# DeepSportradar-v1: Computer Vision Dataset for Sports Understanding with High Quality Annotations

Gabriel Van Zandycke*
g.vanzandycke@sportradar.com
Sportradar AG
Switzerland

Vladimir Somers*
v.somers@sportradar.com
Sportradar AG
Switzerland

Maxime Istasse*
m.istasse@sportradar.com
Sportradar AG
Switzerland

Carlo Del Don
c.deldon@sportradar.com
Sportradar AG
Switzerland

Davide Zambrano*
d.zambrano@sportradar.com
Sportradar AG
Switzerland

## ABSTRACT

With the recent development of Deep Learning applied to Computer Vision, sport video understanding has gained a lot of attention, providing much richer information for both sport consumers and leagues. This paper introduces DeepSportradar-v1, a suite of computer vision tasks, datasets and benchmarks for automated sport understanding. The main purpose of this framework is to close the gap between academic research and real world settings. To this end, the datasets provide high-resolution raw images, camera parameters and high quality annotations. DeepSportradar currently supports four challenging tasks related to basketball: ball 3D localization, camera calibration, player instance segmentation and player re-identification. For each of the four tasks, a detailed description of the dataset, objective, performance metrics, and the proposed baseline method are provided. To encourage further research on advanced methods for sport understanding, a competition is organized as part of the MMSports[1] workshop from the ACM Multimedia 2022 conference, where participants have to develop state-of-the-art methods to solve the above tasks. The four datasets, development kits and baselines are publicly available[2].

## CCS CONCEPTS

• **Computing methodologies → Computer vision tasks**; *Scene understanding*; *Image segmentation*; *Object detection*; *Interest point and salient region detections*; *Reconstruction*; *Matching*.

## KEYWORDS

challenge, competition, dataset, deep learning, computer vision, image understanding, sports, basketball, ball 3D localization, camera calibration, instance segmentation, person re-identification, reid

---

*Authors contributed equally to this work.

[1]http://mmsports.multimedia-computing.de/mmsports2022/index.html

[2]https://github.com/DeepSportRadar

## 1 INTRODUCTION

Individual and professional sports have always had a strong impact on the economic, political, and cultural aspects of our society. When only considering the economical side, this impact is likely going to increase as the global sports market size, including services and goods from sport entities, is expected to grow from $354.96 billion in 2021 to $707.84 billion in 2026[3]. The online live-streaming market alone is going to increase in value from $18.12bn in 2020 to $87.34bn in 2028[4]. A strong contribution to this growth is provided by the recent and rapid technological advancement which has changed the way people watch and enjoy sports. Indeed, Computer Vision (CV) and the recent developments in Deep Learning (DL) [25] provide the opportunity to extract meaningful information from live streamed events resulting in a much richer experience for both consumers and leagues.

The ability of DL-based solutions to be useful and reliable in real world applications strongly depends on the quantity and quality of data on which the model has been trained on in the first place [1]. Specifically for sports, the different disciplines and conditions make them unique in terms of the problems the model has to face. Moreover, the quality of the annotations often determines the model performances overall. In the past few years, the SoccerNet [15, 11, 8, 9] datasets have received increasing attention for the amount of data and the benchmark models provided to the CV community. However, two main issues pertain such initiative: first, considering only soccer as representative of all sports does not allow to extend results to other domains; secondly, and more importantly, the SoccerNet annotations are created out of broadcast videos which bring a series of concerns. These concerns include: a limited spatial and temporal coverage of the game due to, on one side, the frequent camera movements which return a subset of the field and, on the

---

[3]https://www.thebusinessresearchcompany.com/report/sports-global-market-report

[4]https://www.verifiedmarketresearch.com/product/sports-online-live-video-streaming-market

other, the replays or advertisements which interrupt the live stream; a lower image resolution with respect to the original sensor; no access to camera parameters or position; and the overlaying graphics such as scores, teams name, advertisements, game statistics, that obstruct the image. In summary, broadcast video annotations remain distant from the actual sensors and tools used to record the game. In conclusion, while initiatives as SoccerNet represent a strong and valid tool for the computer vision community, the introduction of an high quality dataset with available raw images and camera parameters, will help closing the gap between academic research and real world settings.

This paper introduces two datasets in the basketball domain and four different CV tasks each associated with a task-specific dataset extracted from the first two. The data and annotations are provided by SynergySports[5], a division of Sportradar[6], and have been recorded with the Keemotion/Synergy Automated Camera System™. The proposed four tasks are:

- **Ball 3D localization in calibrated scenes.** This task tackles the estimation of ball size on basketball scenes given the oracle ball position.
- **Camera calibration.** This task aims at predicting the camera calibration parameters from images taken from basketball games.
- **Player instance segmentation.** This task deals with the segmentation of individual humans (players, coaches and referees) on the basketball court.
- **Player re-identification.** In this task, the objective is to re-identify basketball players across multiple video frames captured from the same camera viewpoint at various time instants.

Moreover, a competition around the four tasks has been organized and results will be presented at the 5th International ACM Workshop on Multimedia Content Analysis in Sports[7].

