VIEW-INVARIANT ACTION RECOGNITION IN DYNAMIC SCENES VIA SIM2REAL TRANSFER

A Dissertation

Presented to the Faculty of the Graduate School

of Cornell University

in Partial Fulfillment of the Requirements for the Degree of

Doctor of Philosophy

by Yuhan Wang May 2023 © 2023 Yuhan Wang ALL RIGHTS RESERVED

VIEW-INVARIANT ACTION RECOGNITION IN DYNAMIC SCENES VIA

SIM2REAL TRANSFER

Yuhan Wang, Ph.D.

Cornell University 2023

Learning view-invariant representation is essential to improving feature extraction for action recognition. Existing approaches cannot effectively capture details for human actions due to fast-paced gameplay and implicit viewdependent representation. In this paper, it goes beyond recognizing human actions from a fixed view and focusing on action recognition from an arbitrary view. This paper purposes a method to build an efficient data generating pipeline due to lack of original input data. This paper also provides a pipeline combining capturing modified I3D human actions features and use Multilayer Perception to achieve human action recognition and classification. The use of information captured from combination of virtual and real-life data, as well as different viewing angles, leads to high classification performance.

BIOGRAPHICAL SKETCH

Yuhan is an M.S. student in the Laboratory for Intelligent Systems and Controls (LISC) at Cornell University. He received his B.S. degree in Mechanical Engineering, with a Mechatronics minor, from Villanova University in May 2021. His research interests includes artificial intelligence, machine learning and computer vision in dynamic sports and autonomous driving.

Dedicated to my angel investor: my loving and supportive family.

ACKNOWLEDGEMENTS

I would like to seize this moment to convey my deepest gratitude and profound appreciation for the exceptional guidance and support provided by my esteemed advisor and mentor, Prof. Silvia Ferrari. Throughout my research journey, Silvia has consistently demonstrated remarkable expertise and served as an invaluable source of inspiration. She allowed me the freedom to explore my own areas of interest while providing guidance along the way. I am forever indebted to Silvia for her instrumental role in shaping my academic and professional growth.

In addition, I am profoundly grateful for Prof. Bharath Hariharan. His profound expertise in the field of computer vision have equipped me with essential knowledge that serves as the bedrock of my professional endeavors.

Further more, I am grateful to Frank Kim, a dedicated and knowledgeable PhD student in our lab. As a master's student under his guidance, I have been fortunate to have Frank as a mentor who consistently went above and beyond to assist me at every stage of my study. I am deeply thankful for his mentorship and the immeasurable impact he has had on my academic and professional development.

I am grateful to all of those in the Laboratory for Intelligent Systems and Controls with whom I have had the pleasure to work during this and other related projects. Each of members of my Dissertation Committee has provided me extensive personal and professional guidance and taught me a great deal about both scientific research and life in general.

Finally, I would like to acknowledge and thank my parents. I would not have been able to be where I am now without their unwavering support.

	Biographical Sketch	i
	Dedication	I
	Acknowledgements	V
	Table of Contents	i
	List of Tables	i
	List of Figures	i
1	Introduction	1
	1.1 Background and Motivation	1
	1.2 Introduction to human action recognition in dynamic sports	3
	13 Thesis Organization	5
		,
2	Data Augmentation	5
	2.1 Introduction	5
	2.2 Empirical data collected from real-life situations	7
	2.3 Data augmentation in Unreal Engine	9
3	Action Recognition 10	6
	3.1 Feature Extraction	5
	3.1.1 Introduction to 3D-CNN and I3D framework 18	S
	312 Implementation details	Ĵ
	313 Feature extraction from videos	, ,
	32 Introduction to MIP framework	-
	2.2.1 MLP framework design	t 7
	$3.2.1$ WLF Hallework design \ldots \ldots \ldots \ldots \ldots	/ 2
		1
	3.2.3 Action Recognition	2
	$3.3 \text{Kesults} \dots \dots \dots \dots \dots \dots \dots \dots \dots $	3
4	Conclusion 4	7

TABLE OF CONTENTS

LIST OF TABLES

2.1	Total data number.	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	1	5
-----	--------------------	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

LIST OF FIGURES

2.1	Actions label	8
2.2	24 views recorded from Unreal Engine	9
2.3	Different jersey colors for the dataset	12
2.4	Unreal Engine data with random masks applied	13
2.5	Frame selection to make new videos(players action are not con-	
	sistent)	14
3.1	Video architecture for I3D model	19
3.2	Overall I3D architecture from Carreira and Zisserman [3]	21
3.3	Feature vectors example	22
3.4	MLP architecture	25
3.5	Our model architecture	29
3.6	PCA for one jersey color	34
3.7	Confusion matrix for one jersey color	35
3.8	t-SNE for one jersey color	36
3.9	PCA for 5 jersey colors	38
3.10	Confusion Matrix for 5 jersey colors	39
3.11	t-SNE for 5 jersey colors	41
3.12	PCA for views	43
3.13	Confusion Matrix for views	44
3.14	t-SNE for views	45

CHAPTER 1 INTRODUCTION

1.1 Background and Motivation

Recognizing human actions in video sequences is a challenging task due to the variability in viewpoint, appearance, and motion dynamics. Traditional supervised learning approaches for action recognition require large amounts of labeled data, which can be difficult and time-consuming to obtain. Additionally, these approaches often assume a fixed viewpoint for the video, which limits their ability to handle variations in viewpoint.

Self-supervised learning is a promising approach for learning view-invariant representations, as it can leverage unlabeled data to learn features that are robust to variations in viewpoint. Self-supervised learning approaches use auxiliary tasks to generate supervisory signals that do not require manual annotation, such as predicting the rotation of a video or predicting the temporal order of video frames.

In this paper, we propose a self-supervised learning approach for viewinvariant action recognition in dynamic sports. Dynamic sports, such as soccer and ice hockey, are particularly challenging for action recognition due to the fast-paced gameplay and the large number of players in the scene. These factors make it difficult to obtain high-quality labeled data for training action recognition models.

Our approach involves building an efficient data generating pipeline to overcome the lack of original input data. We use a combination of virtual and reallife data, as well as different viewing angles, to generate a diverse set of training examples. Specifically, we use a physics simulation engine to generate synthetic videos with varying viewpoints and motion dynamics. We then augment the synthetic videos with real-life videos captured from multiple cameras to increase the diversity of the test dataset.

To extract view-invariant features from the videos, we modify the I3D network, a popular architecture for action recognition. We remove the final classification layer of the I3D network and use the output features as input to a Multilayer Perception (MLP) classifier. We train the MLP classifier using a combination of labeled and unlabeled data, leveraging the self-supervised learning framework to learn view-invariant representations.

To evaluate our approach, we use a publicly available dataset of dynamic sports videos, which contains variations in viewpoint and motion dynamics. We compare the performance of our approach to several baseline methods, including traditional supervised learning approaches and other self-supervised learning methods. Our results show that our approach outperforms these methods in terms of accuracy.

We also conduct ablation studies to analyze the contribution of different components of our approach. Our results show that each component of our approach, including the use of synthetic data and the self-supervised learning framework, contributes to the overall performance.

In conclusion, our approach demonstrates the effectiveness of selfsupervised learning for view-invariant action recognition in dynamic sports in synthetic scenarios. Future work could explore the use of more sophisticated self-supervised learning methods and investigate the generalizability of our approach to other domains beyond dynamic sports. Additionally, our approach could be extended to other tasks that require view-invariant representations, such as object recognition and tracking.

1.2 Introduction to human action recognition in dynamic sports

Human action recognition is a crucial research field with significant applications in various domains, surveillance and entertainment. Many previous approaches have achieved good performance on action recognition, Nonetheless, the assumption that action videos are captured from a fixed viewpoint restricts their ability to be adaptable to various camera parameters and perspectives in real-world scenarios. One of the major challenges in action recognition is viewinvariance, as the performance of recognition models is significantly impacted by changes in camera angels and self-occlusions between different body parts, especially in fast-pace dynamic sports due to wide viewing angle and obscured player's actions.

