

# Skateboard Trick Recognition through an AI-based Approach

**Kris Saliba**

Supervisor: Dr. Joseph Bonello

June 2024

*Submitted in partial fulfilment of the requirements  
for the degree of Bachelor of Science in Information Technology (Honours)  
(Software Development).*



**L-Università ta' Malta**  
Faculty of Information &  
Communication Technology

# Abstract

# Acknowledgements

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>Contents</b>	<b>v</b>
<b>List of Figures</b>	<b>vi</b>
<b>List of Tables</b>	<b>vii</b>
<b>List of Abbreviations</b>	<b>1</b>
<b>Glossary of Symbols</b>	<b>1</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	2
1.2 Hypothesis . . . . .	2
1.3 Research Questions . . . . .	2
1.4 Aims and Objectives . . . . .	3
1.4.1 Aims . . . . .	3
1.4.2 Objectives . . . . .	3
1.5 Structure . . . . .	3
<b>2 Background</b>	<b>4</b>
2.1 Skateboard Tricks . . . . .	4
2.2 Machine Learning . . . . .	4
2.3 Activity Recognition . . . . .	5
2.4 Neural Networks . . . . .	5
2.4.1 Artificial Neural Networks . . . . .	5
2.4.2 Convolutional Neural Networks . . . . .	6
2.4.3 Recurrent Neural Networks . . . . .	7
2.5 Optical Flow . . . . .	9
2.6 Dimensionality Reduction . . . . .	9

<b>3</b>	<b>Literature Review</b>	<b>10</b>
3.1	Activity Recognition . . . . .	10
3.1.1	Challenges in this field . . . . .	10
3.1.2	Preprocessing techniques . . . . .	11
3.1.3	Activity Recognition Techniques . . . . .	12
3.2	Advancements in Skateboard Trick Classification . . . . .	14
3.2.1	Accelerometer-based approach . . . . .	14
3.2.2	Computer Vision-based Approaches . . . . .	15
<b>4</b>	<b>Methodology</b>	<b>17</b>
4.1	Class Establishment . . . . .	17
4.2	Data Preparation . . . . .	17
4.2.1	Dataset . . . . .	17
4.2.2	Labelling techniques . . . . .	17
4.3	Preprocessing . . . . .	18
4.3.1	Frame Extraction . . . . .	18
4.3.2	Data Augmentation . . . . .	19
4.3.3	Normalisation . . . . .	19
4.3.4	Feature Extraction with Transfer Learning . . . . .	20
4.3.5	Dimensionality Reduction . . . . .	21
4.4	Adapted Models . . . . .	21
4.5	Architecture . . . . .	22
4.6	Evaluation Methods . . . . .	22
<b>5</b>	<b>Implementation</b>	<b>25</b>
5.1	Development Environment . . . . .	25
5.2	Dataset Split and Configuration . . . . .	25
5.3	Frame Extraction using Optical Flow . . . . .	26
5.3.1	Data Augmentation . . . . .	26
5.4	Feature Extraction and Preprocessing for Training . . . . .	28
5.4.1	PCA Dimensionality Reduction . . . . .	28
5.5	Hyperparameter Optimisation . . . . .	29
5.5.1	Training Process and Hyperameter Tuning . . . . .	29
5.6	Callbacks . . . . .	30
5.6.1	Early Stopping . . . . .	30
5.6.2	Model Checkpoint . . . . .	30
5.7	Training History . . . . .	30
<b>6</b>	<b>Evaluation</b>	<b>32</b>
6.1	Model Training Specifications . . . . .	32

6.2	Experiments . . . . .	32
6.2.1	Choice of Optimiser . . . . .	32
6.2.2	The Effect of Data Augmentation . . . . .	32
6.2.3	The Effect of PCA . . . . .	33
6.2.4	Applicability In Real-Time applications . . . . .	34
6.3	Discussion . . . . .	34
6.3.1	Final Findings . . . . .	34
6.3.2	Classification observations . . . . .	35
6.4	Limitations . . . . .	36
<b>7</b>	<b>Conclusions and Future Work</b>	<b>38</b>
7.1	Future Work . . . . .	38
7.2	Conclusion . . . . .	39
<b>A</b>	<b>Sample A</b>	<b>45</b>
<b>B</b>	<b>Sample B</b>	<b>46</b>

# List of Figures

Figure 1.1	Number of skateboarding participants in the United States from 2010 to 2021 in millions. Reproduced from [4], data sourced from Outdoor Foundation [5]. . . . .	1
Figure 2.1	Schematic Representation of an Artificial Neural Network. Reproduced from López et al. (2022) [14] . . . . .	6
Figure 2.2	CNN architecture comprising 5 layers. Reproduced from O'Shea and Nash (2015) [16] . . . . .	7
Figure 2.3	LSTM architecture. Reproduced from ..... . . . .	8
Figure 2.4	BiLSTM Structure. Reproduced from Du et al. (2020) [22] . . . . .	8
Figure 2.5	Demonstration of Optical Flow between Two Consecutive Frames. Reproduced from K. Host and M. Ivašić-Kos [25] . . . . .	9
Figure 3.1	CNN-LSTM Architecture. Reproduced from Donahue et al. (2015) [46]	13
Figure 3.2	add caption and add reference in paragraph . . . . .	15
Figure 3.3	Tricks selected for Chen's models. Reproduced from Chen (2023) [7] .	16
Figure 4.1	Folder-based labelling. . . . .	18
Figure 4.2	Text-based labelling . . . . .	18
Figure 4.3	Comparison of a Single Frame and Multiple Sequential Frames. . . . .	18
Figure 4.4	Visualisation of optical flow: (a) Optical flow representation (b) Individual frames extracted. . . . .	19
Figure 4.5	Artefact Architecture. . . . .	24
Figure 5.1	Comparison of frame extraction between optical flow method and uniform sampling . . . . .	26
Figure 5.2	Comparison of Augmented sequence against original . . . . .	27
Figure 5.3	Augmentation Parameters for all augmentation . . . . .	27
Figure 5.4	Cumulative Variance plots for the outputs of VGG16 and ResNet50 .	29
Figure 5.5	Loss graphs for each architecture . . . . .	31
Figure 6.1	Model loss graphs for (a) Adam and (b) SGD optimisers. . . . .	32
Figure 6.2	Confusion matrices for each architecture. . . . .	36

# List of Tables

Table 4.1 Characteristics of Transfer Learning Models VGG17 and ResNet50 . . . 20

Table 6.1 Performance comparison of models with and without augmentation . . 33

Table 6.2 Comparison of VGG16-BiLSTM and ResNet50-BiLSTM model  
performances before and after PCA implementation . . . . . 33

Table 6.3 Evaluation time (ss:mm) for skateboard trick classification models. . . . 34

Table 6.4 Performance comparison of models from Hanciao Chen’s study [7] and  
this study . . . . . 35

Table A.1 Model summary for VGG-BiLSTM . . . . . 45



# 1 Introduction

Skateboarding dates back to the late 1940s or early 1950s, evolving from a leisure activity for surfers on flat land to a worldwide phenomenon [1]. This worldwide appeal was further amplified by its inclusion as an official sport in the 2020 Tokyo Olympic Games [2]. This event had a significant impact on the sports popularity globally, but particularly evident in the United States, as demonstrated by the spike in the number of U.S. skateboarding participants between the years 2020 and 2021, illustrated in Figure 1.1.

This dynamic sport encompasses various disciplines and riding styles, each offering unique challenges for skateboarders to explore. Two of the most prominent styles are “vert” and “street.” Vert skateboarding involves riding on specialised obstacles, namely, half-pipes and ramps, focusing on aerial manoeuvres. Street skateboarding takes place in urban environments, utilising various obstacles that can be found outdoors, including stairs, rails, ledges, gaps or flat ground for skaters to showcase their creativity [3].

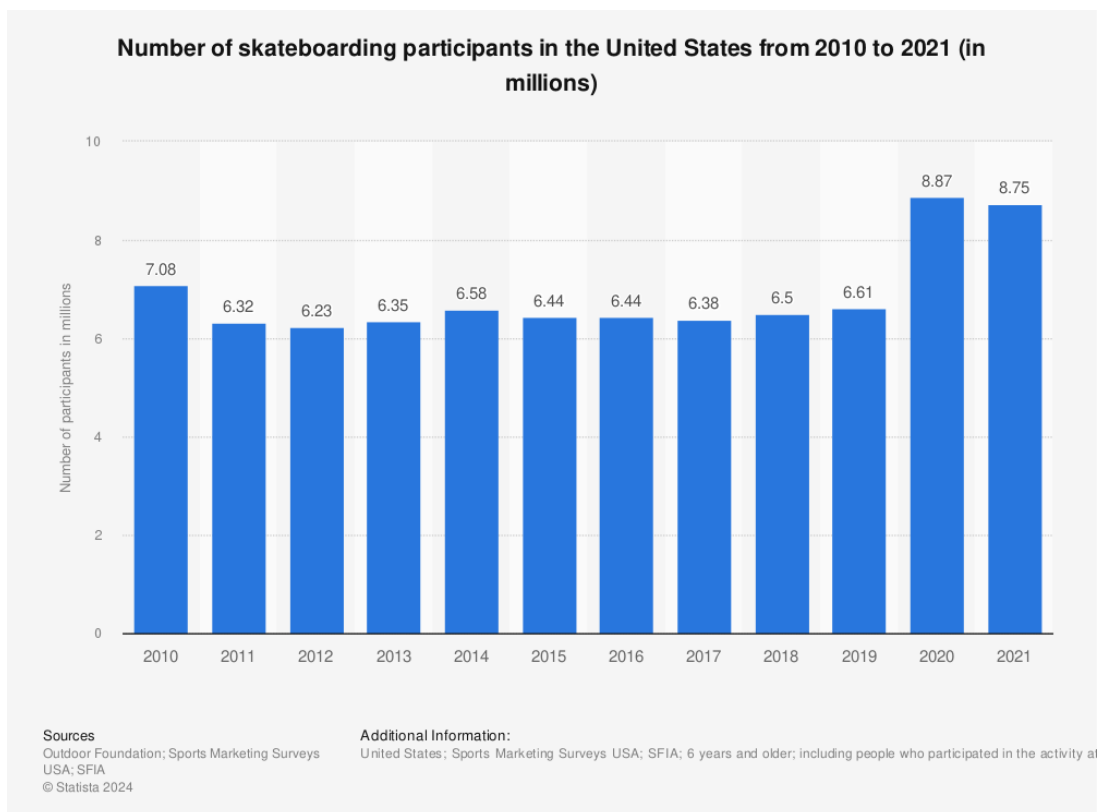


Figure 1.1 Number of skateboarding participants in the United States from 2010 to 2021 in millions. Reproduced from [4], data sourced from Outdoor Foundation [5].

## 1.1 Motivation

The recent surge in skateboarding's popularity across digital platforms like YouTube and Instagram, there is potential for innovative technologies to enhance how viewers engage with skateboarding content. One promising innovation could be the implementation of on-screen trick identification. This feature would allow inexperienced viewers to identify tricks performed by skaters in videos through a digital overlay that labels each manoeuvre.

Traditional methods require manually labelling each trick in a video, a process that is not only time-consuming but also error-prone. Automating this process could save time in this regard, but also offer real-time benefits for live-streamed content, such as skateboard competitions and events. This automation could provide insightful footage to viewers, further enhancing their understanding and appreciation of the sport. An Artificial Intelligence (AI) model capable of accurately classifying skateboard tricks could significantly improve viewer experience in these scenarios.

Despite the popularity of computer vision, particularly in video activity recognition, limited research exists in the domain of skateboard trick classification. While some initial exploration by Shapiee et al. (2020) [6] and Hanxiao Chen (2023) [7] have emerged, a gap in this area persists. This significant lack of research in this area, drives the need for further investigation in this niche field.

## 1.2 Hypothesis

The hypothesis for this research asserts that employing a combination of deep learning strategies and preprocessing techniques can improve the accuracy and robustness of classifying skateboard tricks.

## 1.3 Research Questions

- How effectively can deep learning techniques accurately classify skateboard tricks from video data?
- What is the impact of different video pre-processing techniques on classification accuracy?

## 1.4 Aims and Objectives

### 1.4.1 Aims

- To classify skateboard tricks using images extracted from videos into their respective classes.
- To compare the performance of Deep Learning architectures in the context of skateboard trick classification.

### 1.4.2 Objectives

- To Implement a set of Deep learning architectures and evaluate them using appropriate metrics.
- To Explore and employ suitable frame processing techniques.
- Augment and preprocess the images to determine any performance improvement.

## 1.5 Structure

This study is structured systematically to explore the potential of deep learning in skateboard trick recognition. Chapter 2 begins by establishing the necessary background knowledge required for this study, while Chapter 3 provides an overview of the past research conducted in this domain.

