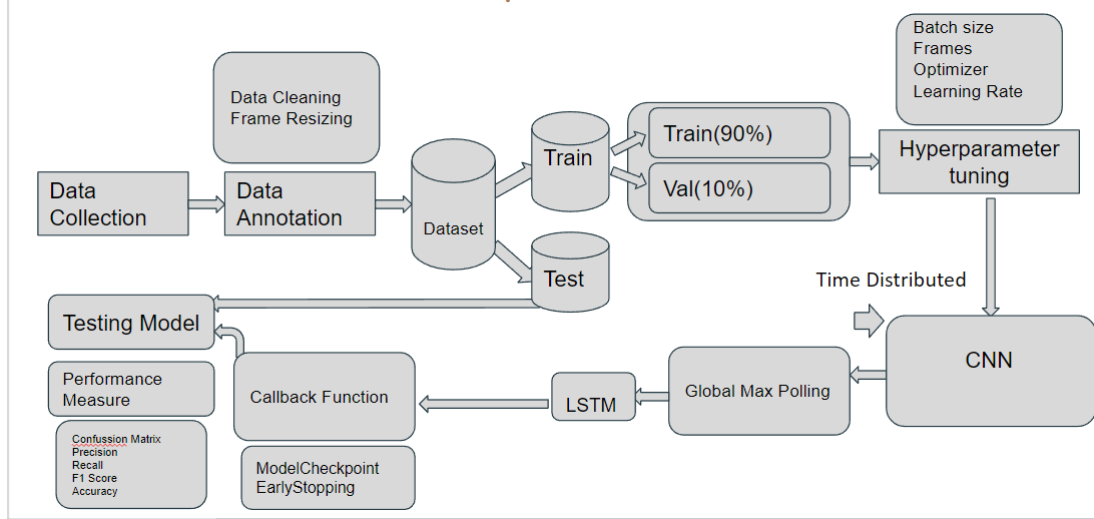


Figure 3.1. Workflow of the proposed human action recognition

3.2 Dataset Collection

To develop a precise classifier capable of accurately identifying student actions such as "Reading," "Writing," and "Using Smartphones," an impeccable dataset is imperative. We curated video clips captured within the library environment to construct this dataset, resulting in a comprehensive collection of 1301 video samples.

In our video samples, each video is 25 FPS. In the model, we extract our video into 10FPS with VideoFrameGenerator Library.

We divided our dataset into two parts train and test. In the training set, we divided it into two parts training and validation test. We split our dataset into 90% and 10%. 90 percent of the dataset is employed to train our classifiers. Subsequently, the remaining 10 percent is for validation tests. The test set is reserved to evaluate the efficiency of our classifiers. This segregation enables a robust assessment of the classifiers' performance on unseen data.

The distribution of samples across each class within our dataset is thoughtfully represented in the accompanying table, highlighting the balanced representation of "Reading," "Writing," and "Using Smartphone" instances. This balanced distribution contributes to a fair and unbiased model development process, ultimately yielding a classifier that is adept at discerning these distinct student activities. 3.1.

Table 3.1. Distribution of Samples of Our Dataset

Classes	Reading	Writing	Using Smartphone
Training Samples	325	325	325
Validation Test Samples	37	37	37
Test Samples	76	52	87

3.3 Overall Working Procedure

3.3.1 VGG16 Model

To detect human actions, a VGG16-based model was selected. To comprehend sequences, four time-distributed layers were added. This methodology is useful when dealing with time-based data or video frames. Instead of using separate input models, a single model can be applied to each input. The initial time-distributed input comprises a VGG16 convolutional neural network, which is commonly known for its 16 layers. The first block of the model includes a 2x2 max-pooling layer followed by two convolution layers with 64 kernels. The second block is similar but includes 128 kernels and a 2x2 max-pooling layer. The third block comprises three convolution layers with 256 kernels followed by a stride 2x2 max-pooling layer. The fourth and fifth blocks each contain three convolution layers with 512 filters and a max-pooling layer. The final two time-distributed layers are dense, one with 512 units and the other with 256 units. The last time-distributed layer flattens the input to match the channel dimensions. The network is trained using gradient descent and dropout to prevent over-fitting. A Bidirectional LSTM layer produces outputs for each time step. The network concludes with three fully connected layers; two with 1024 and 512 units, and the last with 3 units to match the three classes in the dataset.

