

A Survey of Deep Learning in Sports Applications: Perception, Comprehension, and Decision

Zhonghan Zhao*, Wenhao Chai*, Shengyu Hao, Wenhao Hu, Guanhong Wang, Shidong Cao, Mingli Song, *Senior Member, IEEE*, Jenq-Neng Hwang, *Fellow, IEEE*, Gaoang Wang[†], *Member, IEEE*

Abstract—Deep learning has the potential to revolutionize sports performance, with applications ranging from perception and comprehension to decision. This paper presents a comprehensive survey of deep learning in sports performance, focusing on three main aspects: algorithms, datasets and virtual environments, and challenges. Firstly, we discuss the hierarchical structure of deep learning algorithms in sports performance which includes perception, comprehension and decision while comparing their strengths and weaknesses. Secondly, we list widely used existing datasets in sports and highlight their characteristics and limitations. Finally, we summarize current challenges and point out future trends of deep learning in sports. Our survey provides valuable reference material for researchers interested in deep learning in sports applications.

Index Terms—Sports Performance, Internet of Things, Computer Vision, Deep Learning, Survey

I. INTRODUCTION

ARTIFICIAL Intelligence (AI) has found wide-ranging applications and holds a bright future in the world of sports. Its ever-growing involvement is set to revolutionize the industry in myriad ways, enabling new heights of efficiency and precision.

A prominent application of AI in sports is the use of deep learning techniques. Specifically, these advanced algorithms are utilized in areas like player performance analysis, injury prediction, and game strategy formulation [1]. Through capturing and processing large amounts of data, deep learning models can predict outcomes, uncover patterns, and formulate strategies that might not be evident to the human eye. This seamless integration of deep learning and the sports industry [2], [3] exemplifies how technology is enhancing our ability to optimize sporting performance and decision-making.

Although predicting and optimizing athletic performance has numerous advantages, it remains a complex problem. Traditionally, sports experts like coaches, managers, scouts, and sports health professionals have relied on conventional analytical methods to tackle these challenges. However, gathering statistical data and analyzing decisions manually is a demanding and time-consuming endeavor [4]. Consequently,

* Equal contribution.

[†] Corresponding author: Gaoang Wang.

Zhonghan Zhao, Shengyu Hao, Wenhao Hu, Guanhong Wang, Mingli Song are with College of Computer Science and Technology, Zhejiang University.

Shidong Cao is with the Zhejiang University-University of Illinois Urbana-Champaign Institute, Zhejiang University.

Gaoang Wang is with the Zhejiang University-University of Illinois Urbana-Champaign Institute, and College of Computer Science and Technology, Zhejiang University.

Wenhao Chai and Jenq-Neng Hwang are with the University of Washington.

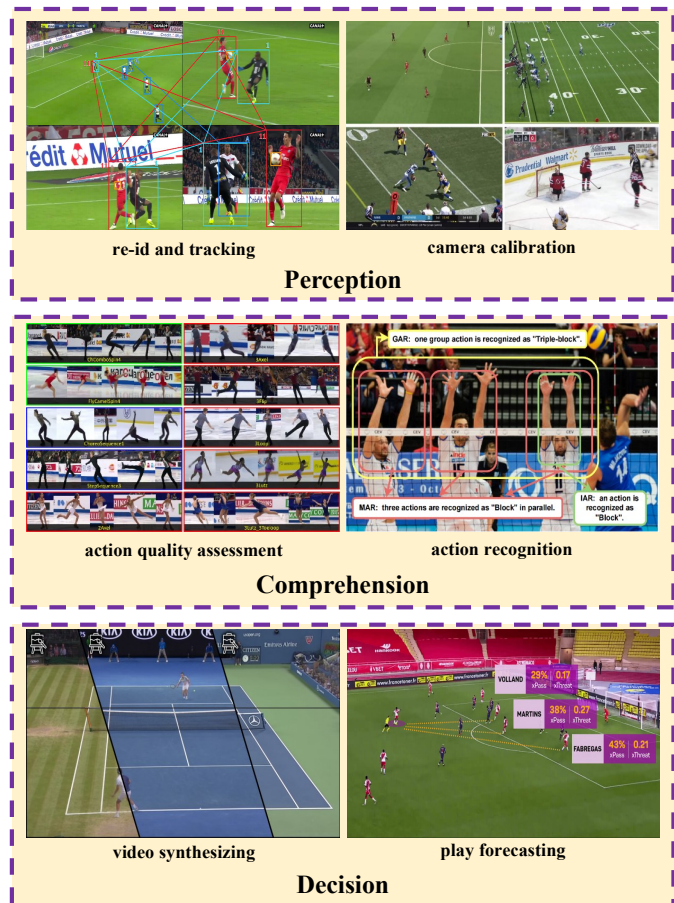


Fig. 1. The examples of the applications in sports performance in perception, comprehension, and decision.

an automated system powered by machine learning emerges as a promising solution that can revolutionize the sports industry by automating the processing of large-scale data.

In recent years, there has been a notable increase in comprehensive surveys exploring the applications of machine learning and deep learning in sports performance. These surveys cover a wide range of topics, including the recognition of sports-specific movements [5], mining sports data [6], and employing AI techniques in team sports [7]. While some surveys focus on specific sports like soccer [7] and badminton [8], others concentrate on particular tasks within computer vision, such as video action recognition [9], video action quality assessment [10], and ball tracking [11]. Furthermore, several studies explore the usage of wearable technology [12], [13]

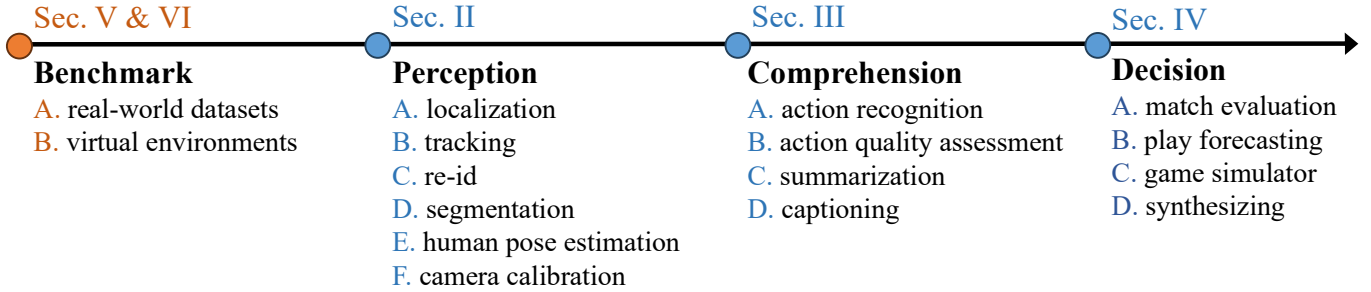


Fig. 2. **Taxonomy.** A hierarchical structure that contains three categories of tasks: Perception, Comprehension, and Decision, as well as Benchmark.

and motion capture systems [14] in sports, with a particular emphasis on the Internet of Things (IoT).

Previous studies [15], [16] have employed a hierarchical approach to analyze sports performance, starting from lower-level aspects and progressing to higher-level components, while also providing training recommendations. In order to comprehend the utilization of deep learning in sports, we have segmented it into three levels: **Perception**, **Comprehension**, and **Decision**. Additionally, we have categorized diverse datasets according to specific sports disciplines and outlined the primary challenges associated with deep learning methodologies and datasets. Furthermore, we have highlighted the future directions of deep learning in motion, based on the current work built upon foundational models.

The contributions of this comprehensive survey of deep learning in sports performance can be summarized in three key aspects.

- We propose a hierarchical structure that systematically divides deep learning tasks into three categories: Perception, Comprehension, and Decision, covering low-level to high-level tasks.
- We provide a summary of sports datasets and virtual environments. Meanwhile, this paper covers dozens of sports scenarios, processing both visual information and IoT sensor data.
- We summarize the current challenges and future feasible research directions for deep learning in various sports fields.

The paper is organized as follows: Section II, III, and IV introduce different tasks with methods for perception, comprehension, and decision tasks in sports. Section V and VI discuss the sports-related datasets and virtual environments. In Section VII and VIII, we highlight the current challenges and future trends of deep learning in sports. Lastly, we conclude the paper in Section IX.

II. PERCEPTION

Perception involves the fundamental interpretation of acquired data. This section presents different deep-learning methodologies tailored to specific sports tasks at the perception level as shown in Figure 3. The subsequent perception segment will encompass tasks such as player tracking, player pose recognition, player instance segmentation, ball localization, camera calibration *etc.*.

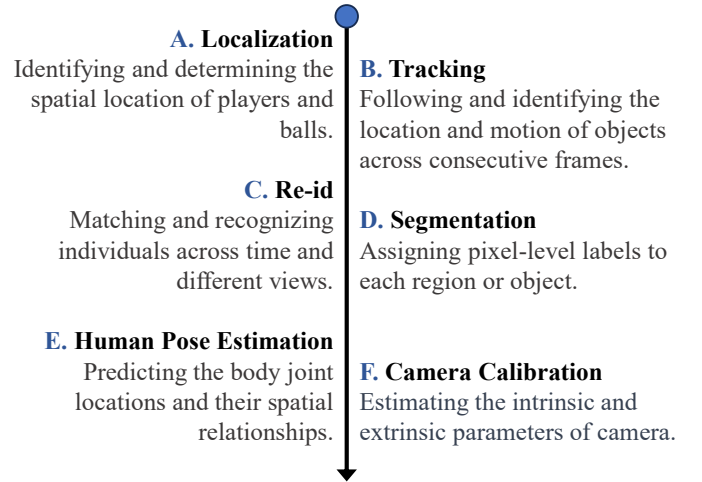


Fig. 3. Taxonomy and description of perception tasks.

A. Player and Ball Localization

Player and ball localization aims at identifying and determining the spatial location of players and balls, which is an essential undertaking in sports video analysis. Precisely identifying these entities can provide valuable insights into team performance, enabling coaches to make well-informed decisions using data. In recent years, numerous deep learning-based techniques have emerged, specifically designed for accurately localizing players and balls in a variety of sports, such as soccer, basketball, and cricket.

1) **Player Localization:** Player localization or detection [17]–[19] serves as a foundation for various downstream applications within the field of sports analysis. These applications include identifying player jersey numbers [20]–[22] and teams [23], [24], predicting movements and intentions [25]–[27]. Some works [28] leverage advancements in generic object detection to enhance the understanding of soccer broadcasts. Others [24] focus on unsupervised methods to differentiate player teams and employ multi-modal and multi-view distillation approaches for player detection in amateur sports [29]. Vandeghen *et al.* [30] introduces a distillation method for semi-supervised learning, which significantly reduces the reliance on labeled data. Moreover, certain studies [31], [32] utilize player localization for action recognition and spotting. Object tracking [33] is also crucial for the temporal localization of players.

2) **Ball Localization:** Ball localization provides crucial 3D positional information about the ball, which offers comprehensive insights into its movement state [11]. This task involves estimating the ball’s diameter in pixels within an image patch centered on the ball, and it finds applications in various aspects of game analytics [34]. These applications include automated offside detection in soccer [35], release point localization in basketball [36], and event spotting in table tennis [37].

Existing solutions often rely on multi-view points [38]–[40] to triangulate the 2D positions of the ball detected in individual frames, providing robustness against occlusions that are prevalent in team sports such as basketball or American football.

However, in single-view ball 3D localization, occlusion becomes a significant challenge. Most approaches resort to fitting 3D ballistic trajectories based on the 2D detections [40], [41], limiting their effectiveness in detecting the ball during free fall when it follows ballistic paths. Nonetheless, in many game situations, the ball may be partially visible or fully occluded during free fall. Van *et al.* [36], [42] address these limitations by deviating from assumptions of ballistic trajectory, time consistency, and clear visibility. They propose an image-based method that detects the ball’s center and estimates its size within the image space, bridging the gap between trajectory predictions offered by ballistic approaches. Additionally, there are also works on reconstructing 3D shuttle trajectories in badminton [43].

B. Player and Ball Tracking

Player and ball tracking is the process of consistently following and identifying the location and motion of objects across consecutive frames. This tracking operation is integral to facilitating an automated understanding of sports activities.

1) **Player Tracking:** Tracking players in the temporal dimension is immensely valuable for gathering player-specific statistics. Recent works [44], [45] utilize the SORT algorithm [46], which combines Kalman filtering with the Hungarian algorithm to associate overlapping bounding boxes. Additionally, Hurault *et al.* [47] employ a self-supervised approach, fine-tuning an object detection model trained on generic objects specifically for soccer player detection and tracking.

In player tracking, a common challenge arises from similar appearances that make it difficult to associate detections and maintain identity consistency. Intuitively, integrating information from other tasks can assist in tracking. Some works [48] explore patterns in jersey numbers, team classification, and pose-guided partial features to handle player identity switches and correlate player IDs using the K-shortest path algorithm. In dance scenarios, incorporating skeleton features from human pose estimation significantly improves tracking performance in challenging scenes with uniform costumes and diverse movements [49].

To address identity mismatches during occlusions, Naik *et al.* [44] utilize the difference in jersey color between teams and referees in soccer. They update color masks in the tracker module from frame to frame, assigning tracker IDs based on jersey color. Additionally, other works [45], [50] tackle occlusion issues using DeepSort [51].

2) **Ball Tracking:** Accurately recognizing and tracking a high-speed, small ball from raw video poses significant challenges. Huang *et al.* [52] propose a heatmap-based deep learning network [53], [54] to identify the ball image in a single frame and learn its flight patterns across consecutive frames. Furthermore, precise ball tracking is essential for assisting other tasks, such as recognizing spin actions in table tennis [55] by combining ball tracking information.

C. Player Re-identification

Player re-identification (ReID) is a task of matching and recognizing individuals across time and different views. In technical terms, this involves comparing an image of a person, referred to as the query, against a collection of other images within a large database, known as the gallery, taken from various camera viewpoints. In sports, the ReID task aims to re-identify players, coaches, and referees across images captured successively from moving cameras [36], [56]. Challenges such as similar appearances and occlusions and the low resolution of player details in broadcast videos make player re-identification a challenging task.

Addressing these challenges, many approaches have focused on recognizing jersey numbers as a means of identifying players [22], [57], or have employed part-based classification techniques [58]. Recently, Teket *et al.* [59] proposed a real-time capable pipeline for player detection and identification using a Siamese network with a triplet loss to distinguish players from each other, without relying on fixed classes or jersey numbers. An *et al.* [60] introduced a multi-granularity network with an attention mechanism for player ReID, while Habel *et al.* [61] utilized CLIP with InfoNCE loss as an objective, focusing on class-agnostic approaches.

To address the issue of low-resolution player details in multi-view soccer match broadcast videos, Comandur *et al.* [56] proposed a model that re-identifies players by ranking replay frames based on their distance to a given action frame, incorporating a centroid loss, triplet loss, and cross-entropy loss to increase the margin between clusters.