A toolkit is provided for each task containing data, annotations and metrics. Moreover, a baseline for each task has been added which serves the purpose of providing an example to consume the data, and as a starting point for the challenge participants' solutions. The next section explains the original datasets, while the subsequent sections will describe each task in detail.

## 2 DATASETS

The four tasks are built on two different datasets. The DeepSport dataset—a multi-labels dataset containing ball 3D annotations, image calibration data and human segmentation masks—is used for the ball 3D localization, court calibration and instance segmentation tasks. The DeepSportradar player ReID dataset is used for the players re-identification task only and will be described further in Section 3.4. Both datasets were acquired during professional basketball matches with the Keemotion/Synergy Automated Camera System™ that offers a sideline view of the court from a position aligned with the court center line. Images were fully annotated and reviewed by human operators, leading to high quality labels, and are made freely available to the research community by Sportradar.

---

[5]https://synergysports.com/

[6]https://sportradar.com/

[7]http://mmsports.multimedia-computing.de/mmsports2022/index.html

**Table 1: The DeepSport dataset was captured in 15 different arenas and three of them are kept for the testing set. It features a variety of angle of views, distance to the court and image resolution.**

| Arena label | Arena name (City) | Number of items |
|---|---|---|
| KS-FR-STCHAMOND | Halle André Boulloche (Saint-Chamond) | 12 |
| KS-FR-FOS | HdS Parsemain (Fos-sur-Mer) | 23 |
| KS-FR-STRASBOURG | Rhénus Sport (Strasbourg) | 8 |
| KS-FR-VICHY | PdS Pierre Coulon (Vichy) | 9 |
| KS-FR-NANTES | la Trocardière (Nantes) | 20 |
| KS-FR-BOURGEB | Ekinox (Bourg-en-Bresse) | 12 |
| KS-FR-GRAVELINES | Sportica (Gravelines) | 129 |
| KS-FR-MONACO | Salle Gaston Médecin (Monaco) | 9 |
| KS-FR-POITIERS | Stade Poitevin (Poitiers) | 5 |
| KS-FR-NANCY | PdS Jean Weille de Gentilly (Nancy) | 40 |
| KS-FR-LEMANS | Antarès (Le Mans) | 16 |
| KS-FR-BLOIS | Le Jeu de Paume (Blois) | 39 |
| **KS-FR-CAEN** | **PdS de Caen (Caen)** | **31** |
| **KS-FR-ROANNE** | **HdS Andre Vacheresse (Roanne)** | **3** |
| **KS-FR-LIMOGES** | **PdS de Beaublanc (Limoges)** | **8** |

## DeepSport dataset

Hereafter, the multi-labels DeepSport dataset, used for three tasks, is described. Originally introduced in [40] with only ball annotations, it was later supplemented with new data and additional annotations. It is now made available publicly on the Kaggle platform [38].

*Description.* The dataset is a collection of *raw-instants*: sets of images captured at the same instant by an array of cameras covering a panorama of the sport field. It features only in-game basketball scenes. Figure 2 shows a *raw-instant* from a two cameras setup. In the DeepSport dataset, camera resolutions range from 2Mpx to 5Mpx. As illustrated in Figure 1, the resulting images have a definition varying between 80px/m and 150px/m, depending on camera resolution, sensor size, lens focal-length and distance to the court.

*Origin.* The dataset was captured in 15 different basketball arenas, each identified by a unique label, during 37 professional games of the French league LNB-Pro A. Figure 1 depicts a cross section of a basketball court where the camera setup height and distance to the court is shown for each arena.

*Split.* The dataset is split in 3 subsets: training, validation and testing. The testing set contains all images from arena labels KS-FR-CAEN, KS-FR-LIMOGES and KS-FR-ROANNE. If not otherwise specified, the remaining instants are split in 15% for the validation set and 85% for the training set. A mapping between arena labels and arena names is given in Table 1 and provides the amount of instants for each arena. An additional challenge set, introduced for the competition, is composed of 35 additional similar instants. They come from a new set of arenas and labels will remain secret.

