

# MultiSports: A Multi-Person Video Dataset of Spatio-Temporally Localized Sports Actions

Yixuan Li    Lei Chen    Runyu He    Zhenzhi Wang    Gangshan Wu    Limin Wang<sup>✉</sup>  
State Key Laboratory for Novel Software Technology, Nanjing University, China

## Abstract

*Spatio-temporal action detection is an important and challenging problem in video understanding. The existing action detection benchmarks are limited in aspects of small numbers of instances in a trimmed video or low-level atomic actions. This paper aims to present a new multi-person dataset of spatio-temporal localized sports actions, coined as MultiSports. We first analyze the important ingredients of constructing a realistic and challenging dataset for spatio-temporal action detection by proposing three criteria: (1) multi-person scenes and motion dependent identification, (2) with well-defined boundaries, (3) relatively fine-grained classes of high complexity. Based on these guidelines, we build the dataset of MultiSports v1.0 by selecting 4 sports classes, collecting 3200 video clips, and annotating 37701 action instances with 902k bounding boxes. Our datasets are characterized with important properties of high diversity, dense annotation, and high quality. Our MultiSports, with its realistic setting and detailed annotations, exposes the intrinsic challenges of spatio-temporal action detection. To benchmark this, we adapt several baseline methods to our dataset and give an in-depth analysis on the action detection results in our dataset. We hope our MultiSports can serve as a standard benchmark for spatio-temporal action detection in the future. Our dataset website is at <https://deeperaction.github.io/multisports/>.*

## 1. Introduction

Spatio-temporal human action detection in untrimmed videos is of great importance for many applications, such as surveillance and sports analysis. Recently, recognizing actions from short trimmed videos has achieved considerable progress [52, 3, 48, 42, 49, 50], but these classification models can not be directly applied for video analysis in a multi-person scene. Meanwhile, although temporal action detection methods [66, 31, 29, 59, 63] for untrimmed videos can

distinguish intervals of human actions from background, they are still unable to spatially detect multiple concurrent human actions, which is important in real-world applications of video analysis.

Current spatio-temporal action detection benchmarks can be mainly classified into two categories: 1) Densely annotated high-level actions such as J-HMDB [20] and UCF101-24 [46]. Their clips only have a single person doing some semantically simple and temporally repeated actions. Typically, the scene context can provide enough cues for recognizing these coarse-grained action categories. Thus, these benchmarks might be impractical for real-world applications such as surveillance, where it is required to deal with more fine-grained actions in a multi-person scene; 2) Sparsely annotated atomic actions such as AVA [15]. They fail to provide clear temporal action boundaries, and simply focus on frame-level spatial localization of atomic actions. This setting removes the requirements of temporal localization for action detection algorithms. Meanwhile, their atomic actions rarely require the complex reasoning over the actors and their surrounding environment.

Based on the analysis above, we argue that a new benchmark is necessary to advance the research of spatio-temporal action detection. The benchmark should satisfy several important requirements to cover the realistic challenges of this task. 1) There should be multiple persons performing different actions concurrently in the same scene, where the background information is not sufficient for action recognition and motion itself of the actor plays a significant role. 2) To address the inherently confusing human action boundaries in time, actions should be both semantically and temporally well-defined with a consensus among humans. 3) Considering the complexity of real-world applications, actions should be fine-grained which requires accurate human pose and motion information, long-term temporal structure, possible interactions between humans, objects and scenes, and reasoning over their relations.

Following the above guidelines, we develop the *MultiSports* dataset, short for *Multi-person Sports Actions*. The dataset is large-scale, high-quality, multi-person, and con-

<sup>✉</sup>: Corresponding author (lmwang@nju.edu.cn).

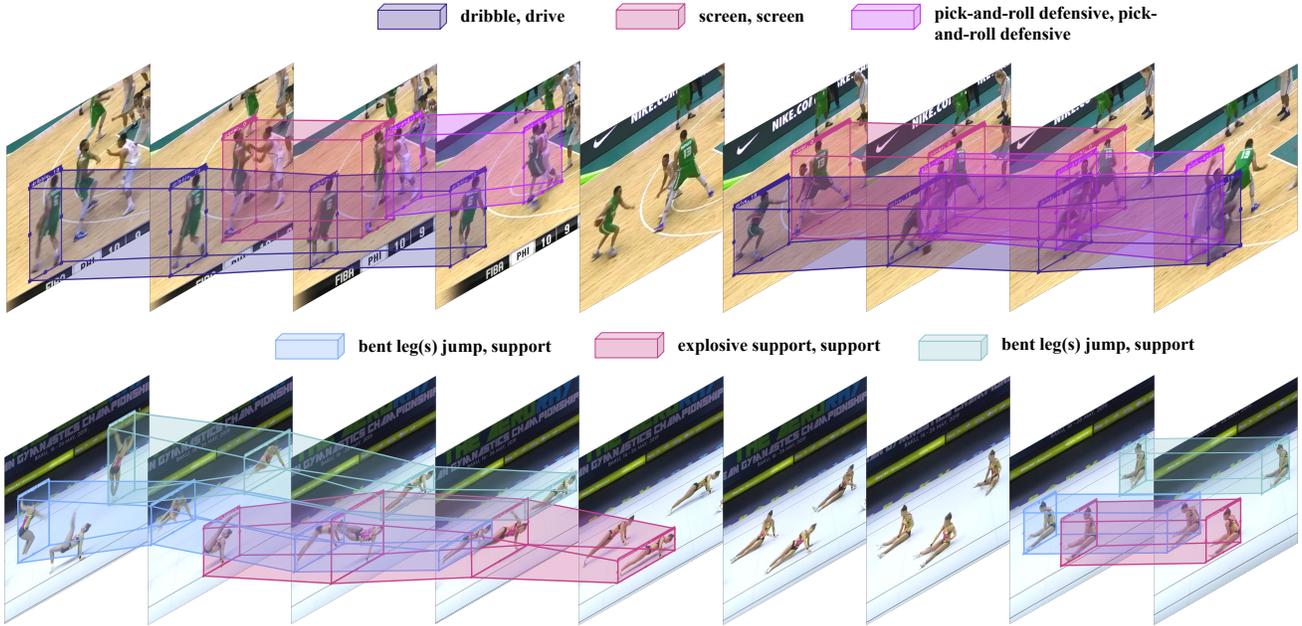


Figure 1. The 25fps tubelets of bounding boxes and fine-grained action category annotations in MultiSports dataset. Multiple concurrent action situations frequently appear in MultiSports with many starting and ending points in the long untrimmed video clips. The frames are cropped and sampled by stride 5 or 7 for visualization propose. Tubes with the same color represent the same person.

tains fine-grained action categories with precise and dense annotations in both spatial and temporal domains. The action vocabulary consists of 66 action classes collected from 4 sports (basketball, volleyball, football and aerobic gymnastics). An example clip has been visualized in Figure 1. We choose these four sports for the following reasons. 1) There are plenty of multiple concurrent action instances in sports competitions. Also, the background is far less characteristic and cannot provide sufficient information for fine-grained action recognition. 2) Sports actions have well-defined categories and boundaries. These boundaries are defined by either professional athletes or official documentations [8]. 3) Due to the complex competition rules, recognizing sports action generally requires to model the long-term structure and the human-object-scene interactions. For example, in football, although the athlete may take only 0.5s to kick the ball, we may need up to 5s context to recognize whether it is pass, long ball, through ball, or cross.

In practice, we conduct exhaustive annotations of 25 fps frame-wise bounding boxes and fine-grained action categories in a two-stage procedure: 1) a team of professional athletes of corresponding sport to annotate the temporal and category labels, and 2) a team of crowd-sourced annotators to finish the bounding boxes with the help of tracking method FCOT [7]. This two-stage annotation procedure as well as careful quality control together can guarantee consistent and clean annotations. To ensure the visual quality, all videos in our dataset are high-resolution records of professional competitions from a diversity of countries and

different performance levels.

Given the well-defined and dense-annotated action instances in *MultiSports v1.0*, we benchmark spatio-temporal action detection on this challenging dataset. We perform empirical studies with several recent state-of-the-art action detector methods. Compared with previous action detection benchmarks such as J-HMDB [20] and UCF101-24 [46], our MultiSports is quite challenging with a much lower frame mAP and video mAP. We also introduce a detailed error analysis on detection results and try to provide more insights on spatio-temporal action detection. According to our analysis on MultiSports benchmark, we figure out several challenges of spatio-temporal action detection that needs to be addressed, such as capturing subtle differences between fine-grained action categories, performing accurate temporal localization, dealing with action occlusion and modeling long-range context. We hope MultiSports could serve as a standard benchmark to advance the area of spatio-temporal action detection in the future. MultiSports spatio-temporal action detection is currently a track of DeeperAction challenge at ICCV 2021 <https://deeperaction.github.io/>.

In summary, our main contribution is twofold. 1) We develop a new benchmark MultiSports of spatio-temporal action detection for well-defined and realistically difficult human actions in a multi-person scene, providing high-quality and 25fps frame-wise annotations from four sports. 2) We conduct extensive studies and systematic error analysis on MultiSports, which reveals the key challenges of spatio-temporal action detection and hopefully can facilitate future

research in this area.

## 2. Related Work

**Action recognition datasets.** Early datasets of action recognition mainly focus on action classification. Those datasets, including KTH [39], Weizmann [2], UCF-101 [46] and HMDB [24], contains manually trimmed short clips to capture semantics of a single action. Their human action cues, however, are overwhelmed by signals of background scenes. Multi-MiT [33] is a multi-label action recognition dataset, which may have several concurrent actions but do not provide temporal duration and spatial annotations. Recently, large-scale video classification datasets such as Sports-1M [22], YouTube-8M [1] and Kinetics [3] have been created for feature representation learning and serve as pre-training in downstream tasks, but appearance cues still play a important role here. Something-something [14] and FineGym [40], with plenty of fine-grained action categories, effectively reduce the influences of background scenes and reveal some key challenges of modeling a single action. They share the similar property of capturing motion cues with *MultiSports*, but only have one concurrent action therefore we address a different need with them.

Temporal action detection datasets such as ActivityNet [16], HACS [64], THUMOS14 [19], MultiTHUMOS [61] and Charades [41] provide temporal action detection annotations for each action of interest in untrimmed videos. But unlike *MultiSports*, they do not provide spatial annotations and could not identify multiple concurrent actions for multiple people.

Previous spatio-temporal action detection datasets, such as UCF Sports [37], UCF101-24 [46] and J-HMDB [20], typically evaluate spatio-temporal action detection for short videos with only a single person and coarse-grained action categories. Our *MultiSports* significantly differs from them in several aspects: multiple concurrent actions by multiple people; less characteristic background scenes; the larger number of action and fine-grained categories; more fast movement and large deformation; and significantly more instances per clip. Recently, a new type of extensions such as DALY [54], AVA [15] and AVA-Kinetics [25] adopt sparse annotations of daily life actions, either in composite or atomic forms, to reduce human labors of annotating and increase the scale of datasets. It may be a good way for evaluating daily life actions without fast movement and large deformation, but unsuitable for areas like sports analysis, since it often requires continuous annotations of all human actions of interest. MEVA [6] is a security dataset, which provides spatial-temporal annotations and some other modality annotations. But our sports actions are more complex and fast-changing than MEVA. Different from previous datasets, our *MultiSports* proposes a more difficult benchmark with multi-person, well-defined bound-

aries, fine-grained setting and frame-by-frame annotations, which focuses on the sports domain.