Table 3.2. Pre-trained VGG16 Architecture for 224×224 Input Shape

Layer Type	Output	Parameters
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0

Table 3.3 says the hyperparameters of our VGG16 model. This model consists of a time distribution (dense) layer, time distribution (flattened), Bidirectional LSTM, and dense layer. The values of time distribution (dense) values consist of 512 and 256. The Bidirectional LSTM parameters are 64 and dense layer hyperparameters consist of 1024 and 512.

Table 3.3. Summary of VGG16 + Bidirectional LSTM Model

Layer (type)	Output Shape	Param #
time_distributed_1	(None, 10, 512)	14714688
bidirectional_1 (Bidirectional)	None, 10, 192)	467712
bidirectional_2 (Bidirectional)	(None, 192)	221952
dropout_1 (Dropout)	(None, 192)	0
dense (Dense)	(None, 256)	49408
dropout_2 (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 256)	65792
dense_2 (Dense)	(None, 96)	24672
dense_3 (Dense)	(None, 3)	291

3.3.2 DenseNet121

Another model DenseNet121 was developed for the Human Action Recognition classification. Similarly, like the first model we used 4 time time-distributed layer and the first-time distributed layer has a DenseNet121 network as its input. The DenseNet121 architecture is also constructed using depth-separable convolution, which is used in the human event recognition model. DenseNet121 consists of 121 layers, including a deep convolution layer, 120 convolution layer, ReLU, average collecting layer, and softmax. Two dense, completely connected layers with 256 and 3 units each make up the following two time-distributed layers. The last time-distributed layer then takes as input a flattened layer that flattens the input's spatial dimensions to their equivalent in channel dimensions. After that, the network is trained by gradient descent and dropout to avoid over-fitting. Then, an output is produced at each time step using a Bidirectional LSTM layer. In the end, three completely linked layers were used. The two dense layers have 256 and 3 units, respectively. Our dataset had three classes, hence the final dense layer contained three units. We utilized "DenseNet121" preserved weights, froze the last 12 layers of the network, and kept the remaining levels trainable such that the final layers' weights would be changed during training.

Table 3.4. Summary of DenseNet121 + Bidirectional LSTM Model

Layer (type)	Output Shape	Param #
time_distributed	(None, 10, 1024)	7037504
bidirectional (Bidirectional)	(None, 10, 128)	557568
bidirectional (Bidirectional)	(None, 128)	98816
dropout (Dropout)	(None, 128)	0
dense (Dense))	(None, 256)	33024
dropout_1 (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 3)	771

Table 3.4 says the hyperparameters of our DenseNet121 model. This model consists of a time distribution (dense) layer, time distribution (flattened), Bidirectional LSTM, and dense layer. The values of time distribution (dense) values consist of 1024 and 512. The Bidirectional LSTM parameters are 98816 and dense layer hyperparameters consist of 256 and 3.

3.3.3 InceptionV3 Model

InceptionV3 is applied as our third model as we know inceptionV3 is made up of a total of 42 layers and is more accurate than its forerunners. Factorization into smaller Convolutions, Spatial Factorization into Asymmetric Convolutions, Utility of Auxiliary Classifiers, and Efficient Grid Size Reduction are the main improvements made to the InceptionV3 model. Traditionally, the grid size of the feature maps was decreased using average and maximum pooling. The activation dimension of the network filters is enhanced in the InceptionV3 model to more effectively minimize the grid size. There are more network filters on InceptionV3 compared to the previous versions of inception. We employed four time-distributed layers and fed InceptionV3 into the first one. The following two time-distributed layers feature two dense layers with nodes of 1024 and 512, similar to the prior ways. The table below shows the final layers and the settings utilized in InceptionV3.

Table 3.5. Summary of InceptionV3 + Bidirectional LSTM Model

Layer (type)	Output Shape	Param
time_distributed	(None, 10, 2048)	21802784
bidirectional (Bidirectional)	(None, 10, 160)	1362560
bidirectional_1 (Bidirectional)	(None, 160)	154240
dropout (Dropout)	(None, 160)	0
dense	(None, 160)	10304
dropout_1 (Dropout)	(None, 64)	0
dense_1	(None, 64)	4160
dense_2	(None, 64)	4160
dense_3	(None, 3)	195

Table 3.5 says the hyperparameters of our InceptionV3 model. This model consists of a time distribution (dense) layer, time distribution (flattened), Bidirectional LSTM, and dense layer. The values of time distribution (dense) values consist of 2048, 1024, and 512. The RNN parameters are 96 and dense layer hyperparameters consist of 1024 and 512.