*Annotations.* The cameras used to capture the *raw-instants* are calibrated, which means that intrinsic and extrinsic parameters are known. The ball 3D annotation was obtained by leveraging the calibration data and clicking two points in the image space: the
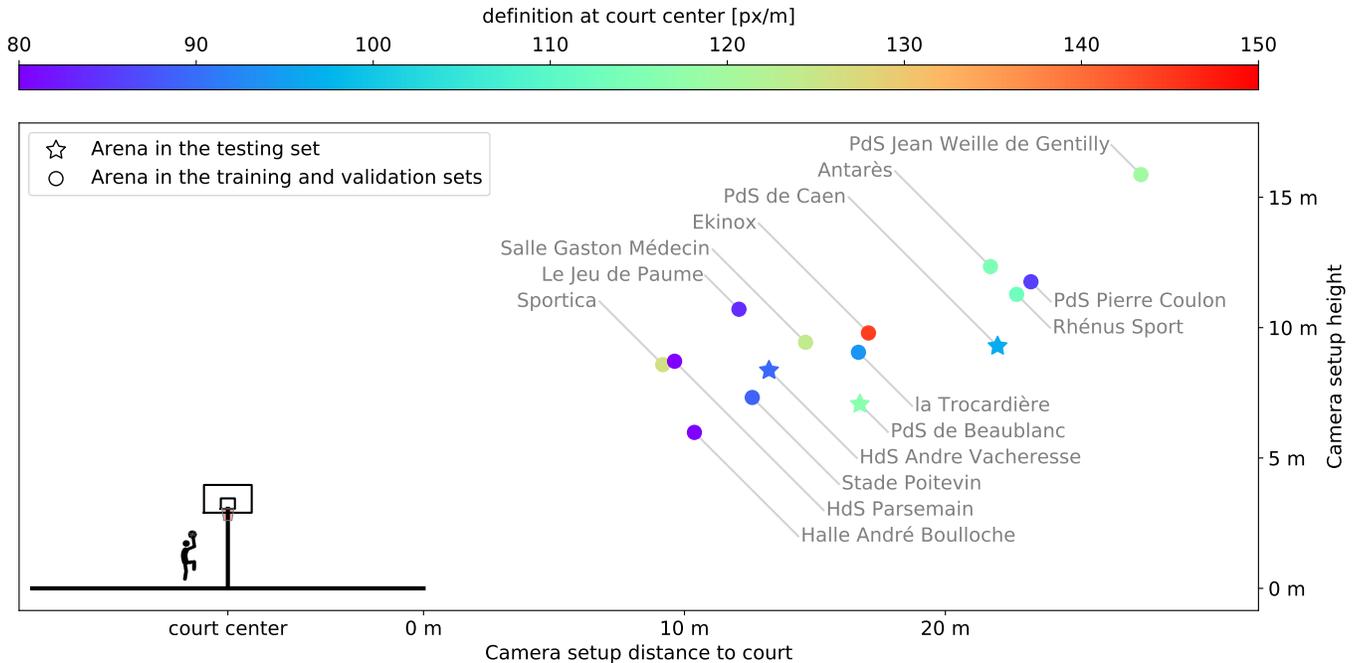
**Figure 1: Cross section showing camera setup height from the ground and distance to the court of the different arenas in which images were acquired. The camera definition depends on camera resolution, sensor size, lens focal length and camera setup distance to the court.**



**Figure 2: A *raw instant* captured by the Keemotion/Synergy Automated Camera System with a two cameras setup.**

ball center and its vertical projection to the ground. This process is described and validated in [39]. The contouring of humans lying near the court were annotated on each image individually following [34], with a special care given to occlusions.

## 3 THE TASKS

This section describes the four tasks, explaining the dataset splits, the metrics used for evaluating and compare results for the competition and the baselines provided to the users.

### 3.1 Ball 3D localization

Automated ball 3D localization in team sport scenes has important applications like assisting referees or feeding game analytics. In sports like basketball where the ball is often occluded and in players hands, an image based approach is required to fill the gap between the trajectories proposed by a ballistic approach. Hence, this task

aims at localizing the ball in 3D on basketball scenes, from a single calibrated image. This problem can be solved by both detecting the ball center and estimating its size in the image space. Indeed, the 3D localization can be recovered using camera calibration information and the knowledge of the real ball size [39]. Since ball detection has been largely studied [22], this task focuses on 3D localization given oracle ball positions. Hence, the task consists in estimating the ball diameter in pixels in the image space from an image patch centered on the ball.

*3.1.1 Dataset.* The task uses $N \times N$ crops around oracle ball positions from the **DeepSport dataset** presented in Section 2, where $N$ is a parameter. Figure 3 shows samples from the dataset with $N = 128$. As shown in Figure 4, ball diameter ranges from 15px to 35px in the dataset.

*3.1.2 Metrics.* The main metric used to evaluate methods is the mean absolute diameter error (MADE) between the prediction and the ground-truth. In addition, the mean absolute projection error (MAPE) and the mean absolute relative error (MARE), described in [39], are used. The MAPE measures the error of a vertical projection of positions on the ground plane. The MARE measures the distance error relative to the camera position.

*3.1.3 Baseline and results.* The baseline proposed with this task is the regression model described in [39]. It is composed of a VGG16 [36] features extractor followed by 3 fully connected layers and is supervised with a Huber loss [21].

The baseline is trained with random scaling and random color-gamma data augmentations on image patches with $N = 64$. The

**Figure 3: Samples of** $128 \times 128$ **crops around oracle ball positions. The dataset features many different scenes with different colors, backgrounds and lighting conditions. Ball is often partly occluded or in players hand and suffers from motion blur.**

supervision is performed using Adam optimizer for 100 epochs with an initial learning rate of $10^{-4}$, exponentially decayed by half during two epochs. This decay is applied every 10 epochs starting at epoch 50.

The baseline reaches a MADE of 2.12 pixels. This corresponds to a MAPE of 3.05 meters and a MARE of 10%.