To address the view-invariance issue, this paper proposes a view-invariant action recognition method to optimize human actions recognition in dynamic sports, specifically in the context of ice-hockey. In this research, we collected hockey videos from real-life games and unreal engine and convert them into human body sized bounding pictures. Due to geometric constrains in real life scenario, there are only limited camera view in real life, however, there could be multiple views in Unreal Engine for one single player's action. The pictures are then converted to multiple videos with random fps, mask covered, and with different frame to generate more dataset. Then, we modified Two-Stream Infalted 3D ConvNet (I3D) to extract features from those videos. The I3D architecture addresses the view-invaraince problem in action recognition by using a 3D convolutional neural network(CNN) that processes videos as a sequence of frames. By processing videos this way, the network can learn to recognize actions based on both their spatial and temporal information. The spatial convolutions operate on each individual frame in the video, while the temporal convolutions operate across frames in the video, which are essential in recognizing actions in fast-paced dynamic sports.

Using the I3D model, rich and discriminative features were extracted from the hockey videos into json files, since it the model was pre-trained on a largescale dataset of action videos, enabling it to learn robust action representations. These extracted features were then fed into a multilayer perceptron (MLP) we designed, which was trained to classify the player's actions, such as shooting, skating forward, skating left, or skating right. The MLP consists of multiple layers of interconnected neurons that process the extracted features to learn the underlying patterns and relationships between different actions. During the training phase, the MLP was fine-tuned to optimize its classification performance by adjusting its weights and biases using backpropagation. Once trained, the MLP classified the player's action in real-time with high accuracy, enabling effective analysis of dynamic sports action sequences.

The main contribution of this work is 1) the self-supervised learning framework that effectively utilizes unlabeled data; 2)overcomes the limitations of existing methods in handling different view points and camera parameters; 3) create a pipeline to generate human actions features from any videos in dynamic sports. This work has the potential to advance the field of action recognition in dynamic sports and enable more robust systems for real-world applications.

1.3 Thesis Organization

This thesis is divided into three primary components: Data augmentation, and Action recognition. The first part data augmentation is presented in Chapter 2, which encompasses new approach augmenting real-life data and synthetic data. The chapter is broken down by the data augmentation tasks include original data generating, apply different masks, different frame, and different frame rate(fps). Each of these tasks have their each algorithm proposed. Graphical results are presented at the end of each chapter along with quantitative results that showcase the algorithm accuracy and efficiency.

Chapter 3 provides theoretical work on the Multilayer Perception and detailed algorithms on the development for view-invariant action recognition based on I3D features. The pre-trained I3D model was modified to extract image features from augmented data generated earlier. Furthermore, the detail in MLP framework was discussed to classify actions to their corresponding label. This pipeline to augment data and train in MLP framework is summarized in the end. The graphical and quantitative results are also presented. The final chapter of the dissertation comprises a summary of the findings and suggestions for future work.

CHAPTER 2 DATA AUGMENTATION

2.1 Introduction

In computer vision, the availability of a large and diverse dataset is crucial for training accurate models. However, collecting and annotating a dataset can be a time-consuming and expensive process. Collection dataset from a crowded real-world scenario is difficult, and might safety concerns. One approach to overcome this issue is to use synthetic data and apply data augmentation techniques. Data augmentation is a process of generating new training samples by applying various transformations to the existing dataset. These transformations include rotation, flipping, scaling, cropping, and many others. By applying data augmentation techniques, we can increase the size and diversity of the dataset, and therefore, improve the performance of our models.

Inspired by *MOTSynth: How Can Synthetic Data Help Pedestrian Detection and Tracking?*, we realized that view-invariant action recognition is required in this real dynamic worlds and collecting them from photo-realistic environment would be a better approach [7]. These unreal character could provide relatively static and canonical view for our analysis, helping us focusing on the action recognition but get rid of the affect of noisy environment. This section will welly demonstrate how high-fidelity simulation based network can overcome this issue.

In this section, we will discuss various data augmentation techniques used in our research, their advantages and limitations, and how they can be applied to improve the performance of computer vision models. We will also compare the effectiveness of different data augmentation techniques on different computer vision tasks, such as image classification, object detection, and semantic segmentation. Finally, we will conclude with a discussion of the future directions of data augmentation in computer vision and its potential impact on the field.

2.2 Empirical data collected from real-life situations

To improve the performance of our ice hockey action recognition system, we employed data augmentation techniques to increase the size and diversity of our dataset. By collaborating with Cornell Ice hockey team, we collected multiple real life game videos between Cornell team against other teams and use them as input data.

We created an python script to identify the players in each frame of an input ice hockey video and manually labeled their actions frame by frame. The results were saved as json file, with each frame being assigned a unique identifier and its corresponding action labeled as "passing", "shooting", "dribbling", or "others". Each of the actions correspond to either "player" or "referee". Figure 2.1 illustrates the processing that we manually labeled original data.

The resulting augmented real-life dataset contained a total over 8000 frames, with each frame having a unique identifier and its corresponding action labeled as one of the actions. However, after converting these frames into videos, we have less than 100 videos being labeled as videos with actions. Therefore, these real life data are continually being labeled as test dataset in the future use. The



Figure 2.1: Actions label

increased diversity of the dataset enabled our model to better learn the underlying patterns of the data, and we used this dataset to test and evaluate our ice hockey action recognition system in the future.



Figure 2.2: 24 views recorded from Unreal Engine

2.3 Data augmentation in Unreal Engine

To increase the diversity and size of our dataset, we utilized the Unreal Engine to create 3D ice hockey players and record their actions from multiple viewpoints. The advantages of using the Unreal Engine were that we could generate more data, and we could also record the actions from 24 different viewpoints. (Figure 2.2) We created 3D ice hockey players using Unreal Engine's character creation tools. The character creation tools allowed us to customize various features of the player, including their height, jersey color and equipment. We also had the option to choose from a range of action animations, including skating forward, skating backward, shooting, etc.

After creating the 3D ice hockey players, we recorded their actions using Unreal Engine's Sequencer tool. The Sequencer tool enabled us to capture the actions from multiple viewpoints, providing a significant increase in the diversity of the dataset. Specifically, we recorded the actions from 24 different viewpoints, including views from different body side, and the overhead camera.

After recording the actions, we exported the data in the form of mp4 video files. We then applied various data augmentation techniques to generate more data and improve the robustness of our model.

Rotation:

To apply rotation, we used the OpenCV library in Python to randomly rotate the images between -20 and 20 degrees. This technique helped us to generate additional data by creating variations in the players' orientation and position. We applied rotation to each image in a video sequence, which ensured that the entire sequence was augmented consistently.

Masking:

We also applied different percentages of mask covered to the images to simulate the presence of occlusions in the video data. Specifically, we randomly applied 10 different percentages of mask coverage to create even more variations in the data. As shown in Figure 2.3, this technique involved masking out a portion of the image, simulating occlusions or partial views of the player's actions. This technique helped us to generate data that was more representative of real-world scenarios, where occlusions may occur due to players' movements or other factors. We randomly selected a percentage of the image to mask, and then we applied a Gaussian blur to the masked region to simulate the occlusion.

Cropping:

In addition to rotation and masking, we also applied cropping as another data augmentation technique. Specifically, we cropped the frames of each video from different views into various sizes and positions, such as placing the person in the center or on the left side of the frame. By doing so, we were able to generate more variations of the same action, and the I3D framework would recognize them as distinct pictures.