3.3.4 A Hybrid Architecture for Human Action Recognition

The CNN model is wrapped in a time-distributed layer, allowing it to be used across all video frames. Within this wrapper, two fully connected layers have 512 and 256 hidden units each. Afterward, a flattening operation occurs in another time-distributed wrapper. The output goes into a Bidirectional LSTM layer with 64 units, and to avoid overfitting, a dropout layer with a 50% rate is applied to the RNN layer. All fully connected layers use the ReLU activation function. Lastly, a 3-unit softmax layer gives the action probabilities.

The entire architecture is visualized in the diagram. in figure 3.2.

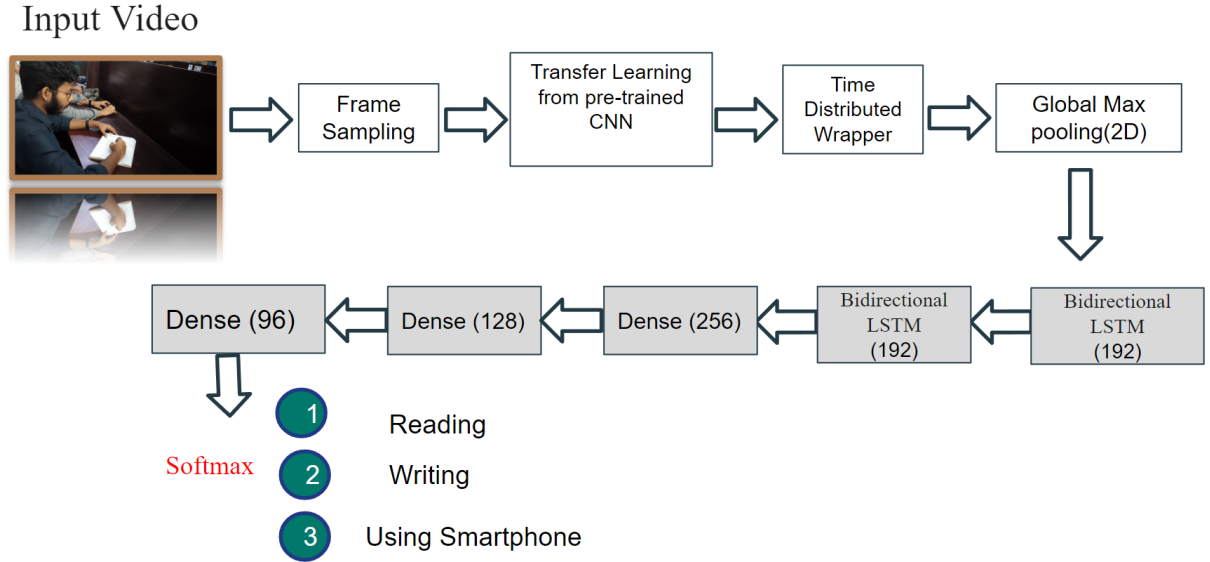


Figure 3.2. Overall Process for Human Action Recognition

3.4 Summary

This section comprehensively presents the development stages of our models designed to identify human action in the library environment and explains their recognition. We delve into data collection methods, our holistic working approach, and the overall procedural framework. Furthermore, we provide insights into the architectures of InceptionV3 and VGG16 models. By amalgamating these discussions, we offer a comprehensive overview of our model creation, operational strategies, data collection practices, and architectural insights into the mentioned models.

CHAPTER 4

EXPERIMENTAL RESULTS AND DISCUSSION

4.1 Introduction

The realm of deep learning involves critical model parameters such as weights and biases, intrinsic to the training data and gradually unveiled throughout the learning journey. A model's efficacy is conventionally appraised via parameters. The cumulative count of parameters within the neural network is the summation of all weights and biases. We also present the parameter calculation method for our most optimal model. A pivotal element in assessing the performance of a classification algorithm is the confusion matrix—a structured table that elucidates the algorithm's effectiveness. By visually summarizing the classification algorithm's outcomes, a confusion matrix aids in performance comprehension. Our analysis entails the examination of the confusion matrices for all three models, employed to derive precision and recall metrics for each.

Our pursuit of identifying the most suitable model for our dataset led us to craft a series of sophisticated deep-learning models. Among these are esteemed architectures like VGG16, InceptionV3, DenseNet121. Each of these models underwent rigorous training using our dataset, accompanied by meticulous hyper-tuning to extract their utmost performance potential.

To gauge performance, we employed a suite of performance metrics including precision, recall, F1-score, and accuracy. With these metrics, we scrutinized the models' performance comprehensively. Feel free to remix, modify, and update this text to suit your needs.