## 3.2 Camera calibration

The Camera calibration task aims at predicting the camera parameters from images taken from basketball games. The toolkit for this task can be found on the main DeepSportRadar GitHub page. Formally, this task objective is to predict the projection matrix, $P_{3\times4}$ that maps a 3D point in homogeneous coordinates (a 4-dimensional vector) to a 2D point in homogeneous coordinates (3-dimensional vector) in the image space[8]. The projection matrix combines intrinsic (sensor and lens) and extrinsic (position and rotation) camera

---

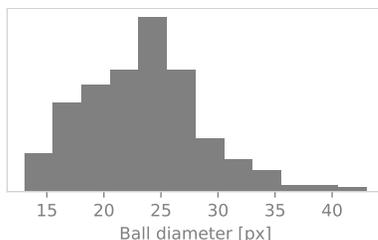[8]https://en.wikipedia.org/wiki/Camera_matrix



**Figure 4: Ball size distribution in the DeepSport dataset.**

parameters as:

$$P := \begin{bmatrix} K_{3\times3}; \mathbf{0}_{3\times1} \end{bmatrix} \begin{bmatrix} R_{3\times3}; T_{3\times1} \\ \mathbf{0}_{1\times3}; 1 \end{bmatrix}$$
$$= K_{3\times3} \begin{bmatrix} R_{3\times3}; T_{3\times1} \end{bmatrix}, \tag{1}$$

where $K$ is the matrix of intrinsic parameters, while $R, T$ are the rotation and translation matrices respectively[9]. For Synergy/Keemotion produced images, the origin of the 3D world is located on the furthest left corner of the basketball court relative to the main camera setup; more precisely in the inner side of the court lines. The unit of length is the centimeter and axis orientation is given by $x$ along the court length, $y$ along the court width and $z$ pointing downward[10]. For simplicity, this task assumes that lenses have no distortion.

The camera calibration parameters are crucial for several CV tasks such as the 3D tracking of players in the field. These parameters can be retrieved on site; however, an automatic method that estimates them is needed when the field and the camera are not accessible anymore. The Camera calibration task falls under the sport-field registration tasks and takes advantage of the known official sizes of the basketball court[11]. Several approaches have been adopted to solve sport-field registrations for different sport domains including tennis, volleyball and soccer [12, 44], relying essentially on keypoints retrieval methods. With the advent of Deep Learning, common approaches tackle the problem as a segmentation task [20, 3, 35, 7]. The baseline introduced for this task adopts this approach.

*3.2.1 Dataset.* This task purpose is to predict the camera calibration parameters from a single frame of a basketball game. The

---

[9]https://ispgroupucl.github.io/calib3d/calib3d/calib.html
[10]https://gitlab.com/deepsport/deepsport_utilities/-/blob/main/calibration.md
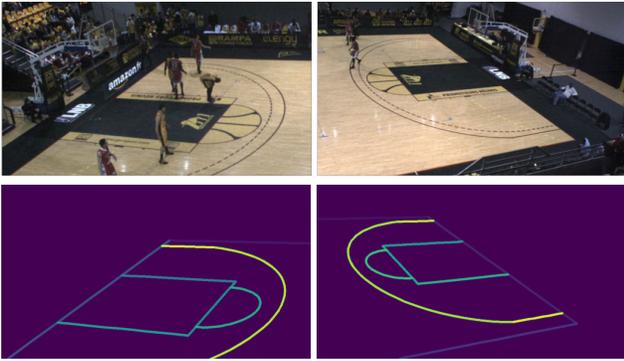[11]https://en.wikipedia.org/wiki/Basketball_court

**Figure 5: Examples of images from the camera calibration task (top). The target camera parameters have been used to generate the court lines (bottom), which can then be used as target for a segmentation model as described for the baseline.**

dataset is made of 728 views (pairs of images and corresponding camera calibration parameters) randomly generated from the **Deep-Sport dataset**. The random view generation process[12] generates a random 3D position within the court and image limits on which the view will be centered. It then samples a pixel density between $\alpha \cdot 20$ px/m and $\alpha \cdot 60$ px/m (see Figure 1), a rotation between $-10°$ and $10°$, and a boolean horizontal flip, to create a crop from the original image of size $\alpha \cdot 480 \times \alpha \cdot 270$. Please note that, the test and challenge sets were provided with $\alpha = 2$, resulting in an output image dimension of $920 \times 540$. The corresponding affine transformation matrix is applied to the intrinsic camera matrix $K$ to produce the camera calibration parameters that correspond to the generated view. These views are divided in train (480), val (164) and test (84) sets. For this challenge, having a validation set on arenas not seen during training is of foremost importance, therefore the arenas of KS-FR-NANTES, KS-FR-BLOIS and KS-FR-FOS are used for the validation set (see Table 1). A final challenge split composed of 84 images is provided for the competition purpose and its camera parameters are kept secret.

A few sport-field registration datasets have been publicly released so far among which the SoccerNet-v2 is the largest [11] (20028 images from 500 games). It is worth noticing that, with our random view generation process, a potentially infinite number of views can be generated from the original dataset images.

*3.2.2 Metrics.* In order to evaluate the proposed methods on this task, the predictions are evaluated based on a Mean Squared Error (MSE) of the projection error of 6 points—left, center and right extremities at the middle and bottom parts of the frame—in the 3D coordinates.