**Spatio-temporal action detection.** Most recent approaches for UCF101-24 and JHMDB can be classified into two categories: frame-level detectors and clip-level detectors. Many efforts have been made to extend an image object detector to the task of spatio-temporal action detection at the frame level [13, 51, 34, 38, 44, 53], where the resulting per-frame detections are then linked to generate final tubes. Although flows could be used to capture motion cues, frame-level detector fails to fully utilize temporal information. To model temporal structures for action detection, some clip-level approaches or action tubelet detectors [18, 26, 21, 60, 27, 65, 45] have been proposed. ACT [21] took several frames as input and detected tubelets regressed from anchor cuboids. STEP [60] progressively refined the proposals by a few steps to solve the large displacement problem and utilized longer temporal information. MOC-detector [27] proposed an anchor-free tubelet detector by treating action instances as trajectories of moving points. For AVA, many methods [11, 12, 47, 55, 56] have been proposed to better make use of spatio-temporal information for atomic action classification.

## 3. The MultiSports Dataset

Our *MultiSports* dataset aims to introduce a new challenging benchmark with high-quality annotations to the area of spatio-temporal action detection, which differs from previous ones in multi-person scene, well-defined temporal boundaries, and fine-grained action categories. Sec. 3.1 introduces our annotation procedure. Statistics and characteristics of *MultiSports* are elaborated in Sec. 3.2 and Sec. 3.3.

### 3.1. Dataset Construction

**Action vocabulary generation.** We select sports of basketball, volleyball, football and aerobic gymnastics, because of their multi-person setting, less ambiguous actions and well-defined temporal boundary. For aerobic gymnastics, we use the official documentations [8]. In practice, we only select *difficulty elements* and discard *movement patterns*. For the remaining ball sports, we use an iterative way to generate our action vocabulary in each sport: we initialize an action list by the suggestions of athletes and write a handbook to clarify the definition of action boundaries. Then we let several annotators try to annotate the data, where inaccurate definitions of action boundaries, ambiguities between action categories and missed action categories will be collected from their feedback. We iteratively adjust our action list and handbook according to the feedback several times before we start massive annotating, which results in the final action hierarchy shown in Figure. 2(a). Note that the annotators of action categories and temporal boundaries are professional athletes of the corresponding sports, so their

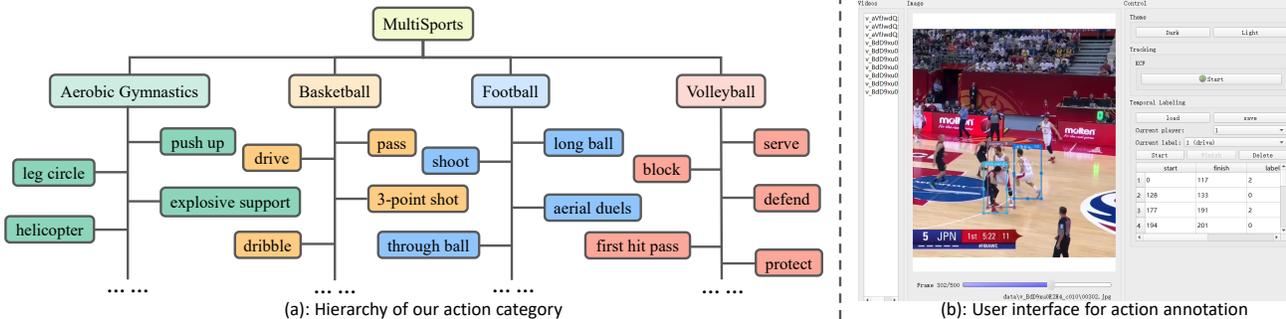


Figure 2. The action vocabulary hierarchy and annotator interface of the MultiSports dataset. (a) Our MultiSports has a two-level hierarchy of action vocabularies, where the actions of each sport are fine-grained. (b) Details of annotations can be found in Sec 3.1.

feedback is important for building a well-defined action vocabulary in practice. To keep action boundaries accurate and make our dataset suitable for spatio-temporal action detection, we do not count common and atomic actions such as run or stand in our action vocabulary. We also exclude foul in ball sports. Because in the 2D video records, we recognize fouls most from the referee’s reaction instead of the actor’s motion. What is worse, it is hard to identify who fouls due to occlusion.

**Data preparation.** After choosing the four sports, we search for their competition videos by querying the name of sports like volleyball and the name of competition levels like Olympics and World Cup on YouTube, and then download videos from top search results. For each video, we only select high-resolution, e.g. 720P or 1080P, competition records and then manually cut them into clips of minutes, with less shot changes in each clip and to be more suitable for action detection. These official records share consistent and rich content, and can guarantee a high-quality dataset.

**Action annotation.** Since our annotations are difficult in labeling fine-grained categories and exhaustive in determining 25fps frame-wise bounding boxes, we naturally decompose our annotation procedure into two stages: 1) A team of professional athletes generate records of the action label, the starting and ending frame, and the person box in the starting frame, which can ensure the efficiency, accuracy and consistency of our annotation results; 2) With the help of FCOT [7] tracking algorithms, a team of crowd-sourced annotators adjust bounding boxes of tracking results at each frame for each record. The ambiguity of spatial human boundaries is much less than that of fine-grained action categories and temporal action boundaries. They use the interface shown in Figure 2(b).

To ensure the consistency of action temporal boundaries, which tends to be ambiguous and remains as a big challenge for most temporal action detection datasets, we write a handbook to clarify the definition of action boundaries as mentioned above. For example, our handbook unifies the annotations of *football pass* as starting from the ball-controlling-leg leaving the ground and ending with this leg

touching the ground again. The annotation handbook is provided in Appendix E.

**Person bounding-box tracking.** As mentioned above, we first tack each record generated by professional athletes and then employ crowd-sourced annotators to refine the bounding boxes at each frame. Specifically, we use FCOT [7] to track the bounding boxes frame-by-frame. We find this tracking-to-refinement labeling process can not only speed up the annotation process, but also increase the annotation quality by enforcing workers to focus on determining precise boundary of each box.

We also evaluate the output of FCOT [7] and results are shown in Table 1. We adopt success and precision metrics proposed in OTB100 [57]. Aerobic turned out the hardest in both success and precision aspects.

	Aerobic gym.	Volleyball	Football	Basketball
Success	0.66	0.72	0.77	0.66
Precision	0.67	0.93	0.92	0.72

Table 1. Tracking results on different sports

**Quality control.** For the first stage of annotation, every clip has at least one annotator with domain knowledge double-checking the annotations. We correct wrong or inaccurate ones and also add missing annotations for a higher recall, e.g., adding missed defence action in football and modifying inconsistent action boundaries. For the second stage, we double-check each instance by playing it in 5fps and manually correct the inaccurate bounding boxes.

### 3.2. Dataset Statistics

Our *MultiSports* v1.0 contains 66 fine-grained action categories from four sports, and has videos selected from 247 competitions. The videos are manually cut into 800 clips per sport to keep data balance between sports. We discard intervals with only background scenes, such as award, and select the highlights of competitions as clips for action detection. Table 2 compares the annotation types and statistics of MultiSports v1.0 with the existing datasets. AVA [15] only has sparse and 1fps annotations of bounding boxes, which fails to provide clear temporal action boundaries and



egories with a long-tailed distribution; 3) the large variance of action instance duration, which makes it difficult to localize the temporal boundary; 4) the fast movement, large deformation and occlusion of actions in sports.

**High Quality.** The videos of MultiSports are with high-resolution (720P or 1080P) competition records, which can preserve details of small humans and objects. Besides, with the help of our annotation team composed of professional athletes, our action categories and their corresponding action boundaries are precisely annotated. The professional annotators and careful quality control is able to provide consistent and clean annotations.

**Diversity.** Our video clips are selected from competitions of different performance levels with diverse countries and genders, making the dataset less biased and good balanced for realistic sports analysis.

**Application.** This task has many application scenarios for sports analysis. Combined with Re-ID techniques, we can automatically perform game commentary, AI referee and technical statistics. It can also assess the player abilities and provide information for developing the training plan and game strategy, and trading players between clubs.

## 4. Experiments and Analysis

### 4.1. Datasets and Metrics

**MultiSports benchmark.** To build a solid action detection benchmark, we manually split the instances into the training set, validation set, and testing set. Due to the long-tailed distribution of action instance numbers, following AVA [15], we only evaluate on 60 classes that have at least 25 instances in validation and test splits to benchmark performance. We resize the whole dataset into 720P. In total, the current version contains 18,422 training instances from 1,574 clips and 6,577 validation instances from 555 clips. We provide the detailed ratio of training and validation instances for each sport in Appendix A. All those instances are selected from 3200 clips covering 247 competition records. Unless otherwise mentioned, we report the results trained on the training set and evaluated on the validation set. The testing set includes 1071 clips and we withhold the annotations in the public release.

**Metrics.** Following the standard practice [53, 21], we utilize frame-mAP and video-mAP to evaluate action detection performance. For video-mAP, we use the 3D IoU, which is defined as the temporal domain IoU of two tracks, multiplied by the average of the IoU between the overlapped frames. The threshold is 0.5 for frame-mAP, 0.2 and 0.5 for video-mAP.

### 4.2. Spatio-temporal Action Detection Results

We evaluate several representative action detection methods on *MultiSports* and compare their performance on

the UCF101-24 [46], JHMDB [20], and AVA [15] in Table 3. For SlowOnly Det. and SlowFast Det., we use the code in MMAAction2 [5]. We use the official released code for ROAD, YOWO and MOC. More details about the methods are provided in Appendix C.

For UCF101-24 [46] and JHMDB [20], which have dense annotations of high-level actions as MultiSports, we find that these methods achieve good performance on them but obtain low performance on MultiSports (frame-mAP of **25.22%**, video-mAP@0.2 of **12.88%** and video-mAP@0.5 of **0.62%** for MOC [27]). In our dataset, the largest performance drop occurs on ROAD [44], which is a frame-level action detector that performs action detection at each frame independently without exploiting temporal information. UCF101-24 [46] and JHMDB [20] have only one category per video. Characteristic visual scenes provide enough cues for predicting their coarse-grained actions. However, MultiSports has a similar background in the same sport, where the background fails to provide sufficient information for fine-grained action recognition. Meanwhile, our temporal boundary annotation is more precise and requires more accurate localization in temporal domain.