4.2 System Configuration

The system configuration for all of our experiments is outlined in [4.1](#)

Table 4.1. System Configuration

Name	Details
Processor	Intel(R) Core(TM) i5-6300HQ CPU @ 2.30GHz
RAM	8 GB
Graphics Card	NVIDIA GeForce 960M
GPU Memory	GPU T4*2
Operating System	Windows 11
Library and Framework	Tensorflow 2.11, Keras 2.3.1
Programming Language	Python 3.8.16
Online Platform	Google colab, Kaggle

4.3 Confusion Matrix

The combination of the actual and anticipated classes is displayed in a confusion matrix. It is a table that is utilized in classification issues to determine where model errors occurred. The rows correspond to the actual courses for which the results were intended. The predictions we have made are represented by the columns. This table makes it simple to identify whose predictions were incorrect. It is a useful indicator of how well models can take into account class property overlap and identify which classes are most susceptible to confusion.

4.4 Classification Report

To evaluate the accuracy of predictions made by a classification algorithm, utilize a classification report. How many of the forecasts came true and how many did not? To be more precise, the metrics of a categorization report are predicted using True Positives, False Positives, True Negatives, and False Negatives. The primary classification metrics precision, recall, and f1-score are displayed in the report on a per-class basis. True and false positives, as well as true and false negatives, are used to determine the metrics. In this instance, "positive" and "negative" are just names for the projected classes. There are four ways to determine if the estimates were accurate or not:

- TN: A case is considered to be true negative(TN) if it was projected to be negative.
- TP: True Positive(TP) refers to a case that was both positive and projected to be positive.
- FN: False Negative, or FN, is when a case was positive but was expected to be negative.
- FP: False Positives(FP) occur when a case was negative but was expected to be positive.

4.4.1 Precision:

Precision is the capacity of a classifier to avoid classifying as positive anything that is negative. It is described for each class as the proportion of true positives to the total of true and false positives. Precision is the accuracy of positive predictions and is calculated using Equation 4.1.

$$Precision = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \quad (4.1)$$

4.4.2 Recall:

The capacity of a classifier to locate every successful instance is known as recall. It is described as the proportion of true positives to the total of true positives and false negatives for each class. Recall is the fraction of positives that were correctly identified and is calculated using Equation 4.2.

$$Recall = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \quad (4.2)$$

4.4.3 F1 score:

The best F1 score is 1.0, while the lowest is 0.0. It is a weighted harmonic mean of precision and recall. F1 scores typically perform worse than accuracy measures because they incorporate precision and recall into their computation. It is often recommended to compare classifier models using the weighted average of F1, rather than overall accuracy. The formula for the F score is given by Equation 4.3.

$$F - Score = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4.3)$$

4.5 Dataset Description

Our classifiers were tested and trained using 1301 videos of human action. The three classifications that make up our dataset are "Reading" "Writing," and "Using Smartphone". In the dataset, we divided our dataset into two parts training and testing set. Training is also divided into two parts train and val test. we split this into 90% training and 10% validation. The figures below show a representation of samples in each class.

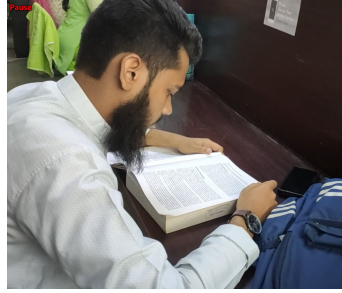


Figure 4.1. Reading



Figure 4.2. Writing



Figure 4.3. Using Smartphone

Table 4.2. Distribution of Samples of Our Dataset

Category	Training Samples	Validation Test Samples	Test Samples
Reading	325	37	76
Writing	325	37	87
Using Smartphone	325	37	52

4.6 Parameter Calculation of Our Best Model VGG16

Table 4.3. Paramter Table Summary of VGG16

Layer Type	Parameters
block1_conv1 (Conv2D)	1792
block1_conv2 (Conv2D)	36928
block1_pool (MaxPooling2D)	0
block2_conv1 (Conv2D)	73856
block2_conv2 (Conv2D)	147584
block2_pool (MaxPooling2D)	0
block3_conv1 (Conv2D)	295168
block3_conv2 (Conv2D)	590080
block3_conv3 (Conv2D))	590080
block3_pool (MaxPooling2D)	0
block4_conv1 (Conv2D)	1180160
block4_conv2 (Conv2D)	2359808
block4_conv3 (Conv2D)	2359808
block4_pool (MaxPooling2D)	0
block5_conv1 (Conv2D)	2359808
block5_conv2 (Conv2D)	2359808
block5_conv3 (Conv2D)	2359808
block5_pool (MaxPooling2D)	0