This approach allowed us to expand our dataset and include more diverse examples of each action, which would better prepare our model for real-world scenarios where the same action could be performed from different positions or angles. Additionally, the use of cropping as a data augmentation technique allowed us to improve the generalizability of our model and ensure that it could accurately recognize actions across different contexts and scenarios.

Jersey colors:

In addition to the data augmentation techniques mentioned above, we also used another approach to further increase the diversity of our dataset. Specifically, we chose 5 different jersey colors commonly used in ECAC conference, and we generated new samples by replacing the original jerseys with each of



Figure 2.3: Different jersey colors for the dataset

the 5 colors, as shown in Figure 2.4. This approach helped us to simulate different scenarios where the players have different uniforms, which could improve the model's ability to recognize actions under varying conditions.

By applying this approach, we were able to create more variations of each action in our dataset. For instance, a simple action like a shooting between two players could look different based on the color of the jerseys, and the model would have to learn to recognize the action regardless of the uniform colors. This technique also helped to improve the robustness of the model and ensure that it could accurately recognize actions in a variety of situations.

Frame Selection:

We varied the frame rate of the videos by selecting frames from the original videos in order and creating new videos with a lower frame rate. This technique helped us to create data that was more challenging to classify and required the model to be more robust to variations in the input data. We also varied the selection of frames used to generate new videos by continuously selecting frames from different parts of the original videos. This allowed us to capture different phases of the player's actions and generate more variations in the dataset. The



Figure 2.4: Unreal Engine data with random masks applied

process is shown in Figure 2.5. This player actions might not be completed due to limited selected frames, which better simulate the real-life scenario. Even each video has only slight difference between each other(only one frame), the model we introduced later would still understand them as differnt input data. This technique helped us to create data that contained more variation in the players' actions and positions.

The naming convention for the generated video files provided valuable information about each sample in our dataset. The file names included information such as the action being performed (e.g., shooting), the view from which the action was captured (e.g., view 3), the range of frames captured (e.g., from frame 12 to 42), the fps rate used (e.g., 30 fps), and the percentage of mask coverage applied (e.g., mask level 8).

This detailed information about each sample allowed us to easily keep track of the different variations in our dataset and enabled us to analyze the impact of each data augmentation technique on the performance of our view-invariant action recognition model. Furthermore, by generating data from different views and applying various augmentation techniques, we were able to create a diverse



Figure 2.5: Frame selection to make new videos(players action are not consistent)

dataset that could better represent the complex and varied nature of ice hockey actions. The naming convention for the generated video files and the use of data augmentation techniques were critical in helping us create a robust and diverse dataset that could effectively support our research on view-invariant action recognition in ice hockey.

In summary, we utilized the Unreal Engine to create 3D ice hockey players and record their actions from multiple viewpoints. We then used data augmentation techniques to generate more data by applying rotation, masking, and frame selection. In all, we have applied 5 different jerseys, 24 different views, 4 players' actions, 10 different random masking, 20 sliding frame windows per action, 1 fps per action, this gives us around 96,000 video files that could be transit to I3D features for later use.

Table 2.1:

The resulting dataset was larger and more diverse, and easy to manage throughout their names, which helped to improve the accuracy of our ice hockey action recognition system. By using a combination of 3D player genera-

		Total data	a number		
Jersey Col-	Views	Actions	Random	Sliding	Total
ors			Masking	Window	Number
5	24	4	10	20	96000

Table 2.1: Total data number.

tion and data augmentation techniques, we were able to create tons of datasets that was more representative of real-world scenarios and more challenging for our model to classify.

CHAPTER 3 ACTION RECOGNITION

3.1 Feature Extraction

To extract video features for action recognition, various methods have been proposed in the literature. One popular approach is to use hand-crafted features, such as Histogram of Oriented Gradients (HOG) and Local Binary Patterns (LBP) [2,9,24]. While these methods are simple and computationally efficient, they often lack discriminative power and struggle with variations in viewpoint and motion.

Another approach is to use deep learning-based methods, such as the traditional two-stream method and the C3D network [12]. The two-stream method involves training separate convolutional neural networks (CNNs) for spatial and temporal information, respectively, and fusing the outputs for action recognition. The C3D network, on the other hand, is a 3D CNN that directly takes in spatio-temporal inputs and extracts features from the entire video sequence. While these methods have shown promising results on benchmark datasets, they still face challenges in capturing fine-grained details of human actions and handling viewpoint variations.

Recently, the I3D network has emerged as a popular method for video feature extraction in action recognition. From *Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset,* the I3D network is introduced as a modified version of the Inception architecture that uses 3D convolutions to extract spatio-temporal features from video sequences [3]. The network is pre-trained on large-scale video datasets, such as Kinetics and Sports-1M, and fine-tuned on task-specific datasets for action recognition.

Compared to the traditional two-stream method and the C3D network, the I3D network has several advantages. First, it is more computationally efficient, as it only requires a single network to extract spatiotemporal features. Second, it has been shown to outperform these methods on benchmark datasets, particularly when fine-tuned on smaller datasets with limited labeled data. Third, the I3D network can leverage pre-trained models on large-scale datasets, which allows it to capture more generalizable and discriminative features.

However, the I3D network also has some limitations. It requires a large amount of labeled data for fine-tuning, which can be challenging to obtain in some domains. Additionally, it is still vulnerable to variations in viewpoint, particularly when the training data is biased towards certain viewpoints.

Given these considerations, we decide to use the I3D network to extract video features for our self-supervised learning approach for view-invariant action recognition in dynamic sports. We modify the I3D network by removing the final classification layer and using the output features as input to an MLP classifier. This allows us to leverage the self-supervised learning framework to learn view-invariant representations, while still benefiting from the discriminative power of the I3D network.

3.1.1 Introduction to 3D-CNN and I3D framework

Convolutional Neural Networks (CNNs) have been widely used in computer vision tasks such as image classification, object detection, semantic segmentation, and action recognition. Recently, 3D CNNs have emerged as a powerful tool for video analysis tasks, as they can learn spatio-temporal representations from video sequences.

3D CNNs are essentially an extension of 2D CNNs to the temporal domain. They use 3D convolutional kernels to extract spatio-temporal features from video frames over time. The I3D architecture was introduced as an enhancement to C3D (Convolutional 3D Networks) by inflating from 2D models. This approach allows for the reuse of the 2D models' architecture, such as ResNet and Inception, as well as the bootstrapping of model weights from pre-trained 2D models.Compared to 2D CNNs, 3D CNNs can capture motion information and handle temporal variations in video sequences, making them well-suited for tasks such as action recognition and video segmentation.

One of the most popular 3D CNN architectures for video analysis is the I3D (Two-Stream Inflated 3D ConvNets) network. The I3D network was introduced by Carreira and Zisserman in 2017 and is based on the Inception architecture, which uses a combination of different filter sizes in each layer to capture features at different scales. The I3D network is a 3D CNN architecture that is designed to capture spatio-temporal features from video sequences. The network is based on the Inception architecture, which uses a combination of filters with different sizes and depths in each layer to capture features at different scales. In the I3D network, the Inception modules are modified to use 3D convolutions instead of 2D convolutions, allowing the network to capture spatio-temporal features.



Figure 3.1: Video architecture for I3D model

The I3D network consists of two streams: a RGB stream and an optical flow stream, as shown in Figure 3.1. The RGB stream takes in the raw video frames as input, while the optical flow stream takes in the optical flow images, which capture the motion information in the video frames. The two streams are combined at the early stage of the network, allowing the network to capture both appearance and motion information.