For AVA [15], which has only sparse annotations of atomic actions, we observe that the performance gap between SlowFast Det. [11] and SlowOnly Det. [11] on MultiSports is more evident than on AVA (frame-mAP gap of **11.02%** vs. **4.54%**). This indicates that the sports actions need a higher frame rate to capture fast motion at a finer temporal granularity. As shown in Figure 5, many aerobic actions gain large absolute improvement, such as aerobic turn (+30 AP) and aerobic horizontal support (+54 AP). We analyze that aerobic actions’ deformation and displacement is the largest among the four sports and benefit more from this finer temporal analysis. We also observe a large performance increase in other sports, such as basketball pass, football clearance and volleyball second attack, which have short temporal duration and intense motion.

### 4.3. Error Analysis

In this section, we analyze the cause of errors to better understand *MultiSports*’ challenges. Based on ACT [21] frame-mAP error analysis, which is designed for the dataset with one action category per video, we propose a new detailed error analysis in video-mAP. We classify the detection errors into 10 mutually exclusive categories to analyze which percentage of the mAP is lost.  $E_R$ : a detection result targets at a ground-truth tube that has already been matched.  $E_N$ : a detection result that has no spatial-temporal intersection with any ground-truth tubes and appears out of thin air.  $E_L$ : a detection result that has the correct action class, accurate temporal localization and inaccurate spatial localization.  $E_C$ : a detection result that has the wrong action class, accurate temporal localization and accurate spatial lo-

Method	Res	MultiSports			UCF101-24			JHMDB			AVA
		F@0.5	V@0.2	V@0.5	F@0.5	V@0.2	V@0.5	F@0.5	V@0.2	V@0.5	F-mAP@0.5
ROAD [44]	300 × 300	3.90	0.00	0.00	70.7	69.8	40.9	-	60.8	59.7	-
YOWO [23]	224 × 224	9.28	10.78	0.87	71.10	72.97	46.42	74.51	88.05	82.57	-
MOC [27] (K=7)	288 × 288	22.51	12.13	0.77	78.0	82.8	53.8	70.8	77.3	77.2	-
MOC [27] (K=11)	288 × 288	25.22	12.88	0.62	-	-	-	-	-	-	-
SlowOnly Det., 4 × 16 [11]	short side 256	16.70	15.71	5.50	-	-	-	-	-	-	20.02
SlowFast Det., 4 × 16 [11]	short side 256	27.72	24.18	9.65	-	-	-	-	-	-	24.56

Table 3. Comparison of the state-of-the-art methods on MultiSports, UCF101-24, JHMDB and AVA.

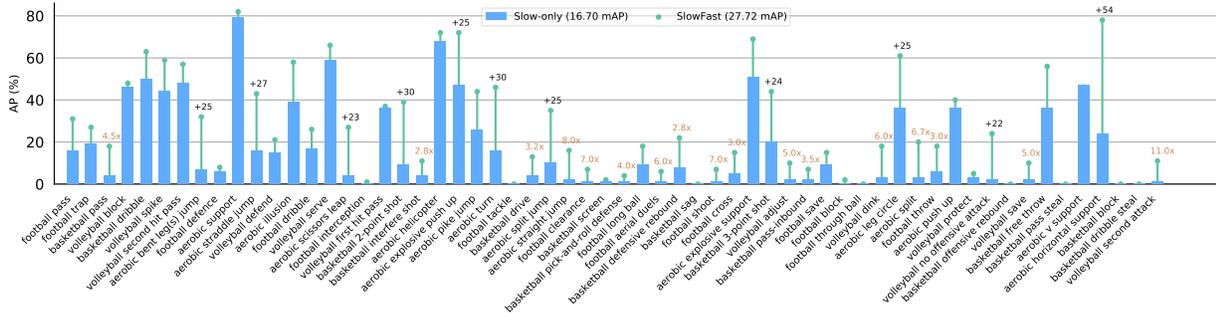


Figure 5. SlowOnly vs. SlowFast frame-mAP. Categories are sorted by descending order on the number of instances.

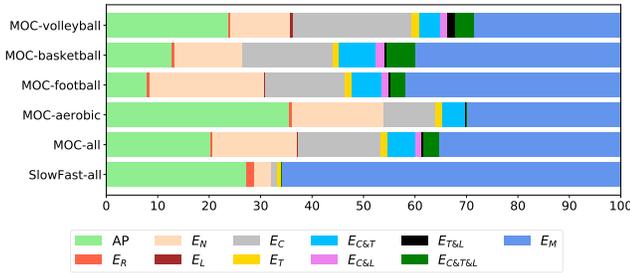


Figure 6. Error Analysis. AP is the correct detection. The threshold for a ground-truth matched by a detection is 0.1. Recall is  $1 - E_M$

calization.  $E_T$ : a detection result that has the correct action class, accurate spatial localization and inaccurate temporal localization.  $E_{C\&T}$ ,  $E_{C\&L}$ ,  $E_{T\&L}$ ,  $E_{C\&T\&L}$ : a detection that is inaccurate in corresponding aspects while acceptable in other aspect (if any). For example,  $E_{C\&T}$  refers to results with wrong action class, inaccurate temporal localization yet accurate spatial localization.  $E_M$ : ground-truth tubes that have not been matched by any detection results. The first nine categories cover the false positive predictions. The partition can be explained with a decision tree which is attached to our Appendix D. The code is provided at <https://github.com/MCG-NJU/MultiSports>.

As shown in Figure 6, despite the relatively low recall, SlowFast Det. achieves higher video-mAP than MOC because it makes much fewer false positive predictions. This can be explained by the fact that SlowFast Det. uses Faster RCNN [36] finetuned on MultiSports as person detector to greatly avoid the person boxes without actions. However, there are still many hard examples missed by Slow-

Fast Det. For MOC,  $E_C$  and  $E_N$  are the most common errors among false positive detection results, indicating the difficulty of our fine-grained action classification. Detection results with  $E_N$  errors means the model indeed detects the person spatio-temporally but unable to identify the action class correctly as the background class.  $E_N$  error is also a result of the training strategy of MOC where only the frames temporally inside action instances are sampled for training, so that although there are negative samples in other spatial location of these frames, the detector does not have enough amount of negative samples for people without doing any sports action. What is more,  $E_{C\&T}$ ,  $E_{C\&T\&L}$  and  $E_T$  are also a large portion of the rest errors (where  $E_{C\&T} > E_{C\&T\&L} > E_T$ ), indicating more temporal errors with inaccurate action boundaries than spatial localization errors for current methods. Therefore we need a more effective way of modeling temporal boundary. Typical error visualization is shown in Figure 7.

#### 4.4. Ablation Study

**How important is temporal information?** The tubelet length  $K$  is important in MOC [27] and we report results on UCF101-24 [46] and *MultiSports* with different  $K$  in Table 4. For frame-mAP, we can find that *MultiSports* can benefit more from longer temporal context than UCF101-24, in spite of the shorter action duration of *MultiSports* than UCF101-24 as shown in Table 2. For video-mAP, the result does not increase as frame-mAP. We analyze there are two reasons. First, predicting movement in MOC is harder with longer input length. What is worse, the categories in *MultiSports* have large deformation and displacement, and



Figure 7. Visualization of typical errors in MultiSports. Green boxes are the ground-truths. Yellow boxes are the detections. Red boxes are the missed ground-truths. 1st and 2nd row: missed detection due to occlusion. 3rd and 4th row:  $E_{C\&T}$ : drive is misclassified as dribble and also has inaccurate action boundary;  $E_M$ : missed detections of screen, pick-and-roll defensive and sag.

MOC Movement Branch can not predict them accurately, which harms the video level detection seriously. Second, Figure 4 shows the variability of action duration. The ratio is 9% for instances duration less than 7 and 23% for less than 11. The fixed clip length  $K$  (e.g. 11) will damage temporal detection ability. So, we need to consider longer temporal context, more accurate movement estimation and flexible temporal detection for MultiSports.

**Which action categories are challenging?** Figure 5 shows that not all categories yield better performance with more training samples. Categories highly correlated with scenes (such as basketball free throw) or aerobics basic categories (such as aerobic horizontal support and V support) can still achieve high performance with fewer samples. Note that aerobics contains basic and complex categories, where complex action combines the motion of basic action and its own core motion, thus longer temporal context is required for these complex actions. In contrast, categories with short temporal duration and intense motion (such as football pass, basketball pass and football interception) achieve low performance even though with lots of training samples. By observing the confusion matrix in Appendix D, we summarize other common challenges: (1) Context modeling, such as basketball 2-point shot vs. 3-point shot (2) Reasoning, such as for volleyball protect vs. defend, we need to focus on whether the ball was blocked back or was spiked by an opponent several frames earlier. (3) Long temporal modeling, such as football long ball vs. pass, they have the similar motion but need to identify how long the ball will be passed.

**Trimmed vs. untrimmed settings.** *MultiSports* has well-defined and high-quality temporal boundaries. We eval-

K	MultiSports			UCF101-24		
	F@0.5	V@0.2	V@0.5	F@0.5	V@0.2	V@0.5
1	14.61	12.53	1.06	68.33	65.47	31.50
3	17.22	11.88	0.76	69.94	75.83	45.94
5	19.29	11.81	<b>0.98</b>	71.63	77.74	49.55
7	22.51	12.13	0.77	<b>73.14</b>	<b>78.81</b>	<b>51.02</b>
9	24.22	11.72	0.57	72.17	77.94	50.16
11	<b>25.22</b>	<b>12.88</b>	0.62	-	-	-
13	24.28	11.23	0.57	-	-	-

Table 4. Exploration study of MOC on the MultiSports and UCF101-24 with different tubelet length  $K$ .

Estimation	MultiSports			AVA
	F@0.5	V@0.2	V@0.5	F-mAP@0.5
Untrimmed	27.72	24.18	9.65	22.57
Trimmed	38.71	24.95	18.34	24.56

Table 5. Test SlowFast Det. on AVA and MultiSports with trimmed way and untrimmed way.

uate the performance of SlowFast Det. under both the untrimmed and trimmed setting on MultiSports and AVA datasets. The results are reported in Table 5. The trimmed setting only evaluates the performance on the frames having annotations and the untrimmed setting reports the performance on all frames. We find that it only drops 2% on AVA while 11% on our dataset, which indicates that temporal localization is really important in our dataset. In addition, video-mAP@0.5 drops far more than video-mAP@0.2. This demonstrates that temporal localization is important for high-quality action tube detection.

## 5. Conclusion

In this paper, we have introduced the *MultiSports* dataset with dense spatio-temporal annotations of actions from four sports. *MultiSports* distinguishes from the existing action detection datasets in many aspects: 1) raising new challenges for recognizing fine-grained action classes; 2) requirement of accurate localization of well-defined boundaries in multiple-person situations; 3) high quality video data and dense annotations; 4) potential applications in sports analysis; 5) less biased dataset with high diversity in competition levels, countries and genders. We have empirically investigated several action detection baseline methods on the *MultiSports* dataset. Our error analysis and ablation studies on the detection results uncover several insightful findings that are beneficial for the future research of spatio-temporal action detection.

