

# Skateboard Trick Recognition through an AI-based Approach

**Kris Saliba**

Supervisor: Prof. 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

# List of Figures

# List of Tables

# 1 Introduction

Skateboarding dates back to the 1940s when handmade skateboards first appeared [1]. It has since developed into a worldwide phenomenon, with its popularity skyrocketing, after gaining recognition as an official sport in the 2020 Tokyo Olympic Games [2]. Skateboarding comprises the dynamic activities of riding a skateboard and skilfully performing a repertoire of tricks, manifesting as a popular and exhilarating “extreme sport”.

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 revolves around riding on specialised obstacles, namely, half-pipes and ramps, emphasising aerial manoeuvres. Street skateboarding transpires 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].

## 2 Background

### 2.1 Skateboard Tricks

Skateboard tricks are the heart and soul of skateboarding. These tricks originate from the dynamic orchestration of rotations and revolutions of a skateboard along various axes emphasising the significance of precise placement of a skateboarder's feet to initiate these rotations. These tricks serve as excellent examples of how the skateboarder's body and skateboard work in perfect harmony. Some common skateboard tricks include:

- **Ollie:** One of the first tricks beginners learn. Where the skateboarder pops the tail of the board while sliding their foot across 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.
- **360 kickflip:** A combination of a kickflip and a  $360^\circ$  board 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 immediately evident or interpretable, and as the volume of available datasets continues to surge so does the demand for machine learning [4].

Morris [5] 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 the algorithm becomes at executing tasks. Training is a multifaceted process that directly influences the overall accuracy of the model. Within this phase, numerous critical factors come into play, shaping the model's performance. These factors encompass dataset quality and diversity, the preprocessing of data, the selection of an appropriate model architecture, the optimal duration of training, and the fine-tuning of essential hyper-parameters [6].

There exist three main categories for ML models[4]:

- **Supervised:** This is a machine learning concept that centres around the development of algorithms to make predictions or classifications using labelled data. The model is trained on a dataset comprising of example input-output pairs.
- **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 algorithms 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 Object Detection

Object detection is a computer vision task that detects instances of objects in images and videos and maps them to a predefined class. For humans, the act of recognising and responding to objects is a trivial task as described in [7], it is an essential feature that enables our performance and communication. Numerous academic and industry researchers have shown a deep interest in the technology, focusing on various applications where object detection plays a major role. These applications include but are not limited to autonomous driving, surveillance systems and face detection [8].

The output of an object detection model yields the instance's location, as the object's centre, a bounding box or even a list of pixels containing the object. The research paper [8] further implies that object detection is consistently defined within the context of a data set that consists of images mapped to a list of relevant object properties, such as their locations and scales, that are specified within each image. This definition makes references to the equation below, where an image is denoted as  $\mathcal{I}$ , and  $O(I)$  represents the collection of object descriptions for objects within the image.

$$O(I) = \{(Y_1^*, Z_1^*), \dots, (Y_i^*, Z_i^*), \dots, (Y_{N^*i}^*, Z_{N^*i}^*)\}$$

In the above equation, each description encompasses two parts,  $Y_i^* \in \mathcal{Y}$  characterises the category or type of an object, and  $Z_{N^*i}^* \in \mathcal{Z}$  represents information about its location, size or shape within the image.  $\mathcal{Z}$  represents the different ways to describe an object, this is typically done by specifying the object's centre  $(x_c, y_c) \in \mathcal{R}^2$  or as a bounding box  $(x_{min}, y_{min}, x_{max}, y_{max}) \in \mathcal{R}^4$ . By utilising these notations, according to [8], object detection can be defined as the operation of combining an image with a set of detections.

## 2.4 Activity Recognition

Activity recognition is the process of identifying and categorizing human activities from video sequences. Human activity involves 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. Similar to object detection, 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 [9],[10].

## 2.5 Neural Networks

### 2.5.1 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 [11]. Figure 2.1 depicts 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 [12].

- **Input Layer:** This is the set of neurons that serve as the initial entry point for external data or features. 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 situated between the Input and Output Layers where the network captures complete nonlinear 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.