In addition, some researchers have explored semi-supervised or weakly supervised methods. Maglo *et al.* [62] developed a semi-interactive system using a transformer-based architecture for player ReID. Similarly, in hockey, Vats *et al.* [63] employed a weakly-supervised training approach with cross-entropy loss to predict jersey numbers as a form of classification.

D. Player Instance Segmentation

Player instance segmentation aims at assigning pixel-level labels to each player. In player instance segmentation, occlusion is the key problem, especially in crowded regions, like basketball [36]. Some works [64], [65] utilize online specific copy-paste method [66] to address the occlusion issue.

Moreover, instance segmentation features can be used to distinguish different players in team sports with different actions [24], [67]. In hockey, Koshkina *et al.* [24] use Mask R-CNN [68] to detect and segment each person on the playing surface. Zhang *et al.* [67] utilize the segmentation task to enhance throw action recognition [67] and event spotting [37].

E. Player Pose Estimation

Player pose estimation contributes to predicting the body joint locations and their spatial relationships. It often serves as a foundational component for various tasks [69], but there are limited works that specifically address the unique characteristics of sports scenes, such as their long processing times, reliance on appearance models, and sensitivity to calibration errors and noisy detections.

Recent approaches have employed OpenPose [70] for action detection or positional predictions of different elements in sports practice [71]–[73]. For sports with rapidly changing player movements, such as table tennis, some works [74] utilize a long short-term pose prediction network [75] to ensure real-time performance. In specific actions analysis of sports videos, certain works [76] use pose estimation techniques. Furthermore, Thilakarathne *et al.* [77] utilize tracked poses as input to enhance group activity recognition in volleyball. In more spatial heavy sports where less action or movement is present but more complexity lies in the poses, researchers focus on providing practitioners with tools to verify the correctness of their poses for more efficient learning, such as in Taichi [78] and Yoga [79].

F. Camera Calibration

Camera calibration in sports, also known as field registration, aims at estimating the intrinsic and extrinsic parameters of cameras. Homography provides a mapping between a planar field and the corresponding visible area within an image. Field calibration plays a crucial role in tasks that benefit from position information within the stadium, such as 3D player tracking on the field. Various approaches have been employed to solve sport-field registrations in different sports domains, including tennis, volleyball, and soccer [80], [81], often relying on keypoint retrieval methods.

With the emergence of deep learning, recent approaches focus on learning a representation of the visible sports field through various forms of semantic segmentation [32], [82]–[84]. These approaches either directly predict or regress an initial homography matrix [85]–[87], or search for the best matching homography in a reference database [84], [88] that contains synthetic images with known homography matrices or camera parameters. In other cases [83], [84], a dictionary of camera views is utilized, connecting an image projection of a synthetic reference field model to a homography. The segmentation is then linked to the closest synthetic view in the dictionary, providing an approximate camera parameter estimate, which is further refined for the final prediction.

III. COMPREHENSION

Comprehension can be defined as the process of understanding and analyzing data. It involves higher-level tasks compared to the perception stage discussed in Section II. In order to achieve a comprehensive understanding of sports, the implementation can utilize raw data and directly or indirectly incorporate the tasks from the perception layer. Namely, it can utilize the outputs obtained from the perception network, such as human skeletons, depth images *etc.*

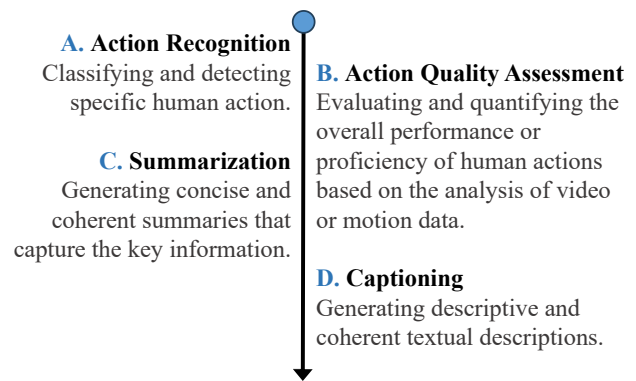


Fig. 4. Taxonomy and description of comprehension tasks.

In this section, we delve into specific tasks related to understanding and analyzing sports as shown in Figure 4. These tasks include individual and group action recognition, action quality assessment, action spotting, sports video summarization, and captioning.

A. Individual Action Recognition

Player action recognition targets classifying and detecting specific human action. Individual action recognition is commonly used for automated statistical analysis of individual sports, such as counting the occurrences of specific actions. Moreover, it plays a crucial role in analyzing tactics, identifying key moments in matches, and tracking player activity, including metrics like running distance and performance. This analysis can assist players and coaches in identifying the essential technical factors required for achieving better results. In team sports, coaches need to monitor all players on the field and their respective actions, particularly how they execute them. Therefore, an automated system capable of tracking all these elements could greatly contribute to the players' success. However, this casts a significant challenge for computers due to the simultaneous occurrence of different actions by multiple players on the sports field, leading to issues such as occlusion and confusing scenes.

While end-to-end models [96], [97], [110] are commonly employed in the literature on video action recognition, they are often better suited for coarse-grained classification tasks [111]–[113], which focus on broader categories like punches or kicks. In contrast, most sports require more fine-grained methods capable of distinguishing between specific techniques within these broader categories [114], [115].

Fine-grained action recognition within a single sport can help mitigate contextual biases present in coarse-grained tasks, making it an increasingly important research area [114], [116], [117]. Skeleton-based methods [118]–[120] have gained popularity for fine-grained action recognition in body-centric sports. These approaches utilize 2D or 3D human pose as input for recognizing human actions. By representing the human skeleton as a graph with joint positions as nodes and modeling the movement as changes in these graph coordinates over time, both the spatial and temporal aspects of the action can be captured. Additionally, some works [121]–[123] focus on fine-

TABLE I
DEEP LEARNING MODELS FOR SPORTS COMPREHENSION. “IAR”, “GAR”, “AQA” STAND FOR INDIVIDUAL ACTION RECOGNITION, GROUP ACTION RECOGNITION, ACTION QUALITY ASSESSMENT.

Task	Method	Venue	Benchmark	Link
IAR	TSM [89]	ICCV-2019	FineGym, P ² A	✓
	CSN [90]	ICCV-2019	Sports 1M	✓
	SlowFast [91]	ICCV-2019	P ² A, Diving48	✓
	G-Blend [92]	CVPR-2020	Sports 1M	✓
	AGCN [93]	TIP-2020	FSD-10	✓
	ResGCN [94]	MM-2020	FSD-10	✓
	MoViNet [95]	CVPR-2021	P ² A	✓
	TimeSformer [96]	ICML-2021	P ² A, Diving48	✓
	ViSwin [97]	arXiv-2021	P ² A	✓
	ORViT [98]	arXiv-2021	Diving48	✓
	BEVT [99]	arXiv-2021	Diving48	✓
VIMPAC [100]	arXiv-2021	Diving48	✓	
CTR-GCN [101]	ICCV-2021	FSD-10	✓	
GAR	DIN [102]	ICCV-2021	Diving48, HierVolleyball-v2	✓
	PoseC3D [103]	CVPR-2022	FineGym, FSD-10, HierVolleyball-v2	✓
AQA	S3D [104]	ICIP-2018	AQA-7	✓
	C3D-LSTM [105]	WACV-2019	AQA-7	✓
	C3D-AVG-MTL [106]	CVPR-2019	MTL-AQA	✓
	C3D-MSLSTM [107]	TCSVT-2020	FisV, MIT-Skate	✓
	I3D-USDL [108]	CVPR-2020	AQA-7, MTL-AQA	✓
	TSA [109]	ACM MM 2021	FR-FS, AQA-7, MTL-AQA	✓

grained action recognition in sports that do not involve body-centric actions.

B. Group Action Recognition

Group activity recognition involves recognizing activities performed by multiple individuals or objects. It plays a significant role in automated human behavior analysis in various fields, including sports, healthcare, and surveillance. Unlike multi-player activity recognition, group / team action recognition focuses on identifying a single group action that arises from the collective actions and interactions of each player within the group. This poses greater challenges compared to individual action recognition and requires the integration of multiple computer vision techniques.

Due to the involvement of multiple players, modeling player interaction relations becomes essential in group action analysis. In general, actor interaction relations can be modeled using graph convolutional networks (GCN) or Transformers in various methods. Transformer-based methods [124]–[129] often explicitly represent spatiotemporal relations and employ attention-based techniques to model individual relations for inferring group activity. GCN-based methods [102], [130] construct relational graphs of the actors and simultaneously explore spatial and temporal actor interactions using graph convolution networks.

Among them, Yan *et al.* [126] construct separate spatial and temporal relation graphs to model actor relations. Gavriluk *et al.* [124] encode temporal information using I3D [111] and establish spatial relations among actors using a vanilla transformer. Li *et al.* [129] introduces a cluster attention mechanism. Dual-AI [131] proposes a dual-path role interaction framework for group behavior recognition, incorporating temporal encoding of the actor into the transformer architecture. Moreover, the use of simple multi-layer perceptrons (MLP)

for feature extraction in group activity analysis [132] is an emerging approach with great potential.

Moreover, some other works focus more on specific action recognition through temporal localization rather than classification. Several automated methods have been proposed to identify important actions in a game by analyzing camera shots or semantic information. Studies [133]–[135] have explored human activity localization in sports videos, salient game action identification [136], [137], and automatic identification and summarization of game highlights [138]–[140]. Recent methods are more on soccer. For instance, Giancola *et al.* [141] introduce the concept of accurately identifying and localizing specific actions within uncut soccer broadcast videos. More recently, innovative methodologies have emerged in this field, aiming to automate the process. Cioppa *et al.* [142] propose the application of a context-aware loss function to enhance model performance. They later demonstrated how integrating camera calibration and player localization features can improve spotting capabilities [32]. Hong *et al.* [143] propose an efficient end-to-end training approach, while Darwish *et al.* [144] utilize spatiotemporal encoders. Alternative strategies, such as graph-based techniques [145] and transformer-based methods [146], offer fresh perspectives, particularly in handling relational data and addressing long-range dependencies. Lastly, Soares *et al.* [147], [148] have highlighted the potential of anchor-based methods in precise action localization and categorization.

C. Action Quality Assessment

Action quality assessment (AQA) is a method used to evaluate and quantify the overall performance or proficiency of human actions based on the analysis of video or motion data. AQA takes into account criteria such as technique, speed, and control to assess the movement and assign a score, which can be used to guide training and rehabilitation programs. AQA

has proven to be reliable and valid for assessing movement quality across various sports. Research in this field primarily focuses on analyzing the actions of athletes in the Olympic Games, such as diving, gymnastics, and other sports mentioned in Section V. Existing methods typically approach AQA as a regression task using various video representations supervised by scores.

Some studies concentrate on enhancing network structures to extract more distinct features. For instance, Xu *et al.* [107] propose self-attentive LSTM and multi-scale convolutional skip LSTM models to predict Total Element Score (TES) and Total Program Component Score (PCS) in figure skating by capturing local and global sequential information in long-term videos. Xiang *et al.* [104] divide the diving process into four stages and employ four independent P3D models for feature extraction. Pan *et al.* [149] develop a graph-based joint relation model that analyzes human node motion using the joint commonality module and the joint difference module. Parisi *et al.* [150] propose a recurrent neural network with a growing self-organizing structure to learn body motion sequences and facilitate matching. Kim *et al.* [151] model the action as a structured process and encode action units using an LSTM network. Wang *et al.* [109] introduce a tube self-attention module for feature aggregation, enabling efficient generation of spatial-temporal contextual information through sparse feature interactions. Yu *et al.* [152] construct a contrastive regression framework based on video-level features to rank videos and predict accurate scores.

Other studies focus on improving the performance of action quality assessment by designing network loss functions. Li *et al.* [153] propose an end-to-end framework that employs C3D as a feature extractor and integrates a ranking loss with the mean squared error (MSE) loss. Parmar *et al.* [106] explore the AQA model in a multi-task learning scenario by introducing three parallel prediction tasks: action recognition, comment generation, and AQA score regression. Tang *et al.* [108] propose an uncertainty-aware score distribution learning approach that takes into account difficulty levels during the modeling process, resulting in a more realistic simulation of the scoring process.

Furthermore, some studies focus on comparing the quality of paired actions. Bertasius *et al.* [154] propose a model for basketball games based on first-person perspective videos, utilizing a convolutional-LSTM network to detect events and evaluate the quality of any two movements.

D. Sports Video Summarization

Sports video summarization aims at generating concise and coherent summaries that capture the key information. It often prioritizes the recognition of player actions [155]. This research field aims to generate highlights of broadcasted sports videos, as these videos are often too lengthy for audiences to watch in their entirety. Given that many sports matches can have durations of 90-180 minutes, it becomes a challenging task to create a summary that includes only the most interesting and exciting events.

Agyeman *et al.* [155] employ a 3D ResNet CNN and LSTM-based deep model to detect five different soccer sports

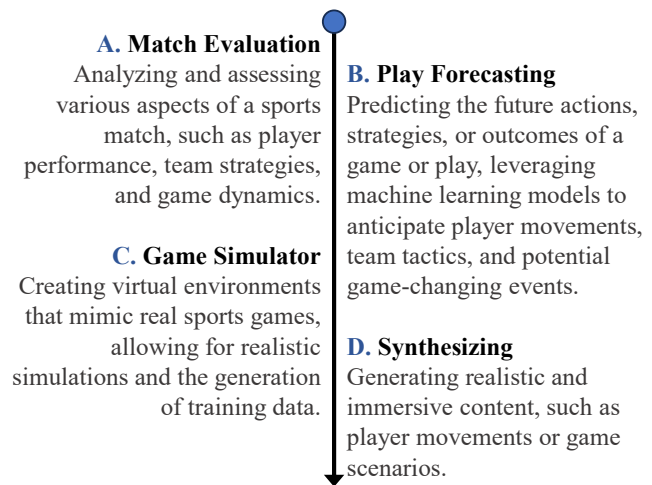


Fig. 5. Taxonomy and description of decision tasks.

action classes. Rafiq *et al.* [156] propose a transfer learning-based classification framework for categorizing cricket match clips [157] into five classes, utilizing a pre-trained AlexNet CNN and data augmentation. Shingrakhia *et al.* [158] present a multimodal hybrid approach for classifying sports video segments, utilizing the hybrid rotation forest deep belief network and a stacked RNN with deep attention for the identification of key events. Li *et al.* [159] propose a supervised action proposal guided Q-learning based hierarchical refinement approach for structure-adaptive summarization of soccer videos. While current research in sports video summarization focuses on specific sports, further efforts are needed to develop a generic framework that can support different types of sports videos.

E. Captioning

Sports video captioning involves generating descriptive and coherent textual descriptions. Sports video captioning models are designed to generate sentences that provide specific details related to a particular sport, which is a multimodal [160] task. For instance, in basketball, Yu *et al.* [161] propose a structure that consists of a CNN model for categorizing pixels into classes such as the ball, teams, and background, a model that captures player movements using optical flow features, and a component that models player relationships. These components are combined in a hierarchical structure to generate captions for NBA basketball videos. Similarly, attention mechanisms and hierarchical recurrent neural networks have been employed for captioning volleyball videos [162].