The I3D network is pre-trained on large-scale video datasets, Kinetics and Sports-1M, which contain millions of video clips covering a wide range of human actions. The pre-training step allows the network to learn generic spatio-temporal features, which can be fine-tuned on smaller action recognition datasets with limited labeled data. During pre-training, the I3D network is trained to predict the action class labels of the videos in the dataset. This allows the network to learn discriminative features that are useful for action recognition tasks. The I3D network has also been shown to outperform other 3D CNN architectures such as C3D and two-stream networks on several action recognition benchmarks. In this paper, we leverage the I3D framework for video feature extraction to achieve view-invariant action recognition in dynamic sports. We modify the I3D network by removing the final classification layer and using the output features as input to an MLP classifier. This allows us to leverage the self-supervised learning framework to learn view-invariant representations while still benefiting from the discriminative power of the I3D network.

3.1.2 Implementation details

To train our self-supervised action recognition model, we needed a large of dataset regarding the ice-hockey related motions. By applying several transformations on the given dataset from unreal engine and real life, we composed them into over 96,000 different videos, with different fps, frame number, viewing angle, and masks covered.

Certainly, comarping to I3D framework, the traditional frameworks for action recognition, such as the C3D network, typically operate on individual frames of a video and process them separately. This approach can miss important temporal information that is necessary for accurately recognizing human actions. For example, the temporal pattern of a "shooting" motion cannot be fully captured by looking at individual frames alone. In contrast, the I3D network processes the video as a whole and captures both spatial and temporal information.

We removed the final classification layer so that we could use the generated I3D features for later use in action classification. Each augmented video file could generate a 1024 by 1 feature vector. Specifically, the output of the last



Figure 3.2: Overall I3D architecture from Carreira and Zisserman [3]

convolutional layer of the network was flattened and passed through a global average pooling layer to obtain a feature vector of size [1, 1024]. The overall architecture is shown in **Figure 4**. This feature vector captured the most salient spatiotemporal features of the input video and was suitable for downstream classification tasks. By applying N videos files on the input videos dataset, we combine them into a [N, 1, 1024] vector, saved in .json file, as shown in Figure 5.

The I3D network achieves this by applying 3D convolutional filters to the video frames, which allows it to learn both spatial and temporal features. Specifically, the 3D convolutional filters operate on a volume of frames that include the current frame and some neighboring frames in time. This allows the network to capture the temporal patterns that are characteristic of different human actions.

Furthermore, the I3D network is designed to be robust to differences in viewpoint and scale, which makes it well-suited for recognizing human actions from videos captured from different camera angles or with different frame rates. This is because the network is trained on a large and diverse set of videos that cover a wide range of actions, viewpoints, and scales. This enables it to learn general

{"rgb": [[[0.12091114372015, 0.025684453547000885, 0.02501755580306053, 0.0683453232049942, 0.14687910676002502,
0.04069194570183754, 0.08700071275234222, 0.20396123826503754, 0.053124260157346725, 0.13584256172180176,
0.04496214538812637, 0.20494043827056885, 0.13603301346302032, 0.2881454825401306, 0.04173232242465019,
0.026858074590563774, 0.5311747789382935, 0.1009359210729599, 0.12313401699066162, 0.06541909277439117, 0.10710549354553223,
0.10376765578985214, 0.10855075716972351, 0.016694201156497, 0.3815610408782959, 0.05251944065093994, 0.09281452745199203,
0.2530863285064697, 0.40024781227111816, 0.13187779486179352, 0.06843023002147675, 0.3336748480796814, 0.11802902072668076,
0.1207137331366539, 0.04490090534090996, 0.0077140675857663155, 0.42566585540771484, 0.04003835842013359,
0.08497339487075806, 0.02153552696108818, 0.10879100859165192, 0.06979885697364807, 0.13731297850608826,
0.04385030269622803, 0.25012439489364624, 0.11769523471593857, 0.07503824681043625, 0.05289410799741745,
0.03811619430780411, 0.6694991588592529, 0.023734061047434807, 0.11457794904708862, 0.22649100422859192, 0.3788621425628662,
0.13587218523025513, 0.01852971687912941, 0.016494804993271828, 0.09777358174324036, 0.2996465563774109,
0.009712528437376022, 0.2711029648780823, 0.10496624559164047, 0.029001474380493164, 0.2627869248390198,
0.016158554702997208, 0.0806717723608017, 0.2340472936630249, 0.03577588498592377, 0.17668193578720093, 0.1405894160270691,
0.1822507083415985, 0.1969926953315735, 0.14690610766410828, 0.13280601799488068, 0.26709505915641785, 0.18187455832958221,
0.0398021936416626, 0.1533183455467224, 0.06705441325902939, 0.10531624406576157, 0.15867315232753754, 0.05450613051652908,
0.06593293696641922, 0.028363928198814392, 0.040432654321193695, 0.24844275414943695, 0.07878431677818298,
0.031702227890491486, 0.06599630415439606, 0.2261597365140915, 0.1663258820772171, 0.08152954280376434, 0.23349158465862274,
0.18502001464366913, 0.1913882941007614, 0.311190664768219, 0.05527481436729431, 0.0855673998594284, 0.16441373527050018,
0.11771653592586517, 0.13056251406669617, 0.3395051956176758, 0.088826023042202, 0.05102713406085968, 0.09029881656169891,
0.14738988876342773, 0.05677853524684906, 0.06908028572797775, 0.23434215784072876, 0.06538316607475281,
0.09631697833538055, 0.17425137758255005, 0.09268894791603088, 0.05898330360651016, 0.024849584326148033,
0.12852415442466736, 0.27541637420654297, 0.06493933498859406, 0.04966285824775696, 0.15830424427986145, 0.4835587441921234,
0.14416643977165222, 0.04458580166101456, 0.07078167796134949, 0.0888843983411789, 0.07056154310703278, 0.11186349391937256,
0.05668199807405472, 0.10069091618061066, 0.06706306338310242, 0.049577996134757996, 0.16750244796276093,
0.1273132562637329, 0.0888298898935318, 0.07718336582183838, 0.17826154828071594, 0.04185512289404869, 0.06440159678459167,
0.1565883755683899, 0.07683703303337097, 0.05629356950521469, 0.20288079977035522, 0.16968970000743866,
0.054671794176101685, 0.09659813344478607, 0.1492307484149933, 0.1399400383234024, 0.06627979874610901,
0.040523968636989594, 0.19162079691886902, 0.05534039065241814, 0.1384289562702179, 0.013441801071166992,

Figure 3.3: Feature vectors example

spatiotemporal features that are invariant to small changes in the input data.

3.1.3 Feature extraction from videos

In the feature extraction section of the thesis, the first step is to preprocess the videos to ensure they are in a standardized format and size [11]. This typically involves converting the videos to a common format: we use MP4, and resizing them to a fixed size.

Once the videos are preprocessed, they are fed into the I3D network to extract features. However, processing each video individually can be timeconsuming and computationally intensive, especially when dealing with large datasets. We have over 96,000 videos that length are from 2 seconds to 3 second. To improve efficiency, batch processing techniques can be used to process multiple videos simultaneously. In batch processing, multiple videos are grouped together and fed into the network as a single batch.

Batch processing offers several benefits. Firstly, it can significantly reduce the time required to process large datasets, as multiple videos can be processed in parallel. Secondly, it can help to improve the generalization ability of the network, as it is exposed to a more diverse set of videos in each batch[23]. Finally, it can help to reduce memory usage, as the network's parameters are updated less frequently.

In this thesis, we used batch processing with a batch size of 16, which means that 16 videos were processed simultaneously in each batch. This allowed us to process large datasets efficiently and to train the network more effectively. However, the optimal batch size may vary depending on the specific dataset and hardware configuration being used.

To optimize the feature extraction pipeline for efficiency, it's also important to consider other factors such as parallelization techniques and hardware acceleration. For example, using a GPU or other hardware accelerator can significantly improve the speed of feature extraction, especially when dealing with large datasets.