Next, Chapter 4 outlines the approach behind the artefact, with implementation specifics detailed in Chapter 5. Chapter 6, discusses the outcomes and evaluates them against other studies and finally, Chapter 7 presents the possibilities of future work and the final conclusions of this research

## 2 Background

### 2.1 Skateboard Tricks

Skateboard tricks can be described as dynamic manoeuvres that involve complex coordination of the skateboard and the skateboarder's body. The key to successfully performing these tricks is appropriate foot placement, which is critical for controlling the skateboard's speed and direction. This control allows skateboarders to manipulate the board in ways that replicate specific tricks, showcasing their technical abilities and creativity. Some of the most common and simple flat ground tricks are listed as the following:

- **Ollie:** One of the first tricks beginners learn. Where the skateboarder pops the tail of the board while simultaneously sliding their foot across the nose of the board, causing the board to level out in the air, used to jump over obstacles.
- **Kickflip:** A trick where the skateboarder flips the board under their feet while jumping, making it spin  $360^\circ$  around the x-axis.
- **Pop-Shuvit:** A trick where the skateboarder scoops the board with their back foot causing a  $180^\circ$  rotation around the y-axis.

Skateboarders continually innovate and come up with new trick combinations, contributing to the dynamic nature of the sport.

### 2.2 Machine Learning

Machine Learning (ML) can be defined as a field of study that explores algorithms and statistical models employed by computer systems to execute tasks without the need to be explicitly programmed. It is particularly applicable in situations where the information we seek from a dataset is not interpretable, and as the volume of available datasets continues to surge so does the demand for machine learning [8].

Morris (2019) [9] characterises ML as the advancement of algorithms that progressively enhance their performance through practice, suggesting that the more training the learning algorithm undergoes, the better it becomes at executing tasks. Numerous critical factors shape a model's performance within this phase, as exemplified by Budach et al. (2022) [10]. Such factors include dataset quality and diversity, data preprocessing, the selection of a suitable model architecture, training time and the fine-tuning of hyper-parameters. The three main categories for ML models are defined as the following [8]:

- **Supervised:** This is a ML concept that involves training a model to make classifications based on input data that has been labelled with the correct label.
- **Unsupervised:** This ML concept concentrates on discovering relationships within data when there are no predefined "correct" answers or labelled examples to guide the learning process. These models are left to autonomously explore and divulge structures in the data.
- **Reinforcement:** This type of learning consists of an agent that interacts with the environment and learns from the continuous feedback it receives in the form of rewards or punishment.

## 2.3 Activity Recognition

Activity recognition is the process of identifying and categorizing human activities from video sequences. Human activities involve a wide range of motions and interactions with objects, varying from simple isolated actions like dancing to more complex activities that engage multiple body parts and external objects such as football matches. The human ability to perceive these behaviours is a trivial task; yet, it is a challenging problem for computers due to the sequential nature and the resemblance of visual content in such activities [11, 12].

## 2.4 Neural Networks

### 2.4.1 Artificial Neural Networks

Artificial Neural Networks (ANNs) are a class of Machine Learning models that are inspired by the interconnected systems of neurons found in the nervous system of living organisms. They consist of connected nodes capable of learning from their environment and adapting to complex patterns in data [13]. Figure 2.1 presents a schematic representation of an ANN. The diagram is organised into three fundamental layers: the Input Layer, the Hidden Layer(s) and the Output Layer [14].

- **Input Layer:** This is the set of neurons that serve as the initial entry point for external data. Each input neuron in this layer corresponds to a specific feature or variable used in the Neural Network model.
- **Hidden Layer(s):** This is the set of neurons that are located between the Input and Output Layers where the network captures complete non-linear behaviours of data and feature transformations.

- **Output Layer:** This is the set of neurons that provide the final predictions produced by the Neural Network. Depending on how the ANN is configured, the final output can be continuous, binary, ordinal, or count.

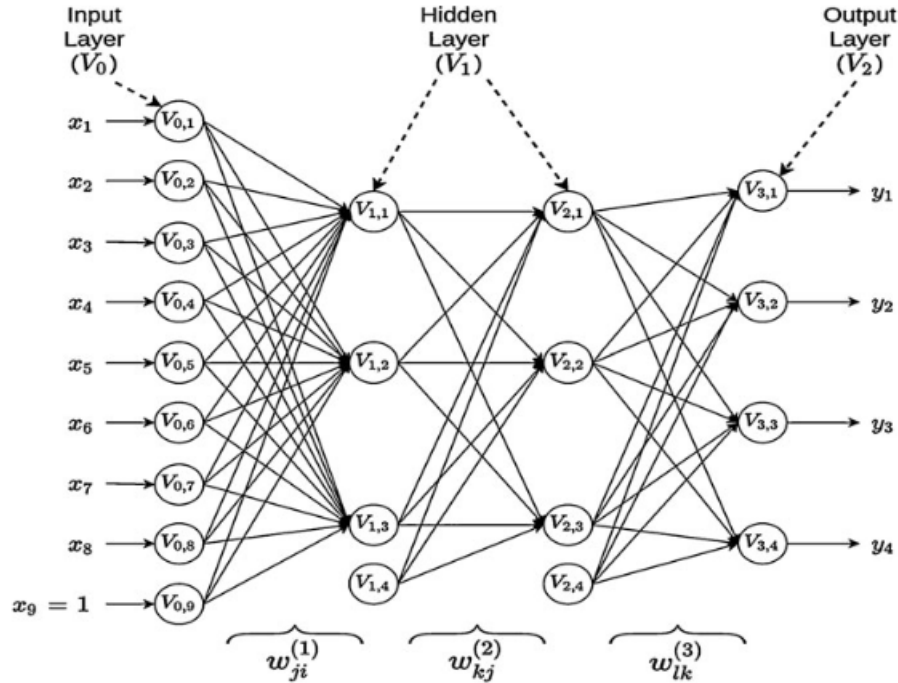


Figure 2.1 Schematic Representation of an Artificial Neural Network. Reproduced from López et al. (2022) [14]

## 2.4.2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs), as described by Gu et al. (2018) [15] are a category of Deep learning architectures with roots in the biological visual perception mechanisms of living organisms. These networks have gained widespread attention for their incredible performance in the field of image recognition and pattern recognition tasks.

CNNs are comprised of three types of layers: convolutional layers, pooling layers and fully-connected layers. The convolutional layer applies a set of filters (or kernels) that slide across the input data performing a localised dot product between their weights and the corresponding values of the input data. This process effectively extracts features from images, critical in understanding complex patterns in images. The results are summed up to generate a single value in the feature map. The pooling layer applies an aggregation function such as max pooling or average pooling to create a downsampled representation of the input data. Finally, the fully-connected attempt to transform the outputs from the previous layer into an output vector that represents a score corresponding to a class label [16, 17].

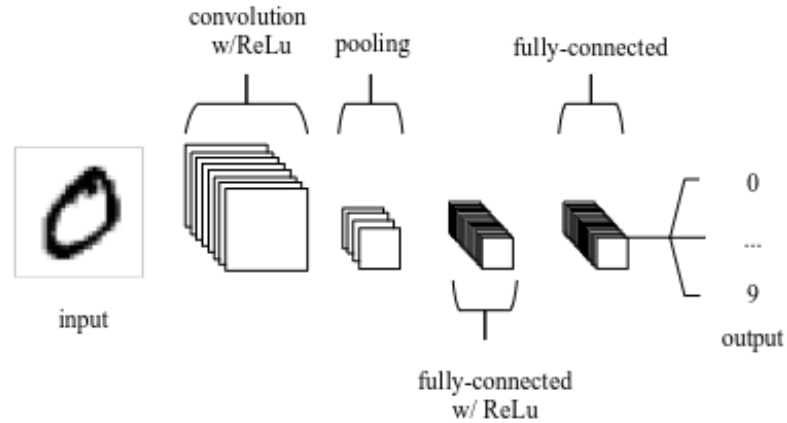


Figure 2.2 CNN architecture comprising 5 layers. Reproduced from O'Shea and Nash (2015) [16]

Figure 2.2 showcases a simplified CNN architecture consisting of 5 layers. Nonetheless, the complexity of CNNs can be scaled by stacking multiple layers, thereby increasing the network's depth to cater for more complex tasks.

### 2.4.3 Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are a subset of Neural networks that are designed for sequential data processing. RNNs are capable of modelling dynamic relationships in sequential data by feeding signals from past time steps back into the network. However, they are limited due to their inability to access long-term data, limited to approximately ten sequential time steps [9]. This constraint arises from the challenge of dealing with vanishing or exploding gradients in the backpropagation procedure while processing extended sequences, as discussed in prior works [18].

#### Long Short-Term Memory Networks

To address the limitation that RNNs encounter in capturing long-term dependencies, Long Short-Term Memory Recurrent Neural Networks (LSTM-RNNs) were introduced as an extension to RNNs in 1997 [19, 20]. LSTM networks automatically the memory of RNNs, allowing them to capture dependencies across more than 1,000 time steps depending on the network's complexity. These models are also able to tackle the vanishing problem associated with RNNs by introducing gates that regulate data flow. These gates maintain a constant error flow through structures called constant error carousels (CECs) to ensure that gradients do not vanish during backpropagation [9].

LSTMs consist of three gates: input, output and forget gates. The input gate, determines whether new information will be uploaded to the network, the output gate manages whether current cell values contribute to the output, and the forget gate

determines whether existing information will be preserved or removed [21]. Figure 2.3 illustrates the architecture of an LSTM network, displaying the interaction between the cell state and the various other gates that regulate the flow of information through the network.

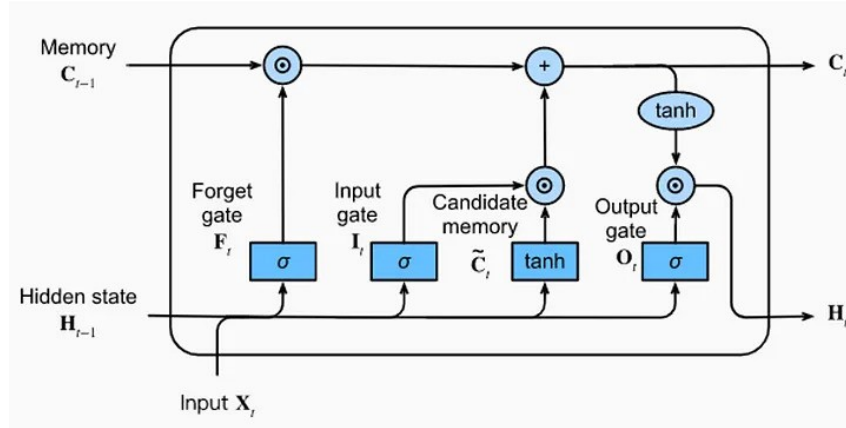


Figure 2.3 LSTM architecture. Reproduced from .....

### Bidirectional LSTMs

Bidirectional LSTMs (BiLSTMs) are variants of LSTMs designed to enhance the ability to capture patterns in sequential data and improve learning long-term dependencies. Unlike LSTMs that only process data in a single direction, BiLSTMs utilise a different approach by applying two LSTM models to the input data. In the first round, the LSTM is applied to the input data, and in the second round, it is applied to the reversed input sequence. Furthermore, Tavakoli et al. (2019) [21], concluded that BiLSTMs reported better accuracies compared to regular LSTMs.

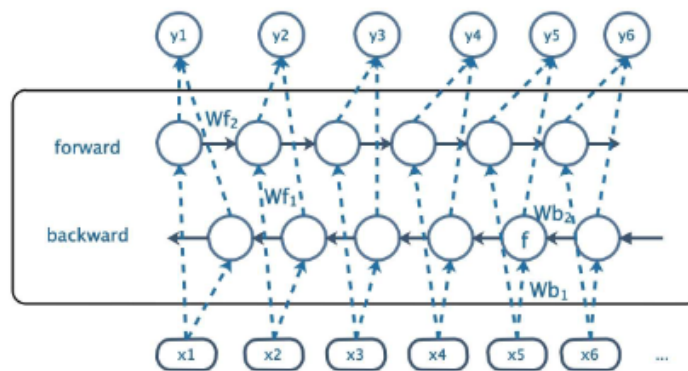


Figure 2.4 BiLSTM Structure. Reproduced from Du et al. (2020) [22]

Figure 2.4 demonstrates the structure of a BiLSTM network, showcasing the flow of data between the two parallel LSTM layers. The first LSTM processes the input data in the forward direction, from  $x_1$  to  $x_6$ , while the other LSTM in the backwards direction from  $x_6$  to  $x_1$ , allowing the network to learn from past and future contexts. The forward



LSTM units are connected through weights  $Wf_1$  and  $Wf_2$ , while the backward LSTM units are connected through weights  $Wb_1$  and  $Wb_2$ . By learning from sequences in both directions, BiLSTMs are particularly effective in scenarios where past and future contexts are crucial.