Furthermore, the utilization of multiple modalities can be extended to explore the creation of detailed captions or narratives for sports videos. Qi *et al.* [163] and Yu *et al.* [164] have successfully generated fine-grained textual descriptions for sports videos by incorporating attention mechanisms that consider motion modeling and contextual information related to groups and relationships.

IV. DECISION

The decision or decision-making process in sports involves the highest level of tasks, where the deployment or implicit perception and understanding of sports are essential before generating more abstract decisions. This section encompasses various tasks such as match evaluation, play forecasting, game simulation, player motion generation, and match generation as shown in Figure 5.

A. Match Evaluation

Match evaluation involves analyzing and assessing various aspects of a sport match, such as player performance, team strategies, and game dynamics. This task requires match modeling, often employing deep reinforcement learning methods. For instance, Wang *et al.* [165] develop a deep reinforcement learning model to study NBA games with the goal of minimizing offensive scores. Luo *et al.* [166] combine Q-function learning and inverse reinforcement learning to devise a unique ranking method and an alternating learning framework for a multi-agent ice hockey Markov game. Liu *et al.* [167] value player actions under different game contexts using Q-function learning and introduce a new player evaluation metric called the Game Impact Metric. Yanai *et al.* [168] model basketball games by extending the DDPG [169] architecture to evaluate the performance of players and teams.

B. Play Forecasting

Play Forecasting aims at predicting the future actions, strategies, or outcomes of a game or play, leveraging machine learning models to anticipate player movements, team tactics, and potential game changing events. The availability of accurate player and ball tracking data in professional sports venues has generated interest in assisting coaches and analysts with data-driven predictive models of player or team behavior [170], [171]. Several studies have utilized multiple years of match data to predict various aspects, such as predicting the ball placement in tennis [172], [173] and the likelihood of winning a point [174]. Le *et al.* [175] focus on predicting how NBA defenses will react to different offensive plays, while Power *et al.* [176] analyze the risk-reward of passes in soccer. In a more recent work, Wang *et al.* [177] delve into the analysis of where and what strokes to return in badminton.

C. Game Simulators

Game simulators typically aim at creating virtual environments that mimic real sports games, allowing for realistic simulations and the generation of training data [178]–[181]. These virtual environments, which are discussed in detail in Section VI, allow agents to move freely based on specific algorithms, simulating real-world sports scenarios. Within such environments, deep reinforcement learning (DRL) algorithms have shown remarkable performance in sport-related tasks. Zhao *et al.* [182] propose a hierarchical learning approach within a multi-agent reinforcement framework to emulate human performance in sports games. Jia *et al.* [183] address the challenges of asynchronous real-time scenarios in a basketball

sports environment, supporting both single-agent and multi-agent training.

The soccer virtual environment GFootball has gained significant attention in recent years [181]. In the 2020 Google Research Football Competition, the winning team, WeKick [184], developed a powerful agent using imitation learning and distributed league training. However, WeKick is specifically designed for single-agent AI and cannot be extended to multi-agent control. To address this limitation, Huang *et al.* [185] propose TiKick, an offline multi-agent algorithm that completes full games in GFootball using replay data generated by WeKick [185]. Another approach, Tizero [186], trains agents from scratch without pre-collected data and employs a self-improvement process to develop high-quality AI for multi-agent control [186].

Although DRL systems have made significant progress, they continue to encounter challenges in several areas, including multi-agent coordination, long-term planning, and non-transitivity [187]–[189]. These challenges highlight the complexity of developing AI systems that can effectively coordinate with multiple agents, make strategic decisions over extended periods, and account for non-transitive relationships in dynamic environments. Further research and advancements in these areas are crucial for enhancing the capabilities of DRL systems.

D. Player Motion Synthesizing

Utilizing video-based sequences to capture and analyze player movements represents a powerful approach to enhancing data diversity in sports. This innovative initiative has the potential to make a positive impact on the development of sports disciplines. Through detailed analysis and reproduction of player movements, we can gain valuable insights that have the potential to improve techniques, elevate athletic performance, and drive progress in the world of sports. This pioneering endeavor holds great promise for advancing the field and benefiting athletes and sports enthusiasts alike.

1) **Auto Choreographer:** Creating choreography involves the creative design of dance movements. However, automating the choreography process computationally is a challenging task. It requires generating continuous and complex motion that captures the intricate relationship with accompanying music.

Music-to-dance motion generation can be approached from both 2D and 3D perspectives. 2D approaches [190]–[192] rely on accurate 2D pose detectors [193] but have limitations in terms of expressiveness and downstream applications. On the other hand, 3D dance generation methods utilize techniques such as LSTMs [194]–[198], GANs [199], [200], transformer encoders with the RNN decoder [201] or transformer decoder [202], and convolutional sequence-to-sequence models [203], [204] to generate motion from audio.

Early works [191], [198], [204] in this field could predict future motion deterministically from audio but struggled when the same audio had multiple corresponding motions. However, recent advancements, such as the work by Li *et al.* [202], have addressed this limitation by formulating the problem with seed

motion. This enables the generation of multiple motions from the same audio, even with a deterministic model. Li *et al.* [202] propose a novel cross-modal transformer-based model that better preserves the correlation between music and 3D motion. This approach results in more realistic and globally translated long human motion.

E. Sport Video Synthesizing

The goal of artificially synthesizing sports videos is to generate realistic and immersive content, such as player movements or game scenarios. Early works in this field train models using annotated videos where each time step is labeled with the corresponding action. However, these approaches use a discrete representation of actions, which make it challenging to define prior knowledge for real-world environments. Additionally, devising a suitable continuous action representation for an environment is also complex. To address the complexity of action representation in tennis, Menapace *et al.* [205] propose a discrete action representation. Building upon this idea, Huang *et al.* [206] model actions as a learned set of geometric transformations. Davtyan *et al.* [207] take a different approach by separating actions into a global shift component and a local discrete action component. More recent works in tennis have utilized a NeRF-based renderer [208], which allows for the representation of complex 3D scenes. Among these works, Menapace *et al.* [209] employ a text-based action representation that provides precise details about the specific ball-hitting action being performed and the destination of the ball.

V. DATASETS AND BENCHMARKS

In the era of deep learning, having access to effective data is crucial for training and evaluating models. In order to facilitate this, we have compiled a list of commonly used public sports datasets, along with their corresponding details, as shown in Table II. Below, we provide a more detailed description of each dataset.

A. Soccer

In soccer, most video-based datasets benefit from active tasks like player tracking and action recognition, while some datasets focus on field localization and registration or player depth maps and meshes.

Some datasets focus more on player detection and tracking. Soccer-ISSIA [240] is an early work and a relatively small dataset with player bounding box annotations. SVPP [241] provides a multi-sensor dataset that includes body sensor data and video data. Soccer Player [242] is specifically designed for player detection and tracking, while SoccerTrack [214] is a novel dataset with multi-view and super high definition.

Other datasets like Football Action [137] and SoccerDB [211] benefit action recognition, and ComprehensiveSoccer [243] and SSET [210] can be used for various video analysis tasks, such as action classification, localization, and player detection. SoccerKicks [212] provides player pose estimation. GOAL [244] supports knowledge-grounded video captioning.

The SoccerNet series [33], [34], [141], [212] is the largest one including annotations for a variety of spatial annotations and cross-view correspondences. It covers multiple vision-based tasks including player understanding like player tracking, re-identification, broadcast video understanding like action spotting, video captioning, and field understanding like camera calibration.

In recent years, the combination of large-scale datasets and deep learning models has become increasingly popular in the field of soccer tasks, raising the popularity of the SoccerNet series datasets [34], [141], [212]. Meanwhile, SoccerDB [211], SSET [210], and ComprehensiveSoccer [243] are more suitable for tasks that require player detection. However, there are few datasets like SoccerKick [213] for soccer player pose estimation. It is hoped that more attention can be paid to the recognition and understanding of player skeletal movements in the future.

B. Basketball

Basketball datasets have been developed for various tasks such as player and ball detection, action recognition, and pose estimation. APIDIS [40], [245] is a challenging dataset with annotations for player and ball positions, and clock and non-clock actions. Basket-1,2 [38] consists of two frame sequences for action recognition and ball detection. NCAA [246] is a large dataset with action categories and bounding boxes for player detection. SPIROUDOME [215] focuses on player detection and localization. BPAD [154] is a first-person perspective dataset with labeled basketball events. SpaceJam [247] is for action recognition with estimated player poses. FineBasketball [248] is a fine-grained dataset with 3 broad and 26 fine-grained categories. NBA [126] is a dataset for group activity recognition, where each clip belongs to one of the nine group activities, and no individual annotations are provided, such as separate action labels and bounding boxes. NPUBasketball [216] contains RGB frames, depth maps, and skeleton information for various types of action recognition models. DeepSportradar-v1 [36] is a multi-label dataset for 3D localization, calibration, and instance segmentation tasks. In Captioning task, NSVA [217] is the largest open-source dataset in the basketball domain. Compared to SVN [249] and SVCDV [162], NSVA is publicly accessible and has the most sentences among the three datasets, with five times more videos than both SVN and SVCDV. Additionally, there are some special datasets that focus on reconstructing the player. NBA2K dataset [250] includes body meshes and texture data of several NBA players.

C. Volleyball

Despite being a popular sport, there are only a few volleyball datasets available, most of which are on small scales. Volleyball-1,2 [38] contains two sequences with manually annotated ball positions. HierVolleyball [251] and its extension HierVolleyball-v2 [252] are developed for team activity recognition, with annotated player actions and positions. Sports Video Captioning Dataset-Volleyball (SVCDV) [162] is a dataset for captioning tasks, with 55 videos from YouTube,

TABLE II

A LIST OF VIDEO-BASED SPORTS-RELATED DATASETS USED IN THE PUBLISHED PAPERS. NOTE THAT SOME OF THEM ARE NOT PUBLICLY AVAILABLE AND “MULTIPLE” MEANS THAT THE DATASET CONTAINS VARIOUS SPORTS INSTEAD OF ONLY ONE SPECIFIC TYPE OF SPORTS. “DET.”, “CLS.”, “TRA.”, “ASS.”, “SEG.”, “LOC.”, “CAL.”, “CAP.” STAND FOR PLAYER/BALL DETECTION, ACTION CLASSIFICATION, PLAYER/BALL TRACKING, ACTION QUALITY ASSESSMENT, OBJECT SEGMENTATION, TEMPORAL ACTION LOCALIZATION, CAMERA CALIBRATION, AND CAPTIONING RESPECTIVELY.

Sport	Dataset	Year	Task	# Videos	Avg. length	Link
Soccer	SoccerNet [141]	2018	loc.& cls.	500	5,400	✓
	SSET [210]	2020	det.&tra.	350	0.8h	✓
	SoccerDB [211]	2020	cls.& loc.	346	1.5h	✓
	SoccerNet-v2 [212]	2021	cls.&loc.	500	1.5h1.5h	✓
	SoccerKicks [213]	2021	pos.	38	-	✓
	SoccerNet-v3 [34]	2022	cls.&tra.	346	1.5h	✓
	SoccerNet-Tracking [33]	2022	cls.&tra.	21	45.5m	✓
	SoccerTrack [214]	2022	tra.&loc.	20	30s	✓
Basketball	BPAD [215]	2017	ass.	48	13m	✓
	NBA [126]	2020	cls.	181	-	✓
	NPUBasketball [216]	2021	cls.	2,169	-	✓
	DeepSportradar-v1 [36]	2022	seq.&cal.	-	-	✓
	NSVA [217]	2022	cls.&cap.	32,019	9.5s	✓
Tennis	PE-Tennis [218]	2022	det.&cal.	14,053	3s	✓
	LGEs-Tennis [209]	2023	cal.&tra.&cap.	7,112	7.8s	✓
Figure Skating	FisV-5 [219]	2020	ass.& cls.	500	2m50s	✓
	FR-FS [220]	2021	ass.& cls.	417	-	✓
Diving	MTL-AQA [221]	2019	ass.	1,412	-	✓
	FineDiving [222]	2022	ass.& cls.	3,000	52s	✓
Dance	GrooveNet [194]	2017	pos.	2	11.5m	✓
	Dance with Melody [195]	2018	pos.	61	92s	✓
	EA-MUD [200]	2020	pos.	17	74s	✓
	AIST++ [202]	2021	det&pos.	1,408	13s	✓
	DanceTrack [49]	2022	tra.	100	52.9s	✓
Golf	GolfDB [223]	2019	cls.	1,400	-	✓
Gymnastics	FineGym [114]	2020	cls.& loc.	-	-	✓
Rugby	Rugby sevens [119]	2022	tra.	346	40s	✓
Baseball	MLB-YouTube [224]	2018	cls.	5,111	-	✓
General	Sports 1M [225]	2014	cls.	1M	36s	✓
	OlympicSports [226]	2014	ass.	309	-	✓
	SVW [227]	2015	det.& cls.	4,100	11.6s	✓
	OlympicScoring [228]	2017	ass.	716	-	✓
	MADS [229]	2017	ass.	30	-	✓
	MultiTHUMOS [230]	2017	cls.	400	4.5m	✓
	AQA-7 [231]	2019	ass.	1,189	-	✓
	C-Sports [232]	2020	cls.&loc.	2,187	-	✓
	MultiSports [233]	2021	cls.&loc.	3,200	20.9s	✓
	ASPset-510 [234]	2021	pos.	510	-	✓
	HAA-500 [235]	2021	cls.	10,000	2.12s	✓
	SMART [236]	2021	cls.	5,000	-	✓
	Win-Fail [237]	2022	cls.	817	3.3s	✓
	SportsPose [238]	2023	pos.	25	11m	✓
	SportsMOT [239]	2023	tra.	240	25s	✓

each containing an average of 9.2 sentences. However, this dataset is not available for download.

D. Hockey

The Hockey Fight dataset [253] contains 1,000 video clips from National Hockey League (NHL) games for binary classification of fight and non-fight. The Player Tracklet dataset [254] consists of 84 video clips from NHL games with annotated bounding boxes and identity labels for players and referees and is suitable for player tracking and identification.

E. Tennis

Various datasets have been constructed for tennis video analysis. ACASVA [255] is designed for tennis action recognition and consists of six broadcast videos of tennis games

with labeled player positions and time boundaries of actions. THETIS [256] includes 1,980 self-recorded videos of 12 tennis actions with RGB, depth, 2D skeleton, and 3D skeleton videos, which can be used for multiple types of action recognition models. Tenniset [257] contains five Olympic tennis match videos with six labeled event categories and textural descriptions, making it suitable for both recognition, localization, and action retrieval tasks.

It should be noted that some recent works focus more on generative tasks, like PVG [258], which obtained a tennis dataset through YouTube videos. PE-Tennis [218] built upon PVG and introduces camera calibration resulting from reconstruction, making it possible to edit the viewpoint. LGEs-Tennis [209] enables generation from text editing on player movement, shot type, and location.