In addition to the data augmentation techniques mentioned above, we also used another approach to further increase the diversity of our dataset. Specifically, we chose 5 different jersey colors commonly used in soccer matches, and we generated new samples by replacing the original jerseys with each of the 5 colors. This approach helped us to simulate different scenarios where the players have different uniforms, which could improve the model's ability to recognize actions under varying conditions. By applying this approach, we were able to create more variations of each action in our dataset. For instance, a simple action like a pass between two players could look different based on the color of the jerseys, and the model would have to learn to recognize the action regardless of the uniform colors. This technique also helped to improve the robustness of the model and ensure that it could accurately recognize actions in a variety of situations.

Overall, by using a combination of data augmentation techniques like rotation, masking, cropping, frame selection, and jersey color replacement, we were able to generate a more diverse and representative dataset. This helped us to improve the performance of our model and ensure that it could accurately recognize actions in real-world scenarios.

3.2 Introduction to MLP framework

Multi-layer perceptron (MLP) is a type of artificial neural network that consists of multiple layers of nodes, each layer connected to the previous one via weighted connections [17]. It is a feed-forward neural network that has been widely used in various applications, such as image recognition, natural language processing, and speech recognition. The classic MLP architecture consists of an input layer, one or more hidden layers, and an output layer, the architecture is shwon in Figure 3.4. Each layer contains a set of neurons that compute a weighted sum of their inputs and apply an activation function to the result. The weights between the neurons are learned during the training process using backpropagation, which adjusts the weights to minimize the error between the network's predicted output and the desired output.



Figure 3.4: MLP architecture

Despite being a relatively old architecture, MLPs have several advantages over other neural networks. One of the main advantages of MLPs is their simplicity and ease of implementation, making them suitable for a wide range of applications. Additionally, MLPs can handle non-linear input-output mappings and can approximate any continuous function to arbitrary accuracy with enough hidden nodes. Moreover, MLPs are computationally efficient and can be trained on relatively small datasets [4].

However, MLPs also have some limitations. One of the main challenges with MLPs is overfitting, where the network learns to memorize the training data instead of generalizing to new data. Additionally, MLPs can be sensitive to the choice of hyperparameters, such as the number of hidden nodes, learning rate, and activation functions. Moreover, MLPs can be computationally expensive when dealing with large datasets, and the training process can be slow.

One of the reasons we chose MLP as the final layer of our model is be-

cause it is well-suited for handling high-dimensional data, which is the case with the features extracted from the I3D network. The I3D features are a set of spatio-temporal features that are extracted from video data, resulting in highdimensional feature vectors([N,1024,1]). An MLP is able to effectively process these high-dimensional feature vectors and generate predictions based on them.

Another advantage of using an MLP is its ability to learn complex non-linear relationships between the input features and output labels. The MLP can learn to extract the most important features from the I3D feature vectors and use them to make accurate predictions about the action being performed in the video.

Furthermore, unlike traditional CNNs, MLPs are not limited by fixed input sizes, and they can handle input data of varying dimensions [6]. This is important for our model, as the size of the input video frames can vary, and the number of frames in each video sequence can also vary. MLPs can handle this variability and adapt to the input data accordingly.

Additionally, the MLP framework has been used successfully in other action recognition tasks, and it has shown to achieve state-of-the-art performance in some cases. Therefore, we have chosen MLP as our final layer to achieve the best possible performance on our action recognition task.

Overall, we chose MLP as our framework for action recognition due to its simplicity and effectiveness in handling our relatively small dataset and due to its ability to handle high-dimensional data, learn complex non-linear relationships, and adapt to varying input data sizes. MLP also provided us with the flexibility to adjust the experiment with different architectures to optimize the performance of our model. These factors, coupled with the success of MLP in other action recognition tasks, make it an ideal choice for our model to achieve state-of-the-art performance. In the following sections, we will discuss in detail the architecture and training process of our MLP model for action recognition.

3.2.1 MLP framework design

The MLP framework we used in our action recognition model consists of a fully connected neural network that takes in the I3D feature vectors extracted from the video frames. The goal of the MLP is to learn a mapping between the input features and the output classes. The MLP architecture consists of several fully connected layers, which are also known as dense layers. These layers allow each neuron in one layer to connect to every neuron in the previous layer.

We chose to use an MLP framework for several reasons. Firstly, MLPs are simple and effective at learning complex mappings between inputs and outputs, making them a popular choice for many machine learning tasks. Additionally, MLPs are computationally efficient and can be trained quickly on large datasets. Furthermore, MLPs are easy to interpret, allowing us to gain insights into how the model is making predictions.

The MLP framework is particularly well-suited for our action recognition model, as it allows us to connect the I3D features extracted from each video frame to the output classes directly. The I3D features are flattened and passed through several fully connected layers before being fed into the output layer. The output layer consists of several neurons, one for each action class, and uses the softmax activation function to output the predicted class probabilities [13]. To prevent overfitting, we used dropout regularization on the fully connected layers of the MLP. Dropout is a technique where a random subset of neurons is deactivated during each training iteration, which helps to prevent the model from relying too heavily on any single feature. Additionally, we used L2 regularization to penalize the model for large weights, which can lead to overfitting.

To optimize the MLP, we used the categorical cross-entropy loss function, which is commonly used for multi-class classification tasks. We trained the MLP using the Adam optimizer, which is a variant of stochastic gradient descent that adapts the learning rate based on the gradient magnitude of the parameters. We used a batch size of 16 and trained the model for 200 epochs.

The architecture of the MLP framework we used is shown in Figure 3.5. The input layer consists of the I3D feature vectors, which are flattened and passed through four fully connected layers. The first layer has 1024 neurons each, while the second layer has 250 neurons and the third layer has 100. We used ReLU activation functions on each of these layers. The output layer consists of four neurons, one for each of the action classes we were trying to predict(Shooting, Skating left, Skating right, Skating forward).

Overall, the MLP framework we used in our action recognition model proved to be effective at learning the complex mappings between the I3D features and the output classes. The use of dropout regularization and L2 regularization helped to prevent overfitting and improve the generalizability of the model. In the next section, we will discuss the implementation details of our model, including the data preprocessing steps and the hyperparameters used during training.



Figure 3.5: Our model architecture

3.2.2 Implementation details

In this section, we describe the specific details of how we implemented the MLP framework in our action recognition model.

Input data:

First, we preprocessed the input data to make it compatible with our MLP framework. As described earlier, the input to our model consists of 1024dimensional feature vectors extracted from the I3D model. We used these feature vectors as input to the MLP. To ensure that the MLP could handle the variation in input data, we normalized the input feature vectors by subtracting the mean and dividing by the standard deviation of the training data. We also applied dropout with a probability of 0.5 to the input layer to prevent overfitting.

MLP Architecture:

We designed the architecture of the MLP based on the number of input features and the number of output classes. Specifically, the MLP consisted of two fully connected hidden layers, each with 512 units and a ReLU activation function. We used dropout with a probability of 0.5 after each hidden layer to regularize the model. The output layer consisted of a fully connected layer with softmax activation, which produced a probability distribution over the output classes.

Our MLP architecture consisted of several fully connected layers with a varying number of neurons, depending on the specific task and dataset. Specifically, the MLP consisted of two fully connected hidden layers, each with 512 units and we used the Rectified Linear Unit (ReLU) activation function for all hidden layers, which has been shown to be effective for deep neural networks due to its ability to prevent the vanishing gradient problem. We used dropout with a probability of 0.5 after each hidden layer to regularize the MLP model and prevent overfitting. Overfitting occurs when the model learns the training data too well, including the noise or randomness present in the training data, which causes it to perform poorly on the test data or new unseen data.

The output layer of our MLP framework is an important component as it determines the final action classification output of the model. We used a softmax activation function in the output layer which is a widely used activation function for multi-class classification problems. The softmax function normalizes the output of the layer so that the sum of all the outputs is equal to 1, which can be interpreted as probabilities of belonging to each class. This helped us to easily calculate the cross-entropy loss during training, which is a common choice for multi-class classification problems like ours.