## 2.5 Optical Flow

Optical flow is an approximation technique used to estimate how objects move across a series of images. It works by analysing the changes in brightness of pixels, which often occur due to the movement of objects. The application of Optical flow generates a two-dimensional vector field where each vector is a displacement vector, representing the movement of a point in the image, showing both the direction and distance of that movement from one frame to the next. [23, 24]. Figure 2.5 demonstrates an optical flow visualisation between two frames, the lines appearing on the image indicate the direction and magnitude of their movement between the frames.



Figure 2.5 Demonstration of Optical Flow between Two Consecutive Frames.  
Reproduced from K. Host and M. Ivašić-Kos [25]

## 2.6 Dimensionality Reduction

Datasets with many features are characterised as high dimensional. Such datasets often include redundant data, including closely related or identical variables. The goal of dimensionality reduction is to remove these unnecessary components and create a simplified lower-dimensional feature space. This helps reduce the impact of redundant data by mapping the valuable data from the original features to a smaller set of features [26].

PCA is a widely used statistical technique that transforms features into a lower-dimensional space by transforming the original features into a new set of variables known as principal components. These principal components are orthogonal to each other, meaning that they are uncorrelated to one another and are reorganised in a way that the first few components represents the maximum variance from the original data [27].

## 3 Literature Review

### 3.1 Activity Recognition

Building on the foundational understanding of activity recognition (AR) defined in the background, it's important to acknowledge the impact it has on various sectors. The recent advancements in human action recognition have not only revolutionised sectors such as healthcare and security as demonstrated in the studies [28], [29], but also hold extensive potential in the realm of sports.

The research conducted by K. Host and M. Ivašić-Kos (2022) in [25] along with Wu et al. (2022) in [30] further elaborate on this potential, such as its numerous applications in categorising complex sports actions, injury prevention methods and refinement of game strategies through video analysis. These studies also propose various methodological advancements, including the utilisation of Deep Learning models to enhance accuracy, the exploration of multimodal data sources for sports activities and an emphasis on real-time analysis capabilities. These insights provide valuable techniques that can be leveraged for the development of AR models in the field of sports.

The study by Beddiar et al. [31], identifies two main streams of Human-Computer Interaction (HCI) technologies: Contact-based and Vision-based systems. The authors categorise contact-based HCI as those technologies that require physical user interaction through mediums such as accelerometers, wearable sensors and multi-touch interfaces. Alternatively, the authors describe vision-based methods as the simplification of HCI due to more natural human communication, eliminating the need for physical contact or equipment. These methods use image and video data to recognise human activities offering an advantage in terms of societal acceptance and usability.

#### 3.1.1 Challenges in this field

The domain of AR comes with many challenges as described by Zhang et al. (2017) [32]. The researchers suggest that while certain scenarios utilise static cameras such as surveillance systems, most situations that benefit from AR adopt dynamic recording devices such as mobile phones and sports event broadcasts. These dynamic devices introduce a significant level of complexity as a result of their tricky dynamic backgrounds present in the video footage. Zhang et al. also point out specific difficulties posed by long-distance and low-quality videos, often encountered in environments like crowded public spaces and sports events. The camera distance, results in smaller subjects, making a detailed analysis of human movements more challenging while lower-quality videos further complicate the task for Human Activity Recognition (HAR) systems

The challenges highlighted by Zhang et al. [32] in the domain of AR are directly relevant to the development of a skateboard trick classifier. In scenarios like televised skateboarding events, skaters may appear relatively small to accommodate the entire skatepark. This factor, along with the complexity of the background as a result of dynamic recording, poses a challenge for trick recognition in live broadcasts. Furthermore, if a skateboard trick classifier is intended for personal development, then it may encounter videos of lower quality, again complicating the task of trick recognition. Addressing these challenges is crucial for the development of a skateboard trick classifier that performs well in real-world conditions.

### 3.1.2 Preprocessing techniques

Preprocessing is a crucial step in the development of ML models, especially in the field of AR. It involves the application of various techniques to enhance raw data before applying Machine Learning algorithms.

#### Data Augmentation

As a result of the limited research in the domain of skateboard trick classification, there is a significant lack of open-source datasets featuring skateboard tricks. This lack of data presents an opportunity to employ data augmentation techniques, especially valuable in studies with limited data. Methods such as flipping, rotating, scaling and colour manipulation not only artificially enhance the size of the original dataset but also lower the likelihood of overfitting as highlighted in prior works, [33, 34]. In the realm of Skateboard trick classifiers, Shapiee et al (2020) [6] effectively employed data augmentation techniques to expand their dataset, demonstrating the application of these methods in improving model performance for trick classification.

#### Feature Extraction

Xudong Jiang (2009) [35], describes feature extraction as the process of capturing the core attributes of an object through the elimination of redundancies, resulting in a set of numerical features ideal for classification. The study by Manjunath Jogin et al. (2018) [36] highlights the effectiveness of CNNs at extracting features due to their capabilities in detecting complex patterns and enhancing features. In this study, a CNN not only performs feature extraction but also performs classification, achieving an accuracy of 86% on the CIFAR-10 [37] dataset.

## Transfer Learning

Following the discussion on CNNs for feature extraction, it is important to highlight the role of transfer learning in enhancing ML models. The concept of transfer learning as detailed by the studies by Weiss et al. (2016) [38] and Soni et al. (2018) [39] involves using a pre-trained ML model to leverage its experience for a new, but related task. This technique is frequently used for feature extraction using pre-trained CNNs, especially in scenarios where there is limited data for training or when training a CNN from scratch is not feasible.

The study by Sargano et al. (2017) [40], employed transfer learning using pre-trained deep CNNs like AlexNet [41] and GoogleNet [42], for HAR. Notably, they utilised these models for feature extraction, followed by an SVM classifier for final classification. Their approach showcases the resource-efficient and time-saving nature of transfer learning, evident in their impressive accuracies of 98.15% and 91.47%.

Moreover, research on transfer learning is not unique within the field of skateboard trick classification. Two other studies by Shapiee et al. (2020) [6] and Hancio Chen (2023)[7] have employed this approach and achieved high accuracies. However, a more in-depth analysis of these papers is detailed in section 3.2. These studies strongly suggest exploring the use of various pre-trained models for feature extraction in the development of a skateboard trick classifier.

### 3.1.3 Activity Recognition Techniques

The accurate classification of skateboard tricks poses a unique challenge for computer vision due to the sport's dynamic and complex nature, characterised by rapid movements, potential background noise and camera angles. This section explores various computer vision techniques and architectures employed in previous literature, exploring their potential impact on this specific task.

Deep Learning (DL) techniques have become increasingly popular over traditional ML methods, for their ability to learn feature representations automatically from raw data, significantly improving performance [43]. One particularly effective DL architecture is the CNN-LSTM. This approach leverages the CNNs strength in extracting spatial features from individual frames and LSTMs ability to capture temporal information across frames, leading to a deeper understanding of video content.

The study by Orozco et al. (2020) [44] adopted this approach to test its effectiveness against three AR datasets: KTH, UCF-11, and HMDB-51, particularly focusing on how the number of LSTM units impacted performance. The authors employed transfer learning, utilising the VGG16 model [45] for feature extraction from videos, followed by an LSTM network for classification. Their findings show that 360

LSTM units achieved an accuracy of 93.86% on the KTH dataset, while 320 units led to an accuracy of 91.93%. However, the performance on the HMDB-51 dataset dropped, with 400 LSTM units resulting in a lower accuracy of 47.36%. These findings show the potential of the CNN-LSTM approach, particularly for simpler datasets, highlighting the need for further investigation and optimisation for more complex datasets like HMDB-51.

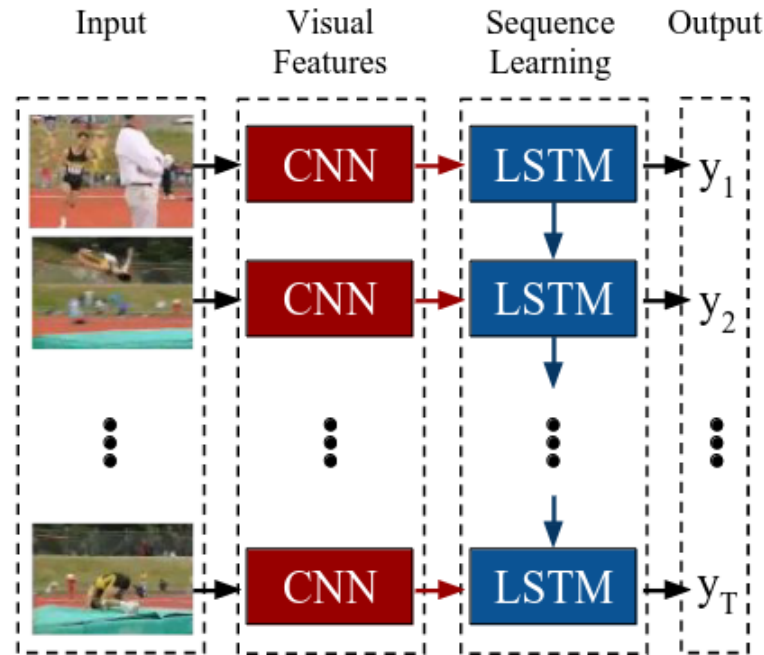


Figure 3.1 CNN-LSTM Architecture. Reproduced from Donahue et al. (2015) [46]

Building upon the foundational CNN-LSTM architecture, a recent study by Saoudi et al. (2023) [47] enhances this model by incorporating three key advancements:

1. **3D Convolutional Neural Networks (3D CNNs):** Unlike regular CNNs, data is processed as 3D volumes, enabling them to capture both spatial and temporal information by performing convolution operations on all three dimensions: width, height and time. The authors opted to use the I3D model, leveraging transfer learning to bypass the resources required to train one from scratch. The I3D model was selected for its proven efficiency and its ability to be fine-tuned for activity recognition tasks.
2. **Bi-directional Long-Short-Term Memory (BiLSTM) network:** The authors chose to use a variant of the LSTM called the Bi-directional LSTM (BiLSTM), to take advantage of its dual-directional approach, enabling it to capture information from both past and future states.
3. **Attention Mechanisms:** Saoudi et al. incorporated attention mechanisms after their BiLSTM, enabling the model to prioritise particular parts of the input data.

This technique allowed the model to learn which temporal features were most applicable to the task, providing a more detailed representation of the input and ultimately, improved performance.

With the integration of these advancements, Saoudi et al. achieved model accuracies of 97.98% and 96.83% on the HMDB51 and UFC101 datasets respectively, demonstrating the potential of 3D CNNs, BiLSTMs and attention mechanisms for activity recognition tasks.

## 3.2 Advancements in Skateboard Trick Classification

In the emergent field of skateboard trick classification, leveraging activity recognition techniques from a video have led to two primary methodologies among researchers. The first technique involves utilising signals obtained from skateboard-mounted accelerometers or artificially generated signals. These signals are then fed into a study-dependent model for classification, as outlined by Abdullah et al. (2021) [48] and Corrêa et al. (2017) [49]. The second approach employs computer vision techniques, leveraging video footage of skateboard tricks to train and refine models for accurate trick identification, as depicted by the studies of Shapiee et al. (2020) [6] and Hanciao Chen (2023) [7].

### 3.2.1 Accelerometer-based approach

The study by Abdullah et al. (2021) [48], makes use of a custom dataset comprising six skateboard tricks most commonly executed in competitive events. Amateur skateboarders performed each trick five times on a modified skateboard equipped with an Inertial Measurement Unit (IMU) to record the signals produced. The researchers capture six signals for each trick; linear accelerations along the  $x$ ,  $y$ , and  $z$  axes ( $a_x$ ,  $a_y$ ,  $a_z$ ) and angular accelerations along the same axes ( $g_x$ ,  $g_y$ ,  $g_z$ ). They then opt for the unique approach of concatenating all six signals onto a single image corresponding to one trick, employing two input image transformations: raw data (RAW) and Continuous Wavelet Transform (CWT).

With the application of six transfer learning models on this data, Abdullah et al. [48] reports exceptionally high accuracies, achieving a 100% test accuracy over multiple models. While these results are remarkable, very high levels of accuracy are rare in ML applications and are typically associated with models that may be overfitting the data.

The study by Corrêa et al. (2017) [49], obtained their sample data by artificially generating 543 signals based on prior research, utilising MATLAB 2015 and Signal Processing Toolbox. These signals were then categorised into five distinct classes

representing different skateboard tricks, each with various samples ranging from 30 to 50 per class, across three axes (X, Y and Z). This study developed and validated individual Artificial Neural Networks (ANNs) for each axis, as well as the combination of the three: ANN XYZ, displaying the potential of Neural Networks to categorise multidimensional skateboard tricks. The ANNs are all multilayer feed-forward neural networks (MFFNNs), structured into three distinct layers. They feature an input layer with 82 neurons, a hidden layer, comprised of 23 neurons utilising a tan-sigmoid transfer function and an output layer consisting of 5 neurons with a softmax function. Finally, the study achieved high accuracies, with ANNs X, Y and Z achieving 94.8%, 96.7% and 98.7%, respectively, while the combined ANN XYZ achieved an accuracy of 92.8%.