F. Table Tennis

Various datasets have been developed for table tennis stroke recognition, such as TTStroke-21 [259], which comprises 129 self-recorded videos of 21 categories, and SPIN [55], which includes 53 hours of self-recorded videos with annotations of ball position and player joints. OpenTTGames [37] consists of 12 HD videos of table tennis games, labeled with ball coordinates and events. Stroke Recognition [260] is similar to TTStroke-21, but much larger, and P²A [261] is one of the largest datasets for table tennis analysis, with annotations of each stroke in 2,721 broadcasting videos.

G. Gymnastics

The FineGym [114] is a recent work developed for gymnastic action recognition and localization. It contains 303 videos with around 708-hour length and is annotated hierarchically, making it suitable for fine-grained action recognition and localization. On the other hand, AFG-Olympics [262] provides challenging scenarios with extensive background, viewpoint, and scale variations over an extended sample duration of up to 2 minutes. Additionally, a discriminative attention module is proposed to embed long-range spatial and temporal correlation semantics.

H. Badminton

The Badminton Olympic [263] provides annotations for player detection, point localization, action recognition, and localization tasks. It comprises 10 YouTube videos of singles badminton matches, each approximately an hour long. The dataset includes annotations for player positions, temporal locations of point wins, and time boundaries and labels of strokes. Meanwhile, Stroke Forecasting [177] contains 43,191 trimmed video clips of badminton strokes categorized into 10 types, which can be used for both action recognition and stroke forecasting.

I. Figure skating

There are 5 datasets proposed for figure skating action recognition in recent years. FineSkating [264] is a hierarchical-labeled dataset of 46 videos of figure skating competitions for action recognition and action quality assessment. FSD-10 [265] comprises ten categories of figure skating actions and provides scores for action quality assessment. FisV-5 [107] is a dataset of 500 figure skating competition videos labeled with scores by 9 professional judges. FR-FS [109] is designed to recognize figure skating falls, with 417 videos containing the movements of take-off, rotation, and landing. MCFS [266] has three-level annotations of figure skating actions and their time boundaries, allowing for action recognition and localization.

J. Diving

There are three diving datasets available for action recognition and action quality assessment. Diving48 [267] contains 18,404 video segments covering 48 fine-grained categories of diving actions, making it a relatively low-bias dataset

suitable for model evaluation. In contrast, MTL-AQA [221] consists of 1,412 samples annotated with action quality scores, class labels, and textural commentary, making it suitable for multiple tasks. In addition, FineDiving [222] is a recent dataset consisting of 3,000 video samples covering 52 types of actions, 29 sub-action types, and 23 difficulty levels, providing fine-grained annotations including action types, sub-action types, coarse and fine time boundaries, and action scores. It is the first fine-grained motion video dataset for the AQA task, filling the gap in fine-grained annotations in AQA and suitable for designing competition strategies and better showcasing athletes' strengths.

K. Dance

The field of deep learning has several research tasks for dance, including music-oriented choreography, dance motion synthesis, and multiple object tracking. Researchers propose several datasets to promote research in this field. GrooveNet [194] consists of approximately 23 minutes of motion capture data recorded at 60 frames per second and four performances by a dancer. Dance with Melody [195] includes 40 complete dance choreographies for four types of dance, totaling 907,200 frames collected with optical motion capture equipment. EA-MUD [200] includes 104 video sequences of 12 dancing genres, while AIST++ [202] is a large-scale 3D human dance motion dataset with frame-level annotations including 9 views of camera intrinsic and extrinsic parameters, 17 COCO-format human joint locations in both 2D and 3D, and 24 SMPL pose parameters. These datasets can be used for tasks such as dance motion recognition, tracking, and quality assessment.

L. Sport Related Datasets for General Purpose

There are several datasets for sports action recognition and assessment tasks, including UCF sports [268], MSR Action3D [269], Olympic [270], Sports 1M [225], SVW [227], MultiSports [233], OlympicSports [226], OlympicScoring [228], and AQA [231]. These datasets cover different sports, including team sports and individual sports, and provide various annotations, such as action labels, quality scores, and bounding boxes.

Additionally, Win-Fail [237] is a dataset specifically designed for recognizing the outcome of actions, while SportPose [238] is the largest markerless dataset for 3D human pose estimation in sports, containing 5 short sports-related activities recorded from 7 cameras, totaling 1.5 million frames. SportsMOT [239] is a large-scale and high-quality multi-object tracking dataset comprising detailed annotations for each player present on the field in diverse sports scenarios. These datasets provide valuable resources for researchers to develop and evaluate algorithms for various sports-related tasks.

M. Others

CVBASE Handball [271] and CVBASE Squash [271] are developed for handball and squash action recognition,

respectively, with annotated trajectories of players and action categories. GolfDB [223] facilitates the analysis of golf swings, providing 1,400 high-quality golf swing video segments, action labels, and bounding boxes of players. Lastly, FenceNet [119] consists of 652 videos of expert-level fencers performing six categories of actions, with RGB frames, 3D skeleton data, and depth data provided. Rugby sevens [62] is a public sports tracking dataset with tracking ground truth and the generated tracks. MLB-YouTube [224] is introduced for fine-grained action recognition in baseball videos.

VI. VIRTUAL ENVIRONMENTS

Researchers can utilize virtual environments for simulation. In a virtual environment that provides agents with simulated motion tasks, multiple data information can be continuously generated and retained in the simulation. For example, Fever Basketball [183] is an asynchronous environment, which supports multiple characters, multiple positions, and both the single-agent and multi-agent player control modes.

There are many virtual soccer games, such as rSoccer [178], RoboCup Soccer Simulator [272], the DeepMind MuJoCo Multi-Agent Soccer Environment [179], [180] and JiDi Olympics Football [273]. rSoccer [178] and JiDi Olympics Football [273] are two toy football games in which plays are just rigid bodies and can just move and push the ball. However, players in GFootball [181] have more complex actions, such as dribbling, sliding, and sprinting. Besides, environments like RoboCup Soccer Simulator [272] and DeepMind MuJoCo Multi-Agent Soccer Environment [179], [180] focus more on low-level control of a physics simulation of robots, while GFootball focuses more on developing high-level tactics. To improve the flexibility and control over environment dynamics, SCENIC [274] is proposed to model and generate diverse scenarios in a real-time strategy environment programmatically.

VII. CHALLENGES

In recent years, deep learning has emerged as a powerful tool in the analysis and enhancement of sports performance. The application of these advanced techniques has revolutionized the way athletes, coaches, and teams approach training, strategy, and decision-making. By leveraging the vast amounts of data generated in sports, deep learning models have the potential to uncover hidden patterns, optimize performance, and provide valuable insights that can inform decision-making processes. However, despite its promising potential, the implementation of deep learning in sports performance faces several challenges that need to be addressed to fully realize its benefits.

a) Task Challenge: The complex and dynamic nature of sports activities presents unique challenges for computer vision tasks in tracking and recognizing athletes and their movements. Issues such as identity mismatch due to similar appearances [48], [49], blurring [52] caused by rapid motion, and occlusion [44], [45] from other players or objects in the scene can lead to inaccuracies and inconsistencies in tracking and analysis. Developing robust and adaptable algorithms that can effectively handle these challenges is essential to improve

the performance and reliability of deep learning models in sports applications.

b) Datasets Standardization: Standardizing datasets for various sports is a daunting task, as each sport has unique technical aspects and rules that make it difficult to create a unified benchmark for specific tasks. For example, taking action recognition tasks as an example, in diving [222], only the movement of the athlete needs to be focused on, and attention should be paid to the details of role actions. However, in team sports such as volleyball [251], more attention is needed to distinguish and identify targets and cluster the same actions after identification. Given the varying emphases of tasks, there are substantial differences in the dataset requirements. To go further, action recognition of the same sport type, involves nuanced differences in label classification, making it challenging to develop a one-size-fits-all solution or benchmark. The creation of standardized, user-friendly, open-source, high-quality, and large-scale datasets is crucial for advancing research and enabling fair comparisons between different models and approaches in sports performance analysis.

c) Data Utilization: The sports domain generates vast amounts of fine-grained data through sensors and IoT devices. However, current data processing methods primarily focus on computer vision and do not fully exploit the potential of end-to-end deep learning approaches. To fully harness the power of these rich data sources, researchers must develop methods that combine fine-grained sensor data with visual information. This fusion of diverse data streams can enable more comprehensive and insightful analysis, leading to significant advancements in the field of sports performance. Some studies have shown that introducing multi-modal data can benefit the analysis of athletic performance. For example, in table tennis, visual and IOT signals can be simultaneously used to analyze athlete performance [275]. In dance, visual and audio signals are both important [202]. More attention is needed on how to utilize diverse data, so as to achieve better fusion. Meanwhile, multi-modal algorithms and datasets [202] are both necessary.

VIII. FUTURE TREND

The integration of deep learning methodologies into sports analytics can empower athletes, coaches, and teams with unprecedented insights into performance, decision-making, and injury prevention. This future work aims to explore the transformative impact of deep learning techniques in sports performance, focusing on data generation methods, multi-modality and multi-task models, foundation models, applications, and practicability.

a) Multi-modality and Multi-task: By harnessing the power of multi-modal data and multi-task learning, robust and versatile models capable of handling diverse and complex sports-related challenges can be fulfilled. Furthermore, we will investigate the potential of large-scale models in enhancing predictive and analytical capabilities. It consists of practical applications and real-world implementations that can improve athlete performance and overall team dynamics. Ultimately, this work seeks to contribute to the growing body of research on deep learning in sports performance, paving the way for

novel strategies and technologies that can revolutionize the world of sports analytics.

b) Foundation Model: The popularity of ChatGPT has demonstrated the power of large language models [276], while the recent segment-anything project showcases the impressive performance of large models in visual tasks [277]. The prompt-based paradigm is highly capable and flexible in natural language processing and even image segmentation, offering unprecedented rich functionality. For example, some recent work has leveraged segment-anything in medical image [278]–[280], achieving promising results by providing point or bounding box prompts for preliminary zero-shot capability assessment, demonstrating that segment anything model (SAM) has good generalization performance in medical imaging. Therefore, the development of large models in the sports domain should consider how to combine existing large models to explore applications, and how to create large models specifically for the sports domain.

Combining large models requires considering the adaptability of the task. Compared to the medical field, sports involve a high level of human participation, inherently accommodating different levels and modalities of methods and data. We believe that both large language models in natural language processing and large image segmentation models in computer vision should have strong compatibility in sports. In short, we believe there is potential for exploring downstream tasks, such as using ChatGPT for performance evaluation and feedback: employ ChatGPT to generate natural language summaries of player or team performance, as well as provide personalized feedback and recommendations for improvement.

Foundation models directly related to the sports domain require a vast amount of data corresponding to the specific tasks. For visual tasks, for example, it is essential to ensure good scalability, adopt a prompt-based paradigm, and maintain powerful capabilities while being flexible and offering richer functionality. It is important to note that large models do not necessarily imply a large number of parameters, but rather a strong ability to solve tasks. Recent work on segment-anything has proven that even relatively simple models can achieve excellent performance when the data volume is sufficiently large. Therefore, creating large-scale, high-quality datasets in the sports domain remains a crucial task.

c) Data Generation: High-quality generated data can significantly reduce manual labor costs while demonstrating the diversity that generative models can bring. Many studies [202], [281] have focused on generating sports videos, offering easily editable, high-quality generation methods, which are elaborated upon in the relevant Section IV-D and IV-E. Meanwhile, by combining large models, additional annotation work can be performed at this stage, and if possible, new usable data can be generated.

d) Applications: Though there are many excellent automatic algorithms for different tasks in the field of sports, they are still insufficient when it comes to deployment for specific tasks. In the daily exercise of ordinary people, who generally lack professional guidance, there should be more applications that make good use of these deep learning algorithms, and use more user-friendly and intelligent methods to promote

sports for everyone. There are already some works [282]–[284] focusing on sports performance analysis, data recording visualization, energy expenditure estimation, and many other aspects. At the same time, in professional sports, there are also some works [16], [275] that focus on combining various data and methods to help improve athletic performance. Broadly speaking, in both daily life and professional fields, there is a need for more applications relating to health and fitness assessments.

e) Practicability: In more challenging, high-level tasks with real-world applications, practicality becomes increasingly important. Many practical challenges remain unexplored or under-explored in applying deep learning to sports performance. In decision-making, for example, current solutions often rely on simulation-based approaches. However, multi-agent decision-making techniques hold great potential for enhancing real-world sports decision-making. Tasks such as ad-hoc teamwork [285] in multi-agent systems and zero-shot human-machine interaction are crucial for enabling effective and practical real-world applications. Further research is needed to bridge the gap between theoretical advancements and their practical implications in sports performance analysis and decision-making. For example, RoboCup [272] aims to defeat human players in the World Cup by 2050. This complex task requires robots to perceive their environment, gather information, understand it, and execute specific actions. Such agents must exhibit sufficient generalization, engage in extensive human-machine interaction, and quickly respond to performance and environmental changes in real-time.

IX. CONCLUSION

In this paper, we present a comprehensive survey of deep learning in sports, focusing on four main aspects: algorithms, datasets, challenges, and future works. We innovatively summarize the taxonomy and divide methods into perception, comprehension, and decision from low-level to high-level tasks. In the challenges and future works, we provide cutting-edge methods and give insights into the future trends and challenges of deep learning in sports.

ACKNOWLEDGMENTS

This work is supported by National Key R&D Program of China under Grant No.2022ZD0162000, and National Natural Science Foundation of China No.62106219.

REFERENCES

- [1] N. Chmait and H. Westerbeek, “Artificial intelligence and machine learning in sport research: An introduction for non-data scientists,” *Frontiers in Sports and Active Living*, p. 363, 2021.
- [2] “Smt,” <https://www.smt.com/>.
- [3] “vizrt,” <https://www.vizrt.com/>.
- [4] A. Duarte, C. Micael, S. Ludovic, S. Hugo, and D. Keith, *Artificial Intelligence in Sport Performance Analysis*, 2021.
- [5] E. E. Cust, A. J. Sweeting, K. Ball, and S. Robertson, “Machine and deep learning for sport-specific movement recognition: a systematic review of model development and performance,” *Journal of sports sciences*, vol. 37, no. 5, pp. 568–600, 2019.
- [6] R. P. Bonidia, L. A. Rodrigues, A. P. Avila-Santos, D. S. Sanches, J. D. Brancher *et al.*, “Computational intelligence in sports: A systematic literature review,” *Advances in Human-Computer Interaction*, vol. 2018, 2018.