**Acknowledgements.** This work is supported by National Natural Science Foundation of China (No. 62076119, No. 61921006), Program for Innovative Talents and Entrepreneur in Jiangsu Province, and Collaborative Innovation Center of Novel Software Technology and Industrialization. Thanks to professional athletes of Nanjing University varsities and MCG students for annotating this dataset.

## Appendix A: More Dataset Details

### A.1 Train split vs. Validation split

In order to guarantee enough instances for each class despite the severely unbalanced distribution, we artificially split the instances into the training set and the validation set in Table 6. To avoid data leakage from the training set to the validation/testing set, we ensure that data from the same match should be used for only one purpose. In other words, clips in the validation set cannot come from the matches covered in the training set. Unless otherwise mentioned, we report the results trained on the training set and evaluated on the validation set.

### A.2 Comparison with other type of Dataset

MEVA [6] is a new security dataset, whose data is from RGB and thermal IR cameras, UAV footage and GPS locations for the actors. It defines 37 activities (66 for *MultiSports*) with 17055 instances (37701 for *MultiSports*), where 29 activities are about person and 8 activities are about vehicle. The categories in this dataset are atomic, such as *person\_close\_trunk* and *person\_stand\_up*, which are different from our fine-grained and complicated sports categories. What’s more, most of the categories in MEVA are daily actions, whose deformation and displacement are not large. Although it is a multi-person dataset, we believe our *MultiSports* can bring new challenges different from MEVA.

## Appendix B: More Ablation Study

**How the well-defined and high quality temporal boundary help?** We add some temporal noise to the train set GT. For a  $L$ -frame length instance, we randomly choose a new length  $\text{new\_L}$  from  $(1, L)$  and then the start point offset from  $(0, L - \text{new\_L})$ . We sample the new annotation from the original. Other settings are kept the same. From the Table 7, we find the performance is much worse without well-defined temporal boundaries. It can conclude that our *MultiSports* has well-defined and high quality temporal annotations, which can help improve the performance and promote the algorithms to localize the boundary more accurately.

## Appendix C: Method Details

**ROAD** [44] is a deep-learning framework for real-time action localisation and classification. It adopts SSD [32] to regress and classify action detection boxes in each frame independently, which does not utilize temporal information. Then, the frame detections are linked into action tubes by an online algorithm. Here we use the python linking code provided by MOC [27] instead of the original MATLAB code. Following the settings of ROAD on UCF101-24 [46],

we use an ImageNet pre-trained VGG16 [43] network. We first try an initial learning rate of  $1e-4$  as their setting on UCF101-24, but the loss diverges into infinity after 20 iterations. The reported experiment on our *MultiSports* adopts an initial learning rate of  $1e-5$ . We use SGD optimizer and the learning rate is reduced to its  $\frac{1}{10}$  after 30000, 60000 iterations, which is the same as their practice on UCF101-24. The maximum iteration number is 120000.

**YOWO** [23] is a frame-level action detector with two branches. A 2D-CNN branch extracts the spatial features of the key frame while a 3D-CNN branch extracts spatio-temporal features of the key frame and the previous  $n$  ( $n=16$ ) frames. Then, the features of two branches are fused by a channel fusion and attention mechanism(CFAM) module and finally passed to a convolution layer to predict the action class and bounding box in Yolov2 [35] manner. Finally, the frame detections are linked into action tubes by a dynamic programming algorithm. Note that the linking algorithm in YOWO is trimmed, thus we use the same linking algorithm as MOC on *MultiSports*. We use 2D Darknet-19 backbone pretrained on PASCAL VOC [10] and 3D ResNeXt-101 backbone pre-trained on Kinetics [3]. To utilize multiple GPUs, we modified the batch size to 80 and the initial learning rate to  $8e-4$ . Following the training strategy of YOWO on UCF101-24 [46], we adopt SGD optimizer and the learning rate is reduced to its  $\frac{1}{2}$  after 30000, 40000, 50000, 60000 iterations. The epoch maximum is set to 5. Note that YOWO only estimates performance on the frames having annotations, thus frame-mAP we report on UCF101-24 is much lower than in the original paper.

**MOC** [27] is an anchor-free tubelet-level action detector with three branches, which firstly takes  $K$  frames as input, then outputs  $K$  frame tubelet results and finally links these tubelets into tubes with a common matching strategy. We use DLA34 [62] as the backbone network, which is pre-trained on COCO [30]. Following the training strategy of MOC on UCF101-24 [46], we use the Adam optimizer with the learning rate  $5e-4$ . The learning rate is reduced to its  $\frac{1}{10}$  after epoch 6 and 8. The epoch maximum is set to 12.

**SlowFast Det.** [11] firstly uses a person detector on the key frame to localize for region proposal. Then, each 2D RoI at the key frame is extended into a 3D RoI by replicating it along the temporal dimension. Finally, it extracts RoI features from the backbone features for predicting category. The person detector is a Faster R-CNN with a ResNeXt-101-FPN [58, 28] backbone, which is pre-trained on ImageNet [9] and the COCO human keypoint images [30]. The backbone is the variant of SlowFast or SlowOnly, which sets the spatial stride of  $res_5$  to 1 and uses a dilation of 2 for its filters. Note that we use the code in MMAAction2 [5]. The results on AVA [15] and our *MultiSports* in the paper are all produced by it. We use the pre-computed proposals for AVA from previous work [11, 55]. Following previous

	Volleyball	Football	Basketball	Aerobic	All
instance ratio	3549:1294	6144:2153	4532:1715	4197:1415	18422:6577
clip ratio	402:130	402:132	379:147	391:146	1574:555
competition ratio	32:11	36:12	34:14	23:8	125:45

Table 6. Train split vs Validation split

Method	GT Noise	MultiSports		
		F@0.5	V@0.2	V@0.5
MOC (K=7)	✓	13.71	8.59	0.63
MOC (K=7)	✗	22.51	12.13	0.77
SlowOnly Det., $4 \times 16$	✓	12.60	8.98	3.05
SlowOnly Det., $4 \times 16$	✗	16.70	15.71	5.50

Table 7. Exploration on the effect of the temporal boundary noise.

work [11, 55], we fine-tune the person detector on our *MultiSports* with MMDetection [4]. We use the SGD optimizer with the learning rate 0.0025 and finetune 2 epochs on our *MultiSports*. The person detector produces 96.16 AR@100 on our *MultiSports* validation set. The detected boxes with confidence of  $> 0.9$  are selected for action detection on both datasets. Our backbones are based on ResNet50, which are pre-trained on Kinetics-400 [3]. The  $T \times \tau$  is set to  $4 \times 16$ . The  $\alpha$  is set to 8 for SlowFast. We use a step-wise learning rate, reducing the learning rate  $10\times$  after epoch 6 and 7. We train for 8 epochs with a linear warm-up for the first 5 epochs, where the result is similar with that of training 20 epochs and a lot of training time is saved. The initial learning rate is set to 0.1125 for SlowFast and 0.2 for SlowOnly. SlowFast and Slowonly Det. use the same link algorithm as MOC.

## Appendix D: Error Analysis

### D.1 Error Tree

To further understand the difficulty in our *MultiSports* dataset, we classify the detection errors into 10 different categories in a tree structure as shown in Figure 8 (code in *VideomAP\_error.py*), which are:

- $E_R$  (Errors of repeated detections): a detection result that has tubelet IoU larger than a threshold and the right action class with some ground-truth tubelets, but the ground-truths have been matched by other detection results before with a confidence score larger than it.
- $E_N$  (Errors of not matched): a detection result that has no intersection with any ground-truth tubelets of any class, indicating there should be no detection results but it appears out of thin air.
- $E_L$  (Errors of spatial localization): a detection result

that has the same action class and temporal IoU larger than a threshold with some ground-truth, but it has a low average spatial bounding box IoU in the area of the temporal intersection of ground-truth tubelets and it so that a lower tubelet IoU than the required threshold.

- $E_C$  (Errors of classification): a detection result that has the tubelet IoU larger than a threshold with a ground-truth, but its action class is not the same with the ground-truth’s class.
- $E_T$  (Errors of temporal localization): a detection result that has the same action class and average spatial bounding box IoU larger than a threshold with some ground-truth in the area of the temporal intersection of ground-truth tubelets and it, but low temporal IoU with ground-truths so that a lower tubelet IoU than the required threshold.
- $E_{C\&T}$  (Errors of both classification and temporal localization): a detection result that has average spatial bounding box IoU larger than a threshold with some ground-truth tubelets in the area of the temporal intersection of ground-truth tubelets and it, but both low temporal IoU and wrong action class.
- $E_{C\&L}$  (Errors of both classification and spatial localization): a detection result that has temporal IoU larger than a threshold with some ground-truth tubelets, but both wrong action class and low average spatial bounding box IoU with some ground-truth in the area of the temporal intersection of ground-truth tubelets and it.
- $E_{T\&L}$  (Errors of both temporal and spatial localization): a detection result in which we first select the ground-truth tubelet from all action classes that has the maximum tubelet IoU with the detection result, then we find they share the same action class, but both temporal IoU and average IoU of spatial bounding boxes lower than a threshold.
- $E_{C\&T\&L}$  (Errors of classification, temporal and spatial localization): a detection result that has some intersection with some ground-truth tubelets, which is different with EN, but wrong action class and both the temporal and average bounding box IoU lower than a threshold.

For each detected tubelet  $d_i$  from a sorted list by descending order of confidence score of class  $c$ .  
 Notation:  $th$ : threshold;  $th_t$ : the square root of  $th$ ;  $th_s$ : the square root of  $th$ ;  $GT(c)$ : set of ground-truths of class  $c$ ;  $dupGT(c)$ : copy of  $GT(c)$ ;  $GT(others)$ : set of all ground-truths that not in class  $c$ ;  $GT(all)$ : set of all ground-truths;  $T\_IoU$ : the temporal domain IoU;  $S\_IoU$ : the average of the IoU between the overlapping frames;  $tubelet\_IoU$ :  $T\_IoU * S\_IoU$ .