*3.2.3 Baseline and results.* The baseline is composed by two models: the first is a segmentation model (DeepLabv3[5]) that predicts the 20 lines of the basketball court (see Figure 5); the second finds the 2D intersections in the image space and matches them with the

visible 3D locations of the court[13]. If enough intersections points are found (>5) the method `cv2.calibrateCamera`[14] predicts the camera parameters. In all the other cases, the model returns an average of the camera parameters in the training set as default. The segmentation model has been fine-tuned on the Camera calibration dataset (with $\alpha = 1$) for 40k steps with AdamW optimizer, base learning rate of 0.001, weight decay of 0.0001, and Amsgrad, reaching an mIoU of 0.46 on the validation set. The current baseline has an MSE error of 592.48 cm on the Test split and 490.31 cm on the Challenge set.

## 3.3 Player instance segmentation

The player instance segmentation task tackles the delineation of individual humans (players, coaches and referees) lying on a basketball court or less than 1 meter away from its borders. Instance segmentation is a pervasive task that can apply to images captured from any domain in which objects can be individually identified and counted. Instance segmentation datasets have been collected, among other domains, in microbiology [24], biology [32], autonomous driving [10] and in everyday life [16, 29]. The set of methods that solve instance segmentation is equally rich. A dichotomy is usually drawn between top-down and bottom-up methods. Top-down methods first propose candidate bounding boxes, and then segment the main object of interest in each of them [17, 2]. Bottom-up methods first label pixels with embedding vectors and then cluster pixels with similar embeddings into instance masks [33, 6]. In both types, main limitations usually arise in crowded regions, where objects are close to or occlude each other. So much so that many of the new state-of-the-art methods explicitly tackle those weaknesses [23, 46].

The dataset we propose here has three key aspects that make it particularly relevant for studying instance segmentation in those challenging cases. Please refer to Figure 6 for visual examples that illustrate those. First, instances only belong to one class. This renders the training and analysis of models less cumbersome, as there is no interference between classes of different frequencies during the training, and no averaging of performance metrics across them. Second, although only one class is present, instances have varied appearances and poses, and are sometimes already tricky to extract from the background. Furthermore, occlusions are frequent, constituting challenging cases. The fact that instances of a same class have high interactions, sometimes leading each other to be split in disconnected parts, stresses greatly current instance segmentation methods. Third, instance masks provided are very precise. Those annotations have been semi-automatically annotated as reported in [34]. All in all, we believe this dataset provides a good compromise that is challenging for state-of-the-art models, yet practical to study models with.

*3.3.1 Dataset.* The task uses the **DeepSport dataset** presented in Section 2, with every images used individually[15], and provided in the COCO format [29]. The *train* and *val* subsets contain respectively 223 and 37 images sampled uniformly from the first set of

---

[12]See implementation at: https://gitlab.com/deepsport/deepsport_utilities/-/blob/main/deepsport_utilities/ds/instants_dataset/views_transforms.py

[13]https://github.com/DeepSportRadar/camera-calibration-challenge/blob/main/utils/intersections.py

[14]https://docs.opencv.org/4.6.0/d9/d0c/group__calib3d.html

[15]Recall that dataset items, the *raw-instants*, are composed of multiple images

**Figure 6: Samples of annotated images from the instance segmentation task (cropped around annotated instances). Annotated instances are highlighted in distinct colors.**

arenas presented in Table 1. They contain respectively 1674 and 344 annotations. The *test* subset contains 64 images coming from the last three arenas of Table 1. It contains 477 annotated humans. For the competition, participants were evaluated on the 84 images coming from the challenge set introduced in Section 2. The number of annotated humans is kept secret.

*3.3.2   Metrics.* Because the mAP metric is well established in instance segmentation [29, 10], and because we are mostly interested about segmentation quality, we focus on it as main metric. The version specific to instance segmentation (sometimes referred to as segm_mAP) is different from that of object detection (bbox_mAP). Indeed, it uses the intersection-over-union (IoU) of predicted and ground-truth masks rather than bounding boxes to compute each intermediate AP curve. This way, good segmentation (IoU $\geq$ 0.80) is strongly rewarded while low-quality segmentation (IoU $\approx$ 0.55) is not.

Because only one class is present, there is no averaging between the metrics of frequent and rare classes. Our mAP simply looks like

$$mAP = \frac{1}{10} \sum_{\tau \in [0.50:0.05:0.95]} AP@\tau \qquad (2)$$

AP@$\tau$ being the area under the precision-recall curve, when considering as true-positives the predicted masks that have IoU > $\tau$ with any ground-truth mask. ($\tau \geq 0.5$ prevents one-to-many mappings)

*3.3.3   Baseline and results.* The code base we provide[16] is built on MMDet [4]. Our baseline consists of a Mask-RCNN model [17], with a ResNeXt 101 backbone [43] and default configuration, trained for 20 epochs. This model reaches a mAP of 0.51.