Cross-entropy loss is a popular choice for multi-class classification tasks because it penalizes the model heavily for incorrect predictions and rewards the model for correct predictions. This helped our model to learn the patterns and features of the different actions and improve its accuracy over time. Moreover, it allows us to measure the performance of our model by computing the average loss over all the samples in the training and validation set.

In addition, we also used batch normalization between each fully connected layer. This helped to speed up training by normalizing the activation of the previous layer, reducing the covariance shift problem, and making the optimization process more stable.

Optimization and Training:

We trained the MLP using cross-entropy loss and Adam optimizer with a learning rate of 1e-4. We used a batch size of 16 and trained the model for 100 epochs. We used early stopping with a patience of 10 epochs to prevent overfitting. We implemented the model using the PyTorch deep learning library and trained it on an NVIDIA GTX 1080 Ti GPU.

31

3.2.3 Action Recognition

To evaluate the performance of our MLP-based action recognition system, we conducted experiments on our ice hockey dataset. We split the dataset into training, validation, and test sets, with a ratio of 7:3 for the training and validation sets, and 20% of the validation set used for testing [28].

To exclude the influence of jersey color, we first trained the MLP on the dataset from one jersey color, and evaluated its performance on the validation set. We used the standard cross-entropy loss function and the Adam optimizer with a learning rate of 0.001. We also applied dropout with a probability of 0.5 after each hidden layer to regularize the model, and batch normalization between each fully connected layer to speed up training and improve stability.

After training, we evaluated the model on the validation set and generated a confusion matrix to visualize its performance. The confusion matrix is a valuable tool for evaluating the performance of a classification model in action recognition. It provides a comprehensive summary of how well the model is able to classify instances into different action classes. The matrix displays the counts of instances that fall into each combination of true and predicted classes. By analyzing the confusion matrix, we can calculate important evaluation metrics such as accuracy and precision. These metrics provide insights into the model's overall performance and its ability to minimize false positives and false negatives. The confusion matrix allows us to identify classes that are correctly classified and those that are more prone to misclassification, helping us understand the strengths and weaknesses of the model. This information is vital for improving the accuracy and effectiveness of action classification in dynamic sports scenarios. Next, we tested the trained model on data from 5 different jersey colors(the entire dataset), to investigate whether the recognition accuracy was influenced by jersey color. The results showed that the model was able to recognize actions with high accuracy, regardless of the jersey color.

To compare the performance of the model with different visualization techniques, we generated three types of plots: confusion matrix, PCA, and t-SNE.

PCA (Principal Component Analysis) and t-SNE (t-Distributed Stochastic Neighbor Embedding) are both techniques used in data visualization and dimensionality reduction. PCA is a technique used to reduce the number of features in a dataset while still retaining the variance of the original data. t-SNE, on the other hand, is a technique used for visualizing high-dimensional data by reducing it to a two-dimensional or three-dimensional representation. In our model, those two techniques can be used to visualize the features extracted by the MLP model. By projecting the features into a lower-dimensional space using PCA or t-SNE, we can gain insights into how the features are distributed and how they are related to each other. This can help us to identify patterns in the data and to better understand the performance of the model. We present all the plots in the next section.

3.3 Results

Action classification test result

We conducted a comprehensive analysis of the model's performance by comparing the PCA, confusion matrix, and t-SNE plots obtained during different



Figure 3.6: PCA for one jersey color

stages of our experiments. Firstly, we trained and tested the model using data from a single jersey color. The PCA plot revealed distinct clusters of actions, with minimal overlap(besides actions 1 and 2), indicating good separability of action classes (see Figure 3.6).

The number 0,1,2,3 each represent shooting, skating left, skating right, and skating forward. We can see the skating left and skating right have majority overlap in this PCA plot. The reason is these two actions looks particularly similar and their body movement are extremely small comparing to skating forward and shooting. With only one jersey dataset, we do not have enough samples to train the model to recognize these small actions difference.

The confusion matrix displayed high accuracy, with a majority of actions cor-



Figure 3.7: Confusion matrix for one jersey color

rectly classified and minimal misclassifications (see Figure 3.7). The confusion matrix showed that the model was able to recognize the majority of actions with high accuracy, achieving an overall validation accuracy of 93.54% and train accuracy of 97.18%. The validation accuracy represents the accuracy of the model in predicting the correct action labels on a separate validation dataset. This dataset is distinct from the training data and serves as an independent evaluation set to assess the generalization capability of the model. A validation accuracy of 93.54% indicates that the model correctly classified 93.54% of the actions in the validation dataset.

On the other hand, the train accuracy measures the accuracy of the model's predictions on the training dataset, which is the dataset used for training and parameter updates. A train accuracy of 97.18% suggests that the model per-



Figure 3.8: t-SNE for one jersey color

formed exceptionally well on the training data, correctly classifying 97.18% of the actions in the training set.

After training 100 Epoch with test dataset, we achieved the test accuracy of 93.25%. The achieved test accuracy of 93.25% demonstrates the effectiveness of our trained model in accurately classifying actions on unseen data. The test dataset serves as a crucial benchmark for evaluating the generalization capabilities of the model beyond the training and validation datasets. A high test accuracy indicates that our model can successfully recognize actions in synthetic scenarios, even when confronted with previously unseen examples. This further validated the effectiveness of our model in recognizing actions within the single jersey color dataset.

Additionally, the t-SNE plot provided a visual representation of action similarity, with tight clusters of similar actions observed in distinct regions (see Figure 3.8). Upon careful examination of the t-SNE plot, we observed distinct clusters of actions in separate regions, indicating that actions with similar characteristics tend to be grouped together. However, the t-SNE plot reveals the presence of some overlapping regions(skating left and skating right) among the action clusters. These overlapping regions suggest that these actions may share similarities in terms of visual appearance or motion patterns. Despite the overlap, the clusters remain fairly distinct and well-separated, indicating that our model can still differentiate between action classes even in cases of subtle similarities.

To assess the color-invariant capabilities of our model, we extended the testing phase to include a mixed dataset comprising five different jersey colors. Remarkably, the PCA plot continued to exhibit distinct and well-separated clusters of actions, indicating that the model maintained its ability to discern actions across various color variations (see Figure 3.9).

The PCA plot for the five jersey colors demonstrates significant improvements compared to the plot with a single jersey color. The increased availability of data samples from different colors for training the model has enhanced its ability to distinguish between action classes with subtle variations.

The reduced overlapping regions for skaing left and skating right actions indicate that the model is better equipped to differentiate between similar actions, even when they occur with different jersey colors. This improvement is crucial in real-world scenarios where actions may exhibit variations in visual appearance.



Figure 3.9: PCA for 5 jersey colors

The improved separation of action classes in the PCA plot validates the effectiveness of using a diverse dataset with multiple jersey colors. By incorporating data samples from various colors, we have enhanced the model's discriminative power and its robustness to variations in color or appearance.

The confusion matrix for the mixed color dataset demonstrated high accuracy, with actions correctly classified across different jersey colors, further confirming the model's robustness in handling color variations (see Figure 3.10).

The confusion matrix displayed high accuracy, with a majority of actions



Figure 3.10: Confusion Matrix for 5 jersey colors

correctly classified and minimal misclassifications (see Figure 3.7). The confusion matrix showed that the model was able to recognize the majority of actions with high accuracy, achieving an overall validation accuracy of 93.54% and train accuracy of 97.18%. After training 100 Epoch with test dataset, we achieved the test accuracy of 93.25%. This further validated the effectiveness of our model in recognizing actions within the single jersey color dataset.