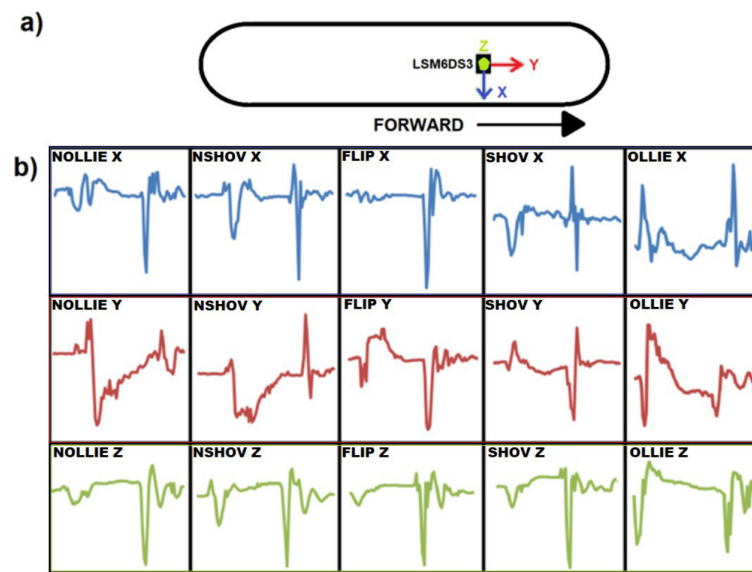


Figure 3.2 add caption and add reference in paragraph

### 3.2.2 Computer Vision-based Approaches

The study by Shapiee et al. (2020) [6] leverages a custom data set comprising videos capturing the execution of five distinct skateboard tricks, each attempted five times. Each video spans two to three seconds, yielding a total of 750 images by extracting 30 frames per video. This study made use of data augmentation techniques to expand their dataset further. Consequently, they introduced an additional 2,250 images, achieving 3,000 images in their data set. Shapiee et al. utilised Transfer Learning (TL) combined with  $k$ -Nearest Neighbour ( $k$ -NN) fusion pipelines. They opted for a pre-trained models like MobileNet, NasNetMobile and NasNetLarge for feature extraction followed by a  $k$ -NN classification algorithm to identify the tricks. The researchers. As a result, this technique demonstrated high classification accuracies, with MobileNet achieving 95%, NASNetMobile 92% and NASNetLarge 90%.

On the other hand, Chen (2023) [7] compiled a comprehensive data set by collecting videos from multiple platforms, including YouTube, Twitter and Instagram. Furthermore, Chen trained the model using 15 fundamental tricks commonly observed in competitive settings, illustrated in Figure 3.3. The researcher collected 50 videos per trick, summing up to a total number of 750 videos. Of these, 45 videos per trick were allocated for training, and the remaining 5 were reserved for validation.

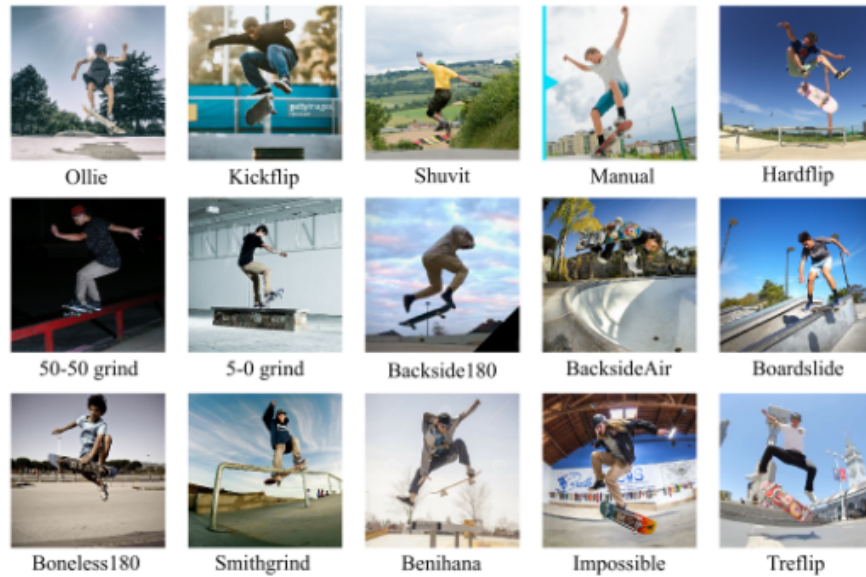


Figure 3.3 Tricks selected for Chen’s models. Reproduced from Chen (2023) [7]

In the abstract by Hanciao Chen [7], extensive experimentation is conducted using diverse models, exploring various combinations of CNN-LSTM and CNN-BiLSTM architectures. The study also incorporated attention mechanisms and explored transfer-based methods for activity recognition. This study further documents and analyses important metrics such as training time, training accuracy and validation accuracy for each model experimented on. Among these, the top three models that stood out in terms of validation accuracy were the ResNet50 with Attention and BiLSTM (84%), ResNet50 with BiLSTM (81%) and ResNet50 with LSTM (80%). Chen’s study provides valuable insight into the application of diverse models in activity recognition in skateboarding.



## 4 Methodology

### 4.1 Class Establishment

This study pursued a multi-class classification strategy, targeting three fundamental skateboard tricks: ollie, pop-shuvit and kickflip. These tricks were selected based on two primary criteria. Firstly, they are often associated with the first tricks learnt by beginners, highlighting their role in foundational skateboarding skills. Secondly, their popularity within the skateboarding community, often performed in competitions emphasises their relevance, making them highly relevant for analysing and evaluating competitive performance.

### 4.2 Data Preparation

#### 4.2.1 Dataset

This study utilised video recordings of skateboarders performing tricks as its primary data source. To ensure model robustness, the final dataset consisted of videos across diverse environmental conditions and varying skateboarder skill levels.

The initial dataset was sourced from the publicly available "SkateboardML" repository on GitHub [50], comprising 200 video clips corresponding to two common tricks: the ollie and the kickflip. To expand the dataset's diversity, additional data was obtained through direct communication with Hanxiao Chen, the author of the SkateboardAI paper [7]. This communication yielded a dataset containing 750 videos covering 15 tricks. However, since the majority of tricks were beyond the scope of this study, only a subset of this data was included into the final dataset. Ultimately, the final dataset comprised of 128 videos per class, totalling 384 videos.

#### 4.2.2 Labelling techniques

This study explores two primary labelling techniques: the folder-based and the text-based approach as illustrated in Figures 4.1 and 4.2. The folder-based method categorises videos into folders named after their corresponding class label offering a simple organisation method. On the other hand, the text-based approach, lists each video's path and corresponding label in a text file, providing more flexibility. Given the limited number of classes and manageable dataset size, this study selected the folder-based approach, considering the extra complexity of the text-based method unnecessary for this project.



Figure 4.1 Folder-based labelling.



Figure 4.2 Text-based labelling

## 4.3 Preprocessing

### 4.3.1 Frame Extraction

Recognising complex human actions requires the examination of sequential data as opposed to relying on single frames or images [11]. To illustrate this point, consider the example of a skateboarder executing a skateboard trick like the "Kickflip", as depicted in Figure 4.3. By feeding an AI model a single image of this trick as shown in Figure 4.3a, it may misinterpret the manoeuvre as another trick due to its subjective nature. Whereas, by providing the model with a sequence of frames, as illustrated in Figure 4.3b, the entire trick sequence is captured including the wind-up, the trick and the landing, capturing actions such as foot placement and board rotation, which together characterise the full skateboard trick.



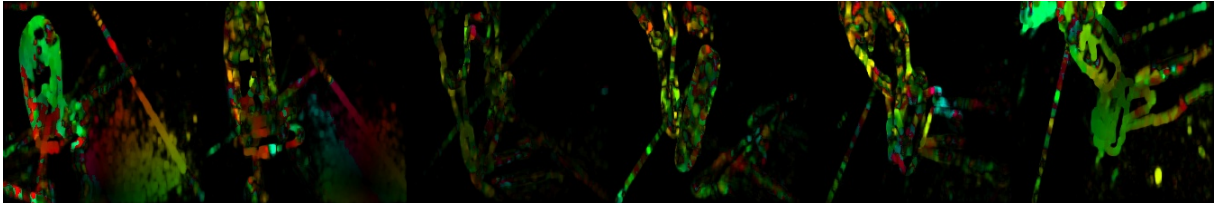
(a) A Single Frame of a "Kickflip".



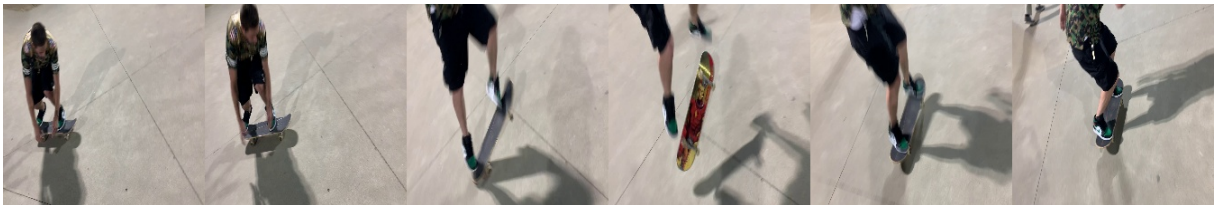
(b) Multiple Frames of a "Kickflip".

Figure 4.3 Comparison of a Single Frame and Multiple Sequential Frames.

This study explored two primary frame extraction techniques: uniform sampling and optical flow method. The former technique evenly splits a video into its corresponding frames by using a pre-calculated step size determined by the number of frames of a video. Whereas, the second technique leverages optical flow to select the frames with the most significant movement. This method calculates the motion of pixels between successive frames and assigns weights to each frame determined by the magnitude of this motion. Figure 4.4 illustrates a visualisation of this technique being applied to "kickflip" trick.



(a) Optical flow visualisation of a kickflip sequence, highlighting regions with significant motion.



(b) Extracted frames from the sequence based on the greatest weightings from the optical flow representation.

Figure 4.4 Visualisation of optical flow: (a) Optical flow representation (b) Individual frames extracted.

### 4.3.2 Data Augmentation

The limited number of data available for this research necessitated the implementation of data augmentation techniques to enlarge the dataset's size before model training. Techniques such as rotating, flipping and adding noise to images also served to add more diversity to the dataset mitigating the risk of overfitting. This risk arises from the over-parameterised nature of the chosen models, whose trainable parameter count surpasses the size of the dataset, causing them to be vulnerable to overfitting.

### 4.3.3 Normalisation

During the preprocessing stage, normalisation played a crucial role in preparing the video frame data before further processing. In particular this study utilised min-max normalisation to scale the pixel intensities between 0 and 1. Normalising data in this manner ensured that the model's input values data was within a standardised range, calculated by Equation 4.1.

$$\text{Normalized Value} = \frac{\text{Pixel Value} - \text{Min Value}}{\text{Max Value} - \text{Min Value}} \quad (4.1)$$

#### 4.3.4 Feature Extraction with Transfer Learning

This research leveraged two well-established CNNs for this task: VGG16 and ResNet50. Both models were pre-trained on the large ImageNet dataset [51], leveraging extensive image recognition knowledge that were then fine-tuned to the domain of skateboard trick recognition. Table 4.1 summarises the key characteristics of VGG16 and ResNet50.

**VGG16:** VGG16, introduced in 2014 by Simonyan and Zisserman [45], is a CNN architecture known for its success in image classification and feature extraction for Action Recognition tasks. This model is relatively simple made up of 16 convolutional and fully connected (FC) layers with repeated use of 3x3 convolutional filters to preserve spatial information. Furthermore, VGG16 has proved to be effective for feature extraction in previous literature [7], supporting its suitability in this study.

**Resnet50:** Resnet50, introduced in 2016 by He et al. [52], is another powerful CNN architecture that is well known for its advancements in image classification and frame extraction capabilities. Unlike VGG16, this model has a deeper architecture employing 50 convolutional layers allowing it to learn more intricate feature representations of the data. Furthermore, it incorporates residual connections to counteract the vanishing gradient problem that is often caused by many layers.

Model	VGG16	ResNet50
Published	2014	2016
Layers	16 (Convolutional and FC layers)	50 (Convolutional layers)
Parameters	~138 million	~25 million
Input Image Size	224x224	224x224

Table 4.1 Characteristics of Transfer Learning Models VGG17 and ResNet50

The final four fully-connected (FC) layers of these models, which are typically used for classification were removed. This was crucial to adapt these models for feature extraction, as the aim was not to classify each frame of the video, but to extract a collection of features that could help sequence models gain a better understanding. These features in the form of vectors not only encapsulate the visual information present in each frame but also more abstract ideas such as low-level information (like edges/shapes) to higher-level features like objects and even actions themselves.