Convolutional Layer 1 of Block 1

Shape of filter = 3×3

Number of filters= 64

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

$$= ((3 \times 3 \times 3 + 1) \times 64)$$

$$= 1792$$

Convolutional Layer 2 of Block 1

Shape of filter = 3×3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

$$= ((3 \times 3 \times 64 + 1) \times 64)$$

$$= 36928$$

Maxpooling Layer

Pooling layer 1 has no parameters.

Convolutional Layer 1 of Block 2

Shape of filter = 3×3

Number of filters= 128

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)
= ((3×3×64+1) ×128)
= 73856

Convolutional Layer 2 of Block 1

Shape of filter = 3 × 3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)
= ((3×3×128+1) ×128)
= 147584

Maxpooling Layer

Pooling layer 1 has no parameters.

Convolutional Layer 1 of Block 3

Shape of filter = 3 × 3

Number of filters= 256

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)
= ((3×3×128+1) ×256)
= 295168

Convolutional Layer 2 of Block 3

Shape of filter = 3 × 3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)
= ((3×3×256+1) ×256)
= 590080

Convolutional Layer 3 of Block 3

Shape of filter = 3 × 3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)
= ((3×3×256+1) ×256)
= 590080

Maxpooling Layer

Pooling layer 1 has no parameters.

Convolutional Layer 1 of Block 4

Shape of filter = 3 × 3

Number of filters= 512

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)
= ((3×3×256+1) ×512)

= 1180160

Convolutional Layer 2 of Block 4

Shape of filter = 3×3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

= $((3 \times 3 \times 512 + 1) \times 512)$

= 2,359,808

Convolutional Layer 3 of Block 4

Shape of filter = 3×3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

= $((3 \times 3 \times 512 + 1) \times 512)$

= 2,359,808

Maxpooling Layer

Pooling Layer 1 has no parameters.

Convolutional Layer 1 of Block 5

Shape of filter = 3×3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

= $((3 \times 3 \times 512 + 1) \times 512)$

= 2,359,808

Convolutional Layer 2 of Block 5

Shape of filter = 3×3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

= $((3 \times 3 \times 512 + 1) \times 512)$

= 2,359,808

Convolutional Layer 3 of Block 5

Shape of filter = 3×3

Parameters = ((shape of width of filter*shape of height filter*number of filters in the previous layer+1)*number of filters)

= $((3 \times 3 \times 512 + 1) \times 512)$

= 2,359,808

Maxpooling Layer

Pooling Layer 1 has no parameters.

4.6.1 Results and Experiments VGG16 of Approach

Table 4.4. Experiments of the VGG16 Approach

Convenet	Dense	Dense 1	LSTM	Dense 2	Dense 3	Learning Rate	Optimizer	Test Accuracy
VGG16	256	256	96,96	96	3	0.00001	Adam	84%
VGG16	256	256	64,64	96	3	0.00001	Adam	83%
VGG16	256	256	96,96	96	3	0.00001	Adam	74%
VGG16	256	256	96,96	96	3	0.00001	RmsProp	80%
VGG16	256	256	96,96	96	3	0.00001	RmsProp	35%

Table 4.4 describes all the hyper tuning and experiments of our VGG16 models. The best result we got from the VGG16 model is an 84% accuracy level with the first two dense layers consisting of 256 and 256 hyperparameters values, Bidirectional LSTM value is 96, 96 and the other two dense layer values are 96 and 3 which is dense 2 and dense 3. We use Adam as our optimizer. In our many experiments, the dense and dense 2 hyperparameters values vary from 256 to 96. Dense 1 and Dense 3 hyperparameters values vary from 256 to 3. Bidirectional LSTM values vary from 64 to 256. We use two different types of optimizers RMSProp, and Adam.