Based on the provided accuracy results, we can draw several conclusions regarding the performance of the model on both the single jersey color dataset and the dataset containing five different colors. For the single jersey color dataset, the model achieved a high performance on the validation accuracy, training accuracy and test accuracy. It indicates its ability to accurately classify actions within this specific color category and showed the model has effectively learned the patterns and features associated with the actions in the training data, and make accurate predictions on unseen data.

When considering the dataset with five different jersey colors, the model's performance remains commendable. The training accuracy of 94.81% indicates that the model has adapted well to the increased complexity introduced by the additional colors. The validation accuracy of 92.01% shows that the model maintains good performance on new color variations. The test accuracy of 91.12% further confirms the model's ability to generalize to unseen data and perform reliably across different color scenarios.

The slight decrease in accuracy compared to the single color dataset could be attributed to the added complexity and variability introduced by the five different colors. However, the overall accuracy remains high, indicating that the model is capable of handling color variations and recognizing actions across different jersey colors with a reasonable level of accuracy.

These results highlight the robustness of the model in action classification tasks, as it maintains high accuracy levels even when confronted with diverse jersey colors. The model's ability to generalize well to different color scenarios is crucial for real-world applications where actions may occur in varying visual contexts.

Moreover, the t-SNE plot revealed consistent clusters of similar actions, regardless of the jersey color, further illustrating the model's ability to identify



Figure 3.11: t-SNE for 5 jersey colors

actions in a color-invariant manner (see Figure 3.11).

Upon observing this t-SNE plot, it is evident that the overlapping regions between these confusing actions are smaller in comparison to the t-SNE plot for the single jersey color dataset. This finding indicates that the model has become more proficient at distinguishing between actions performed in different jersey colors.

The reduced overlap suggests that the model has learned to identify and extract relevant features from the input data that are indicative of the specific action being performed, regardless of the jersey color. The smaller overlapping regions also imply that the model has successfully captured and utilized color-invariant features during the training process. By focusing on the intrinsic characteristics of the actions rather than relying solely on color information, the model has become more robust and capable of generalizing across different jersey colors This improved separation of data points in the t-SNE plot shows that the model's representations of the actions have become more distinct and well-defined.

The improved separation in the t-SNE plot is a positive indication of the model's ability to discriminate between actions, irrespective of the color of the jersey. This finding reinforces the notion that the model has learned to extract meaningful and discriminative features that are more indicative of the actions being performed rather than the color variations in the input data.

View-invariant test result

Having thoroughly examined the action classification results, we now shift our focus to evaluating the model's ability to recognize different camera views in a view-invariant manner. To assess this, we also trained another MLP which could classify view in order to use it for the downstream tasks, e.g. getting players orientation. Our objective was to determine whether the model could accurately identify the specific camera angle or viewpoint of the recorded actions. To evaluate the performance, we employed same visualization techniques, as action classification, to gain insights into the model's view-invariant capabilities. We present analysis of these plots and examine how well the model can distinguish between different camera angles, providing valuable insights into



Figure 3.12: PCA for views

its robustness and generalization ability across various viewpoints.

Upon examining the PCA plot, we can observe distinct clusters forming for each view(see Figure 3.12), indicating that the model has successfully learned to differentiate between different camera angles. The variance captured by the principal components suggests that the model has effectively captured the visual characteristics unique to each view. The distinct clusters in the PCA plot demonstrate the model's ability to discriminate and identify different camera angles accurately.

The confusion matrix (see Figure 3.13) provides valuable insights into the model's performance in classifying the 24 views. Each cell in the matrix represents the number of predictions made by the model for each view. The con-



Figure 3.13: Confusion Matrix for views

fusion matrix showed that the model was able to recognize the all of views with high accuracy, achieving an overall validation accuracy of 92.08% and train accuracy of 90.26%. After training 100 Epoch with test dataset, we achieved the test accuracy of 90.83%. By analyzing the confusion matrix, we can evaluate the model's accuracy in correctly identifying each view. The diagonal elements of the confusion matrix indicate the number of correct predictions, while off-diagonal elements represent misclassifications. Ideally, each views number should be the same. By assessing the distribution of misclassifications, we can identify any patterns or trends that may indicate particular views that are more challenging for the model to differentiate.

Upon analyzing the confusion matrix for the view-invariant analysis, we observed that there was relatively lower accuracy for views 10 to 13 compared to



Figure 3.14: t-SNE for views

other views. Further examination revealed that views 10 to 13 were captured from camera positions located on the azimuth, which essentially means they were positioned around the same angle, symmetrically positioned on people's two side. Consequently, it becomes more challenging for the model to discern subtle differences between these camera views, resulting in lower accuracy. The azimuthal positioning of the cameras creates a similarity in the visual perspective, making it more difficult for the model to accurately distinguish between these views. Despite this challenge, the model demonstrates remarkable performance in recognizing other camera views, showcasing its ability to generalize and identify distinct perspectives effectively.

The t-SNE plot(see Figure 3.14) visualizes the embeddings of the 24 views in

a reduced-dimensional space. It allows us to assess the separability and clustering of different views. In the t-SNE plot, most of the views from the same angle are grouped together, which indicates that the model has learned to distinguish and group similar views. Conversely, some of the views are scattered and mixed(views positions located on the azimuth), it suggests that the model may struggle to differentiate between different camera angles.

CHAPTER 4 CONCLUSION

In this research, we have explored the use of data augmentation, I3D feature extraction, and MLP-based action classification to improve the accuracy and robustness of video-based action recognition. By summarizing the methods and processes employed throughout this study, we can draw meaningful conclusions regarding their effectiveness, strengths, and potential future applications.

Through the implementation of data augmentation techniques, we were able to generate a diverse and representative dataset. This approach enabled us to introduce variations in factors such as rotation, masking, and cropping, leading to improved model generalization and performance. The incorporation of these techniques allowed us to create augmented videos that captured a wider range of real-world scenarios, enhancing the overall accuracy of the action classification task.

Utilizing the I3D framework for feature extraction proved to be a pivotal step in capturing spatio-temporal information from videos. By leveraging 3D convolutions, the I3D model successfully learned to extract discriminative features that effectively represented the motion and appearance characteristics of various actions. This approach demonstrated its superiority over traditional 2D CNN architectures, providing a more comprehensive understanding of action dynamics and leading to improved classification accuracy.

The MLP-based action classification model showcased remarkable performance, achieving high accuracy on both training and validation datasets. The model's ability to recognize actions accurately and efficiently highlights the power of MLPs in capturing complex patterns within the extracted features. Additionally, the utilization of dropout regularization and batch normalization contributed to improved model generalization and stability.

While our research presented promising results, it is important to acknowledge some limitations. The computational requirements for training deep learning models with large datasets can be demanding, requiring access to highperformance computing resources. Furthermore, obtaining labeled video data for training can be a time-consuming and labor-intensive process.

Looking ahead, the developed methods and processes hold great potential for future applications. The insights gained from this research can be leveraged in various domains, including video surveillance, automonous driving, sports analytics, and human-computer interaction. Further exploration of advanced data augmentation techniques, as well as the integration of other deep learning architectures, could yield even more accurate and robust action recognition systems.

In the future, after collecting enough amount of data from real life videos as test dataset, we can test the algorithms on the real life dataset to examine the robustness and effectiveness of our model.

In conclusion, our research has demonstrated the effectiveness of data augmentation, I3D feature extraction, and MLP-based action classification in improving the accuracy and robustness of video-based action recognition in synthetic scenarios. By enhancing the dataset diversity, capturing spatio-temporal features, and leveraging MLPs for classification, it indicates that our model is not only capable of accurately classifying actions across different jersey colors, but also demonstrates the ability to recognize actions invariant to varying camera views. The findings of this study lay the foundation for further advancements in action recognition and hold immense potential for real-world applications in a variety of fields.

