# DeepSportLab: a Unified Framework for Ball Detection, Player Instance Segmentation and Pose Estimation in Team Sports Scenes

Seyed Abolfazl Ghasemzadeh[1], Gabriel Van Zandycke[1,2], Maxime Istasse[1,2], Niels Sayez[1], Amirafshar Moshtaghpour[1], and Christophe De Vleeschouwer[1]

[1]UCLouvain ICTEAM/ELEN Belgium
`<firstname>.<lastname>@uclouvain.be`
[2]SportRadar AG
`<f>.<lastname>@sportradar.com`

## Abstract

*This paper presents a unified framework to (i) locate the ball, (ii) predict the pose, and (iii) segment the instance mask of players in team sports scenes. Those problems are of high interest in automated sports analytics, production, and broadcast. A common practice is to individually solve each problem by exploiting universal state-of-the-art models,* e.g.*, Panoptic-DeepLab for player segmentation. In addition to the increased complexity resulting from the multiplication of single-task models, the use of the off-the-shelf models also impedes the performance due to the complexity and specificity of the team sports scenes, such as strong occlusion and motion blur. To circumvent those limitations, our paper proposes to train a single model that simultaneously predicts the ball and the player mask and pose by combining the part intensity fields and the spatial embeddings principles. Part intensity fields provide the ball and player location, as well as player joints location. Spatial embeddings are then exploited to associate player instance pixels to their respective player center, but also to group player joints into skeletons. We demonstrate the effectiveness of the proposed model on the DeepSport basketball dataset, achieving comparable performance to the SoA models addressing each individual task separately.*

## 1 Introduction

The automation of team sports analytics and broadcasting [5, 4, 13], relies on detailed scene interpretation, which itself depends on the ability to detect the ball, segment the players (e.g. for improved tracking and recognition), and predict their pose (e.g. for action recognition). Our work leverages Convolutional Neural Networks (CNNs) to tackle those tasks.

The most natural approach to solve the three tasks above with CNNs is to use a pre-trained state-of-the-art model for each individual problem, *e.g.*, Mask R-CNN [17], PifPaf [25], and Panoptic-DeepLab [6] for ball localization, player pose estimation, and player instance segmentation, respectively. Such approach, however, re-

sults in poor performance since team sports scenes – especially those of indoor sports – are more complex compared with in-the-wild scenes. First, they involve strong player occlusions, *e.g.*, a defending player blocking an attacking one. Such occlusions are usually considered as "crowd" in in-the-wild datasets, such as CityScapes; hence, excluded from training. Second, they contain fast moving players and balls causing motion blur. Third, dealing with sports players is subject to deformation, since the players often jump, run, or stretch their bodies. Fourth, due to its frequent interactions with players, the ball is often partially occluded. Fifth, indoor team sports scenes suffer from weak contrast between the ball and field and from the reflection of the players in the field.

An improved approach consists in fine-tuning the weights of pre-trained models with task-specific datasets [42]. Nevertheless, that solution has its shortcomings. General purpose models are often demanding in terms of the memory and computation, which prevents their real-time application, especially if multiple models have to run in parallel. Moreover, by tackling each task individually, the CNN model ignores the correlation between them, which might hamper their performance [21].