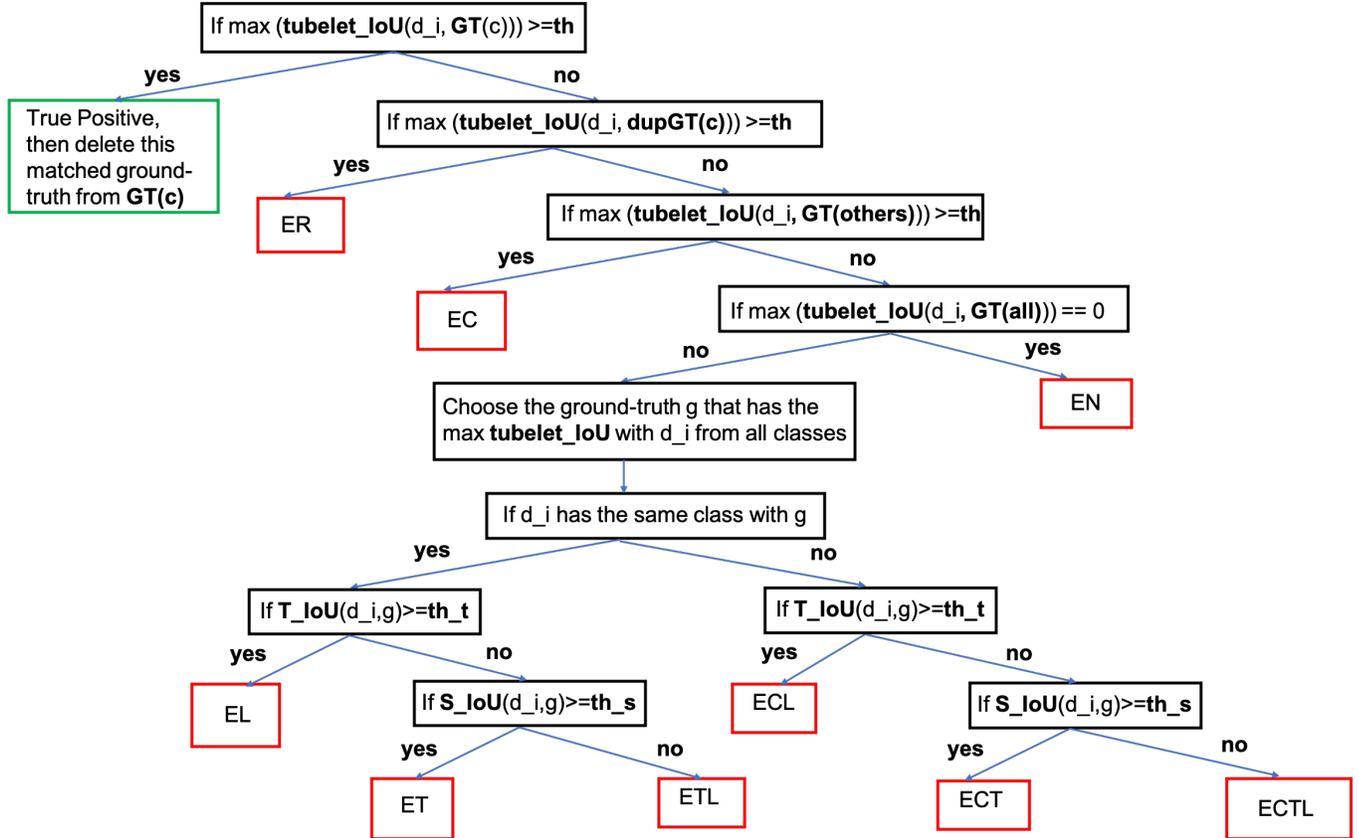


Figure 8. Error Tree

- $E_M$  (Errors of missed detections): ground truth tubelets that have not been matched by any detection results.

## D.2 More Visualization of Error Analysis

As shown in Figure 9, we collect more visualizations of MOC(K=11) as a supplementary of Figure 7 in our paper.

## D.3 Confusion Matrix

We draw the confusion matrices of the predictions which are classified into  $AP$  and  $E_C$  in Figure 10. We observe that the aerobic performs best because its categories relate only to individual actors. Actions having similar motions but different spatio-temporal contexts tend to confuse. For example, 1) drive vs. dribble in basketball, drive emphasizes on breaking through defender and being closer to the basket, which needs to model person-person interaction and spatial localization; 2) through ball vs. pass in football, through ball will break through the opponent’s line of defense and

be passed in front of the teammate, which needs long-term temporal modeling and reasoning. 3) offensive rebound vs. defensive rebound, the difference is whether the offensive player or defensive player gains control of the ball; 4) defend vs. protect in volleyball, we need to focus on whether the ball was blocked back or was spiked by an opponent several frames earlier.

## Appendix E: Annotation Documentation

### E.1 Aerobic Gymnastics

There are four groups of difficulty elements in aerobic gymnastics, namely dynamic strength, static strength, jumps & leaps, and balance & flexibility. We pick out 21 elements to form the aerobic categories of our *MultiSports*. The following is a detailed definition of these categories, a simplified version of the definition in [8].

**Group A: Dynamic Strength.** All elements in Group A ending in a split position, must have both hands on each

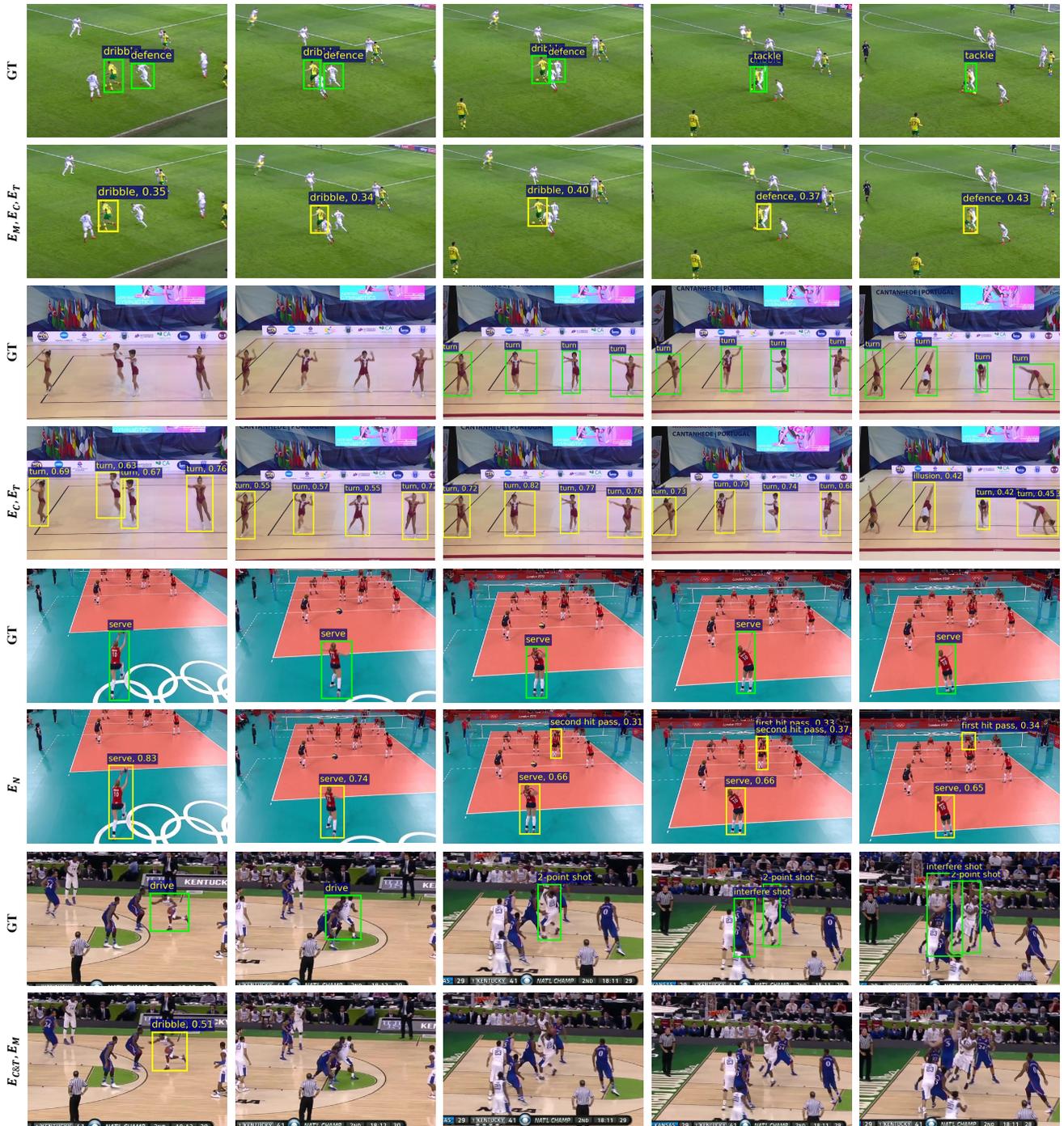
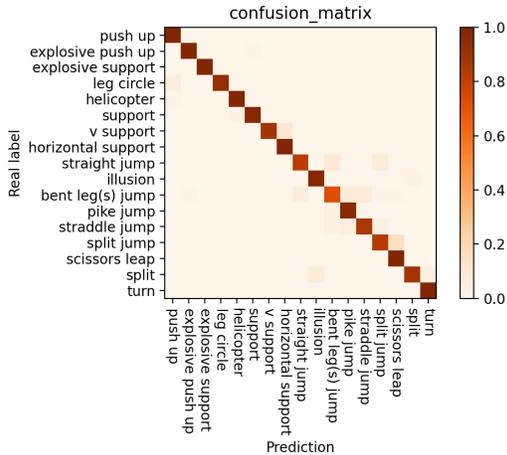


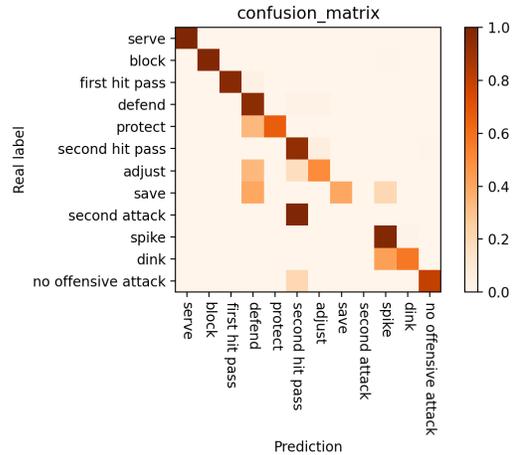
Figure 9. More detailed visualizations on our *MultiSports* dataset with our novel error categories of video-mAP. Green boxes are the ground-truths. Yellow boxes are the detections. 1st and 2nd row:  $E_M$ : missed detection of defence;  $E_C$ : tackle is misclassified as defence;  $E_T$ : dribble has inaccurate action ending boundary. 3rd and 4th row:  $E_C$ : turn is misclassified as illusion in the last picture in 4th row;  $E_T$ : turn has inaccurate action boundary. 5rd and 6th row:  $E_N$ : detection results contain that athletes actually doing none of sports actions but the model identifies first hit pass and second hit pass for them. 7rd and 8th row:  $E_C \& T$ : drive is misclassified as dribble and also has inaccurate action boundary;  $E_M$ : missed detections of interfere shot and 2-point shot.

side of the body on the floor.