Training Accuracy vs. Test Accuracy of The Best VGG16 Model From The Experiment Table

<matplotlib.legend.Legend at 0x7fa69fc57d90>

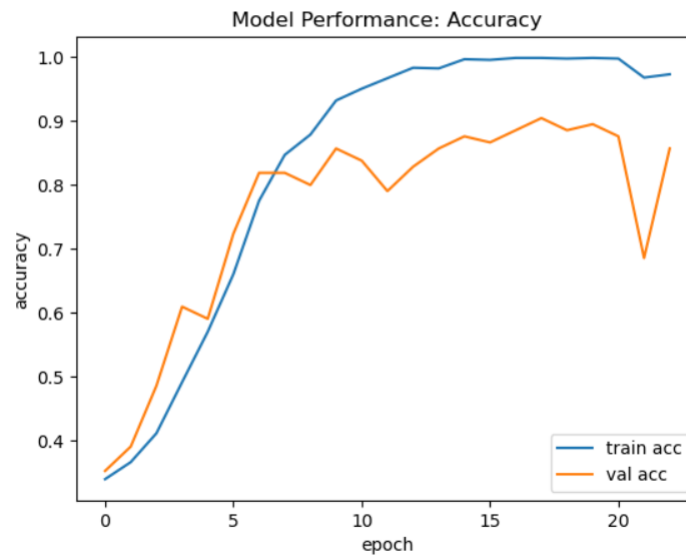


Figure 4.4. Training Accuracy vs Test Accuracy of The Best VGG16 Model

The figure 4.4 shows that while our model performs well during training, it yields the greatest test accuracy of 84% when applied to new data, which it overfits.

Precision and Recall Calculation From The Confusion Matrix

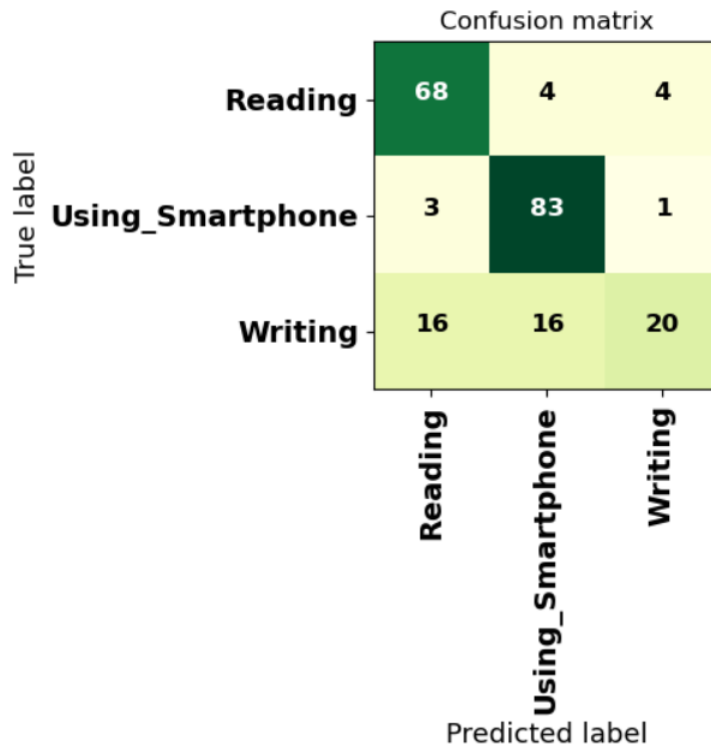


Figure 4.5. Confusion Matrix of The Best VGG16 Model + Bidirectional LSTM Model

Table 4.5. Calculations of Confusion Matrix from

Class	TP	FN	FP
Reading	68	8	19
Writing	20	32	5
Using Smartphone	83	4	20

Table 4.7 Shows us the calculations of the confusion matrix of the best VGG16 model. Now we calculate Precision and Recall below:

True Positive of Class Reading (TP): 68

False Negative of Class Reading (FN): 8

False Positive of Class Reading (FP): 19

Precision=(68/68+8)=0.78

Recall=(68/68+19)=0.89

True Positive of Class Writing (TP): 20

False Negative of Class Writing (FN): 32

False Positive of Class Writing (FP): 5

Precision=(20/20+5)=0.8

Recall=(20/20+32)=0.38

True Positive of Class Using Smartphone (TP):83

False Negative of Class Using Smartphone (FN):4

False Positive of Class Using Smartphone (FP):20

Precision=(83/83+20)=0.80

Recall=(83/83+4)=0.95

4.6.2 Results and Experiments InceptionV3 of Approach

Table 4.6. Experiments of the InceptionV3 Approach

Convenet	Dense	Dense 1	LSTM	Dense 2	Dense 3	Learning Rate	Optimizer	Test Accuracy
InceptionV3	128	128	80,80	128	3	0.00001	Adam	69%
InceptionV3	128	128	64,64	128	3	0.00001	Adam	75%
InceptionV3	128	128	64,64	128	3	0.00001	Adam	78%
InceptionV3	128	128	64,96	128	3	0.00001	Adam	73%
InceptionV3	128	128	80,80	128	3	0.00001	Adam	67%