### 4.3.5 Dimensionality Reduction

The use of pre-trained models used in this study generate high-dimensional feature vectors. While these features capture valuable information, a large number of dimensions can lead to challenges. Firstly, it increases the computational costs of training models. Secondly, it could potentially lead to the "curse of dimensionality" as explained by Venkat (2018) [53], where models may struggle to learn effectively with high dimensional data. To address these issues, dimensionality reduction is employed after feature extraction to reduce the number of features while still retaining valuable information for classification.

## 4.4 Adapted Models

This research investigates the performance of four different deep learning architectures for skateboard trick classification. These models leverage pre-trained CNNs for feature extraction followed by sequence modelling techniques to capture the temporal dynamics of skateboard tricks. The architectures investigated are as follows.

1. **VGG16-LSTM:** This architecture leverages the pre-trained VGG16 model to extract features, which are then connected to an LSTM for sequence modelling.
2. **ResNet50-LSTM:** This architecture utilises the deeper ResNet50 pre-trained CNN for feature extraction, followed by an LSTM for sequence modelling. The increased depth of ResNet50 potentially allows it to learn more complex feature representations than VGG16.
3. **VGG16-BiLSTM:** This architecture employs the VGG16 model for feature extraction and a BiLSTM network for sequence modelling. BiLSTMs are known to be able to learn dependencies in both directions of a sequence, which could be useful for capturing subtle details in skateboard tricks.
4. **ResNet50-BiLSTM:** This architecture combines ResNet50 for feature extraction with a BiLSTM network. This combination leverages the potential advantages of deeper feature learning and the bi-directional capabilities of BiLSTMs.

The effectiveness of these models in prior research for video-based Action Recognition tasks motivated their selection for this study. For instance, Orozco et al. (2020) achieved 91.93% accuracy using a VGG-LSTM architecture [44]. Additionally, Chen's work [7] explored all four of these models in the context of skateboard trick classification and demonstrated the potential for competitive results.

## 4.5 Architecture

The Deep Learning Architecture, as illustrated in Figure 4.5 demonstrates the data flow pipeline for the skateboard trick classifier. The process begins with preprocessing the labeled video data including frame extraction, augmenting the training set for variability and reducing dimensionality for efficiency. Pre-trained CNNs such as VGG16 and ResNet50, then extract features from the preprocessed frames, before feeding them into their respective sequence models.

The accuracy of these models are then evaluated on a held-out test set, to ensure the models performance on unseen data. Finally, the evaluation utilises performance plots such as confusion matrices to visualise classification performance, model epoch loss and accuracy graphs to track training history.

## 4.6 Evaluation Methods

To thoroughly evaluate the performance of the proposed models, this study employed various evaluation metrics, including accuracy, precision, recall and F1-score. These metrics provided valuable insights into the model's ability to correctly classify different tricks.

The above mentioned metrics depend on the following definitions, described in the context of the "kickflip" class and tabulated in a confusion matrix.

- **True Positive (TP):** The number of instances that the model correctly identified a video as containing a kickflip.
- **True Negative (TN):** The number of instances that the model incorrectly identified a video as containing a kickflip.
- **False Positive (FP):** The number of instances that the model correctly identified a video as not containing a kickflip.
- **False Negative (FN):** The number of instances that the model incorrectly identified a video as not containing a kickflip.

**Accuracy:** Measures the proportion of correctly classified instances, over the total number of predictions made by the model. This metric provides a generic indicator of the model's reliability, however this specific metric can be sensitive in scenarios with class imbalance [54].

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (4.2)$$

**Precision:** Measures the proportion of predicted positive instances that are truly real positives [55]. In the context of kickflip detection, it represents the proportion of true kickflips against all predicted kickflips.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (4.3)$$

**Recall:** Measures the proportion of positive instances that are truly predicted positive [55]. In the context of kickflip detection, it represents the proportion of true kickflips that are identified correctly.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (4.4)$$

**F1-score:** This metric provides a balanced measure of performance by combining both precision and recall. It is calculated using the harmonic mean on both values and outputs a value ranging from zero to one, with values closer to one, indicating better performance.

$$\text{F1-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4.5)$$

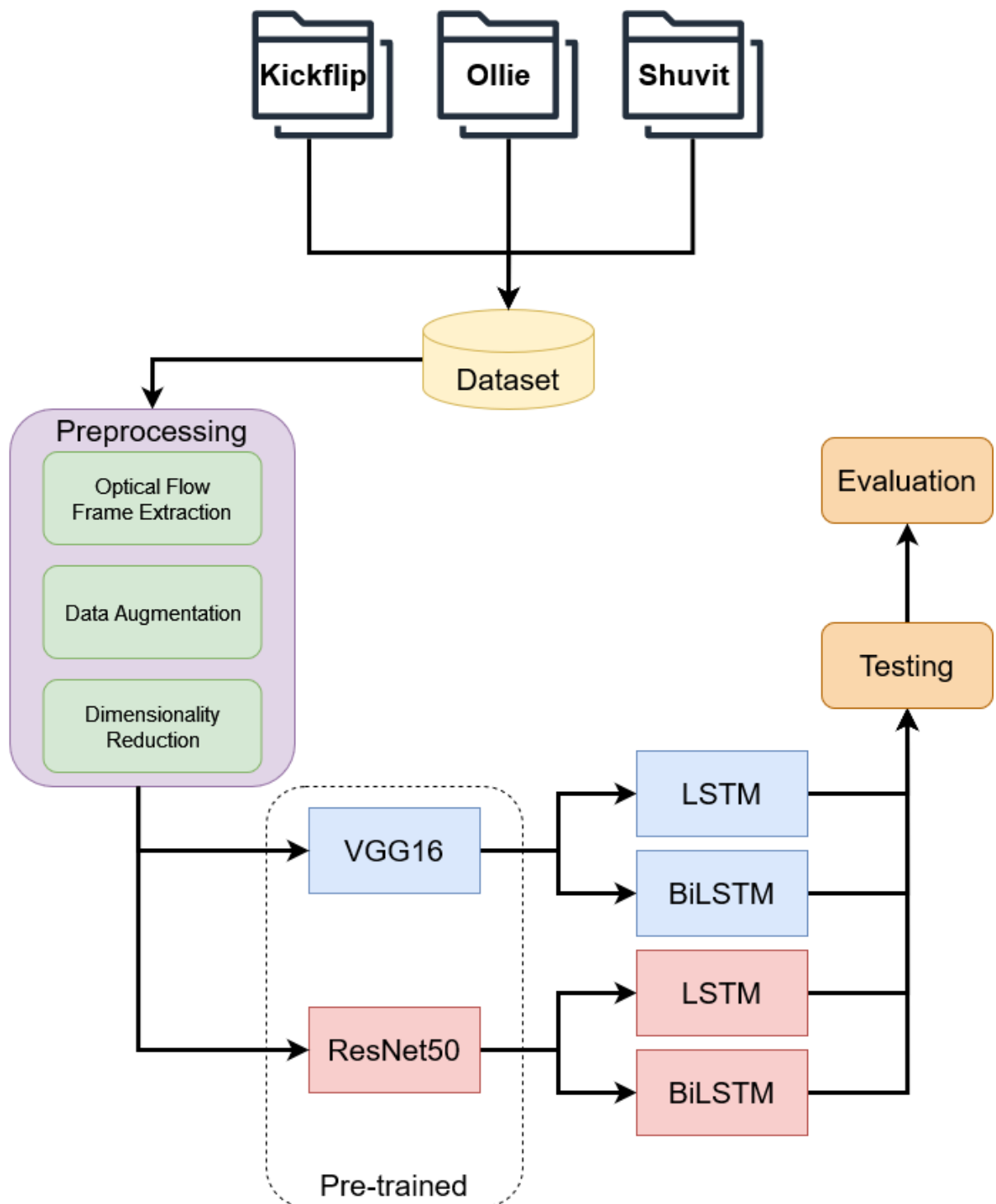


Figure 4.5 Artefact Architecture.



# 5 Implementation

## 5.1 Development Environment

**Python:** Python's large database of libraries, along with its wide use in Machine Learning, made it the ideal choice for developing this artefact. Furthermore, its large and active community made tackling problems and troubleshooting issues simpler. Prior experience using Python also contributed to this decision.

**Tensorflow:** Tensorflow is an open-source library, powerful in numerical computation and Machine Learning applications. It provides a rich toolset that allows for efficient development, training and deployment of Machine Learning models. Notably, Tensorflow comes with the Keras API, further simplifying the creation development by offering wide range of tools such as access pre-trained models, pre-defined layers and training and evaluation abstractions.

**OpenCv (cv2):** This open source library played a crucial role in the computer vision components of development, offering essential tools related to video processing, particularly during the frame extraction phase.

**Matplotlib:** This research made use of Matplotlib for its plotting and data visualisation functions, heavily used during the evaluation stage to properly visualise results or any other insights gained throughout.

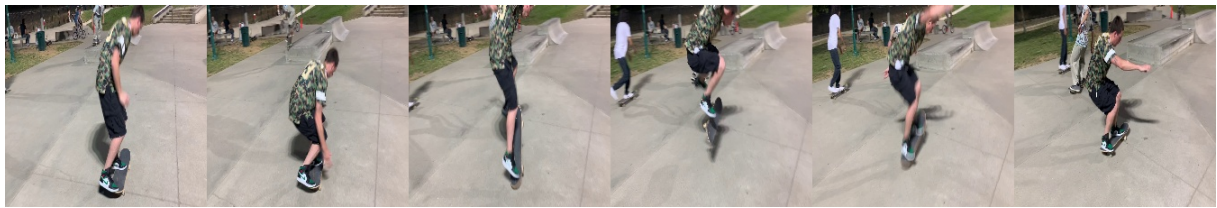
## 5.2 Dataset Split and Configuration

This study opted for a training-validation-testing split on the original dataset allocating 80% for training, 10% for validation and 10% for testing. The training set allowed the models to learn the patterns and relationships within the input data. The validation enabled a portion of the data to be used to evaluate the performance after every epoch. Finally, testing set provided a final assessment of the model's generalizability, by using data the model's have not seen before. The decision for this split stemmed from prioritising training data due to limited dataset size.

In consideration of Chen's (2023) [7] research, which advocated for a sequence length of 45 frames per video, with each frame sized at  $299 \times 299$  pixels, these parameters were initially assessed. However, experimentation with such parameters revealed significant computational costs with negligible improvement in results. Consequently, this research struck a balance between costs and performance, by opting for a sequence length of 20 and a frame size of  $224 \times 224$  pixels.

### 5.3 Frame Extraction using Optical Flow

The exploration of two frame extraction techniques aimed to find the most effective way to capture the crucial motions in a video through a series of frames. The optical flow method was hypothesised to capture frames with the most significant motion, unlike uniform sampling which often missed crucial parts of the skateboard trick, capturing more frames before or after the trick execution. An example of this behaviour at a smaller scale can be observed in Figure 5.1.



(a) Extracted frames using optical flow method.



(b) Extracted frames using uniform sampling.

Figure 5.1 Comparison of frame extraction between optical flow method and uniform sampling

This study employed Farneback's algorithm [56] using the OpenCV (cv2) library, to estimate the optical flow magnitudes on each frame. The selection process involved reading pairs of successive frames from a video, resized for faster processing and converted to greyscale to minimise noise caused by colour variations. The cv2 function `cv2.calcOpticalFlowFarneback()` performed dense optical flow estimation between these frames, using parameters such as pyramid scale, levels, winsize, iterations that can be found in Appendix values.

The function `cv2.cartToPolar()` converted the Cartesian flow vectors returned by the optical flow function into polar coordinates, discarding the directional information, while retaining only the magnitudes. Finally, the code computed the average magnitude and appended it to a list. This process iterated through all frames to select those with the highest average to be extracted from the video.

#### 5.3.1 Data Augmentation

This research explored a number of augmentation techniques using the `ImageDataGenerator` library provided by Tensorflow. This tool allowed for the definition

of specific parameters that influenced how each frame was augmented. To expand the number of training samples, this study created two copies of the original dataset, each augmented with its own set of unique parameters controlling image shifts, brightness, scaling and other transformations. The two unique augmentation parameters utilised in this study are illustrated in Figure 5.3.

Figure 5.2 showcases an augmented trick sequence against its original, split into six frames for illustrative purposes. Experimentation revealed that certain augmentation parameters such as image flipping and large rotational shifts disrupted the visualisation of the trick, these were removed or adjusted to a smaller scale from the augmentation parameters.



(a) Augmented frames.