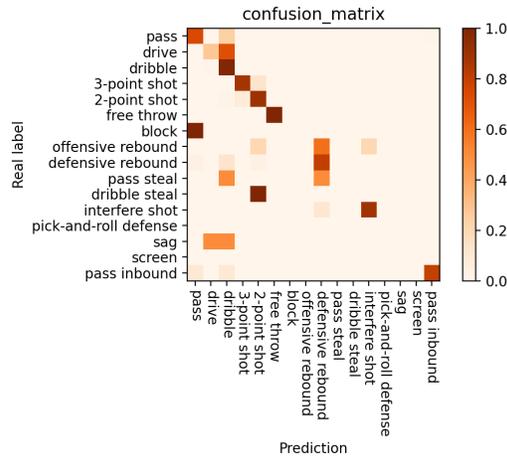
- **Push up:** Starting and/or finishing: one or both hands are in contact with the floor, elbows extended. Shoul-



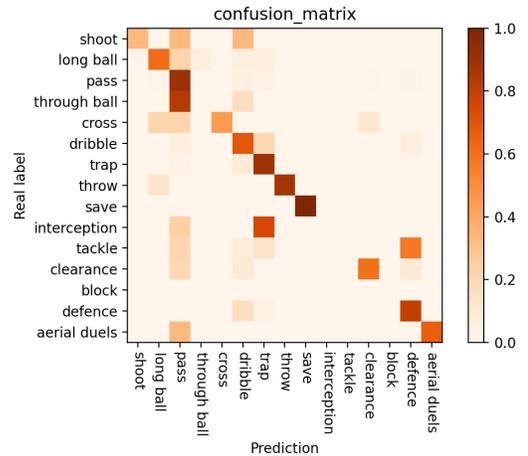
(a) Aerobic



(b) Volleyball



(c) Basketball



(d) Football

Figure 10. Confusion Matrix of SlowFast Det. on different sports.

ders must be parallel to the floor; head in line with the spine and pelvis tucked with abdominal muscles contracted. Flexion of elbows: All push-ups must have, at the end of the downwards phase, a maximum distance of 10cm from the chest to the floor. The downward and/or the upward phase of a push up must be controlled with shoulders parallel to the floor. Lateral and Hinge push up, 4 phases have to be shown. Wenson push up: one leg on the upper part of the arm (Triceps) of the same side.

- **Explosive push up:** 1) A Frame: Pike position in the airborne phase (60° between trunk and legs). 2) Cut: While airborne, the legs straddle sideways and forward to land extended in rear support, feet lifted off the floor during the skill.
- **Explosive support:** Back support on the floor, back parallel to the floor, extending the legs upward and for-

ward with a flight phase. Impulse from High V support position, airborne phase and landing to push up or split position.

- **Leg circle:** The starting position must be from free front support on both hands; the hips must be lifted and extended during the full rotation. Feet may not touch the floor before the completion of the circle. 1) Leg circle: the hips must be lifted and extended. 2) Flair: legs straddle, the hips must be lifted and extended during the full rotation. Feet may not touch the floor before the completion of the circle.
- **Helicopter:** After alternative leg circles, legs close to the chest, body alignment on the upper back (feet off the floor). The legs are extended upward and forward. ½ twist initiated from the feet is made to land in push up or wenson or split.

**Group B: Static Strength.** These elements demonstrate

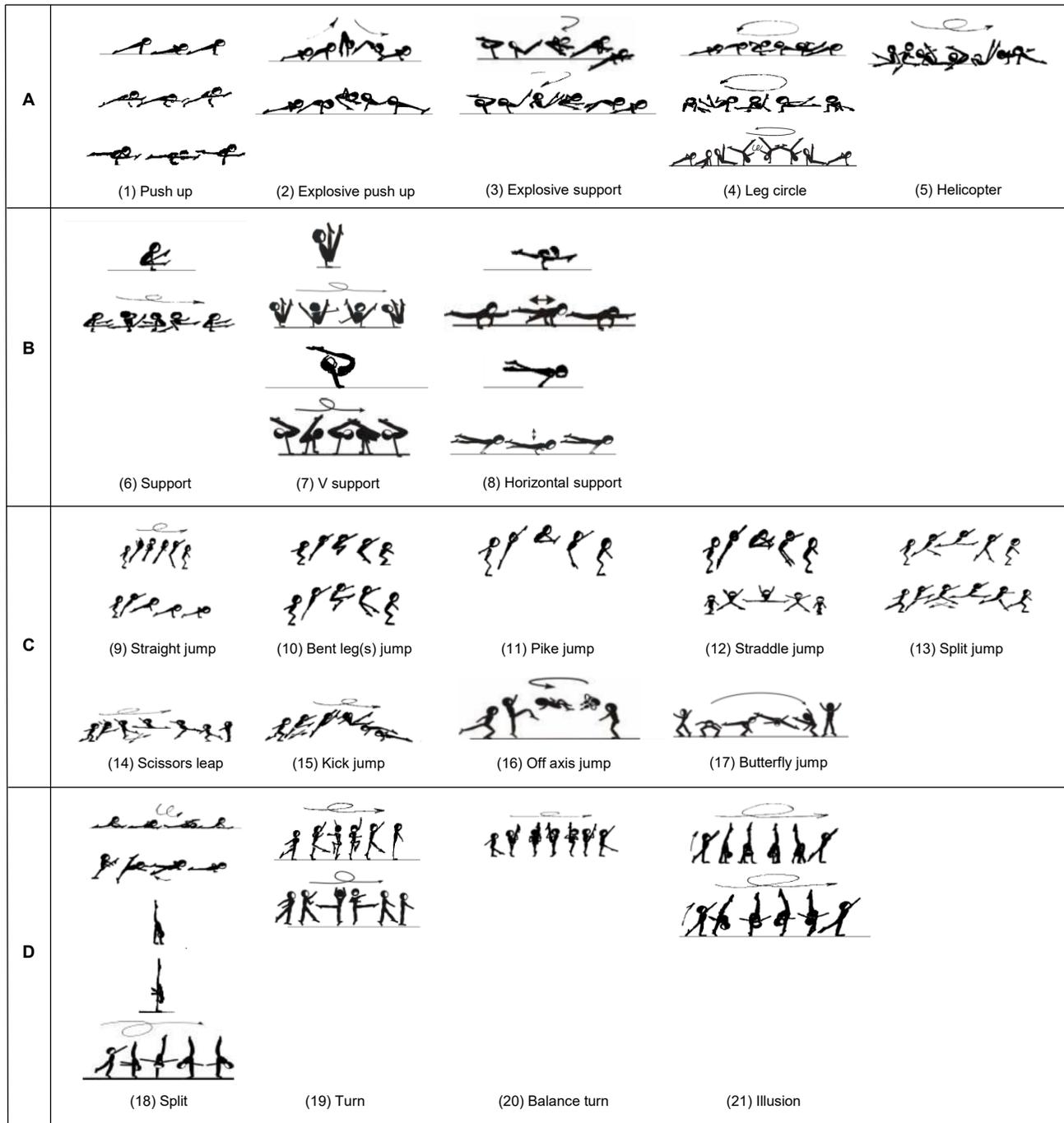


Figure 11. Diagrams of each difficulty element in aerobic gymnastics.

isometric strength and must be held for 2 seconds. In the case of turns in support, the support must be held for 2 seconds either at the start, during or end of the turn. The body is fully supported by one or both arms and only the hands are in contact with the floor. Feet and/or hips must not touch the floor during the whole skill. While in support, the hands must be flat on the floor.

- **Support:** 1) Straddle support: Legs must be straight parallel to the floor in Straddle position (90° minimum). 2) L support: Legs must be straight together and parallel to the floor.
- **V support:** 1) Straddle V support: Hips are flexed and legs straddled 90° open and vertical, minimum width 90°. 2) V support: Hips are flexed and legs are together

vertical. 3) High V support: The back is parallel to the floor.

- **Horizontal support:** 1) Wenson support: the body is extended parallel to the floor, one leg supported on the upper part of the Triceps. 2) Planche: the body is supported on both hands with straight arms, not more than 20° above parallel.

**Group C: Jumps & Leaps.** All jumps and leaps must demonstrate explosive power and maximum amplitude. All jumps that can be performed from 1 foot or two feet will be considered as the same element and will receive the same value. This applies also for landing. Take off preparation: head, shoulder, chest, hips, knees, feet must in the same direction. Body shape while airborne must be clearly recognizable. Body and legs must be tight and straight, with head in line with the spine. **Landing Positions:** 1) Standing: Landing on one foot or two feet must be in a vertical position, with bend leg(s) before finishing in perfect alignment. 2) In push up: both hands and supporting feet must land at the same time in a controlled manner. 3) In split: must land from airborne phase to split form with both hands on each side of the body on the floor. 4) In frontal split: must land from airborne phase to frontal split form, both hands in front of the body.

- **Straight jump:** The body is in extended alignment, the pelvis is fixed – 2 different kinds of jumps and leaps: 1) Vertical: All air turns, Free fall. 2) Vertical to Horizontal: Gainer.
- **Bent leg(s) jump:** 1) Tuck: Both legs are lifted close to the chest with knees bent. 2) Cossack: After takeoff, the body shows a pike shape, legs together parallel to the floor or higher, one leg straight, one leg bent. The angle between the trunk and legs: not be more than 60°. The angle at the knee joint may not be more than 60°.
- **Pike jump:** After takeoff, the body shows a pike shape, legs together and straight, parallel to the floor or higher. The angle between the trunk and legs may not be more than 60°.
- **Straddle jump:** 1) Straddle: Legs are lifted in straddle position (minimum 90° angle), parallel to the floor or higher, arms and trunk extended over them. The angle between the trunk and legs may not be more than 60°. 2) Frontal split: Legs are fully abducted laterally (right and left) frontal (180°) with the upright upper body.
- **Split jump:** 1) Split: Legs are fully stretched front and back in sagittal split (180°) with the upright upper body. 2) Switch: After takeoff, the leading leg must be parallel to the floor and switch with the rear leg to show a split (180°) in the air.

- **Scissors leap:** The leading leg must be parallel to the floor and switches forward with 1/2 turn (180°).
- **Kick jump:** The leading leg must be parallel to the floor and switches forward.
- **Off axis jump:** A one-foot take off, kicking the free leg (bend or straight) upward diagonally. While airborne, the body inclines backward to be out of axis with a longitudinal rotation(s) in tuck or straight position, arms close to the chest. Landing in 1 foot/feet together or in split.
- **Butterfly jump:** A one-foot take off, kicking the free leg backward to lift the body upward. While airborne, legs fly open in straddle (or feet together) with the body in a horizontal position (with or without longitudinal rotation(s)). Landing on one leg.

#### **Group D: Balance & Flexibility.**

- **Split:** Legs must be straight, in line, showing 180°. In Vertical Split: supporting leg must be in vertical position.
- **Turn:** All exercises requiring turns must demonstrate complete rotations on the ball of the foot. Turns are completed when the heel of the turning foot touches the floor.
- **Balance turn:** A Balance turn where one leg is lifted to either in sagittal or frontal balance and is supported by one hand.
- **Illusion:** Starting position of illusion: head, shoulder, chest, hips, knees, toes must be in alignment. A full split (180°) must be shown during the movement.

For the temporal definition, strictly follow the diagrams in Figure 11 (quoted from [8]) to determine the starting and ending of actions, except for the following situations: when an athlete's action is not in place or is completely blocked by other athletes.