### 2.5.2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs), as described by Gu et al. (2018) [13] 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

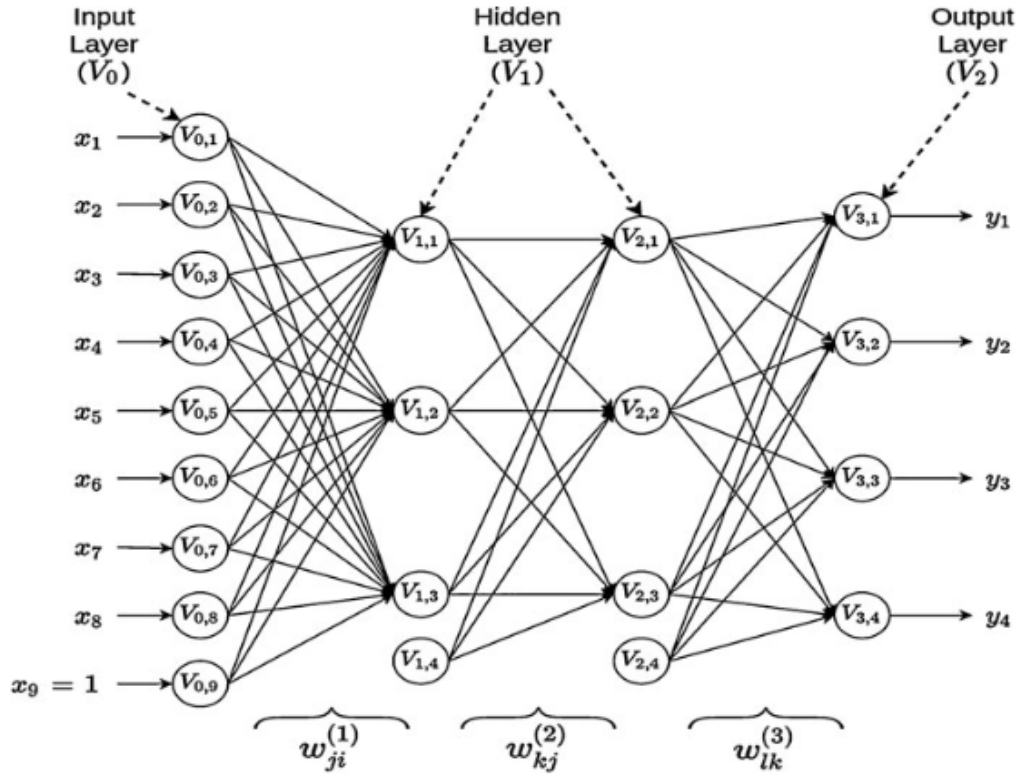


Figure 2.1 Schematic Representation of an Artificial Neural Network. Reproduced from López et al. (2022) [12]

their incredible performance in various fields such as visual recognition, speech recognition and natural language processing.

CNNs, incorporate multiple layers and are capable of extracting effective representations

### 2.5.3 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 by their inability to access long-term data, limited to approximately ten sequential time steps. This constraint arises from the challenge of dealing with vanishing or exploding gradients during the backpropagation procedure during the processing of extended sequences, as discussed in prior works [5],[14].

To address the limitation that RNNs encounter in capturing long-term dependencies, Long Short-Term Memory Recurrent Neural Networks (LSTM-RNNs) were introduced [15]. LSTM networks unlike RNNs are capable of capturing dependencies across more than 1,000 time steps depending on the network's complexity and are more biologically credible [5].

## 3 Literature Review

### 3.1 Activity Recognition

Building on the foundational understanding of activity recognition defined in the background, it's important to acknowledge the impact it has on various sectors. The recent advancements in recognising human actions in videos have not only revolutionalised sectors such as healthcare and security as demonstrated in the studies [16], [17], but also hold extensive potential in the realm of sports.

The research conducted by K. Host and M. Ivašić-Kos in [18] along with Wu et al. in [19] 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 activity recognition models in the field of sports.

The study by Beddiar et al. [20], 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.

While both contact and vision-based systems have their merits, this study will specifically focus on vision-based techniques and their application in the development of a skateboard trick classifier. These methods, which utilise image and video data to recognise human activities are selected due to their societal acceptability and their applicability in sports broadcasts.

#### 3.1.1 Challenges in this field

The domain of Activity recognition comes with many challenges as depicted by Zhang et al. (2017) [21]. The researchers suggest that while certain scenarios utilise static cameras such as surveillance systems, most situations that benefit from Activity

Recognition 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 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 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. [21] in the domain of activity recognition 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 home use, then it may encounter videos of lower quality further 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 activity recognition. It involves the application of various techniques to enhance raw data before feeding it to Machine Learning algorithms. —continue or leave out

As a result of the limited research on the emergence of a skateboard trick classifier, 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 [22], [23]. In the realm of Skateboard trick classifiers, Shapiee et al (2020) [24] effectively employed data augmentation techniques to expand their dataset, demonstrating the application of these methods in improving model performance for trick classification.