(b) Original frames.

Figure 5.2 Comparison of Augmented sequence against original

```

1  {
2      "width_shift_range": 0.1,
3      "height_shift_range": 0.1,
4      "shear_range": 0.1,
5      "rotation_range": 10,
6      "zoom_range": 0.2,
7      "brightness_range": [0.7, 1.3],
8      "channel_shift_range": 20.0,
9      "fill_mode": "nearest"
10 }
```

(a) Augmentation Parameters 1

```

1  {
2      "width_shift_range": 0.2,
3      "height_shift_range": 0.1,
4      "shear_range": 0.8,
5      "rotation_range": 15,
6      "zoom_range": 0.2,
7      "brightness_range": [0.2, 0.1],
8      "channel_shift_range": 10.0,
9      "fill_mode": "nearest"
10 }
```

(b) Augmentation Parameters 2

Figure 5.3 Augmentation Parameters for all augmentation

## 5.4 Feature Extraction and Preprocessing for Training

After successfully extracting the most significant frames from each video using optical flow, the next step involved preparing the data for further processing. The main algorithm consisted of an iterative process responsible for resizing and normalising all frames before input to the pre-trained CNNs for feature extraction. The `cv2.resize()` function resized all frames to 224x224 to satisfy the VGG and ResNet input image requirements, before performing min-max normalisation on each frame by dividing the pixel intensities values by 255.

Once resized and normalised, the pre-trained CNNs extracted features from each frame. With the final classification layers removed, the extracted features of a singular image using VGG returned a shape of  $(1 \times 7 \times 7 \times 512)$ , while ResNet50, resulted in a shape of  $(1 \times 7 \times 7 \times 2048)$ . In both cases, the outputs were flattened using the Tensorflow `flatten()` function to transform the extracted features into a one-dimensional vector suitable to be fed into the dense layers of the network. Thus, VGG generates a final feature shape of  $(20 \times 25088)$  for an entire video, while ResNet50 generates  $(20 \times 100352)$ .

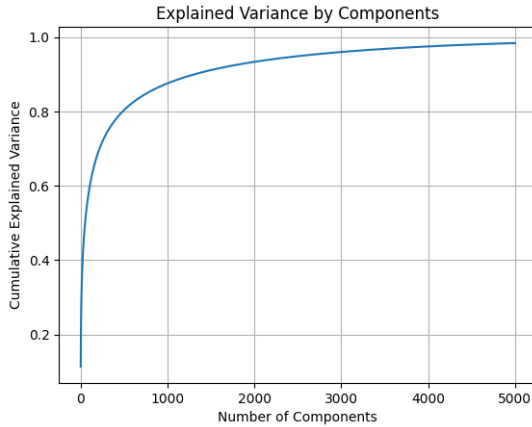
The NumPy function `np.save()` saved these extracted features along with their respective labels for training, validation, and test sets. This allowed for easy retrieval of the pre-processed data, eliminating the need to re-extract features with every new session.

### 5.4.1 PCA Dimensionality Reduction

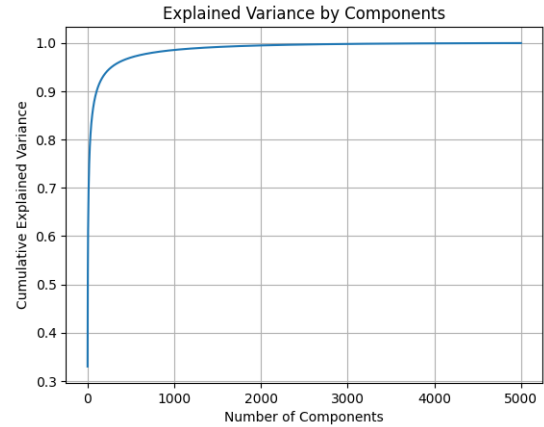
This research, leveraged Principal Component Analysis (PCA) dimensionality reduction, a technique, crucial in refining the feature sets produced by the pre-trained models.

It was aimed to achieve around 95% variance retention when determining the optimal number of components for reduction. This involved plotting two explained variance plots: one for the output of pre-training with VGG16 and another for the output of ResNet50, as illustrated in Figure 5.4. These plots display the cumulative explained variance against the number of components, allowing approximate selection of the number of components by examining the 95% percentile from both plots. Upon examining both plots, 2514 components were chosen for VGG, while 278 were selected for ResNet50. Consequently, this results in final feature shapes of  $(20 \times 2514)$  and  $(20 \times 278)$  respectively.

The application of PCA dimensionality reduction significantly reduced training times. This efficiency significantly sped up the process of hyperparameters tuning by reducing waiting times associated with training each iteration.



(a) Cumulative Variance plot of the output of VGG16.



(b) Cumulative Variance plot of the output of ResNet50.

Figure 5.4 Cumulative Variance plots for the outputs of VGG16 and ResNet50

## 5.5 Hyperparameter Optimisation

The architectures described in this study represent the best versions of themselves after a number of iterations and optimisations. To speed up the time-consuming task of hyperparameter tuning, this research leveraged the power of Optuna, a hyperparameter optimisation framework [57]. Optuna leverages a Bayesian optimisation algorithm that runs through a number of iterations, observing the performance of past configurations and strategically selecting parameter combinations such as learning rate, dropout rate and neuron count.

This iterative approach enables Optuna to navigate the search space and converge on optimal hyperparameters for a specific mode. This library was integrated within a custom-built development environment specifically designed to identify the optimal hyperparameter configurations for multiple model architectures and their corresponding input data.

### 5.5.1 Training Process and Hyperparameter Tuning

Throughout the training phase, Optuna took the responsibility of initialising the model's hyperparameter configurations. Leveraging its Bayesian optimisation algorithm, Optuna provided initial configurations, which served as starting points for the iterative training process.

In this iterative process, each model's hyperparameters underwent continuous evaluation and refinement after every iteration at an attempt to improve accuracy. Performance metrics such as loss curves, confusion matrices and F-scores, were monitored closely to assess the effectiveness of each hyperparameter configuration.

Common pitfalls observed, included overfitting, indicated in model loss graphs when validation loss begins to rise, and fluctuations in performance metrics such as accuracy and F-scores. These fluctuations often signify an unstable learning process, suggesting the need for hyperparameters reconfiguration. Finally, with every iteration, the hyperparameters and performance scores were saved for future analysis and comparison, aiding in the refinement of upcoming iterations.

## 5.6 Callbacks

Callbacks served as essential tools, used to monitor and influence the training process dynamically. They provided ways to automate certain tasks at different phases of training, such as saving checkpoints of the model or stopping the model early.

### 5.6.1 Early Stopping

An early stopping callback reduced overfitting and saved computational resources. This method monitored the validation loss at every epoch and halted training if improvement stopped after a predetermined patience value. All experiments investigated a patience of 10 epochs, based on the observation that the models were unlikely to improve after 10 epochs, with no validation loss advancements.

### 5.6.2 Model Checkpoint

This study incorporated model check pointing in the training process to save intermediate models after every epoch. This implementation was configured to monitor the validation loss and only save the model when it showed an improvement, allowing a seamless resumption of training in case of interruption.

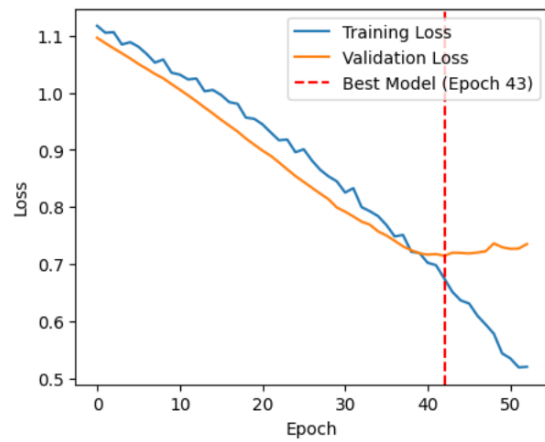
## 5.7 Training History

Figure 5.5 presets the training and validation loss curves for the models that achieved the best test accuracy within each architecture. The vertical dashed lines indicate the epoch at which the early stopping callback restored the model's weights, preventing overfitting.

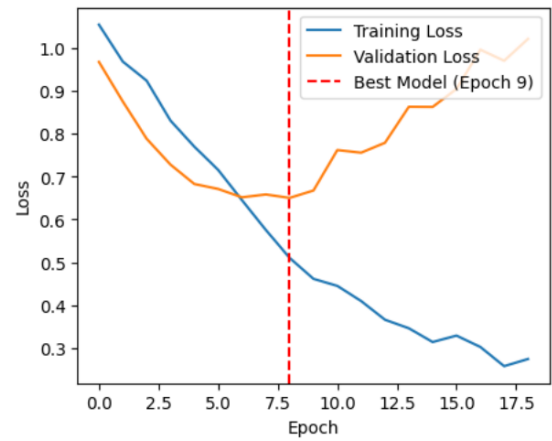
The loss curves across the four architectures all exhibit distinct patterns, highlighting the differences in learning dynamics specific to each architecture. Notably, smaller learning rates were applied to the VGG16-LSTM and ResNet50-BiLSTM models, with the intention to enhance the stability of training, thereby improving performance.

Additionally, the validation loss for the BiLSTM configurations generally exhibited more volatility, likely due to the complexity of their bidirectional layers.

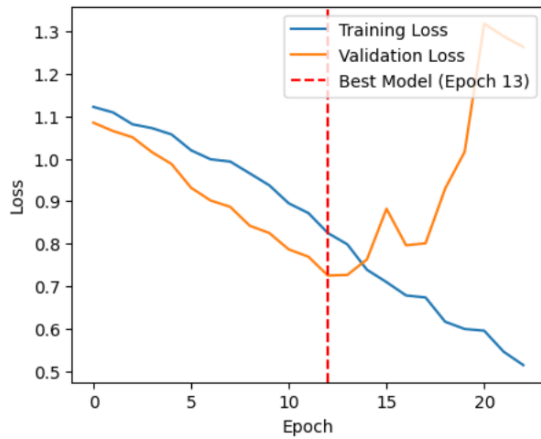
All models employed the Stochastic Gradient Descent (SGD) optimiser throughout the training process due to excessive overfitting observed with Adam in Section 6.1.1. While SGD performed better than Adam in the context of classifying skateboard tricks, it required longer training times to converge.



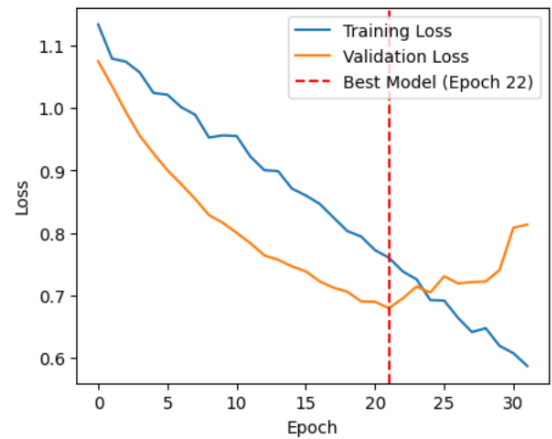
(a) VGG16-LSTM



(b) VGG16-BiLSTM



(c) ResNet50-LSTM



(d) ResNet50-BiLSTM

Figure 5.5 Loss graphs for each architecture

## 6 Evaluation

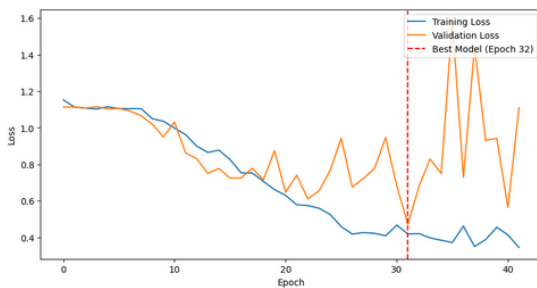
### 6.1 Model Training Specifications

Due to the absence of a dedicated graphics card, training exclusively relied on the CPU, featuring a Ryzen 5 3600 6-core processor, operating at a base clock speed of 3.59 GHz. Despite lacking GPU acceleration, the Ryzen 5 3600 processor reliably handled the training of several deep learning architectures and experiments conducted.