#### **E.2 Volleyball**

- **Serve:** Send the ball over the net from behind the end line to start a new round. **Start:** The ball leaves the player's hand. **End:** If the player takes off, any foot touches the ground. Otherwise, the upper arm of the serving arm is below the horizontal plane.
- **Block:** Deflect the ball coming from an attacker on the net. The one that doesn't take off is not considered a block. The one that takes off but doesn't touch the ball is considered a block. **Start:** Any foot leaves the ground. **End:** Any foot touches the ground.

- **First Hit Pass:** Receive the serve. The player can receive the ball overhand, one-hand or underhand. **Start:** If overhand, the player raises any hand over the chest. If one-hand, the player's arm begins to stretch out. If underhand, the player begins to hold hands together. **End:** If overhand, the player puts any hand below the chest. If one-hand, the player's hitting-ball arm relaxes. If underhand, the player's hands loose.
- **Defend:** Receive the ball from the opposite side except for the serve. The player can receive the ball overhand, one-hand or underhand. **Start:** If overhand, the player raises any hand over the chest. If one-hand, the player's arm begins to stretch out. If underhand, the player begins to hold hands together. **End:** If overhand, the player puts any hand below the chest. If one-hand, the player's hitting-ball arm relaxes. If underhand, the player's hands loose.
- **Protect:** Receive the ball returned by the block. The player can receive the ball overhand, one-hand or underhand. **Start:** If overhand, the player raises any hand over the chest. If one-hand, the player's arm begins to stretch out. If underhand, the player begins to hold hands together. **End:** If overhand, the player puts any hand below the chest. If one-hand, the player's hitting-ball arm relaxes. If underhand, the player's hands loose.
- **Second Hit Pass:** The second overhand pass to organize the offense. **Start:** The player raises any hand over the chest. **End:** The player puts any hand below the chest.
- **Adjust:** For the second touch, due to the inadequacy of first hit pass or defending or protecting, the player has to adjust the ball underhand to organize offense. **Start:** The player begins to hold hands together. **End:** The player's hands loose.
- **Save:** Due to the poor first hit pass or defending or protecting, the route of the ball changes dramatically. The actor can't second hit pass overhand or adjust underhand to organize offense but uses one hand or both hands to reach the ball to prevent the ball from landing directly. **Start:** If one-hand, the player's arm begins to stretch out. Otherwise, the player begins to hold hands together. **End:** If one-hand, the player's hitting-ball arm relaxes. Otherwise, the player's hands loose.
- **Second Attack:** For the second touch, a direct attack by the setter. **Start:** Any foot leaves the ground. **End:** Any foot touches the ground.
- **Spike:** Slam the ball over the net into the opposing court. **Start:** Any foot leaves the ground. **End:** Any foot touches the ground.

- **Dink:** Lightly tap the ball over the net to an area on the opponent's side of the court that is not being guarded or occupied by a defensive player. **Start:** Any foot leaves the ground. **End:** Any foot touches the ground.
- **No Offensive Attack:** For the second or third touch, the ball is passed over the net non-aggressively, because of the bad first hit pass or defending or protecting. The actor can push the ball overhand, pass the ball underhand or tap the ball from a position below the net with one hand, where the actor doesn't take off. **Start:** If overhand, the player raises any hand over the chest. If underhand, the player begins to hold hands together. Otherwise, the upper arm of hitting the ball arm is above the horizontal plane. **End:** If overhand, the player puts any hand below the chest. If underhand, the player's hands loose. Otherwise, the upper arm of hitting the ball arm is below the horizontal plane.

### E.3 Football

- **Shoot:** Hit the ball in an attempt to score a goal. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the ball-controlling foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the ball-controlling foot touches the ground.
- **Long Ball:** Middle and long distance (over 30 meters) pass. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the ball-controlling foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the ball-controlling foot touches the ground.
- **Pass:** Short distance (within 30 meters) pass. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the ball-controlling foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the ball-controlling foot touches the ground.

- **Through Ball:** A pass that can clearly break through the opponent's line of defense and has a penetrating effect. At least one defensive player is passed. The ball is passed in front of the player's teammate. In other words, the player should pass the ball to where his running teammate is going to be. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the ball-controlling foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the ball-controlling foot touches the ground.
- **Cross:** A medium-to-long-range pass from a wide area of the field towards the centre of the field near the opponent's goal. Provide direct or indirect shooting opportunities for offensive players. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the ball-controlling foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the ball-controlling foot touches the ground.
- **Dribble:** Have control over the ball for a period of time and distance. **Start:** At the first touch with the ball, the ball-controlling foot leaves the ground. **End:** At the last touch with the ball, the ball-controlling foot touches the ground.
- **Trap:** Use effective parts of the body to adjust the ball, including the speed and position of the ball. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the ball-controlling foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the ball-controlling foot touches the ground.
- **Throw:** The player throws the ball from out of the field and the goalkeeper throws the ball. **Start:** Upper arms swing forward. **End:** Upper arms are below the horizontal plane.
- **Save:** The goalkeeper uses his body parts (except his feet) to destroy the ball that is threatening to the goal. **Start:** After the ball is shot, the goalkeeper begins to move. **End:** Any part of the body touches the ground.
- **Interception:** The defensive player consciously destroys the ball on the opponent's pass route. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the interception foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the interception foot touches the ground.
- **Tackle:** The defensive player snatches the ball under the control of the offensive player. **Start:** the tackling foot leaves the ground. **End:** the tackling foot touches the ground.
- **Clearance:** The defensive player destroys the ball in the backfield in order to gain the initiative in time and space. The main difference between clearance and long ball/ball is that long ball/ball aims at some player but clearance is aimless. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the clearance foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the clearance foot touches the ground.
- **Block:** Intentionally destroy the opponent's threatening shot or block the opponent's shooting angle. The goalkeeper blocked the ball with his foot. Feet, torso and head are all allowed. **Start:** If using torso and head: if the player takes off, any part of the body leaves the ground (such as a foot); if the player does not take off, the player stands firmly and prepares to touch the ball. If using feet, the blocking foot leaves the ground. **End:** If using torso and head: if the player takes off, any part of the body touches the ground (such as a foot); if the player does not take off, stand firmly after touching the ball. If using feet, the blocking foot touches the ground.
- **Defence:** The defensive player approaches the player, of whom the ball is under the control, to make restriction and interference. **Start:** The defender is shorter

than 1.2 meters from the offensive player who is controlling the ball. **End:** 1) the offensive player passes the ball out or the ball is gained by other defensive players. 2) this defender begins to tackle. 3) this defender is longer than 1.2 meters from the offensive player who is controlling the ball.

- **Aerial duels:** Two or more people compete for the high-altitude ball in order to obtain the ball, where all people are annotated. If the player does not take off, it is not considered aerial duels. Note that the player who has an obvious purpose, such as clearance and pass, is annotated that action. **Start:** Any part of the body leaves the ground (such as a foot). **End:** Any part of the body touches the ground (such as a foot).

#### E.4 Basketball

- **Pass:** The player moves the ball to the teammate. **Start:** The player begins to push the ball outwards with his arms. **End:** The ball leaves both hands of the player.
- **Drive:** The player, who controls the ball, gets rid of the defense by passing the defensive player or stopping suddenly during the movement. The aim is to get closer to the basket and create a space that is conducive to shooting. The next step is usually to shoot, layup, or pass the ball to teammates. **Start:** At the first touch with the ball, the hand presses the ball down. **End:** At the last touch with the ball, the ball is bounced into the hand.
- **Dribble:** The player slaps the ball bounced from the ground continuously while on the spot or on the move. **Start:** At the first touch with the ball, the hand presses the ball down. **End:** At the last touch with the ball, the ball is bounced into the hand.
- **3-point Shot:** Shoot from beyond the three-point line. **Start:** The player raises the shooting hand over the chest. **End:** If the player takes off, any part of the body touches the ground (such as a foot). Otherwise, the player puts the shooting hand below the chest.
- **2-point Shot:** Shoot from within the three-point line. **Start:** The player raises the shooting hand over the chest. **End:** If the player takes off, any part of the body touches the ground (such as a foot). Otherwise, the player puts the shooting hand below the chest.
- **Free Throw:** Unopposed attempts to score points by shooting from behind the free throw line. **Start:** The player raises the shooting hand over the chest. **End:** If the player takes off, any part of the body touches the ground (such as a foot). Otherwise, the player puts the shooting hand below the chest.
- **Block:** When the offense shoots, the defender successfully knocks the ball out as the ball goes up. **Start:** The defender raises the blocking hand over the chest. **End:** If the defender takes off, any part of the body touches the ground (such as a foot). Otherwise, the defender puts the blocking hand below the chest.
- **Offensive Rebound:** After a missed shot, the two sides compete for a rebound and the offensive player grabs it. **Start:** The player raises the grabbing-ball hand over the chest. **End:** If the player takes off, any part of the body touches the ground (such as a foot). Otherwise, the player catches the ball firmly.
- **Defensive Rebound:** After a missed shot, the two sides compete for a rebound and the defensive player grabs it. **Start:** The player raises the grabbing-ball hand over the chest. **End:** If the player takes off, any part of the body touches the ground (such as a foot). Otherwise, the player catches the ball firmly.
- **Pass Steal:** The defensive player intercept the ball in the process of passing, which is not under the control of offensive player. **Start:** The defender's stealing-ball hand begins to stretch out. **End:** 1) Route of the ball changes; 2) The defender catches the ball firmly.
- **Dribble Steal:** The defensive player steals the ball under the control of offensive player. **Start:** The defender's stealing-ball hand begins to stretch out. **End:** The ball is out of the control of the offensive player who had the control.
- **Interfere Shot:** The defender interferes with the shot but does not touch the ball. **Start:** The defender raises the interfering hand over the chest. **End:** If the defender takes off, any part of the body touches the ground (such as a foot). Otherwise, the defender puts the interfering hand below the chest.
- **Pick-and-roll Defense:** In pick-and-roll, the defender of the offensive ball-controlling player is blocked by the teammate of this offensive player. **Start:** The defender has physical contact with the offensive screening player. **End:** The defender does not have physical contact with the offensive screening player.
- **Sag:** The defender gives up the offensive player he is responsible for and turns to defend the offensive ball-controlling player. **Start:** The defender consciously approach the offensive ball-controlling player. **End:** 1) the offensive player passes or shoots the ball; 2) this defender is broken through; 3) this defender gives up.
- **Screen:** In pick-and-roll, the offensive player uses his body to set a pick for his ball-controlling teammate.

**Start:** Both feet of the offensive player touches the ground. **End:** Any foot of the offensive player is ready to leave the ground completely. Small range movement is not considered the end.

- **Pass-inbound:** The player passes the ball from the boundary lines to restart the play. **Start:** The player begins to push the ball outwards with his arms. **End:** The ball leaves both hands of the player.
- **Save:** The player gets back the ball that is about to go out of bounds. **Start:** The player begins to push the ball outwards with his arms. **End:** The ball leaves both hands of the player.
- **Jump Ball:** A method used to begin or resume the play. Two opposing players attempt to gain control of the ball after an official tosses it into the air between them, where both players are annotated. **Start:** The player raises the grabbing-ball hand over the chest. **End:** Any part of the body touches the ground (such as a foot).