In this work, we propose to rely on a single CNN to jointly localize the ball and predict the player poses and player instance masks, given a single input image. An overview of our model is presented in Fig. 1. As an important feature of the proposed model, we define an extended set of 19 keypoint-types including 17 human body parts, ball centroid, and player centroid; hence, treating the player centroids and ball detection tasks as keypoint detection problems. Our CNN is illustrated in more details in Fig. 2, where it outputs two sets of predictions. Inspired by Panoptic-DeepLab [6], one head network predicts pixel-wise semantic classes and offset vectors, providing the information required to associate player pixels to their centroid. As shown in Fig. 1, the other three head networks predict the Part Intensity Field (PIF) of the Pif-Paf framework [25] for the 19 keypoint-types, *i.e.*, a collection of confidence score, localization vectors, and the size and scale of that keypoint. Our instance segmentation decoder fuses the semantic classes, offset vectors, and
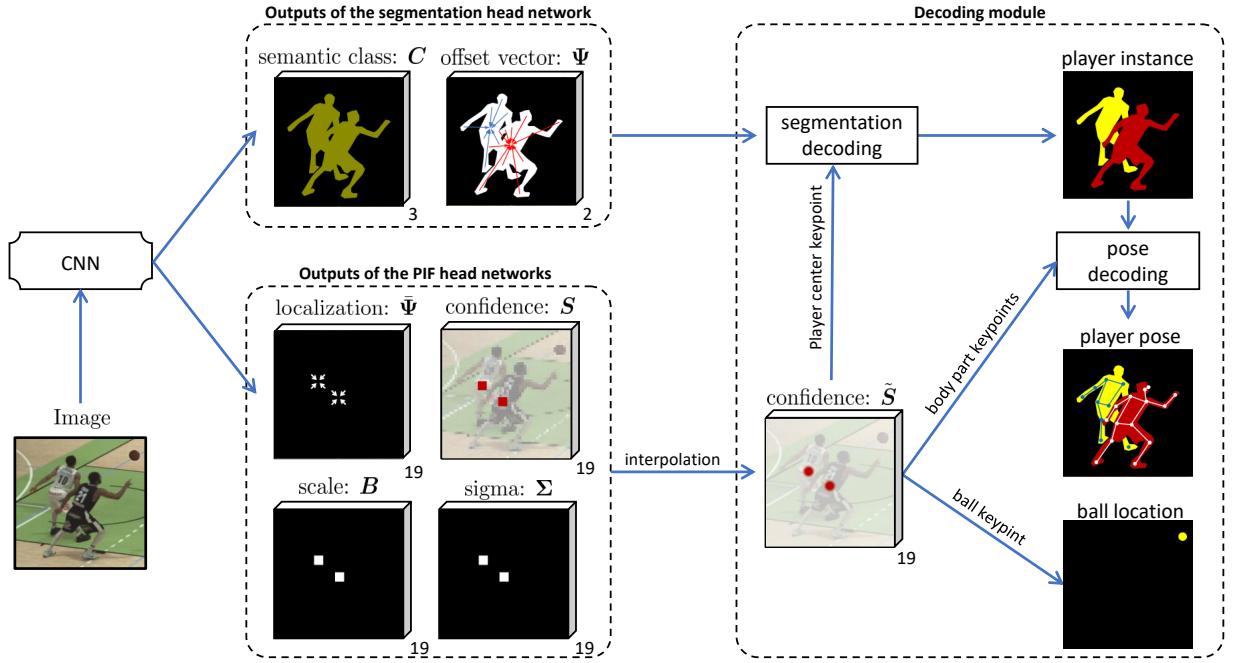
Figure 1: **An overview of DeepSportLab.** Our CNN outputs: *(i)* semantic class scores, *(ii)* offset vectors to associate pixels to their instance center, and *(iii)* a set of low resolution outputs used to generate the high-resolution keypoint confidence maps. This set includes the keypoints confidence map, the keypoints localization vectors, the size (or width) of the keypoints, and the scale parameters, which is used for scaling the per-keypoint localization loss (see Eq. (2)). The segmentation decoding module fuses the first two outputs and the high-resolution confidence map (*i.e.*, only the map of the center-of-the-player keypoint) in order to predict the player masks. The confidence maps of the body part keypoints and the resulting player masks are then fused by the pose decoding module and yield the player poses. Finally, the location of the ball is extracted from the high-resolution map of the ball keypoint. Such pipeline results in a light-weight decoding process. Note that the images in this figure are only for visualization purposes and are neither the ground truth annotations nor the predicted maps. In the outputs of the PIF head networks, we only show the maps associated with the center-of-the-player keypoint.

player centroids. The resulting player mask is then leveraged to associate the predicted body parts to the individual players. Therefore, unlike [25], our specialized model does not require the Part Affinity Field (PAF) for building the player skeletons; resulting in a light-weight decoding procedure. We train our model on the combination of the COCO [28] and DeepSport [1] datasets. The software is made freely available at https://github.com/ispgroupucl/DeepSportLab.

## 2   Related works

**Instance segmentation.**   Instance segmentation systems can be grouped into *proposal-based* [16, 15, 17] and *proposal-free* [31, 6] methods. The common idea in the former type of approaches is to predict a number of object proposals (or bounding-boxes) and then to classify the objects within each of those proposals. In contrast, most of the works in the latter category (including the proposed DeepSportLab) rely on predicting pixel-wise embedding vectors such that pixels belonging to an instance receive similar embeddings. In that context, a clustering algorithm is exploited to group the pixels based on the similarity between their embedding vectors, and in turn, to create the instance masks. The first proposal-free method

is reported in [27], where the CNN outputs the pixel-wise coordinates of the top-left and bottom-right corners of the corresponding object bounding-box. Those coordinates can be thought as 4-D embeddings. Newell *et al.* [31] introduced associative embeddings, which can be thought as a vector representing each pixel's cluster. In particular, the embedding vectors in [31] are related to an (non-physical) abstract space and are supervised without ground truth references – unlike the formalism in [27]. Other intermediate approaches have been proposed in [12, 24, 9] with different clustering algorithms and loss functions. Novotny *et al.* [32] extended the concept of associative embedding to spatial embedding: the network predicts pixel-wise offset vectors such that pixels of an instance point to that instance's center. That idea later inspired several works [21, 41, 30], specially Panoptic-DeepLab [6], which is one of the current top-performing instance segmentation models. Inspired by [6], our network predicts offset vectors from each pixel to its corresponding player instance center, which will be used to form the player instances.