Table 4.6 describes all the hyper-tuning and experiments of our InceptionV3 models. The best result we got from the InceptionV3 model is a 78% accuracy level with the first two dense layers consisting of 128 and 128 hyperparameters values, Bidirectional LSTM value is 64, and the other two dense layer values are 128 and 3 which is dense 2 and dense 3. We use Adam as our optimizer. In our many experiments, the dense and dense 2 hyperparameters values vary from 1024 to 128. Dense 1 and Dense 2 hyperparameters values vary from 128 to 128. Bidirectional LSTM values vary from 64 to 256. We use optimizers, Adam.

Training Accuracy vs. Test Accuracy of The Best InceptionV3 Model From The Experiment Table

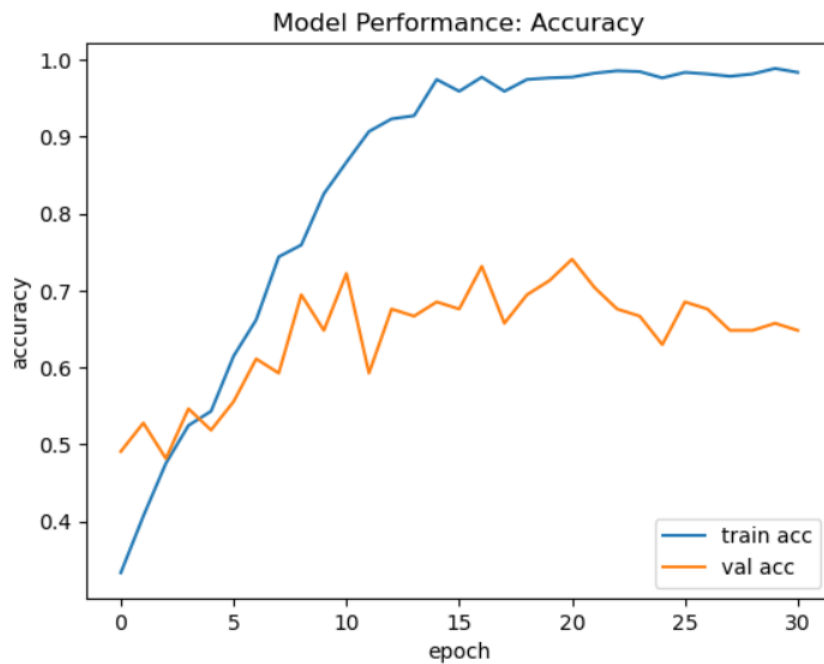


Figure 4.6. Training Accuracy vs Val Accuracy of The Best InceptionV3 Model

The figure 4.6 shows that while our model performs well during training, it yields the greatest accuracy of 78% when applied to new data, which it over-fits.

Precision and Recall Calculation From The Confusion Matrix

		Confusion matrix		
True label	Reading	63	7	6
	Using_Smartphone	7	76	4
	Writing	13	10	29
		Reading	Using_Smartphone	Writing
		Predicted label		

Figure 4.7. Confusion Matrix of The Best InceptionV3 + Bidirectional LSTM Model

Table 4.7. Calculations of Confusion Matrix from

Class	TP	FN	FP
Reading	63	13	20
Writing	29	23	10
Using Smartphone	76	11	17

Table 4.7 Shows us the calculations of the confusion matrix of the best InceptionV3 model. Now we calculate Precision and Recall below:

True Positive of Class Reading (TP): 63

False Negative of Class Reading (FN): 13

False Positive of Class Reading (FP): 20

Precision=(63/63+20)=0.76

Recall=(63/63+13)=0.83

True Positive of Class Writing (TP): 29

False Negative of Class Writing (FN): 23

False Positive of Class Writing (FP): 10

Precision=(29/29+10)=0.74

Recall=(29/29+23)=0.56

True Positive of Class Using Smartphone (TP):76

False Negative of Class Using Smartphone (FN):11

False Positive of Class Using Smartphone (FP):17

Precision=(76/76+17)=0.82

Recall=(76/76+11)=0.87

CHAPTER 5

CONCLUSION

Human Action Recognition is the most interest field for researchers over last ten years. Video-based human activity recognition area is growing day by day. This is the system to automatically recognize and interpret human actions from visual data, from video. In every environment this action varies from man to man. So the application of deep learning method need to implement all the environment. For educational institute, library environment is most important zone where no regarding work.