## References

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *CoRR*, abs/1609.08675, 2016. 3
- [2] Moshe Blank, Lena Gorelick, Eli Shechtman, Michal Irani, and Ronen Basri. Actions as space-time shapes. In *ICCV*, pages 1395–1402, 2005. 3
- [3] João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *CVPR*, pages 4724–4733, 2017. 1, 3, 5, 9, 10
- [4] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Mmdetection: Open mmlab detection toolbox and benchmark. *CoRR*, abs/1906.07155, 2019. 10
- [5] MMAAction2 Contributors. Openmmlab’s next generation video understanding toolbox and benchmark. <https://github.com/open-mmlab/mmaaction2>, 2020. 6, 9
- [6] Kellie Corona, Katie Osterdahl, Roderic Collins, and Anthony Hoogs. MEVA: A large-scale multiview, multimodal video dataset for activity detection. In *WACV*, pages 1059–1067, 2021. 3, 9
- [7] Yutao Cui, Cheng Jiang, Limin Wang, and Gangshan Wu. Fully convolutional online tracking. *CoRR*, abs/2004.07109, 2020. 2, 4
- [8] Federation Internationale de Gymnastique. Aerobic gymnastics-code of points. *FIG Aerobic Gymnastics FIG Executive Committee*, 2017. 2, 3, 11, 15
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. 9
- [10] Mark Everingham, S. M. Ali Eslami, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *Int. J. Comput. Vis.*, pages 98–136, 2015. 9
- [11] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, pages 6201–6210, 2019. 3, 6, 7, 9, 10
- [12] Rohit Girdhar, João Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *CVPR*, pages 244–253, 2019. 3
- [13] Georgia Gkioxari and Jitendra Malik. Finding action tubes. In *CVPR*, pages 759–768, 2015. 3
- [14] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fründ, Peter Yianilos, Moritz Mueller-Freitag, Florian Hoppe, Christian Thureau, Ingo Bax, and Roland Memisevic. The “something something” video database for learning and evaluating visual common sense. In *ICCV*, pages 5843–5851, 2017. 3
- [15] Chunhui Gu, Chen Sun, David A. Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A video dataset of spatio-temporally localized atomic visual actions. In *CVPR*, pages 6047–6056, 2018. 1, 3, 4, 5, 6, 9
- [16] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015. 3
- [17] Grant Van Horn and Pietro Perona. The devil is in the tails: Fine-grained classification in the wild. *CoRR*, abs/1709.01450, 2017. 5
- [18] Rui Hou, Chen Chen, and Mubarak Shah. Tube convolutional neural network (T-CNN) for action detection in videos. In *ICCV*, pages 5823–5832, 2017. 3
- [19] Haroon Idrees, Amir Roshan Zamir, Yu-Gang Jiang, Alex Gorban, Ivan Laptev, Rahul Sukthankar, and Mubarak Shah. The THUMOS challenge on action recognition for videos “in the wild”. *Comput. Vis. Image Underst.*, pages 1–23, 2017. 3
- [20] Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J. Black. Towards understanding action recognition. In *ICCV*, pages 3192–3199, 2013. 1, 2, 3, 5, 6
- [21] Vicky Kalogeiton, Philippe Weinzaepfel, Vittorio Ferrari, and Cordelia Schmid. Action tubelet detector for spatio-temporal action localization. In *ICCV*, pages 4415–4423, 2017. 3, 6
- [22] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Fei-Fei Li. Large-scale video classification with convolutional neural networks. In *CVPR*, pages 1725–1732, 2014. 3
- [23] Okan Köpüklü, Xiangyu Wei, and Gerhard Rigoll. You only watch once: A unified CNN architecture for real-time

- spatiotemporal action localization. *CoRR*, abs/1911.06644, 2019. 7, 9
- [24] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso A. Poggio, and Thomas Serre. HMDB: A large video database for human motion recognition. In *ICCV*, pages 2556–2563, 2011. 3
- [25] Ang Li, Meghana Thotakuri, David A. Ross, João Carreira, Alexander Vostrikov, and Andrew Zisserman. The ava-kinetics localized human actions video dataset. *CoRR*, abs/2005.00214, 2020. 3, 5
- [26] Dong Li, Zhaofan Qiu, Qi Dai, Ting Yao, and Tao Mei. Recurrent tubelet proposal and recognition networks for action detection. In *ECCV*, pages 306–322, 2018. 3
- [27] Yixuan Li, Zixu Wang, Limin Wang, and Gangshan Wu. Actions as moving points. In *ECCV*, pages 68–84, 2020. 3, 6, 7, 9
- [28] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 936–944, 2017. 9
- [29] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. BMN: boundary-matching network for temporal action proposal generation. In *ICCV*, pages 3888–3897, 2019. 1
- [30] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, pages 740–755, 2014. 9
- [31] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. BSN: boundary sensitive network for temporal action proposal generation. In *ECCV*, pages 3–21, 2018. 1
- [32] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: single shot multibox detector. In *ECCV*, pages 21–37, 2016. 9
- [33] Mathew Monfort, Kandan Ramakrishnan, Alex Andonian, Barry A. McNamara, Alex Lascelles, Bowen Pan, Quanfu Fan, Dan Gutfreund, Rogério Schmidt Feris, and Aude Oliva. Multi-moments in time: Learning and interpreting models for multi-action video understanding. *CoRR*, abs/1911.00232, 2019. 3
- [34] Xiaojiang Peng and Cordelia Schmid. Multi-region two-stream R-CNN for action detection. In *ECCV*, pages 744–759, 2016. 3
- [35] Joseph Redmon and Ali Farhadi. YOLO9000: better, faster, stronger. In *CVPR*, pages 6517–6525, 2017. 9
- [36] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6):1137–1149, 2017. 7
- [37] Mikel D. Rodriguez, Javed Ahmed, and Mubarak Shah. Action MACH a spatio-temporal maximum average correlation height filter for action recognition. In *CVPR*, 2008. 3
- [38] Suman Saha, Gurkirt Singh, Michael Sapienza, Philip H. S. Torr, and Fabio Cuzzolin. Deep learning for detecting multiple space-time action tubes in videos. In *BMVC*, 2016. 3
- [39] Christian Schüldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: A local SVM approach. In *ICPR*, pages 32–36, 2004. 3
- [40] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. Finegym: A hierarchical video dataset for fine-grained action understanding. In *CVPR*, pages 2613–2622, 2020. 3, 5
- [41] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *ECCV*, pages 510–526, 2016. 3
- [42] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, pages 568–576, 2014. 1
- [43] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 9
- [44] Gurkirt Singh, Suman Saha, Michael Sapienza, Philip H. S. Torr, and Fabio Cuzzolin. Online real-time multiple spatiotemporal action localisation and prediction. In *ICCV*, pages 3657–3666, 2017. 3, 6, 7, 9
- [45] Lin Song, Shiwei Zhang, Gang Yu, and Hongbin Sun. Tacnet: Transition-aware context network for spatio-temporal action detection. In *CVPR*, pages 11987–11995, 2019. 3
- [46] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *CoRR*, abs/1212.0402, 2012. 1, 2, 3, 5, 6, 7, 9
- [47] Jiajun Tang, Jin Xia, Xinzhi Mu, Bo Pang, and Cewu Lu. Asynchronous interaction aggregation for action detection. In *ECCV*, pages 71–87, 2020. 3
- [48] Du Tran, Lubomir D. Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, pages 4489–4497, 2015. 1
- [49] Limin Wang, Wei Li, Wen Li, and Luc Van Gool. Appearance-and-relation networks for video classification. In *CVPR*, pages 1430–1439, 2018. 1
- [50] Limin Wang, Yu Qiao, and Xiaoou Tang. Action recognition with trajectory-pooled deep-convolutional descriptors. In *CVPR*, pages 4305–4314, 2015. 1
- [51] Limin Wang, Yu Qiao, Xiaoou Tang, and Luc Van Gool. Actionness estimation using hybrid fully convolutional networks. In *CVPR*, pages 2708–2717, 2016. 3
- [52] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, pages 20–36, 2016. 1
- [53] Philippe Weinzaepfel, Zaïd Harchaoui, and Cordelia Schmid. Learning to track for spatio-temporal action localization. In *ICCV*, pages 3164–3172, 2015. 3, 6
- [54] Philippe Weinzaepfel, Xavier Martin, and Cordelia Schmid. Towards weakly-supervised action localization. *CoRR*, abs/1605.05197, 2016. 3
- [55] Chao-Yuan Wu, Christoph Feichtenhofer, Haoqi Fan, Kaiming He, Philipp Krähenbühl, and Ross B. Girshick. Long-term feature banks for detailed video understanding. In *CVPR*, pages 284–293, 2019. 3, 9, 10
- [56] Jianchao Wu, Zhanghui Kuang, Limin Wang, Wayne Zhang, and Gangshan Wu. Context-aware RCNN: A baseline for action detection in videos. In *ECCV*, pages 440–456, 2020. 3

- [57] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Object tracking benchmark. *IEEE Trans. Pattern Anal. Mach. Intell.*, 37(9):1834–1848, 2015. 4
- [58] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, pages 5987–5995, 2017. 9
- [59] Huijuan Xu, Abir Das, and Kate Saenko. R-C3D: region convolutional 3d network for temporal activity detection. In *ICCV*, pages 5794–5803, 2017. 1
- [60] Xitong Yang, Xiaodong Yang, Ming-Yu Liu, Fanyi Xiao, Larry S. Davis, and Jan Kautz. STEP: spatio-temporal progressive learning for video action detection. In *CVPR*, pages 264–272, 2019. 3
- [61] Serena Yeung, Olga Russakovsky, Ning Jin, Mykhaylo Andriluka, Greg Mori, and Li Fei-Fei. Every moment counts: Dense detailed labeling of actions in complex videos. *Int. J. Comput. Vis.*, pages 375–389, 2018. 3
- [62] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *CVPR*, pages 2403–2412, 2018. 9
- [63] Runhao Zeng, Wenbing Huang, Chuang Gan, Mingkui Tan, Yu Rong, Peilin Zhao, and Junzhou Huang. Graph convolutional networks for temporal action localization. In *ICCV*, pages 7093–7102, 2019. 1
- [64] Hang Zhao, Antonio Torralba, Lorenzo Torresani, and Zhicheng Yan. HACS: human action clips and segments dataset for recognition and temporal localization. In *ICCV*, pages 8667–8677, 2019. 3, 5
- [65] Jiaojiao Zhao and Cees G. M. Snoek. Dance with flow: Two-in-one stream action detection. In *CVPR*, pages 9935–9944, 2019. 3
- [66] Yue Zhao, Yuanjun Xiong, Limin Wang, Zhirong Wu, Xiaoou Tang, and Dahua Lin. Temporal action detection with structured segment networks. In *ICCV*, pages 2933–2942, 2017. 1