- [7] R. Beal, T. J. Norman, and S. D. Ramchurn, "Artificial intelligence for team sports: a survey," *The Knowledge Engineering Review*, vol. 34, p. e28, 2019.
- [8] D. Tan, H. Ting, and S. Lau, "A review on badminton motion analysis," in *2016 International Conference on Robotics, Automation and Sciences (ICORAS)*. IEEE, 2016, pp. 1–4.
- [9] F. Wu, Q. Wang, J. Bian, N. Ding, F. Lu, J. Cheng, D. Dou, and H. Xiong, "A survey on video action recognition in sports: Datasets, methods and applications," *IEEE Transactions on Multimedia*, 2022.
- [10] S. Wang, D. Yang, P. Zhai, Q. Yu, T. Suo, Z. Sun, K. Li, and L. Zhang, "A survey of video-based action quality assessment," in *2021 International Conference on Networking Systems of AI (INSAI)*. IEEE, 2021, pp. 1–9.
- [11] P. R. Kamble, A. G. Keskar, and K. M. Bhurchandi, "Ball tracking in sports: a survey," *Artificial Intelligence Review*, vol. 52, no. 3, pp. 1655–1705, 2019.
- [12] Y. Adesida, E. Papi, and A. H. McGregor, "Exploring the role of wearable technology in sport kinematics and kinetics: A systematic review," *Sensors*, vol. 19, no. 7, p. 1597, 2019.
- [13] M. Rana and V. Mittal, "Wearable sensors for real-time kinematics analysis in sports: a review," *IEEE Sensors Journal*, vol. 21, no. 2, pp. 1187–1207, 2020.
- [14] E. Van der Kruk and M. M. Reijne, "Accuracy of human motion capture systems for sport applications; state-of-the-art review," *European journal of sport science*, vol. 18, no. 6, pp. 806–819, 2018.
- [15] A. M. Turing, *Computing machinery and intelligence*. Springer, 2009.
- [16] J. Wang, K. Qiu, H. Peng, J. Fu, and J. Zhu, "Ai coach: Deep human pose estimation and analysis for personalized athletic training assistance," in *Proceedings of the 27th ACM international conference on multimedia*, 2019, pp. 374–382.
- [17] U. Rao and U. C. Pati, "A novel algorithm for detection of soccer ball and player," in *2015 International Conference on Communications and Signal Processing (ICCSPP)*. IEEE, 2015, pp. 0344–0348.
- [18] Y. Yang, M. Xu, W. Wu, R. Zhang, and Y. Peng, "3d multiview basketball players detection and localization based on probabilistic occupancy," in *2018 Digital Image Computing: Techniques and Applications (DICTA)*. IEEE, 2018, pp. 1–8.
- [19] M. Şah and C. Direkçöğlü, "Evaluation of image representations for player detection in field sports using convolutional neural networks," in *13th International Conference on Theory and Application of Fuzzy Systems and Soft Computing—ICAIFS-2018 13*. Springer, 2019, pp. 107–115.
- [20] S. Gerke, A. Linnemann, and K. Müller, "Soccer player recognition using spatial constellation features and jersey number recognition," *Computer Vision and Image Understanding*, vol. 159, pp. 105–115, 2017.
- [21] G. Li, S. Xu, X. Liu, L. Li, and C. Wang, "Jersey number recognition with semi-supervised spatial transformer network," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 1783–1790.
- [22] H. Liu and B. Bhanu, "Pose-guided r-cnn for jersey number recognition in sports," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [23] M. Istasse, J. Moreau, and C. De Vleeschouwer, "Associative embedding for team discrimination," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [24] M. Koshkina, H. Pidaparthi, and J. H. Elder, "Contrastive learning for sports video: Unsupervised player classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4528–4536.
- [25] M. Manafifard, H. Ebadi, and H. A. Moghaddam, "A survey on player tracking in soccer videos," *Computer Vision and Image Understanding*, vol. 159, pp. 19–46, 2017.
- [26] R. Theagarajan, F. Pala, X. Zhang, and B. Bhanu, "Soccer: Who has the ball? generating visual analytics and player statistics," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1749–1757.
- [27] A. Arbues-Sanguesa, A. Martín, J. Fernández, C. Ballester, and G. Haro, "Using player's body-orientation to model pass feasibility in soccer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 886–887.
- [28] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 28, 2015, pp. 91–99.
- [29] A. Cioppa, A. Deliege, M. Istasse, C. De Vleeschouwer, and M. Van Droogenbroeck, "Arthus: Adaptive real-time human segmentation in sports through online distillation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [30] R. Vandeghen, A. Cioppa, and M. Van Droogenbroeck, "Semi-supervised training to improve player and ball detection in soccer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3481–3490.
- [31] R. Sanford, S. Gorji, L. G. Hafemann, B. Pourbabaee, and M. Javan, "Group activity detection from trajectory and video data in soccer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 898–899.
- [32] A. Cioppa, A. Deliege, F. Magera, S. Giancola, O. Barnich, B. Ghanem, and M. Van Droogenbroeck, "Camera calibration and player localization in soccernet-v2 and investigation of their representations for action spotting," in *Proceedings of the IEEE/CVF Conference on CVPR*, 2021, pp. 4537–4546.
- [33] A. Cioppa, S. Giancola, A. Deliege, L. Kang, X. Zhou, Z. Cheng, B. Ghanem, and M. Van Droogenbroeck, "Socccernet-tracking: Multiple object tracking dataset and benchmark in soccer videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3491–3502.
- [34] A. Cioppa, A. Deliege, S. Giancola, B. Ghanem, and M. Van Droogenbroeck, "Scaling up soccernet with multi-view spatial localization and re-identification," *Scientific Data*, vol. 9, no. 1, p. 355, 2022.
- [35] I. Uchida, A. Scott, H. Shishido, and Y. Kameda, "Automated offside detection by spatio-temporal analysis of football videos," in *Proceedings of the 4th International Workshop on Multimedia Content Analysis in Sports*, 2021, pp. 17–24.
- [36] G. Van Zandycke, V. Somers, M. Istasse, C. D. Don, and D. Zambrano, "Deepsporthead-v1: Computer vision dataset for sports understanding with high quality annotations," in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 1–8.
- [37] R. Voeikov, N. Falaleev, and R. Baikulov, "Ttnet: Real-time temporal and spatial video analysis of table tennis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 884–885.
- [38] A. Maksai, X. Wang, and P. Fua, "What players do with the ball: A physically constrained interaction modeling," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 972–981.
- [39] X. Cheng, N. Ikoma, M. Honda, and T. Ikenaga, "Simultaneous physical and conceptual ball state estimation in volleyball game analysis," in *2017 IEEE Visual Communications and Image Processing (VCIP)*. IEEE, 2017, pp. 1–4.
- [40] P. Parisot and C. De Vleeschouwer, "Consensus-based trajectory estimation for ball detection in calibrated cameras systems," *Journal of Real-Time Image Processing*, vol. 16, no. 5, pp. 1335–1350, 2019.
- [41] J. Sköld, "Estimating 3d-trajectories from monocular video sequences," 2015.
- [42] G. Van Zandycke and C. De Vleeschouwer, "3d ball localization from a single calibrated image," in *Proceedings of the IEEE/CVF Conference on CVPR*, 2022, pp. 3472–3480.
- [43] P. Liu and J.-H. Wang, "Monotrack: Shuttle trajectory reconstruction from monocular badminton video," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3513–3522.
- [44] B. T. Naik, M. F. Hashmi, Z. W. Geem, and N. D. Bokde, "Deepplayer-track: player and referee tracking with jersey color recognition in soccer," *IEEE Access*, vol. 10, pp. 32 494–32 509, 2022.
- [45] B. T. Naik and M. F. Hashmi, "Yolov3-sort: detection and tracking player/ball in soccer sport," *Journal of Electronic Imaging*, vol. 32, no. 1, pp. 011 003–011 003, 2023.
- [46] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in *2016 IEEE international conference on image processing (ICIP)*. IEEE, 2016, pp. 3464–3468.
- [47] S. Hurault, C. Ballester, and G. Haro, "Self-supervised small soccer player detection and tracking," in *Proceedings of the 3rd international workshop on multimedia content analysis in sports*, 2020, pp. 9–18.
- [48] R. Zhang, L. Wu, Y. Yang, W. Wu, Y. Chen, and M. Xu, "Multi-camera multi-player tracking with deep player identification in sports video," *Pattern Recognition*, vol. 102, p. 107260, 2020.
- [49] P. Sun, J. Cao, Y. Jiang, Z. Yuan, S. Bai, K. Kitani, and P. Luo, "Dancetrack: Multi-object tracking in uniform appearance and diverse motion," *arXiv preprint arXiv:2111.14690*, 2021.

- [50] M. Buric, M. Ivacic-Kos, and M. Pobar, "Player tracking in sports videos," in *2019 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, 2019, pp. 334–340.
- [51] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," in *2017 IEEE international conference on image processing (ICIP)*. IEEE, 2017, pp. 3645–3649.
- [52] Y.-C. Huang, I.-N. Liao, C.-H. Chen, T.-U. Ik, and W.-C. Peng, "Tracknet: A deep learning network for tracking high-speed and tiny objects in sports applications," in *2019 16th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 2019, pp. 1–8.
- [53] V. Belagiannis and A. Zisserman, "Recurrent human pose estimation," in *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*. IEEE, 2017, pp. 468–475.
- [54] T. Pfister, J. Charles, and A. Zisserman, "Flowing convnets for human pose estimation in videos," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1913–1921.
- [55] S. Schwarcz, P. Xu, D. D'Ambrosio, J. Kangaspunta, A. Angelova, H. Phan, and N. Jaitly, "Spin: A high speed, high resolution vision dataset for tracking and action recognition in ping pong," *arXiv preprint arXiv:1912.06640*, 2019.
- [56] B. Comandur, "Sports re-id: Improving re-identification of players in broadcast videos of team sports," *arXiv preprint arXiv:2206.02373*, 2022.
- [57] A. Nady and E. E. Hemayed, "Player identification in different sports," in *VISIGRAPP*, 2021.
- [58] A. Senocak, T.-H. Oh, J. Kim, and I. S. Kweon, "Part-based player identification using deep convolutional representation and multi-scale pooling," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018, pp. 1813–18137.
- [59] O. M. Teket and I. S. Yetik, "A fast deep learning based approach for basketball video analysis," in *Proceedings of the 2020 4th International Conference on Vision, Image and Signal Processing*, ser. ICVISP 2020. New York, NY, USA: Association for Computing Machinery, 2020. [Online]. Available: <https://doi.org/10.1145/3448823.3448882>
- [60] Q. An, K. Cui, R. Liu, C. Wang, M. Qi, and H. Ma, "Attention-aware multiple granularities network for player re-identification," in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 137–144.
- [61] K. Habel, F. Deuser, and N. Oswald, "Clip-reident: Contrastive training for player re-identification," in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 129–135.
- [62] A. Maglo, A. Orcesi, and Q.-C. Pham, "Efficient tracking of team sport players with few game-specific annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3461–3471.
- [63] K. Vats, W. McNally, P. Walters, D. A. Clausi, and J. S. Zelek, "Ice hockey player identification via transformers," *arXiv preprint arXiv:2111.11535*, 2021.
- [64] B. Yan, Y. Li, X. Zhao, and H. Wang, "Dual data augmentation method for data-deficient and occluded instance segmentation," in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 117–120.
- [65] B. Yan, F. Qi, Z. Li, Y. Li, and H. Wang, "Strong instance segmentation pipeline for mmsports challenge," *arXiv preprint arXiv:2209.13899*, 2022.
- [66] G. Ghiasi, Y. Cui, A. Srinivas, R. Qian, T.-Y. Lin, E. D. Cubuk, Q. V. Le, and B. Zoph, "Simple copy-paste is a strong data augmentation method for instance segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 2918–2928.
- [67] C. Zhang, M. Wang, and L. Zhou, "Recognition method of basketball players' throwing action based on image segmentation," *International Journal of Biometrics*, vol. 15, no. 2, pp. 121–133, 2023.
- [68] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.
- [69] W. Chai, Z. Jiang, J.-N. Hwang, and G. Wang, "Global adaptation meets local generalization: Unsupervised domain adaptation for 3d human pose estimation," *arXiv preprint arXiv:2303.16456*, 2023.
- [70] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "Openpose: real-time multi-person 2d pose estimation using part affinity fields," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 1, pp. 172–186, 2021.
- [71] N. Promrit and S. Waijanya, "Model for practice badminton basic skills by using motion posture detection from video posture embedding and one-shot learning technique," in *Proceedings of the 2019 2nd artificial intelligence and cloud computing conference*, 2019, pp. 117–124.
- [72] S. Suda, Y. Makino, and H. Shinoda, "Prediction of volleyball trajectory using skeletal motions of setter player," in *Proceedings of the 10th Augmented Human International Conference 2019*, 2019, pp. 1–8.
- [73] T. Shimizu, R. Hachiuma, H. Saito, T. Yoshikawa, and C. Lee, "Prediction of future shot direction using pose and position of tennis player," in *Proceedings Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports*, 2019, pp. 59–66.
- [74] E. Wu and H. Koike, "Futurepong: Real-time table tennis trajectory forecasting using pose prediction network," in *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, 2020, pp. 1–8.
- [75] H. Sak, A. W. Senior, and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling," 2014.
- [76] M. Einfalt, C. Dampayrou, D. Zecha, and R. Lienhart, "Frame-level event detection in athletics videos with pose-based convolutional sequence networks," in *Proceedings Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports*, 2019, pp. 42–50.
- [77] H. Thilakarathne, A. Nibali, Z. He, and S. Morgan, "Pose is all you need: The pose only group activity recognition system (pogars)," *Machine Vision and Applications*, vol. 33, no. 6, p. 95, 2022.
- [78] A. Tharatipyakul, K. T. Choo, and S. T. Perrault, "Pose estimation for facilitating movement learning from online videos," in *Proceedings of the International Conference on Advanced Visual Interfaces*, 2020, pp. 1–5.
- [79] E. W. Trejo and P. Yuan, "Recognition of yoga poses through an interactive system with kinect based on confidence value," in *2018 3rd international conference on advanced robotics and mechatronics (ICARM)*. IEEE, 2018, pp. 606–611.
- [80] D. Farin, S. Krabbe, W. Effelsberg *et al.*, "Robust camera calibration for sport videos using court models," in *Storage and Retrieval Methods and Applications for Multimedia 2004*, vol. 5307. SPIE, 2003, pp. 80–91.
- [81] Q. Yao, A. Kubota, K. Kawakita, K. Nonaka, H. Sankoh, and S. Naito, "Fast camera self-calibration for synthesizing free viewpoint soccer video," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 1612–1616.
- [82] N. Homayounfar, S. Fidler, and R. Urtasun, "Sports field localization via deep structured models," in *Proceedings of the IEEE Conference on CVPR*, 2017, pp. 5212–5220.
- [83] J. Chen and J. J. Little, "Sports camera calibration via synthetic data," in *Proceedings of the IEEE/CVF conference on CVPR workshops*, 2019, pp. 0–0.
- [84] L. Sha, J. Hobbs, P. Felsen, X. Wei, P. Lucey, and S. Ganguly, "End-to-end camera calibration for broadcast videos," in *Proceedings of the IEEE/CVF conference on CVPR*, 2020, pp. 13627–13636.
- [85] X. Nie, S. Chen, and R. Hamid, "A robust and efficient framework for sports-field registration," in *Winter Conference on Applications of Computer Vision, WACV*. IEEE, 2021, pp. 1935–1943. [Online]. Available: <https://doi.org/10.1109/WACV48630.2021.00198>
- [86] F. Shi, P. Marchwica, J. C. G. Higuera, M. Jamieson, M. Javan, and P. Siva, "Self-supervised shape alignment for sports field registration," in *Winter Conference on Applications of Computer Vision, WACV*. IEEE, 2022, pp. 3768–3777. [Online]. Available: <https://doi.org/10.1109/WACV51458.2022.00382>
- [87] Y.-J. Chu, J.-W. Su, K.-W. Hsiao, C.-Y. Lien, S.-H. Fan, M.-C. Hu, R.-R. Lee, C.-Y. Yao, and H.-K. Chu, "Sports field registration via keypoints-aware label condition," in *Conference on Computer Vision and Pattern Recognition Workshops, CVPRW*. IEEE/CVF, 2022, pp. 3523–3530. [Online]. Available: <https://doi.org/10.1109/CVPRW56347.2022.00396>
- [88] N. Zhang and E. Izquierdo, "A high accuracy camera calibration method for sport videos," in *International Conference on Visual Communications and Image Processing, VCIP*. IEEE, 2021, pp. 1–5. [Online]. Available: <https://doi.org/10.1109/VCIP53242.2021.9675379>
- [89] J. Lin, C. Gan, and S. Han, "Tsm: Temporal shift module for efficient video understanding," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7083–7093.
- [90] D. Tran, H. Wang, L. Torresani, and M. Feiszli, "Video classification with channel-separated convolutional networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5552–5561.
- [91] C. Feichtenhofer, H. Fan, J. Malik, and K. He, "Slowfast networks for video recognition," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6202–6211.