To address this issue, we implement deep learning to recognize human action in library environment. we have developed a dataset of 1301 video clips in this article. We classify three human action including Reading, Writing and Using Smartphone.

To determine the optimal model for our dataset, we have created deep learning models, including VGG16, DenseNet, and InceptionV3. These models have been trained on the dataset and then hyper-tuned to optimize their performance. We have used performance assessments, including precision, recall, F1-score, and accuracy, to examine the performance.

The model based on VGG16 with a Bidirectional LSTM has demonstrated the highest level of test accuracy among the evaluated approaches which is 84%. This result indicates that our proposed model can effectively predict human action properly. The model has been trained on a diverse dataset to predict correctly.

To implement human action recognition in libraries, it's important to consider ethical and

privacy concerns when implementing such technology, ensuring that user consent and data protection are properly addressed. This technology can be applied to various scenarios, such as monitoring library spaces, understanding user behavior, enhancing security, and improving services.

In conclusion, the proposed dataset of 1301 video clips and the associated deep learning model have the potential to recognize how human act in library environment.

5.1 Limitation

We generated and experimented our model on a small set of dataset. Our model can't work with multi action recognition. Our experimented model got low accuracy, which requires more experiment on large dataset. Some actions share similar visual patterns specially writing, leading to misclassifications.

5.2 Future Work

In future we will include more human actions in other models like MobileNet, Image dataset improving accuracy, especially in the writing category handling more complex data, better handling of long-term dependencies. We would like to work with multi action label also.

BIBLIOGRAPHY

- [1] M. Koperski, “Human action recognition in videos with local representation,” Ph.D. dissertation, COMUE Université Côte d’Azur (2015-2019), 2017.
- [2] M. V. d. Silva, “Human action recognition based on spatiotemporal features from videos,” 2020.
- [3] A. Bux, *Vision-based human action recognition using machine learning techniques*. Lancaster University (United Kingdom), 2017.
- [4] C. Liang, J. Lu, and W. Q. Yan, “Human action recognition from digital videos based on deep learning,” in *Proceedings of the 5th International Conference on Control and Computer Vision*, 2022, pp. 150–155.
- [5] S. Umakanthan, “Human action recognition from video sequences,” Ph.D. dissertation, Queensland University of Technology, 2016.
- [6] M. M. Ullah, “Supervised statistical representations for human action recognition in video,” Ph.D. dissertation, Université Européenne de Bretagne, 2012.
- [7] P. T. Biliński, “Human action recognition in videos,” Ph.D. dissertation, Université Nice Sophia Antipolis, 2014.
- [8] L. Romeo, R. Marani, T. D’Orazio, and G. Cicirelli, “Video based mobility monitoring of elderly people using deep learning models,” *IEEE Access*, vol. 11, pp. 2804–2819, 2023.
- [9] M. Marszalek, “A spatio-temporal descriptor based on 3d-gradients,” in *British Machine Vision Conference, 2008*, 2008.

-
- [10] Y. Yang, R. A. Sarkis, R. El Atrache, T. Loddenkemper, and C. Meisel, "Video-based detection of generalized tonic-clonic seizures using deep learning," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 8, pp. 2997–3008, 2021.
 - [11] C. Feichtenhofer, H. Fan, J. Malik, and K. He, "Slowfast networks for video recognition," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6202–6211.
 - [12] K. R. Campbell, S. W. Marshall, J. F. Luck, *et al.*, "Head impact telemetry system's video-based impact detection and location accuracy," *Medicine and science in sports and exercise*, vol. 52, no. 10, p. 2198, 2020.
 - [13] J. Sasadai, N. Maeda, R. Shimizu, *et al.*, "Analysis of team-sport wheelchair falls during the rio 2016 summer paralympic games: A video-based cross-sectional observational study," *BMJ open*, vol. 10, no. 3, e033088, 2020.
 - [14] M. Abdellaoui and A. Douik, "Human action recognition in video sequences using deep belief networks.," *Traitement du Signal*, vol. 37, no. 1, 2020.
 - [15] S. Z. Gurbuz and M. G. Amin, "Radar-based human-motion recognition with deep learning: Promising applications for indoor monitoring," *IEEE Signal Processing Magazine*, vol. 36, no. 4, pp. 16–28, 2019.
 - [16] G. Rajendran, O. T. Lee, A. Gopi, N. Gautham, *et al.*, "Study on machine learning and deep learning methods for human action recognition," 2020.