Relying on MMDet gives contestants the possibility to use a wide range of renowned and top-performing models off-the-shelf.
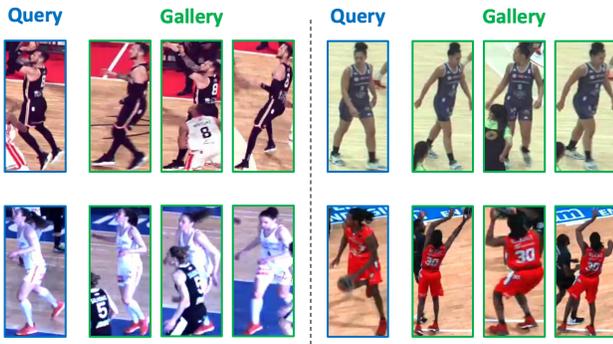
## 3.4   Player re-identification

Person re-identification [45], or simply ReID, is a person retrieval task which aims at matching an image of a person-of-interest, called *the query*, with other person images within a large database, called *the gallery*, captured from various camera viewpoints. ReID has important applications in smart cities, video-surveillance and sports analytics, where it is used to perform person retrieval or tracking. The objective of the DeepSportradar player ReID task is to re-identify players, coaches and referees across images captured successively from the same moving camera during a basketball game, as illustrated in Figure 7. Compared to traditional street surveillance type re-identification datasets [31, 48, 49, 27, 47], the DeepSportradar ReID dataset is challenging because players from the same team have very similar appearance, which makes it hard to tell them apart. However, as opposed to standard ReID datasets, all images are captured by the same camera, from the same point of view.

ReID has gained more and more attention recently, with several works proposing state-of-the-art methods based on global [13, 19,

---

[16]https://github.com/DeepSportRadar/instance-segmentation-challenge

**Table 2: The player re-identification dataset in numbers.**

| Subset | # Sequences | Split | # Ids | # Images |
|--------|-------------|-------|-------|----------|
| Train | 45 | - | 436 | 8569 |
| Test | 5 | Query | 50 | 50 |
| | | Gallery | 50 | 910 |
| Challenge | 49 | Query | (undisclosed) | 468 |
| | | Gallery | (undisclosed) | 8703 |



**Figure 7: The player re-identification task: illustration of some correctly retrieved gallery samples for four players of interest given as queries.**

30, 41], or part-based [37, 28] feature extractor. Other works introduced alternative ReID tasks, such as occluded ReID [31] or video-based ReID [47]. Multiple frameworks were also open-sourced to support further research on supervised [50, 19] or unsupervised ReID [14].

*3.4.1 Dataset.* The DeepSportradar player ReID dataset was built using 99 short video sequences from 97 different professional basketball games of the LNB proA league played in 29 different arenas. This dataset of short basketball video sequences was originally introduced for the VIPriors 2021 workshop [26], and is different from the DeepSport dataset described in Section 2. These sequences are 20 frames long, with a frequency of 10 frames per second (FPS), and they contain on average 10 different tracklets, i.e. identities. Therefore, the dataset contains a wide variety of players and sportswear appearance, within multiple arenas with different illumination and court appearance. Image crops[17] from players, coaches and referees have been extracted within each of these frames. The resulting re-identification dataset is composed of 18.703 thumbnails divided into three subsets: train, test and challenge set, as summarized in Table 2. Similar to other ReID datasets, the subset used for evaluating model performance, namely the test and challenge set, are split into a query and a gallery set. For these sets, we chose thumbnails from the $1^{st}$ frame of the sequence as queries, and remaining thumbnails

from the $2^{nd}$ to the $20^{th}$ frame as galleries. Labels from the challenge are kept secret to avoid any cheating in the DeepSportradar Player ReID challenge.

*3.4.2 Metrics.* Two standard retrieval evaluation metrics are used to compare different ReID models: the mean average precision [48] (mAP), and the cumulative matching characteristics (CMC) [42] at Rank-1 and Rank-5. The mAP is used to assess the average retrieval performance considering all ranked gallery samples. The Rank-K accuracy is the probability that at least one correct match appears in the top-K ranked retrieved results. Participants to the DeepSportradar Player ReID challenge are ranked according to their mAP score on the challenge set.

*3.4.3 Baseline and results.* Person re-identification is generally formulated as a metric learning task [45]. Firstly, a feature vector, also called "embedding", is extracted for each image in the dataset using a feature extractor. Secondly, the query to gallery similarity scores are measured as the pairwise euclidean distance of these features vectors in the embedding space. To address the DeepSportradar ReID challenge, we provide a simple CNN-based feature extractor as a baseline. This feature extractor was implemented using the Open-ReID[18] framework, a lightweight library of person re-identification, open-sourced for research purpose. Open-ReID aims to provide a uniform interface for different datasets, a full set of models and evaluation metrics. The baseline employed a ResNet-50 [18] CNN as backbone and is trained with a classification objective: the model tries to predict each sample identity among the 436 identities in the training set. The model is trained for 50 epochs with an SGD optimizer and a cross-entropy loss. Training batches are made of 64 players thumbnails, all resized to $256 \times 128$. We refer readers to our open-source toolbox on GitHub[19] for more details about the baseline architecture and training setup.

The baseline achieves **65% mAP**, **90% Rank-1** and **96% Rank-5** on the testing set of the DeepSportradar Player ReID dataset.