After establishing the role of data augmentation in addressing the lack of data available, another important step in computer vision is data normalisation. This technique is essential for standardising the range and distribution of pixel values in an image and has shown an increase in model performance, [25], [26]. Among methods like Z-score, min-max normalisation is particularly prevalent in normalising pixel intensities to a 0-1 scale as follows:

$$\text{Normalized Value} = \frac{\text{Pixel Value} - \text{Min Value}}{\text{Max Value} - \text{Min Value}} \quad (3.1)$$

While less common in computer vision due to its normality assumption, Z-score normalisation is defined as:

$$\text{Z-score} = \frac{x - \mu}{\sigma} \quad (3.2)$$

Here,  $x$  represents the individual pixel value,  $\mu$  is the mean of all pixel values, and  $\sigma$  is their standard deviation.

The research by Pei et al. (2023) [25], explored the impact of normalisation on classification accuracy. They demonstrated that, for 8-bit images, min-max normalisation outperformed Z-score, significantly enhancing classification accuracy. On the other hand, the study by de Raad et al. [27] suggests that the impact of normalisation on model performance varies depending on certain dataset characteristics. These contrasting conclusions presented by the studies [25] and [27] suggest that while normalisation is an important preprocessing step, its application should be carefully adapted to the characteristics of the dataset.

Having optimised the distribution of pixel values through the normalisation process, the next crucial step is feature extraction. Xudong Jiang [28], describes this technique as the process of capturing the core attributes of an object by eliminating redundancies, resulting in a set of numerical features ideal for classification. CNNs excel at this due to their capabilities in detecting complex patterns and enhancing features. Used primarily in image recognition, as an initial for classification algorithms, their effectiveness is demonstrated by studies like Manjunath Jogin et al. (2018) [29] where they achieved 86% accuracy rate through the use of CNNs for feature extraction alongside several classifiers. Beyond this primary function, CNN's ability to extract complex features from images makes them very effective when coupled with other models for more complex tasks such as activity recognition. In the context of skateboard trick classification, CNNs can efficiently extract detailed spatial features like board rotations, foot positioning and limb movements. This involves CNNs operating on individual frames to extract spatial features, which can then be passed through a sequence analysis model like an LSTM, proven to be efficient at understanding temporal information.

Following the discussion on CNN for feature extraction, it is important to highlight the role of transfer learning in enhancing ML models, especially in fields with limited data like skateboard trick classification. The concept of transfer learning as detailed by the studies [30] and [31] involves using a pre-trained ML model to leverage its experience for a new, but related task. The study by Sargano et al. (2017) [32], employed transfer learning with pre-trained deep CNNs like AlexNet [33] and GoogleNet [34], for human activity recognition. Notably, they utilised the pre-trained

models for feature extraction, followed by an SVM classifier for final recognition. 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 realm of skateboard trick classification. Two other studies [24], [35] have utilised this approach and achieved high accuracies, however, a more in-depth analysis of these papers can be found in the 'Advancements in Skateboard Trick Classification' section. These studies strongly suggest the exploration 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 occlusions 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 [36]. One particularly effective DL architecture is the CNN-LSTM. It combines the strength of CNN's for extracting spatial features from individual frames, with LSTM networks, which capture the temporal information across frames.

The study by Orozco et al. (2020) [37] adopted this approach to test its effectiveness against three activity recognition datasets: KTH, UCF-11, HMDB-51, particularly focusing on how the number of LSTM units impacted performance. The authors employed transfer learning, utilising the VGG16 model [38] 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.

Building upon the foundational CNN-LSTM architecture, a recent study by Saoudi et al. (2023) [39] enhances this model by incorporating three key advancements:

- 1. 3D Convolutional Neural Networks (3D CNNs):** Unlike regular CNNs, these process video data 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 extensive resources required to build 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 LSTM called Bi-directional Long Short-Term Memory (BiLSTM). Whereas LSTMs only process data in one direction, BiLSTMs can analyse data in both directions, allowing them to understand relationships between preceding and subsequent actions.
- 3. Attention Mechanisms:** This technique focuses the model processing on the most relevant parts of the data, improving performance by important information for the task. By integrating an attention layer after their BiLSTM, Saoudi et al. enabled the model to focus on temporal features within the input, providing a more detailed representation of the input, and improving overall 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 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 signals artificially generated based on the findings of prior studies. These signals are then fed into a study-dependent model for classification, as outlined in [40] and [41]. 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 [24] and [35].