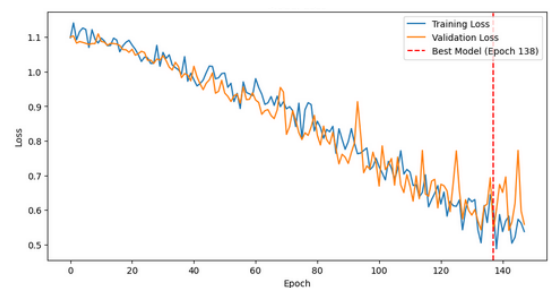
### 6.2 Experiments

#### 6.2.1 Choice of Optimiser

This study compared the results of Stochastic Gradient Descent (SGD) and Adam [58] as optimisation algorithms during training. The results suggest that SGD may be more effective at reducing overfitting compared to Adam. This improvement is illustrated in Figure 6.1, using an iteration of the VGG-BiLSTM architecture. While, Adam initially shows promise with good generalisation, it quickly encounters a pitfall, where the validation loss begins to climb, indicating overfitting. In contrast, the SGD-trained model exhibits a more stable validation loss, likely due to its algorithmic stability and its capability to achieve small generalisation errors as explained in the study by Hardt et al. (2016) [59]. The results observed with the VGG-BiLSTM aligns with the insights discussed by Hardt et al. suggesting that models trained using SGD are less vulnerable to overfitting.



(a) Training and Validation Loss using Adam



(b) Training and Validation Loss using SGD

Figure 6.1 Model loss graphs for (a) Adam and (b) SGD optimisers.

#### 6.2.2 The Effect of Data Augmentation

This research, also compared the effects of Augmenting data before feeding it into the sequence models, as illustrated in Table 6.1. Experiments revealed that each of the



four models demonstrated improvements in accuracies following the application of data augmentation on the input data. Enhancements ranged from a 5% increase in Model 3, to a 9% increase in Models 1 and 4. On average, model accuracy increased by 7.75% highlighting its effectiveness in enhancing the model's ability to generalise on unseen data by adding more variability to the training data.

Model	Unenhanced		Augmented	
	Accuracy	Training Time	Accuracy	Training Time
1. VGG16-LSTM	77%	00:07:47	86%	00:11:38
2. VGG16-BiLSTM	80%	00:07:45	88%	00:18:05
3. ResNet50-LSTM	75%	00:01:19	80%	00:02:19
4. ResNet50-BiLSTM	77%	00:01:41	86%	00:03:01

Table 6.1 Performance comparison of models with and without augmentation

### 6.2.3 The Effect of PCA

Another experiment conducted, involved evaluating the impact of PCA on model performance. The analysis compared VGG16-BiLSTM and ResNet50-BiLSTM as they were the top-performing models from each category of pre-trained models. This comparison focused on their accuracies before and after the application of PCA, as well as differences in training times. To ensure fair results, both models were trained using the same augmented data, with features reduced to their respective pre-calculated number of components discussed in Section 5.4.1. Table 6.2 presents a side-by-side comparison of these models with their training times.

Model	Pre-PCA		Post-PCA	
	Accuracy	Training Time	Accuracy	Training Time
VGG16-BiLSTM	58%	04:29:00	88%	00:18:05
ResNet50-BiLSTM	33%	03:29:00	86%	00:03:01

Table 6.2 Comparison of VGG16-BiLSTM and ResNet50-BiLSTM model performances before and after PCA implementation

The application of PCA shows significant results, not only improving accuracies but also reducing training times drastically. These results highlight the benefits of applying dimensionality reduction to combat the 'curse of dimensionality', often encountered with high-dimensional data, noted in Section 4.3.5. The observed improvements are a result of the removal of redundant features in the dataset, further enhancing the model's learning capabilities by allowing them to focus on the most informative features.

## 6.2.4 Applicability In Real-Time applications

For real-time applications like analysing skateboarding events, the speed at which video data is processed and evaluated is crucial. To evaluate the model's applicability in such applications, this study measured the time taken to classify a single 2-second video at  $\sim 30$ fps for each architecture.

Model	Evaluation Time (ss:ms)
Model 1 (VGG16-LSTM)	09:10
Model 2 (VGG16-BiLSTM)	09:90
Model 3 (ResNet50-LSTM)	08:20
Model 4 (ResNet50-BiLSTM)	9:00

Table 6.3 Evaluation time (ss:mm) for skateboard trick classification models.

The slow processing speeds outlined in Table 6.3 is caused by the computational expensive video analysis pipeline required before evaluation. While this pipeline negatively affected evaluation times, it was responsible for achieving high accuracies and significantly accelerating training time.

## 6.3 Discussion

### 6.3.1 Final Findings

This study's extensive testing on various Deep Learning architectures and the application of various preprocessing techniques, have demonstrated advancements in accuracies and training times in trick classification. These results highlight the effectiveness of applying tailored preprocessing methods in optimising model performance.

As a result of the extensive preprocessing techniques employed, the VGG-BiLSTM model emerged as the top performer with a final accuracy of 88%. The full table of final model accuracies is presented in Table 6.4.

Across all configurations, BiLSTM variants outperformed their equivalent LSTM counterparts. While the LSTM models still demonstrated favourable accuracies, the BiLSTM models improved by an average of 4%, likely due to their ability to process input data from both forward and backward directions. This dual-direction processing enhances the ability of these models to recognise specific temporal patterns in the data that might be missed by LSTM's.

Despite the common preference for ResNet50 for feature extraction due to its deeper and more complex structure, in this specific scenario ResNet50 models slightly underperformed compared to VGG16 models. This could suggest that in the context

of skateboard trick classification, ResNet50 does not necessarily translate to better performance. It is possible that the simpler and lighter capabilities of VGG16 are more effective for specific characteristics found in skateboard tricks.

Model	Study by Chen		This Study	
	Accuracy	Training Time	Accuracy	Training Time
1. VGG16-LSTM	70%	00:40:08	86%	00:11:38
2. VGG16-BiLSTM	69%	00:39:11	88%	00:18:05
3. ResNet50-LSTM	80%	00:59:60	80%	00:02:19
4. ResNet50-BiLSTM	81%	00:59:80	86%	00:03:01

Table 6.4 Performance comparison of models from Hanciao Chen’s study [7] and this study

### Comparative Analysis

When comparing this study’s results to that of Chen (2023) [7], Chen employed a wide range of architectures, including the ones evaluated in this study. Despite this, this study achieved a higher accuracy of 88% with the VGG16-BiLSTM model, in contrast to Chen’s highest outlined accuracy of 84% using the ResNet50 + Attention + BiLSTM model. In addition, compared to Chen, this study exhibits remarkable improvements in training time with an average reduction of 41:04 minutes, compared to their models. Interestingly, Chen also demonstrated higher accuracies with models based on the ResNet50 feature extractor as apposed to VGG, which contrasts the findings of this study where VGG16 based models outperformed ResNet50 variants.

The differences in model performance could be due to several factors. Firstly, the two studies employed different hyperparameter configurations such as learning rate, epochs, and dropout rates during training. Secondly, compared to Chen, this research included more training data per class. Thirdly, the use of optical flow and PCA in this research, could have provided a more concentrated feature set for training the models, leading to higher accuracies.

### 6.3.2 Classification observations

Trained models often encountered difficulty in differentiating between the ollie and kickflip compared to the pop shuvit. This challenge likely arises due to their shared visual characteristics. Both the ollie and the kickflip involve common elements such as the initial pop of the skateboard and the absence of rotation around the y-axis, which distinguishes them to the pop shuvit. The confusion matrices illustrated in Figure 6.2 demonstrate a favourable classification rate on pop-shuvits across all models compared

to ollies and kickflips. Interestingly, in the classification of pop shuvits, BiLSTM models demonstrated 100% classification accuracy, while LSTM models, while still performing well, slightly lagged behind with classifying 11 out of 12 instances.

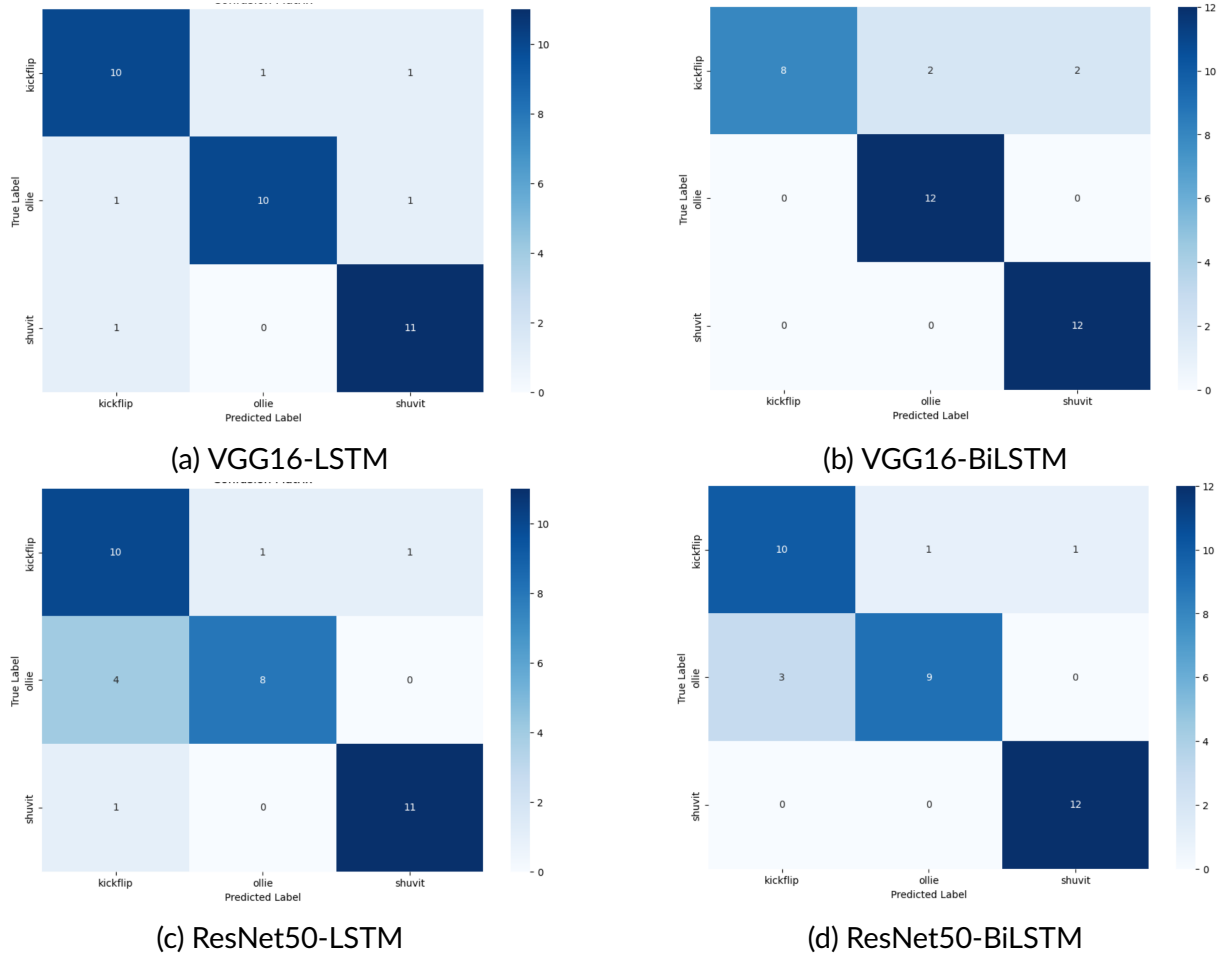


Figure 6.2 Confusion matrices for each architecture.

## 6.4 Limitations

**Limited Data:** The insufficient amount of data available for training the models constrained the ability to train more robust and accurate models. This limitation potentially prohibited the models from capturing the full variability of skateboard tricks filmed by video. Thus, likely not accounting for various lighting conditions, camera angles, video quality and environmental conditions. The result of not considering these aspects, is the unsuitability of the application in real-world scenarios.

**Processing efficiency:** The current implementation of video analysis pipeline, while effective in improving accuracies and reducing training times, suffers from processing delays. This limitation stems from the computationally expensive analysis steps required before trick evaluation. These steps, including optical-flow frame extraction, feature

extraction and PCA dimensionality reduction, delay the processing time, making it less suitable for real-time applications.

## 7 Conclusions and Future Work

### 7.1 Future Work

**Classification Extension:** As outlined in Section 4.1, this study adopted a multi-class classification approach, focusing initially on three fundamental skateboard tricks. While this served as an adequate foundation for exploring the underlying relationships between these tricks and their application in Deep Learning, its scope limited its usefulness in real-world scenarios.

To enhance its applicability, particularly in skateboarding events, the inclusion of more complex tricks more commonly observed in competitions is crucial. Additionally, incorporating tricks that interact with external entities, such as rails, ledges and gaps would enhance its real-world usefulness.

**Automated Judging Systems:** Currently, skateboarding competitions rely on a panel of judges to evaluate skater performance live. While judges possess expertise and understanding, this system is heavily subjective. Personal perspectives and preferences can influence their decisions, raising fairness issues.

An automated judging system would address these limitations by removing the element of human bias, facilitating a more consistent and fair scoring system. However, this requires a more complex classification strategy, as the true challenge lies in translating all the subjective aspects of style and creativity into an objective score.