## 4 CONCLUSIONS

This paper has introduced two new datasets, namely the DeepSport dataset and the Basketball ReID dataset both acquired during professional basketball games with the Keemotion/Synergy Automated Camera System™. Together with these datasets, four CV tasks have been set up: the Ball 3D localization, Camera calibration, Player instance segmentation and Player re-identification. For each task, the dataset, the metrics and the baseline have been specified. The aim of this contribution was to provide a high-quality sports dataset framework where images, camera parameters and annotations are available and built close to the actual game recording setup, therefore providing an unique tool to experiment methods and solutions on real world settings.

---

[17]We refer to these image crops as "player thumbnails" for conciseness, without mentioning coaches and referees, because the large majority of these thumbnails actually depicts players.

[18]https://github.com/Cysu/open-reid
[19]https://github.com/DeepSportRadar/player-reidentification-challenge

# REFERENCES

[1] Md Zahangir Alom et al. 2019. A state-of-the-art survey on deep learning theory and architectures. *Electronics*, 8, 3, 292.

[2] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. 2019. Yolact: real-time instance segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 9157–9166.

[3] Jianhui Chen and James J Little. 2019. Sports camera calibration via synthetic data. In *Proceedings of the IEEE/CVF conference on CVPR workshops*, 0–0.

[4] Kai Chen et al. 2019. MMDetection: open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*.

[5] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*.

[6] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. 2020. Panoptic-deeplab: a simple, strong, and fast baseline for bottom-up panoptic segmentation. In *Proceedings of the IEEE/CVF conference on CVPR*, 12475–12485.

[7] Anthony Cioppa, Adrien Deliege, Floriane Magera, Silvio Giancola, Olivier Barnich, Bernard Ghanem, and Marc Van Droogenbroeck. 2021. 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*, 4537–4546.

[8] Anthony Cioppa, Adrien Deliège, Silvio Giancola, Bernard Ghanem, and Marc Van Droogenbroeck. 2022. Scaling up soccernet with multi-view spatial localization and re-identification. *Scientific Data*, 9, 1, 1–9.

[9] Anthony Cioppa, Silvio Giancola, Adrien Deliege, Le Kang, Xin Zhou, Zhiyu Cheng, Bernard Ghanem, and Marc Van Droogenbroeck. 2022. Soccernet-tracking: multiple object tracking dataset and benchmark in soccer videos. In *Proceedings of the IEEE/CVF Conference on CVPR*, 3491–3502.

[10] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. 2016. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on CVPR*, 3213–3223.

[11] Adrien Deliege, Anthony Cioppa, Silvio Giancola, Meisam J Seikavandi, Jacob V Dueholm, Kamal Nasrollahi, Bernard Ghanem, Thomas B Moeslund, and Marc Van Droogenbroeck. 2021. Soccernet-v2: a dataset and benchmarks for holistic understanding of broadcast soccer videos. In *Proceedings of the IEEE/CVF Conference on CVPR*, 4508–4519.

[12] Dirk Farin, Susanne Krabbe, Wolfgang Effelsberg, et al. 2003. Robust camera calibration for sport videos using court models. In *Storage and Retrieval Methods and Applications for Multimedia 2004*. Vol. 5307. SPIE, 80–91.

[13] Dengpan Fu, Dongdong Chen, Jianmin Bao, Hao Yang, Lu Yuan, Lei Zhang, Houqiang Li, and Dong Chen. 2021. Unsupervised pre-training for person re-identification. *Proceedings of the IEEE conference on CVPR*.

[14] Yixiao Ge, Feng Zhu, Dapeng Chen, Rui Zhao, and Hongsheng Li. 2020. Self-paced contrastive learning with hybrid memory for domain adaptive object re-id. In *Advances in Neural Information Processing Systems*.

[15] Silvio Giancola, Mohieddine Amine, Tarek Dghaily, and Bernard Ghanem. 2018. Soccernet: a scalable dataset for action spotting in soccer videos. In *Proceedings of CVPR workshops*, 1711–1721.

[16] Agrim Gupta, Piotr Dollar, and Ross Girshick. 2019. Lvis: a dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on CVPR*, 5356–5364.

[17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *2017 IEEE International Conference on Computer Vision (ICCV)*, 2980–2988. DOI: 10.1109/ICCV.2017.322.

[18] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. *2016 IEEE Conference on CVPR*, 770–778.

[19] Lingxiao He, Xingyu Liao, Wu Liu, Xinchen Liu, Peng Cheng, and Tao Mei. 2020. Fastreid: a pytorch toolbox for general instance re-identification. *arXiv preprint arXiv:2006.02631*.

[20] Namdar Homayounfar, Sanja Fidler, and Raquel Urtasun. 2017. Sports field localization via deep structured models. In *Proceedings of the IEEE Conference on CVPR*, 5212–5220.

[21] Peter J. Huber. 1964. Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35, 1, (Mar. 1964), 73–101. DOI: 10.1214/aoms/11777037 32.

[22] Paresh R Kamble, Avinash G Keskar, and Kishor M Bhurchandi. 2019. Ball tracking in sports: a survey. *Artificial Intelligence Review*, 52, 3, 1655–1705.

[23] Lei Ke, Yu-Wing Tai, and Chi-Keung Tang. 2021. Deep occlusion-aware instance segmentation with overlapping bilayers. In *Proceedings of the IEEE/CVF conference on CVPR*, 4019–4028.