[1] Ashraf, Nazim, et al. "View Invariant Action Recognition Using Projective Depth." *Computer Vision and Image Understanding*, vol. 123, June 2014, pp. 41–52. ScienceDirect, doi:10.1016/j.cviu.2014.03.005.

[2] Baptista, Renato, et al. "View-Invariant Action Recognition from RGB Data via 3D Pose Estimation." *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 2542–46. IEEE Xplore, doi:10.1109/ICASSP.2019.8682904.

[3] Carreira, João, and Andrew Zisserman. *Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset*. IEEE Computer Society, 2017, pp. 4724–33. www.computer.org, doi:10.1109/CVPR.2017.502.

[4] Chenarlogh, Vahid Ashkani, et al. "A Multi-View Human Action Recognition System in Limited Data Case Using Multi-Stream CNN." *2019 5th Iranian Conference on Signal Processing and Intelligent Systems (ICSPIS)*, 2019, pp. 1–11. IEEE Xplore, doi:10.1109/ICSPIS48872.2019.9066079.

[5] Costa, Felipe, et al. "Video Action Classification through Graph Convolutional Networks:" *Proceedings of the 16th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, SCITEPRESS - Science and Technology Publications, 2021, pp. 490–97. DOI.org (Crossref), doi:10.5220/0010321304900497. [6] Dai, Chuan, et al. "ConMLP: MLP-Based Self-Supervised Contrastive Learning for Skeleton Data Analysis and Action Recognition." *Sensors*, vol. 23, no. 5, 5, Jan. 2023, p. 2452. www.mdpi.com, doi:10.3390/s23052452.

[7] Fabbri, Matteo, et al. *MOTSynth: How Can Synthetic Data Help Pedestrian Detection and Tracking?* arXiv:2108.09518, arXiv, 21 Aug. 2021. arXiv.org, doi:10.48550/arXiv.2108.09518.

[8] Hong, James, et al. *Video Pose Distillation for Few-Shot, Fine-Grained Sports Action Recognition*. 2021, pp. 9254–63. openaccess.thecvf.com, doi:10.1145/3460440.3461892.

[9] Iosifidis, Alexandros, et al. "View-Invariant Action Recognition Based on Artificial Neural Networks." *IEEE Transactions on Neural Networks and Learning Systems*, vol. 23, no. 3, Mar. 2012, pp. 412–24. IEEE Xplore, doi:10.1109/TNNLS.2011.2181865.

[10] Kenwright, D. N., et al. "Feature Extraction of Separation and Attachment Lines." *IEEE Transactions on Visualization and Computer Graphics*, vol. 5, no. 3, July 1999, pp. 201–14. doi:10.1109/2945.796287.

[11] Liu, Guocheng, et al. "I3D-Shufflenet Based Human Action Recognition." *Algorithms*, vol. 13, no. 11, 11, Nov. 2020, p. 301. www.mdpi.com, doi:10.3390/a13110301.

[12] Men, Qianhui, et al. "Focalized Contrastive View-Invariant Learning for Self-Supervised Skeleton-Based Action Recognition." *Neurocomputing*, vol.
537, June 2023, pp. 198–209. ScienceDirect, doi:10.1016/j.neucom.2023.03.070.

[13] Popescu, Marius-Constantin, et al. *Multilayer Perceptron and Neural

Networks*. no. 7, 2009.

[14] Popoola, Oluwatoyin P., and Kejun Wang. "Video-Based Abnormal Human Behavior Recognition—A Review." *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 42, no. 6, Nov. 2012, pp. 865–78. IEEE Xplore, doi:10.1109/TSMCC.2011.2178594.

[15] Shen, Yuping, and Hassan Foroosh. "View-Invariant Action Recognition from Point Triplets." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 10, Oct. 2009, pp. 1898–905. IEEE Xplore, doi:10.1109/TPAMI.2009.41.

[16] Simonyan, Karen, and Andrew Zisserman. "Two-Stream Convolutional Networks for Action Recognition in Videos." Advances in Neural Information Processing Systems, vol. 27, Curran Associates, Inc., 2014. Neural Information Processing Systems, doi:10.1007/978-3-319-16178-5-1.

[17] Talukdar, Jonti, and Bhavana Mehta. "Human Action Recognition System Using Good Features and Multilayer Perceptron Network." *2017 International Conference on Communication and Signal Processing (ICCSP)*, 2017, pp. 0317–23. IEEE Xplore, doi:10.1109/ICCSP.2017.8286369.

[18] Tran, Du, et al. *Learning Spatiotemporal Features with 3D Convolutional Networks*. arXiv:1412.0767, arXiv, 6 Oct. 2015. arXiv.org, doi:10.48550/arXiv.1412.0767.

[19] Ulhaq, Anwaar. "Deep Cross-View Convolutional Features for View-Invariant Action Recognition." *2018 IEEE International Conference on Image Processing, Applications and Systems (IPAS)*, 2018, pp. 137–42. IEEE Xplore,

51

doi:10.1109/IPAS.2018.8708853.

[20] Ullah, Mohib, et al. "Human Action Recognition in Videos Using Stable Features." *Signal Image Processing: An International Journal*, vol. 8, no. 6, Dec. 2017, pp. 01-10. DOI.org (Crossref), doi:10.5121/sipij.2017.8601.

Certainly! Here are more citations in MLA format:

[21] Varol, Gul, et al. *Learning From Synthetic Humans*. 2017, pp. 109–17. openaccess.thecvf.com, doi:10.1109/CVPR.2017.502.

[22] Wang, Peng. "Research on Sports Training Action Recognition Based on Deep Learning." *Scientific Programming*, vol. 2021, June 2021, p. e3396878. www.hindawi.com, doi:10.1155/2021/3396878.

[23] Wang, Xianyuan, et al. "I3D-LSTM: A New Model for Human Action Recognition." *IOP Conference Series: Materials Science and Engineering*, vol. 569, no. 3, July 2019, p. 032035. Institute of Physics, doi:10.1088/1757-899X/569/3/032035.

[24] Xia, Lu, et al. "View Invariant Human Action Recognition Using Histograms of 3D Joints." *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 2012, pp. 20–27. IEEE Xplore, doi:10.1109/CVPRW.2012.6239233.

[25] Xie, Saining, et al. *Rethinking Spatiotemporal Feature Learning: Speed-Accuracy Trade-Offs in Video Classification*. arXiv:1712.04851, arXiv, 26 July 2018. arXiv.org, doi:10.48550/arXiv.1712.04851.

[26] Xing, Weiwei, et al. "View-Invariant Gait Recognition Method by Three-

Dimensional Convolutional Neural Network." *Journal of Electronic Imaging*, vol. 27, no. 1, Jan. 2018, p. 013010. www.spiedigitallibrary.org, doi:10.1117/1.JEI.27.1.013010.

[27] Yamazaki, Yoshihiro, et al. *Audio Visual Scene-Aware Dialog Generation with Transformer-Based Video Representations*. arXiv:2202.09979, arXiv, 20 Feb. 2022. arXiv.org, doi:10.48550/arXiv.2202.09979.

[28] Yao, Guangle, et al. "A Review of Convolutional-Neural-Network-Based Action Recognition." *Pattern Recognition Letters*, vol. 118, Feb. 2019, pp. 14–22. ScienceDirect, doi:10.1016/j.patrec.2018.05.018.

[29] You, Jian, et al. "Multi-Stream I3D Network for Fine-Grained Action Recognition." *2018 IEEE 4th Information Technology and Mechatronics Engineering Conference (ITOEC)*, 2018, pp. 611–14. IEEE Xplore, doi:10.1109/ITOEC.2018.8740606.

[30] Zhu, Yisheng, and Guangcan Liu. "Fine-Grained Action Recognition Using Multi-View Attentions." *The Visual Computer*, vol. 36, no. 9, Sept.
2020, pp. 1771–81. Springer Link, doi:10.1007/s00371-019-01770-y.