- [92] W. Wang, D. Tran, and M. Feiszli, "What makes training multi-modal classification networks hard?" in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 695–12 705.
- [93] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Skeleton-based action recognition with multi-stream adaptive graph convolutional networks," *IEEE Transactions on Image Processing*, vol. 29, pp. 9532–9545, 2020.
- [94] Y.-F. Song, Z. Zhang, C. Shan, and L. Wang, "Stronger, faster and more explainable: A graph convolutional baseline for skeleton-based action recognition," in *proceedings of the 28th ACM international conference on multimedia*, 2020, pp. 1625–1633.
- [95] D. Kondratyuk, L. Yuan, Y. Li, L. Zhang, M. Tan, M. Brown, and B. Gong, "Movinets: Mobile video networks for efficient video recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 020–16 030.
- [96] G. Bertasius, H. Wang, and L. Torresani, "Is space-time attention all you need for video understanding," *arXiv preprint arXiv:2102.05095*, vol. 2, no. 3, p. 4, 2021.
- [97] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video swin transformer," *arXiv preprint arXiv:2106.13230*, 2021.
- [98] R. Herzig, E. Ben-Avraham, K. Mangalam, A. Bar, G. Chechik, A. Rohrbach, T. Darrell, and A. Globerson, "Object-region video transformers," *arXiv preprint arXiv:2110.06915*, 2021.
- [99] R. Wang, D. Chen, Z. Wu, Y. Chen, X. Dai, M. Liu, Y.-G. Jiang, L. Zhou, and L. Yuan, "Bevt: Bert pretraining of video transformers," *arXiv preprint arXiv:2112.01529*, 2021.
- [100] H. Tan, J. Lei, T. Wolf, and M. Bansal, "Vimpac: Video pre-training via masked token prediction and contrastive learning," *arXiv preprint arXiv:2106.11250*, 2021.
- [101] Y. Chen, Z. Zhang, C. Yuan, B. Li, Y. Deng, and W. Hu, "Channel-wise topology refinement graph convolution for skeleton-based action recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13 359–13 368.
- [102] H. Yuan, D. Ni, and M. Wang, "Spatio-temporal dynamic inference network for group activity recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 7476–7485.
- [103] H. Duan, Y. Zhao, K. Chen, D. Shao, D. Lin, and B. Dai, "Revisiting skeleton-based action recognition," *arXiv preprint arXiv:2104.13586*, 2021.
- [104] X. Xiang, Y. Tian, A. Reiter, G. D. Hager, and T. D. Tran, "S3d: Stacking segmental p3d for action quality assessment," in *ICIP*, 2018, pp. 928–932.
- [105] P. Parmar and B. T. Morris, "Action quality assessment across multiple actions," in *WACV*, 2018.
- [106] —, "What and how well you performed? a multitask learning approach to action quality assessment," in *CVPR*, 2019.
- [107] C. Xu, Y. Fu, B. Zhang, Z. Chen, and X. Xue, "Learning to score figure skating sport videos," *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, vol. PP, no. 99, pp. 1–1, 2019.
- [108] Y. Tang, Z. Ni, J. Zhou, D. Zhang, and J. Zhou, "Uncertainty-aware score distribution learning for action quality assessment," in *CVPR*, 2020.
- [109] S. Wang, Y. D., Z. P., C. C., and Z. L., "Tsa-net: Tube self-attention network for action quality assessment," in *ACM MM*, 2021.
- [110] Z. Qi, R. Zhu, Z. Fu, W. Chai, and V. Kindratenko, "Weakly supervised two-stage training scheme for deep video fight detection model," *arXiv preprint arXiv:2209.11477*, 2022.
- [111] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [112] M. Monfort, A. Andonian, B. Zhou, K. Ramakrishnan, S. A. Bargal, T. Yan, L. Brown, Q. Fan, D. Gutfreund, C. Vondrick *et al.*, "Moments in time dataset: one million videos for event understanding," *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 2, pp. 502–508, 2019.
- [113] G. Wang, K. Lu, Y. Zhou, Z. He, and G. Wang, "Human-centered prior-guided and task-dependent multi-task representation learning for action recognition pre-training," in *2022 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2022, pp. 1–6.
- [114] D. Shao, Y. Zhao, B. Dai, and D. Lin, "Finegym: A hierarchical video dataset for fine-grained action understanding," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2616–2625.
- [115] S. Sun, F. Wang, Q. Liang, and L. He, "Taichi: A fine-grained action recognition dataset," in *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, 2017, pp. 429–433.
- [116] J. Choi, C. Gao, J. C. Messou, and J.-B. Huang, "Why can't i dance in the mall? learning to mitigate scene bias in action recognition," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [117] P. Weinzaepfel and G. Rogez, "Mimetics: Towards understanding human actions out of context," *International Journal of Computer Vision*, vol. 129, no. 5, pp. 1675–1690, 2021.
- [118] Z. Liu, H. Zhang, Z. Chen, Z. Wang, and W. Ouyang, "Disentangling and unifying graph convolutions for skeleton-based action recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 143–152.
- [119] K. Zhu, A. Wong, and J. McPhee, "Fencenet: Fine-grained footwork recognition in fencing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3589–3598.
- [120] J. Hong, M. Fisher, M. Gharbi, and K. Fatahalian, "Video pose distillation for few-shot, fine-grained sports action recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9254–9263.
- [121] Y. Ben-Shabat, X. Yu, F. Saleh, D. Campbell, C. Rodriguez-Opazo, H. Li, and S. Gould, "The ikea asm dataset: Understanding people assembling furniture through actions, objects and pose," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 847–859.
- [122] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price *et al.*, "Scaling egocentric vision: The epic-kitchens dataset," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 720–736.
- [123] R. Goyal, S. Ebrahimi Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag *et al.*, "The" something something" video database for learning and evaluating visual common sense," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5842–5850.
- [124] K. Gavriluyk, R. Sanford, M. Javan, and C. G. Snoek, "Actor-transformers for group activity recognition," in *CVPR*, 2020, pp. 839–848.
- [125] G. Hu, B. Cui, Y. He, and S. Yu, "Progressive relation learning for group activity recognition," in *CVPR*, 2020, pp. 980–989.
- [126] R. Yan, L. Xie, J. Tang, X. Shu, and Q. Tian, "Social adaptive module for weakly-supervised group activity recognition," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*. Springer, 2020, pp. 208–224.
- [127] M. Ehsanpour, A. Abedin, F. Saleh, J. Shi, I. Reid, and H. Rezatofighi, "Joint learning of social groups, individuals action and sub-group activities in videos," in *ECCV*. Springer, 2020, pp. 177–195.
- [128] R. R. A. Pramono, Y. T. Chen, and W. H. Fang, "Empowering relational network by self-attention augmented conditional random fields for group activity recognition," in *ECCV*. Springer, 2020, pp. 71–90.
- [129] S. Li, Q. Cao, L. Liu, K. Yang, S. Liu, J. Hou, and S. Yi, "Groupformer: Group activity recognition with clustered spatial-temporal transformer," *ICCV*, 2021.
- [130] J. Wu, L. Wang, L. Wang, J. Guo, and G. Wu, "Learning actor relation graphs for group activity recognition," in *CVPR*, 2019, pp. 9964–9974.
- [131] M. Han, D. J. Zhang, Y. Wang, R. Yan, L. Yao, X. Chang, and Y. Qiao, "Dual-ai: dual-path actor interaction learning for group activity recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 2990–2999.
- [132] G. Xu and J. Yin, "Mlp-air: An efficient mlp-based method for actor interaction relation learning in group activity recognition," *arXiv preprint arXiv:2304.08803*, 2023.
- [133] V. Bettadapura, C. Pantofaru, and I. Essa, "Leveraging contextual cues for generating basketball highlights," in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 908–917.
- [134] F. C. Heilbron, W. Barrios, V. Escorcia, and B. Ghanem, "Sec: Semantic context cascade for efficient action detection," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 3175–3184.
- [135] P. Felsen, P. Agrawal, and J. Malik, "What will happen next? forecasting player moves in sports videos," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 3342–3351.
- [136] A. Cioppa, A. Deliege, and M. Van Droogenbroeck, "A bottom-up approach based on semantics for the interpretation of the main camera stream in soccer games," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1765–1774.
- [137] T. Tsunoda, Y. Komori, M. Matsugu, and T. Harada, "Football action recognition using hierarchical lstm," in *Proceedings of the IEEE*

- conference on computer vision and pattern recognition workshops, 2017, pp. 99–107.
- [138] Z. Cai, H. Neher, K. Vats, D. A. Clausi, and J. Zelek, “Temporal hockey action recognition via pose and optical flows,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [139] M. Sanabria, F. Precioso, and T. Menguy, “A deep architecture for multimodal summarization of soccer games,” in *Proceedings Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports*, 2019, pp. 16–24.
- [140] F. Turchini, L. Seidenari, L. Galteri, A. Ferracani, G. Becchi, and A. Del Bimbo, “Flexible automatic football filming and summarization,” in *Proceedings Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports*, 2019, pp. 108–114.
- [141] S. Giancola, M. Amine, T. Dghaily, and B. Ghanem, “Socccernet: A scalable dataset for action spotting in soccer videos,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 1711–1721.
- [142] A. Cioppa, A. Deliege, S. Giancola, B. Ghanem, M. V. Droogenbroeck, R. Gade, and T. B. Moeslund, “A context-aware loss function for action spotting in soccer videos,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 126–13 136.
- [143] J. Hong, H. Zhang, M. Gharbi, M. Fisher, and K. Fatahalian, “Spotting temporally precise, fine-grained events in video,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*. Springer, 2022, pp. 33–51.
- [144] A. Darwish and T. El-Shabrawy, “Ste: Spatio-temporal encoder for action spotting in soccer videos,” in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 87–92.
- [145] A. Cartas, C. Ballester, and G. Haro, “A graph-based method for soccer action spotting using unsupervised player classification,” in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 93–102.
- [146] H. Zhu, J. Liang, C. Lin, J. Zhang, and J. Hu, “A transformer-based system for action spotting in soccer videos,” in *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*, 2022, pp. 103–109.
- [147] J. V. Soares and A. Shah, “Action spotting using dense detection anchors revisited: Submission to the socccernet challenge 2022,” *arXiv preprint arXiv:2206.07846*, 2022.
- [148] J. V. Soares, A. Shah, and T. Biswas, “Temporally precise action spotting in soccer videos using dense detection anchors,” in *2022 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2022, pp. 2796–2800.
- [149] J. H. Pan, J. Gao, and W. S. Zheng, “Action assessment by joint relation graphs,” in *ICCV*, 2019.
- [150] G. I. Parisi, S. Magg, and S. Wermter, “Human motion assessment in real time using recurrent self-organization,” in *IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, 2016.
- [151] S. T. Kim and M. R. Yong, “Evaluationnet: Can human skill be evaluated by deep networks?” *arXiv:1705.11077*, 2017.
- [152] X. Yu, Y. Rao, W. Zhao, J. Lu, and J. Zhou, “Group-aware contrastive regression for action quality assessment,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 7919–7928.
- [153] Y. Li, X. Chai, and X. Chen, “End-to-end learning for action quality assessment,” in *Advances in Multimedia Information Processing – PCM*, 2018.
- [154] G. Bertasius, H. S. Park, S. X. Yu, and J. Shi, “Am i a baller? basketball performance assessment from first-person videos,” in *ICCV*, 2019.
- [155] R. Agyeman, R. Muhammad, and G. S. Choi, “Soccer video summarization using deep learning,” in *2019 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*, 2019, pp. 270–273.
- [156] M. Rafiq, G. Rafiq, R. Agyeman, G. S. Choi, and S.-I. Jin, “Scene classification for sports video summarization using transfer learning,” *Sensors*, vol. 20, no. 6, p. 1702, Mar 2020. [Online]. Available: <http://dx.doi.org/10.3390/s20061702>
- [157] A. A. Khan, J. Shao, W. Ali, and S. Tumrani, “Content-aware summarization of broadcast sports videos: An audio–visual feature extraction approach,” *Neural Processing Letters*, pp. 1–24, 2020.
- [158] H. Shingrakhia and H. Patel, “Sgrnn-am and hrf-dbn: A hybrid machine learning model for cricket video summarization,” *Vis. Comput.*, vol. 38, no. 7, p. 2285–2301, jul 2022. [Online]. Available: <https://doi.org/10.1007/s00371-021-02111-8>
- [159] W. Li, G. Pan, C. Wang, Z. Xing, and Z. Han, “From coarse to fine: Hierarchical structure-aware video summarization,” *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 18, no. 1s, jan 2022. [Online]. Available: <https://doi.org/10.1145/3485472>
- [160] W. Chai and G. Wang, “Deep vision multimodal learning: Methodology, benchmark, and trend,” *Applied Sciences*, vol. 12, no. 13, p. 6588, 2022.
- [161] H. Yu, S. Cheng, B. Ni, M. Wang, J. Zhang, and X. Yang, “Fine-grained video captioning for sports narrative,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6006–6015.
- [162] M. Qi, Y. Wang, A. Li, and J. Luo, “Sports video captioning via attentive motion representation and group relationship modeling,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 8, pp. 2617–2633, 2019.
- [163] —, “Sports video captioning via attentive motion representation and group relationship modeling,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 8, pp. 2617–2633, 2020.
- [164] H. Yu, S. Cheng, B. Ni, M. Wang, J. Zhang, and X. Yang, “Fine-grained video captioning for sports narrative,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6006–6015.
- [165] J. Wang, I. Fox, J. Skaza, N. Linck, S. Singh, and J. Wiens, “The advantage of doubling: a deep reinforcement learning approach to studying the double team in the nba,” *arXiv preprint arXiv:1803.02940*, 2018.
- [166] Y. Luo, “Inverse reinforcement learning for team sports: Valuing actions and players,” 2020.
- [167] G. Liu and O. Schulte, “Deep reinforcement learning in ice hockey for context-aware player evaluation,” *arXiv preprint arXiv:1805.11088*, 2018.
- [168] C. Yanai, A. Solomon, G. Katz, B. Shapira, and L. Rokach, “Q-ball: Modeling basketball games using deep reinforcement learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 8, 2022, pp. 8806–8813.
- [169] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *arXiv preprint arXiv:1509.02971*, 2015.
- [170] “statsperform-optical-tracking,” <https://www.statsperform.com/team-performance/football/optical-tracking>.
- [171] “secondspectrum,” <https://www.secondspectrum.com>.
- [172] X. Wei, P. Lucey, S. Morgan, and S. Sridharan, “Forecasting the next shot location in tennis using fine-grained spatiotemporal tracking data,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 11, pp. 2988–2997, 2016.
- [173] T. Fernando, S. Denman, S. Sridharan, and C. Fookes, “Memory augmented deep generative models for forecasting the next shot location in tennis,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 9, pp. 1785–1797, 2019.
- [174] X. Wei, P. Lucey, S. Morgan, M. Reid, and S. Sridharan, “The thin edge of the wedge: Accurately predicting shot outcomes in tennis using style and context priors,” in *Proceedings of the 10th Annu MIT Sloan Sport Anal Conf, Boston, MA, USA*, 2016, pp. 1–11.
- [175] H. M. Le, P. Carr, Y. Yue, and P. Lucey, “Data-driven ghosting using deep imitation learning,” 2017.
- [176] P. Power, H. Ruiz, X. Wei, and P. Lucey, “Not all passes are created equal: Objectively measuring the risk and reward of passes in soccer from tracking data,” in *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, 2017, pp. 1605–1613.
- [177] W.-Y. Wang, H.-H. Shuai, K.-S. Chang, and W.-C. Peng, “ShuttleNet: Position-aware fusion of rally progress and player styles for stroke forecasting in badminton,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- [178] F. B. Martins, M. G. Machado, H. F. Bassani, P. H. M. Braga, and E. S. Barros, “rsoccer: A framework for studying reinforcement learning in small and very small size robot soccer,” 2021.
- [179] S. Liu, G. Lever, J. Merel, S. Tunyasuvunakool, N. Heess, and T. Graepel, “Emergent coordination through competition,” *arXiv preprint arXiv:1902.07151*, 2019.
- [180] S. Liu, G. Lever, Z. Wang, J. Merel, S. Eslami, D. Hennes, W. M. Czarnecki, Y. Tassa, S. Omidshafiei, A. Abdolmaleki *et al.*, “From motor control to team play in simulated humanoid football,” *arXiv preprint arXiv:2105.12196*, 2021.
- [181] K. Kurach, A. Raichuk, P. Stańczyk, M. Zajac, O. Bachem, L. Espeholt, C. Riquelme, D. Vincent, M. Michalski, O. Bousquet *et al.*, “Google research football: A novel reinforcement learning environment,” in

- Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 4501–4510.
- [182] Y. Zhao, I. Borovikov, J. Rupert, C. Somers, and A. Beirami, “On multi-agent learning in team sports games,” *arXiv preprint arXiv:1906.10124*, 2019.
- [183] H. Jia, Y. Hu, Y. Chen, C. Ren, T. Lv, C. Fan, and C. Zhang, “Fever basketball: A complex, flexible, and asynchronous sports game environment for multi-agent reinforcement learning,” *arXiv preprint arXiv:2012.03204*, 2020.
- [184] F. Z. Ziyang Li, Kaiwen Zhu, “Wekick,” <https://www.kaggle.com/c/google-football/discussion/202232>, 2020.
- [185] S. Huang, W. Chen, L. Zhang, Z. Li, F. Zhu, D. Ye, T. Chen, and J. Zhu, “Tikick: Towards playing multi-agent football full games from single-agent demonstrations,” *arXiv preprint arXiv:2110.04507*, 2021.
- [186] F. Lin, S. Huang, T. Pearce, W. Chen, and W.-W. Tu, “Tizero: Mastering multi-agent football with curriculum learning and self-play,” *arXiv preprint arXiv:2302.07515*, 2023.
- [187] C. Yu, A. Velu, E. Vinitzky, Y. Wang, A. Bayen, and Y. Wu, “The surprising effectiveness of mappo in cooperative, multi-agent games,” *arXiv preprint arXiv:2103.01955*, 2021.
- [188] M. Wen, J. G. Kuba, R. Lin, W. Zhang, Y. Wen, J. Wang, and Y. Yang, “Multi-agent reinforcement learning is a sequence modeling problem,” *arXiv preprint arXiv:2205.14953*, 2022.
- [189] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune, “First return, then explore,” *Nature*, vol. 590, no. 7847, pp. 580–586, 2021.
- [190] P. Tendulkar, A. Das, A. Kembhavi, and D. Parikh, “Feel the music: Automatically generating a dance for an input song,” *arXiv preprint arXiv:2006.11905*, 2020.
- [191] X. Ren, H. Li, Z. Huang, and Q. Chen, “Self-supervised dance video synthesis conditioned on music,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 46–54.
- [192] J. P. Ferreira, T. M. Coutinho, T. L. Gomes, J. F. Neto, R. Azevedo, R. Martins, and E. R. Nascimento, “Learning to dance: A graph convolutional adversarial network to generate realistic dance motions from audio,” *Computers & Graphics*, vol. 94, pp. 11–21, 2021.
- [193] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, “Realtime multi-person 2d pose estimation using part affinity fields,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7291–7299.
- [194] O. Alemi, J. Françoise, and P. Pasquier, “Groovenet: Real-time music-driven dance movement generation using artificial neural networks,” *networks*, vol. 8, no. 17, p. 26, 2017.
- [195] T. Tang, J. Jia, and H. Mao, “Dance with melody: An lstm-autoencoder approach to music-oriented dance synthesis,” in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1598–1606.
- [196] N. Yalta, S. Watanabe, K. Nakadai, and T. Ogata, “Weakly-supervised deep recurrent neural networks for basic dance step generation,” in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.
- [197] W. Zhuang, Y. Wang, J. Robinson, C. Wang, M. Shao, Y. Fu, and S. Xia, “Towards 3d dance motion synthesis and control,” *arXiv preprint arXiv:2006.05743*, 2020.
- [198] H.-K. Kao and L. Su, “Temporally guided music-to-body-movement generation,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 147–155.
- [199] H.-Y. Lee, X. Yang, M.-Y. Liu, T.-C. Wang, Y.-D. Lu, M.-H. Yang, and J. Kautz, “Dancing to music,” *Advances in neural information processing systems*, vol. 32, 2019.
- [200] G. Sun, Y. Wong, Z. Cheng, M. S. Kankanhalli, W. Geng, and X. Li, “Deepdance: music-to-dance motion choreography with adversarial learning,” *IEEE Transactions on Multimedia*, vol. 23, pp. 497–509, 2020.
- [201] R. Huang, H. Hu, W. Wu, K. Sawada, M. Zhang, and D. Jiang, “Dance revolution: Long-term dance generation with music via curriculum learning,” *arXiv preprint arXiv:2006.06119*, 2020.
- [202] R. Li, S. Yang, D. A. Ross, and A. Kanazawa, “Ai choreographer: Music conditioned 3d dance generation with aist++,” 2021.
- [203] H. Ahn, J. Kim, K. Kim, and S. Oh, “Generative autoregressive networks for 3d dancing move synthesis from music,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3501–3508, 2020.
- [204] Z. Ye, H. Wu, J. Jia, Y. Bu, W. Chen, F. Meng, and Y. Wang, “Choreonet: Towards music to dance synthesis with choreographic action unit,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 744–752.
- [205] W. Menapace, S. Lathuiliere, S. Tulyakov, A. Siarohin, and E. Ricci, “Playable video generation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 10061–10070.
- [206] J. Huang, Y. Jin, K. M. Yi, and L. Sigal, “Layered controllable video generation,” in *Proceedings of the European Conference of Computer Vision (ECCV)*, S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds., 2022.
- [207] A. Davtyan and P. Favaro, “Controllable video generation through global and local motion dynamics,” in *Proceedings of the European Conference of Computer Vision (ECCV)*, 2022.
- [208] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *Proceedings of the European Conference of Computer Vision (ECCV)*, 2020.
- [209] W. Menapace, A. Siarohin, S. Lathuiliere, P. Achlioptas, V. Golyanik, E. Ricci, and S. Tulyakov, “Plotting behind the scenes: Towards learnable game engines,” *arXiv preprint arXiv:2303.13472*, 2023.
- [210] N. Feng, Z. Song, J. Yu, Y.-P. P. Chen, Y. Zhao, Y. He, and T. Guan, “Sset: a dataset for shot segmentation, event detection, player tracking in soccer videos,” *Multimedia Tools and Applications*, vol. 79, pp. 28 971–28 992, 2020.
- [211] Y. Jiang, K. Cui, L. Chen, C. Wang, and C. Xu, “Soccerdb: A large-scale database for comprehensive video understanding,” in *Proceedings of the 3rd International Workshop on Multimedia Content Analysis in Sports*, 2020, pp. 1–8.
- [212] A. Deliege, A. Cioppa, S. Giancola, M. J. Seikavandi, J. V. Dueholm, K. Nasrollahi, B. Ghanem, T. B. Moeslund, and M. Van Droogenbroeck, “Socccnet-v2: A dataset and benchmarks for holistic understanding of broadcast soccer videos,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4508–4519.
- [213] N. M. Lessa, E. L. Colombini, and A. D. S. Simões, “Soccerkicks: a dataset of 3d dead ball kicks reference movements for humanoid robots,” in *2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE, 2021, pp. 3472–3478.
- [214] A. Scott, I. Uchida, M. Onishi, Y. Kameda, K. Fukui, and K. Fujii, “Soccertrack: A dataset and tracking algorithm for soccer with fish-eye and drone videos,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3569–3579.
- [215] P. Parisot and C. De Vleeschouwer, “Scene-specific classifier for effective and efficient team sport players detection from a single calibrated camera,” *Computer Vision and Image Understanding*, vol. 159, pp. 74–88, 2017.
- [216] C. Ma, J. Fan, J. Yao, and T. Zhang, “Npu rgb+d dataset and a feature-enhanced lstm-dgcn method for action recognition of basketball players,” *Applied Sciences*, vol. 11, no. 10, p. 4426, 2021.
- [217] D. Wu, H. Zhao, X. Bao, and R. P. Wildes, “Sports video analysis on large-scale data,” in *ECCV*, Oct. 2022.
- [218] W. Menapace, S. Lathuiliere, A. Siarohin, C. Theobalt, S. Tulyakov, V. Golyanik, and E. Ricci, “Playable environments: Video manipulation in space and time,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3584–3593.
- [219] C. Xu, Y. Fu, B. Zhang, Z. Chen, Y.-G. Jiang, and X. Xue, “Learning to score figure skating sport videos,” *IEEE transactions on circuits and systems for video technology*, vol. 30, no. 12, pp. 4578–4590, 2019.
- [220] S. Wang, D. Yang, P. Zhai, C. Chen, and L. Zhang, “Tsa-net: Tube self-attention network for action quality assessment,” in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 4902–4910.
- [221] P. Parmar and B. T. Morris, “What and how well you performed? a multitask learning approach to action quality assessment,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 304–313.
- [222] J. Xu, Y. Rao, X. Yu, G. Chen, J. Zhou, and J. Lu, “Finediving: A fine-grained dataset for procedure-aware action quality assessment,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2949–2958.
- [223] W. McNally, K. Vats, T. Pinto, C. Dulhanty, J. McPhee, and A. Wong, “GolfdB: A video database for golf swing sequencing,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [224] A. Piergiovanni and M. S. Ryoo, “Fine-grained activity recognition in baseball videos,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 1740–1748.
- [225] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, “Large-scale video classification with convolutional neural networks,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2014, pp. 1725–1732.