**Data limitation:** A major limitation encountered during development was the insufficient amount of data available to train with. Future work in this area, would greatly benefit from additional data to expand the model's robustness and generalisation capabilities. Moreover, if additional data is not accessible, it could be worth exploring techniques to generate synthetic data using models such as Generative Adversarial Networks (GANs) or Variational Autoencoders (VAEs).

**Improvements in evaluation time:** Future work would benefit from simplifying the expensive video analysis pipeline utilised in this study. Improvements in processing time will greatly increase its potential in real-world scenarios, targeting broadcasted skateboard competitions. Possible approaches include, frame extraction algorithm optimisations, model simplifications such as using lighter pre-trained feature extractors and utilising hardware acceleration.

## 7.2 Conclusion

This study aimed to classify skateboard tricks from video and compare the performance of various deep learning architectures in accomplishing this task. This research aimed at enhancing live broadcasts and social media content, in order to contribute to the possibilities for innovation in the skateboarding industry.

Four architectures were implemented: VGG16-LSTM, VGG16-BiLSTM, ResNet50-LSTM and ResNet50-BiLSTM to compare performances of two pre-trained feature extractors and sequence models. Evaluations revealed that the selected models were capable of matching or exceeding the performance of relevant literature.

Additionally, this research explored the impact of Principle Component Analysis (PCA) and data augmentation on model performance. These techniques were employed to assess their potential in improving accuracies and efficiencies of the selected model architectures. As detailed in Chapter 6, these techniques not only enhanced models' accuracy but also reduced training times. presents these experiments with relevant performance metrics to demonstrate the benefits of these approaches.

# Bibliography

- [1] *Encyclopaedia britannica*, 2023. [Online]. Available: <https://www.britannica.com/sports/skateboarding>.
- [2] *One year on: How skateboarding's olympic debut changed the games 2022*, 2022. [Online]. Available: <https://olympics.com/en/news/one-year-on-skateboarding-olympic-debut-feature>.
- [3] Z. Foley, *Vert, street, park - what are the different styles of skateboarding?* 2021. [Online]. Available: <https://www.goskate.com/top/skateboarding-styles-full-guide/>.
- [4] Statista. "Skateboarding participation in the u.s. 2010-2021." (), [Online]. Available: <https://www.statista.com/statistics/191308/participants-in-skateboarding-in-the-us-since-2006/>.
- [5] O. Foundation. "2022 outdoor participation trends report." (), [Online]. Available: <https://outdoorindustry.org/wp-content/uploads/2023/03/2022-Outdoor-Participation-Trends-Report.pdf>.
- [6] M. Shapiee *et al.*, "The classification of skateboarding tricks: A transfer learning and machine learning approach," *MEKATRONIKA*, vol. 2, pp. 1–12, Oct. 2020. DOI: 10.15282/mekatronika.v2i2.6683.
- [7] H. Chen, "Skateboardai: The coolest video action recognition for skateboarding (student abstract)," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, pp. 16 184–16 185, Jun. 2023. DOI: 10.1609/aaai.v37i13.26952.
- [8] B. Mahesh, "Machine learning algorithms -a review," Jan. 2019, ISSN: 1662-5161. DOI: 10.3389/fnhum.2016.00066.
- [9] R. C. Staudemeyer and E. R. Morris, *Understanding lstm – a tutorial into long short-term memory recurrent neural networks*, 2019. arXiv: 1909.09586 [cs.NE].
- [10] L. Budach *et al.*, *The effects of data quality on machine learning performance*, 2022. arXiv: 2207.14529 [cs.DB].
- [11] A. Ullah, J. Ahmad, K. Muhammad, M. Sajjad, and S. W. Baik, "Action recognition in video sequences using deep bi-directional lstm with cnn features," *IEEE Access*, vol. 6, pp. 1155–1166, 2018. DOI: 10.1109/ACCESS.2017.2778011.
- [12] M. Vrigkas, C. Nikou, and I. Kakadiaris, "A review of human activity recognition methods," *Frontiers in Robotics and Artificial Intelligence*, vol. 2, Nov. 2015. DOI: 10.3389/frobt.2015.00028.
- [13] A. Thakur and A. Konde, "Fundamentals of neural networks," *International Journal for Research in Applied Science and Engineering Technology*, vol. 9, pp. 407–26, 2021.



- [14] O. Montesinos-López, A. Montesinos, and J. Crossa, "Fundamentals of artificial neural networks and deep learning," in Jan. 2022, pp. 379–425, ISBN: 978-3-030-89009-4. DOI: 10.1007/978-3-030-89010-0\_10.
- [15] J. Gu *et al.*, "Recent advances in convolutional neural networks," *Pattern Recognition*, vol. 77, pp. 354–377, 2018, ISSN: 0031-3203. DOI: <https://doi.org/10.1016/j.patcog.2017.10.013>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0031320317304120>.
- [16] K. O'shea and R. Nash, "An introduction to convolutional neural networks," *arXiv preprint arXiv:1511.08458*, 2015.
- [17] C.-C. J. Kuo, "Understanding convolutional neural networks with a mathematical model," *Journal of Visual Communication and Image Representation*, vol. 41, pp. 406–413, 2016.
- [18] J. Donahue *et al.*, "Long-term recurrent convolutional networks for visual recognition and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2625–2634.
- [19] F. A. Gers, N. N. Schraudolph, and J. Schmidhuber, "Learning precise timing with lstm recurrent networks," *Journal of machine learning research*, vol. 3, no. Aug, pp. 115–143, 2002.
- [20] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [21] S. Siامي-Nامینی, N. Tavakoli, and A. S. Namin, "The performance of lstm and bilstm in forecasting time series," in *2019 IEEE International conference on big data (Big Data)*, IEEE, 2019, pp. 3285–3292.
- [22] J. Du, Y. Cheng, Q. Zhou, J. Zhang, X. Zhang, and G. Li, "Power load forecasting using bilstm-attention," in *IOP Conference Series: Earth and Environmental Science*, IOP Publishing, vol. 440, 2020, p. 032 115.
- [23] P. O'Donovan, "Optical flow: Techniques and applications," *International Journal of Computer Vision*, vol. 1, p. 26, 2005.
- [24] S. S. Beauchemin and J. L. Barron, "The computation of optical flow," *ACM computing surveys (CSUR)*, vol. 27, no. 3, pp. 433–466, 1995.
- [25] K. Host and M. Ivašić-Kos, "An overview of human action recognition in sports based on computer vision," *Heliyon*, 2022.
- [26] W. Jia, M. Sun, J. Lian, and S. Hou, "Feature dimensionality reduction: A review," *Complex & Intelligent Systems*, vol. 8, no. 3, pp. 2663–2693, 2022.
- [27] J. Ma and Y. Yuan, "Dimension reduction of image deep feature using pca," *Journal of Visual Communication and Image Representation*, vol. 63, p. 102 578, 2019.

- [28] J. Yin *et al.*, "Mc-Istm: Real-time 3d human action detection system for intelligent healthcare applications," *IEEE Transactions on Biomedical Circuits and Systems*, vol. 15, no. 2, pp. 259–269, 2021. DOI: 10.1109/TBCAS.2021.3064841.
- [29] J. Yin, Q. Yang, and J. J. Pan, "Sensor-based abnormal human-activity detection," *IEEE Transactions on Knowledge and Data Engineering*, vol. 20, no. 8, pp. 1082–1090, 2008. DOI: 10.1109/TKDE.2007.1042.
- [30] F. Wu *et al.*, "A survey on video action recognition in sports: Datasets, methods and applications," *IEEE Transactions on Multimedia*, 2022.
- [31] D. R. Beddiar, B. Nini, M. Sabokrou, and A. Hadid, "Vision-based human activity recognition: A survey," *Multimedia Tools and Applications*, vol. 79, no. 41-42, pp. 30 509–30 555, 2020.
- [32] S. Zhang, Z. Wei, J. Nie, L. Huang, S. Wang, Z. Li, *et al.*, "A review on human activity recognition using vision-based method," *Journal of healthcare engineering*, vol. 2017, 2017.
- [33] S.-A. Rebuffi, S. Gowal, D. A. Calian, F. Stimberg, O. Wiles, and T. A. Mann, "Data augmentation can improve robustness," in *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., vol. 34, Curran Associates, Inc., 2021, pp. 29 935–29 948. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/fb4c48608ce8825b558ccf07169a3421-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/fb4c48608ce8825b558ccf07169a3421-Paper.pdf).
- [34] X. Ying, "An overview of overfitting and its solutions," in *Journal of physics: Conference series*, IOP Publishing, vol. 1168, 2019, p. 022 022.
- [35] X. Jiang, "Feature extraction for image recognition and computer vision," in *2009 2nd IEEE international conference on computer science and information technology*, IEEE, 2009, pp. 1–15.
- [36] M. Jogin, M. Madhulika, G. Divya, R. Meghana, S. Apoorva, *et al.*, "Feature extraction using convolution neural networks (cnn) and deep learning," in *2018 3rd IEEE international conference on recent trends in electronics, information & communication technology (RTEICT)*, IEEE, 2018, pp. 2319–2323.
- [37] A. Krizhevsky, G. Hinton, *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [38] K. Weiss, T. M. Khoshgoftaar, and D. Wang, "A survey of transfer learning," *Journal of Big data*, vol. 3, no. 1, pp. 1–40, 2016.
- [39] A. N. Soni, "Application and analysis of transfer learning-survey," *International Journal of Scientific Research and Engineering Development*, vol. 1, no. 2, pp. 272–278, 2018.

- [40] A. B. Sargano, X. Wang, P. Angelov, and Z. Habib, "Human action recognition using transfer learning with deep representations," in *2017 International joint conference on neural networks (IJCNN)*, IEEE, 2017, pp. 463–469.
- [41] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [42] C. Szegedy et al., "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [43] M. Al-Faris, J. Chiverton, D. Ndzi, and A. I. Ahmed, "A review on computer vision-based methods for human action recognition," *Journal of imaging*, vol. 6, no. 6, p. 46, 2020.
- [44] C. I. Orozco, E. Xamena, M. E. Buemi, and J. J. Berlles, "Human action recognition in videos using a robust cnn lstm approach," *Ciencia y Tecnología*, pp. 23–36, 2020.
- [45] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [46] J. Donahue et al., "Long-term recurrent convolutional networks for visual recognition and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2625–2634.
- [47] E. M. Saoudi, J. Jaafari, and S. J. Andaloussi, "Advancing human action recognition: A hybrid approach using attention-based lstm and 3d cnn," *Scientific African*, vol. 21, e01796, 2023.
- [48] M. A. Abdullah et al., "The classification of skateboarding tricks via transfer learning pipelines," *PeerJ Computer Science*, vol. 7, e680, 2021.
- [49] N. K. Corrêa, J. C. M. d. Lima, T. Russomano, and M. A. d. Santos, "Development of a skateboarding trick classifier using accelerometry and machine learning," *Research on Biomedical Engineering*, vol. 33, pp. 362–369, 2017.
- [50] LightningDrop and C. Fitzgerald, *Skateboardml 1.0*, version 1.0, Aug. 2020. DOI: 10.5281/zenodo.3986905. [Online]. Available: <https://doi.org/10.5281/zenodo.3986905>.
- [51] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*, IEEE, 2009, pp. 248–255.
- [52] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

- [53] N. Venkat, "The curse of dimensionality: Inside out," *Pilani (IN): Birla Institute of Technology and Science, Pilani, Department of Computer Science and Information Systems*, vol. 10, 2018.
- [54] A. Tharwat, "Classification assessment methods," *Applied computing and informatics*, vol. 17, no. 1, pp. 168–192, 2020.
- [55] D. M. Powers, "Evaluation: From precision, recall and f-measure to roc, informedness, markedness and correlation," *arXiv preprint arXiv:2010.16061*, 2020.
- [56] G. Farnebäck, "Two-frame motion estimation based on polynomial expansion," in *Image Analysis: 13th Scandinavian Conference, SCIA 2003 Halmstad, Sweden, June 29–July 2, 2003 Proceedings 13*, Springer, 2003, pp. 363–370.
- [57] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2623–2631.
- [58] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [59] M. Hardt, B. Recht, and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," in *International conference on machine learning*, PMLR, 2016, pp. 1225–1234.

## Appendix A Sample A

Layer (type)	Output Shape	Param #
bidirectional_22 (Bidirectional)	(None, 20, 512)	5675008
dropout_105 (Dropout)	(None, 128)	0
dense_90 (Dense)	(None, 64)	8256
dropout_106 (Dropout)	(None, 64)	0
dense_91 (Dense)	(None, 64)	4160
dropout_107 (Dropout)	(None, 64)	0
dense_92 (Dense)	(None, 3)	195
Total params: 5983043 (22.82 MB)		
Trainable params: 5983043 (22.82 MB)		
Non-trainable params: 0 (0.00 Byte)		

Table A.1 Model summary for VGG-BiLSTM

## **Appendix B   Sample B**