[24] Neeraj Kumar et al. 2019. A multi-organ nucleus segmentation challenge. *IEEE transactions on medical imaging*, 39, 5, 1380–1391.

[25] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature*, 521, 7553, 436–444.

[26] Attila Lengyel, Robert-Jan Bruintjes, Marcos Baptista Rios, Osman Semih Kayhan, Davide Zambrano, Nergis Tomen, and Jan van Gemert. 2022. Vipriors 2: visual inductive priors for data-efficient deep learning challenges. *arXiv preprint arXiv:2201.08625*.

[27] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. 2014. Deepreid: deep filter pairing neural network for person re-identification. In *CVPR*.

[28] Yulin Li, Jianfeng He, Tianzhu Zhang, Xiang Liu, Yongdong Zhang, and Feng Wu. 2021. Diverse part discovery: occluded person re-identification with part-aware transformer. *2021 IEEE/CVF Conference on CVPR*, 2897–2906.

[29] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: common objects in context. In *European conference on computer vision*. Springer, 740–755.

[30] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. 2019. Bag of tricks and a strong baseline for deep person re-identification. *2019 IEEE/CVF Conference on CVPR Workshops*, 1487–1495.

[31] Jiaxu Miao, Yu Wu, Ping Liu, Yuhang Ding, and Yi Yang. 2019. Pose-guided feature alignment for occluded person re-identification. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 542–551.

[32] Massimo Minervini, Andreas Fischbach, Hanno Scharr, and Sotirios A Tsaftaris. 2016. Finely-grained annotated datasets for image-based plant phenotyping. *Pattern recognition letters*, 81, 80–89.

[33] Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. 2019. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. In *Proceedings of the IEEE/CVF Conference on CVPR*, 8837–8845.

[34] Niels Sayez and Christophe De Vleeschouwer. 2022. Accelerating the creation of instance segmentation training sets through bounding box annotation. (2022). DOI: 10.48550/ARXIV.2205.11563.

[35] Long Sha, Jennifer Hobbs, Panna Felsen, Xinyu Wei, Patrick Lucey, and Sujoy Ganguly. 2020. End-to-end camera calibration for broadcast videos. In *Proceedings of the IEEE/CVF conference on CVPR*, 13627–13636.

[36] Karen Simonyan and Andrew Zisserman. 2015. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. International Conference on Learning Representations, ICLR, (Sept. 2015). https://arxiv.org/abs/1409.1556v6 arXiv: 1409.1556.

[37] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. 2018. Beyond part models: person retrieval with refined part pooling. In *ECCV*.

[38] Gabriel Van Zandycke. 2021. Deepsport dataset. https://www.kaggle.com/datasets/gabrielvanzandycke/deepsport-dataset. (2021).

[39] Gabriel Van Zandycke and Christophe De Vleeschouwer. 2022. 3d ball localization from a single calibrated image. In *Proceedings of the IEEE/CVF Conference on CVPR*, 3472–3480.

[40] Gabriel Van Zandycke and Christophe De Vleeschouwer. 2019. Real-time CNN-based segmentation architecture for ball detection in a single view setup. *MMSports 2019 - Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports, co-located with MM 2019*, 51–58. ISBN: 9781450369114. arXiv: 2007.11876. DOI: 10.1145/3347318.3355517.

[41] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. 2018. Learning discriminative features with multiple granularities for person re-identification. *Proceedings of the 26th ACM international conf. on Multimedia*.

[42] Xiaogang Wang, Gianfranco Doretto, Thomas Sebastian, Jens Rittscher, and Peter Tu. 2007. Shape and appearance context modeling. In Rio de Janeiro, Brazil, 1–8. DOI: 10.1109/iccv.2007.4409019.

[43] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on CVPR*, 1492–1500.

[44] Qiang Yao, Akira Kubota, Kaoru Kawakita, Keisuke Nonaka, Hiroshi Sankoh, and Sei Naito. 2017. Fast camera self-calibration for synthesizing free viewpoint soccer video. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1612–1616.

[45] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven C. H. Hoi. 2022. Deep learning for person re-identification: a survey and outlook. 44, 6, (June 2022), 2872–2893. DOI: 10.1109/tpami.2021.3054775.

[46] Xiaoding Yuan, Adam Kortylewski, Yihong Sun, and Alan Yuille. 2021. Robust instance segmentation through reasoning about multi-object occlusion. In *Proceedings of the IEEE/CVF Conference on CVPR*, 11141–11150.

[47] Liang Zheng, Zhi Bie, Yifan Sun, Jingdong Wang, Chi Su, Shengjin Wang, and Qi Tian. 2016. Mars: a video benchmark for large-scale person re-identification. In *ECCV*.

[48] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable person re-identification: a benchmark. *2015 IEEE International Conference on Computer Vision (ICCV)*, 1116–1124.

[49] Zhedong Zheng, Liang Zheng, and Yi Yang. 2017. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. *2017 IEEE International Conference on Computer Vision (ICCV)*, 3774–3782.

[50] Kaiyang Zhou and Tao Xiang. 2019. Torchreid: a library for deep learning person re-identification in pytorch. *arXiv preprint arXiv:1910.10093*.