- [226] H. Pirsiavash, C. Vondrick, and A. Torralba, "Assessing the quality of actions," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*. Springer, 2014, pp. 556–571.
- [227] S. M. Safdarnejad, X. Liu, L. Udpa, B. Andrus, J. Wood, and D. Craven, "Sports videos in the wild (svw): A video dataset for sports analysis," in *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, vol. 1. IEEE, 2015, pp. 1–7.
- [228] P. Parmar and B. Tran Morris, "Learning to score olympic events," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 20–28.
- [229] W. Zhang, Z. Liu, L. Zhou, H. Leung, and A. B. Chan, "Martial arts, dancing and sports dataset: A challenging stereo and multi-view dataset for 3d human pose estimation," *Image and Vision Computing*, vol. 61, pp. 22–39, 2017.
- [230] S. Yeung, O. Russakovsky, N. Jin, M. Andriluka, G. Mori, and L. Fei-Fei, "Every moment counts: Dense detailed labeling of actions in complex videos," *International Journal of Computer Vision*, vol. 126, pp. 375–389, 2018.
- [231] P. Parmar and B. Morris, "Action quality assessment across multiple actions," in *2019 IEEE winter conference on applications of computer vision (WACV)*. IEEE, 2019, pp. 1468–1476.
- [232] C. Zalluhoglu and N. Ikizler-Cinbis, "Collective sports: A multi-task dataset for collective activity recognition," *Image and Vision Computing*, vol. 94, p. 103870, 2020.
- [233] Y. Li, L. Chen, R. He, Z. Wang, G. Wu, and L. Wang, "Multisports: A multi-person video dataset of spatio-temporally localized sports actions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13 536–13 545.
- [234] A. Nibali, J. Millward, Z. He, and S. Morgan, "Aspset: An outdoor sports pose video dataset with 3d keypoint annotations," *Image and Vision Computing*, vol. 111, p. 104196, 2021.
- [235] J. Chung, C.-h. Wu, H.-r. Yang, Y.-W. Tai, and C.-K. Tang, "Haa500: Human-centric atomic action dataset with curated videos," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13 465–13 474.
- [236] X. Chen, A. Pang, W. Yang, Y. Ma, L. Xu, and J. Yu, "Sportscap: Monocular 3d human motion capture and fine-grained understanding in challenging sports videos," *International Journal of Computer Vision*, Aug 2021. [Online]. Available: <https://doi.org/10.1007/s11263-021-01486-4>
- [237] P. Parmar and B. Morris, "Win-fail action recognition," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 161–171.
- [238] C. K. Ingwersen, C. Mikkelsen, J. N. Jensen, M. R. Hannemose, and A. B. Dahl, "Sportspose: A dynamic 3d sports pose dataset," in *Proceedings of the IEEE/CVF International Workshop on Computer Vision in Sports*, 2023.
- [239] Y. Cui, C. Zeng, X. Zhao, Y. Yang, G. Wu, and L. Wang, "Sportsmot: A large multi-object tracking dataset in multiple sports scenes," *arXiv preprint arXiv:2304.05170*, 2023.
- [240] T. D'Orazio, M. Leo, N. Mosca, P. Spagnolo, and P. L. Mazzeo, "A semi-automatic system for ground truth generation of soccer video sequences," in *2009 Sixth IEEE International Conference on Advanced Video and Signal Based Surveillance*. IEEE, 2009, pp. 559–564.
- [241] S. A. Pettersen, D. Johansen, H. Johansen, V. Berg-Johansen, V. R. Gaddam, A. Mortensen, R. Langseth, C. Griwodz, H. K. Stensland, and P. Halvorsen, "Soccer video and player position dataset," in *Proceedings of the 5th ACM Multimedia Systems Conference*, 2014, pp. 18–23.
- [242] K. Lu, J. Chen, J. J. Little, and H. He, "Light cascaded convolutional neural networks for accurate player detection," *arXiv preprint arXiv:1709.10230*, 2017.
- [243] J. Yu, A. Lei, Z. Song, T. Wang, H. Cai, and N. Feng, "Comprehensive dataset of broadcast soccer videos," in *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*. IEEE, 2018, pp. 418–423.
- [244] J. Qi, J. Yu, T. Tu, K. Gao, Y. Xu, X. Guan, X. Wang, Y. Dong, B. Xu, L. Hou *et al.*, "Goal: A challenging knowledge-grounded video captioning benchmark for real-time soccer commentary generation," *arXiv preprint arXiv:2303.14655*, 2023.
- [245] C. De Vleeschouwer, F. Chen, D. Delannay, C. Parisot, C. Chaudy, E. Martrou, A. Cavallaro *et al.*, "Distributed video acquisition and annotation for sport-event summarization," *NEM summit*, vol. 8, no. 10.1016, 2008.
- [246] V. Ramanathan, J. Huang, S. Abu-El-Haija, A. Gorban, K. Murphy, and L. Fei-Fei, "Detecting events and key actors in multi-person videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3043–3053.
- [247] S. Francia, S. Calderara, and D. F. Lanzi, "Classificazione di azioni cestistiche mediante tecniche di deep learning," URL: https://www.researchgate.net/publication/330534530_Classificazione_di_Azioni_Cestistiche_mediante_Tecniche_di_Deep_Learning, 2018.
- [248] X. Gu, X. Xue, and F. Wang, "Fine-grained action recognition on a novel basketball dataset," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 2563–2567.
- [249] Y. Yan, N. Zhuang, B. Ni, J. Zhang, M. Xu, Q. Zhang, Z. Zhang, S. Cheng, Q. Tian, Y. Xu *et al.*, "Fine-grained video captioning via graph-based multi-granularity interaction learning," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 2, pp. 666–683, 2019.
- [250] L. Zhu, K. Rematas, B. Curless, S. Seitz, and I. Kemelmacher-Shlizerman, "Reconstructing nba players," in *Proceedings of the European Conference on Computer Vision (ECCV)*, August 2020.
- [251] M. S. Ibrahim, S. Muralidharan, Z. Deng, A. Vahdat, and G. Mori, "A hierarchical deep temporal model for group activity recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1971–1980.
- [252] Ibrahim, Mostafa S and Muralidharan, Srikanth and Deng, Zhiwei and Vahdat, Arash and Mori, Greg, "Hierarchical deep temporal models for group activity recognition," *CoRR*, vol. abs/1607.02643, 2016. [Online]. Available: <http://arxiv.org/abs/1607.02643>
- [253] E. Bermejo Nievas, O. Deniz Suarez, G. Bueno García, and R. Sukthankar, "Violence detection in video using computer vision techniques," in *Computer Analysis of Images and Patterns: 14th International Conference, CAIP 2011, Seville, Spain, August 29-31, 2011, Proceedings, Part II 14*. Springer, 2011, pp. 332–339.
- [254] K. Vats, P. Walters, M. Fani, D. A. Clausi, and J. Zelek, "Player tracking and identification in ice hockey," *arXiv preprint arXiv:2110.03090*, 2021.
- [255] T. De Campos, M. Barnard, K. Mikolajczyk, J. Kittler, F. Yan, W. Christmas, and D. Windridge, "An evaluation of bags-of-words and spatio-temporal shapes for action recognition," in *2011 IEEE Workshop on Applications of Computer Vision (WACV)*. IEEE, 2011, pp. 344–351.
- [256] S. Gourgari, G. Goudelis, K. Karpouzis, and S. Kollias, "Thetis: Three dimensional tennis shots a human action dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2013, pp. 676–681.
- [257] H. Faulkner and A. Dick, "Tennisnet: a dataset for dense fine-grained event recognition, localisation and description," in *2017 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*. IEEE, 2017, pp. 1–8.
- [258] W. Menapace, S. Lathuiliere, S. Tulyakov, A. Siarohin, and E. Ricci, "Playable video generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 061–10 070.
- [259] P.-E. Martin, J. Benois-Pineau, R. Péteri, and J. Morlier, "Sport action recognition with siamese spatio-temporal cnns: Application to table tennis," in *2018 International Conference on Content-Based Multimedia Indexing (CBMI)*. IEEE, 2018, pp. 1–6.
- [260] K. M. Kulkarni and S. Shenoy, "Table tennis stroke recognition using two-dimensional human pose estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4576–4584.
- [261] J. Bian, Q. Wang, H. Xiong, J. Huang, C. Liu, X. Li, J. Cheng, J. Zhao, F. Lu, and D. Dou, "P2a: A dataset and benchmark for dense action detection from table tennis match broadcasting videos," *arXiv preprint arXiv:2207.12730*, 2022.
- [262] S. Zahan, G. M. Hassan, and A. Mian, "Learning sparse temporal video mapping for action quality assessment in floor gymnastics," *arXiv preprint arXiv:2301.06103*, 2023.
- [263] A. Ghosh, S. Singh, and C. Jawahar, "Towards structured analysis of broadcast badminton videos," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 296–304.
- [264] Z. T. L. Shan, "Fineskating: A high-quality figure skating dataset and multi-task approach for sport action," *Peng Cheng Laboratory Communications*, vol. 1, no. 3, p. 107, 2020.
- [265] S. Liu, X. Liu, G. Huang, L. Feng, L. Hu, D. Jiang, A. Zhang, Y. Liu, and H. Qiao, "Fsd-10: a dataset for competitive sports content analysis," *arXiv preprint arXiv:2002.03312*, 2020.

- [266] S. Liu, A. Zhang, Y. Li, J. Zhou, L. Xu, Z. Dong, and R. Zhang, "Temporal segmentation of fine-gained semantic action: A motion-centered figure skating dataset," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 3, 2021, pp. 2163–2171.
- [267] Y. Li, Y. Li, and N. Vasconcelos, "Resound: Towards action recognition without representation bias," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 513–528.
- [268] M. D. Rodriguez, J. Ahmed, and M. Shah, "Action mach a spatio-temporal maximum average correlation height filter for action recognition," in *2008 IEEE conference on computer vision and pattern recognition*. IEEE, 2008, pp. 1–8.
- [269] W. Li, Z. Zhang, and Z. Liu, "Action recognition based on a bag of 3d points," in *2010 IEEE computer society conference on computer vision and pattern recognition-workshops*. IEEE, 2010, pp. 9–14.
- [270] J. C. Niebles, C.-W. Chen, and L. Fei-Fei, "Modeling temporal structure of decomposable motion segments for activity classification," in *Computer Vision—ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part II 11*. Springer, 2010, pp. 392–405.
- [271] J. Pers, "Cvbase 06 dataset: a dataset for development and testing of computer vision based methods in sport environments," *SN, Ljubljana*, 2005.
- [272] H. Kitano, M. Asada, Y. Kuniyoshi, I. Noda, and E. Osawa, "Robocup: The robot world cup initiative," in *Proceedings of the first international conference on Autonomous agents*, 1997, pp. 340–347.
- [273] JiDi, "Jidi olympics football," https://github.com/jidi/ai_lib/blob/master/env/olympics_football.py, 2022.
- [274] A. S. Azad, E. Kim, Q. Wu, K. Lee, I. Stoica, P. Abbeel, A. Sangiovanni-Vincentelli, and S. A. Seshia, "Programmatic modeling and generation of real-time strategic soccer environments for reinforcement learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 6, 2022, pp. 6028–6036.
- [275] J. Wang, J. Ma, K. Hu, Z. Zhou, H. Zhang, X. Xie, and Y. Wu, "Tac-trainer: A visual analytics system for iot-based racket sports training," *IEEE Transactions on Visualization and Computer Graphics*, vol. 29, no. 1, pp. 951–961, 2022.
- [276] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu, Z. Wu, D. Zhu, X. Li, N. Qiang, D. Shen, T. Liu, and B. Ge, "Summary of chatgpt/gpt-4 research and perspective towards the future of large language models," 2023.
- [277] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," *arXiv preprint arXiv:2304.02643*, 2023.
- [278] R. Deng, C. Cui, Q. Liu, T. Yao, L. W. Remedios, S. Bao, B. A. Landman, L. E. Wheless, L. A. Coburn, K. T. Wilson *et al.*, "Segment anything model (sam) for digital pathology: Assess zero-shot segmentation on whole slide imaging," *arXiv preprint arXiv:2304.04155*, 2023.
- [279] S. Roy, T. Wald, G. Koehler, M. R. Rokuss, N. Disch, J. Holzschuh, D. Zimmerer, and K. H. Maier-Hein, "Sam.md: Zero-shot medical image segmentation capabilities of the segment anything model," 2023.
- [280] Y. Liu, J. Zhang, Z. She, A. Kheradmand, and M. Armand, "Samm (segment any medical model): A 3d slicer integration to sam," 2023.
- [281] J. Z. Wu, Y. Ge, X. Wang, W. Lei, Y. Gu, W. Hsu, Y. Shan, X. Qie, and M. Z. Shou, "Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation," *arXiv preprint arXiv:2212.11565*, 2022.
- [282] J. Liu, N. Saquib, Z. Chen, R. H. Kazi, L.-Y. Wei, H. Fu, and C.-L. Tai, "Posecoach: A customizable analysis and visualization system for video-based running coaching," *IEEE Transactions on Visualization and Computer Graphics*, 2022.
- [283] Z. Zhao, S. Lan, and S. Zhang, "Human pose estimation based speed detection system for running on treadmill," in *2020 International Conference on Culture-oriented Science & Technology (ICCST)*. IEEE, 2020, pp. 524–528.
- [284] T. Perrett, A. Masullo, D. Damen, T. Burghardt, I. Craddock, M. Mirmehdi *et al.*, "Personalized energy expenditure estimation: Visual sensing approach with deep learning," *JMIR Formative Research*, vol. 6, no. 9, p. e33606, 2022.
- [285] D. Radke and A. Orchard, "Presenting multiagent challenges in team sports analytics," *arXiv preprint arXiv:2303.13660*, 2023.



Zhonghan Zhao received the BE degree from Communication University of China. He is currently working toward the PhD degree with Zhejiang University - University of Illinois Urbana-Champaign Institute, Zhejiang University. His research interests include machine learning, reinforcement learning and computer vision.



Wenhao Chai received the BE degree from Zhejiang University, China. He is currently working toward the Master degree with University of Washington. His research interests include 3D human pose estimation, generative models, and multi-modality learning.



Shengyu Hao received the MS degree from Beijing University of Posts and Telecommunications, China. He is currently working toward the PhD degree with Zhejiang University - University of Illinois Urbana-Champaign Institute, Zhejiang University. His research interests include machine learning and computer vision.



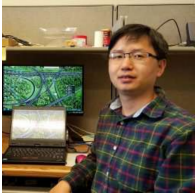
Wenhao Hu received the BS degree from Zhejiang University, China. He is currently working toward the PhD degree with Zhejiang University - University of Illinois Urbana-Champaign Institute, Zhejiang University. His research interests include generative models and 3D reconstruction.



Guanhong Wang received the MS degree from Huaqiao University, China. He is currently working toward the PhD degree with Zhejiang University - University of Illinois Urbana-Champaign Institute, Zhejiang University. His research interests include deep learning and computer vision.



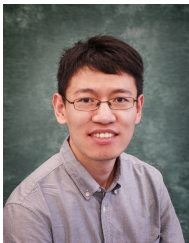
Shidong Cao received the BE degree from Beijing University of Posts and Telecommunications, China. He is currently working toward the MS degree with Zhejiang University - University of Illinois Urbana-Champaign Institute, Zhejiang University. His research interests include machine learning and computer vision.



Dr. Mingli Song received the Ph.D. degree in computer science from Zhejiang University, China, in 2006. He is currently a Professor with the Microsoft Visual Perception Laboratory, Zhejiang University. His research interests include face modeling and facial expression analysis. He received the Microsoft Research Fellowship in 2004.



Dr. Jenq-Neng Hwang received the BS and MS degrees, both in electrical engineering from the National Taiwan University, Taipei, Taiwan, in 1981 and 1983 separately. He then received his Ph.D. degree from the University of Southern California. In the summer of 1989, Dr. Hwang joined the Department of Electrical and Computer Engineering (ECE) of the University of Washington in Seattle, where he has been promoted to Full Professor since 1999. He is the Director of the Information Processing Lab. (IPL), which has won several AI City Challenges and BMTT Tracking awards in the past years. Dr. Hwang served as associate editors for IEEE T-SP, T-NN and T-CSVT, T-IP and Signal Processing Magazine (SPM). He was the General Co-Chair of 2021 IEEE World AI IoT Congress, as well as the program Co-Chairs of IEEE ICME 2016, ICASSP 1998 and ISCAS 2009. Dr. Hwang is a fellow of IEEE since 2001.



Dr. Gaoang Wang joined the international campus of Zhejiang University as an Assistant Professor in September 2020. He is also an Adjunct Assistant Professor at UIUC. Gaoang Wang received a B.S. degree at Fudan University in 2013, a M.S. degree at the University of Wisconsin-Madison in 2015, and a Ph.D. degree from the Information Processing Laboratory of the Electrical and Computer Engineering department at the University of Washington in 2019. After that, he joined Megvii US office in July 2019 as a research scientist working on multi-frame fusion. He then joined Wyze Labs in November 2019 working on deep neural network design for edge-cloud collaboration. His research interests are computer vision, machine learning, artificial intelligence, including multi-object tracking, representation learning, and active learning. Gaoang Wang published papers in many renowned journals and conferences, including IEEE T-IP, IEEE T-MM, IEEE T-CSVT, IEEE T-VT, CVPR, ICCV, ECCV, ACM MM, IJCAI.