### 3.2.1 Accelerometer-based approaches

The study by Abdullah et al (2021) [40], 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, including 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. [40] 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. Recognising the rarity of such high accuracies, this study will consider these findings and efforts will be made to ensure a robust model by employing techniques to avoid overfitting such as early-stopping and the use of a diverse dataset [23].

The study by Corrêa et al (2017) [41], obtained their sample data by artificially generating 543 signals based on prior research, utilising tools such as 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%.

### 3.2.2 Computer Vision-based Approaches

The paper by Shapiee et al (2020) [24] 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. On the other hand, Chen (2023) [35] 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. 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.

The paper by Shapiee et al. [24] utilises data augmentation techniques and applies three rotation augmentation techniques: horizontal rotation, positive 90° rotation and negative 90° rotation. The researchers experimented on three Transfer learning models: MobileNet, NASNetMobile and NASNetLarge, each evaluated using a k-Nearest Neighbor (k-NN) classifier. As a result, the models demonstrated impressive classification accuracies, with MobileNet achieving 95%, NASNetMobile 92% and NASNetLarge 90%.

Recurrent Neural Networks (RNNs) and Long Short-Term Memory Networks (LSTMs) are popular architectures due to their capabilities in modelling the dynamic relationships in sequential data [5]. In the student abstract by Hanciao Chen [35], 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.

# 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] B. Mahesh, "Machine learning algorithms -a review," Jan. 2019, ISSN: 1662-5161. DOI: 10.3389/fnhum.2016.00066.
- [5] 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].
- [6] L. Budach et al., *The effects of data quality on machine learning performance*, 2022. arXiv: 2207.14529 [cs.DB].
- [7] Watson, Rebecca and Huis in 't Veld, Elisabeth M. J. and de Gelder, Beatrice, "The neural basis of individual face and object perception," *Frontiers in Human Neuroscience*, vol. 10, 2016, ISSN: 1662-5161. DOI: 10.3389/fnhum.2016.00066. [Online]. Available: <https://www.frontiersin.org/articles/10.3389/fnhum.2016.00066>.
- [8] S. Agarwal, J. O. du Terrail, and F. Jurie, "Recent advances in object detection in the age of deep convolutional neural networks," *ArXiv*, vol. abs/1809.03193, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:52183570>.
- [9] 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.
- [10] 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.
- [11] 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.

- [12] 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.
- [13] 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>.
- [14] 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.
- [15] 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.
- [16] J. Yin *et al.*, "Mc-lstm: 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.
- [17] 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.
- [18] K. Host and M. Ivašić-Kos, "An overview of human action recognition in sports based on computer vision," *Heliyon*, 2022.
- [19] F. Wu *et al.*, "A survey on video action recognition in sports: Datasets, methods and applications," *IEEE Transactions on Multimedia*, 2022.
- [20] 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.
- [21] 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.
- [22] 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).

- [23] X. Ying, "An overview of overfitting and its solutions," in *Journal of physics: Conference series*, IOP Publishing, vol. 1168, 2019, p. 022 022.
- [24] 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.
- [25] X. Pei *et al.*, "Robustness of machine learning to color, size change, normalization, and image enhancement on micrograph datasets with large sample differences," *Materials & Design*, vol. 232, p. 112 086, 2023.
- [26] X. Zhong, B. Gallagher, K. Eves, E. Robertson, T. N. Mundhenk, and T. Y.-J. Han, "A study of real-world micrograph data quality and machine learning model robustness," *npj Computational Materials*, vol. 7, no. 1, p. 161, 2021.
- [27] K. De Raad *et al.*, "The effect of preprocessing on convolutional neural networks for medical image segmentation," in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, IEEE, 2021, pp. 655–658.
- [28] 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.
- [29] 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.
- [30] 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.
- [31] 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.
- [32] 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.
- [33] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [34] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

- [35] 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.
- [36] 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.
- [37] 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.
- [38] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [39] 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.
- [40] M. A. Abdullah *et al.*, "The classification of skateboarding tricks via transfer learning pipelines," *PeerJ Computer Science*, vol. 7, e680, 2021.
- [41] 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.

## 4 Sample A

## **5 Sample B**