**Human pose estimation.**   Human pose estimation approaches can be divided into two groups: *top-down* [34, 17] and *bottom-up* [25, 23, 33, 31, 2]. Top-down meth-

ods first perform a human detector and then locate the body parts within every detected bounding-box. Bottom-up methods (*i.e.*, the context of this paper), on the contrary, first estimate body parts followed by forming the skeleton.

OpenPose [2, 3] model revolutionized the bottom-up approaches by showing that a non-parametric representation (*i.e.*, PAF) can be learned to associate body parts with individual human instances. Newell *et al.* [31] used the notion of associative embedding, which can be thought as a vector representing each keypoint's cluster. In particular, keypoints with similar embedding vectors are assigned to the same skeleton. PersonLab [33] integrates the human pose estimation with instance segmentation. It learns to predict relative displacement of body parts, *i.e.*reminiscent of the idea of spatial embedding [32], allowing to group body parts into human skeleton. By combining and extending the use of short-range offset vectors (as in PersonLab) and PAF (as in OpenPose), Kreiss *et al.*[25] presented PifPaf framework, which reaches excellent performance for crowded scenes. Our work adopts the PIF part of that framework for predicting the location of the keypoints. Regarding the association of the body part keypoints to the poses, we follow the Panoptic-DeepLab's formalism. Unlike PersonLab, where human instance masks are computed from the human poses (in combination with predicted long-range offset vectors and semantic segmentation), in this work, we leverage the predicted player instance masks to group the keypoints into player poses.

**Deep learning applications in sports.** Deep learning has recently offered promising solutions for sports production, such as jersey number recognition [14], segmentation of the field, players, and lines [8], player discrimination [29, 20], event segmentation [10], player detection [7, 38, 35], swimming stroke rate detection [43], player pose estimation [46, 45, 19, 11], and ball localization [42, 36, 40]. In particular, the classifier in [36] decides whether a patch of an image contains a tennis ball. The authors in [40] and [42] propose to predict the position of the ball by formulating the problem as, respectively, a regression and segmentation task. In this paper, however, we treat ball localization as a keypoint detection problem. Unlike those works on player pose estimation, which are designed for scenes containing a single swimmer [46, 45] or athlete [19, 11], our proposed method aims at multi-player pose estimation in team sports scenes.

# 3 Proposed Method: DeepSportLab

**Notations.** Domain dimensions are represented by capital letters, *e.g.*$P$. Tensors are denoted by upper case bold symbols. Ground truth data or associated parameters are distinguished by an asterisk, *e.g.*, $N_{\mathrm{ply}}^*$ and $N_{\mathrm{ply}}$ are, respectively, the number of annotated and predicted players. The set of keypoint-types is denoted by $\mathcal{K} := \mathcal{K}_{\mathrm{part}} \cup \{\mathrm{ball, ply}\}$, where "ply" denotes the center-of-the-player keypoint and $\mathcal{K}_{\mathrm{part}} := \{\mathrm{left\ eye}, \cdots, \mathrm{right\ ankle}\}$ is the set of 17 body parts.

## 3.1 Principle

Given an input image with $P$ pixels and $N_{\mathrm{ply}}$ players, the goal of DeepSportLab is to predict *(i)* the location of the ball; *(ii)* the instance mask (or set of pixels) of each player $\mathcal{I}_i \subset \{1, \cdots, P\}$ for $i \in \{1, \cdots, N_{\mathrm{ply}}\}$; and *(iii)* the skeleton of players.

As illustrated in Fig. 1, our CNN (detailed in Sec. 3.2) outputs two groups of parameters. One group consists of the predictions from segmentation head network: semantic class scores $\boldsymbol{C} \in [0,1]^P$ (*i.e.*player vs. non-player) and offset (or spatial embedding) vectors $\boldsymbol{\Psi} \in \mathbb{R}^{P \times 2}$, *i.e.*, displacement of each pixel from the center of a player it belongs to (similarly to Panoptic-DeepLab [6]). The second group of outputs includes PIF predictions for each keypoint-type $k \in \mathcal{K}$, *i.e.*, a (low-resolution) pixel-wise confidence map $\boldsymbol{S}(k) \in [0,1]^{\bar{P}}$, a vector component (or localization vector) $\bar{\boldsymbol{\Psi}}(k) \in \mathbb{R}^{\bar{P} \times 2}$ pointing to the closest keypoint, size of the keypoint $\boldsymbol{\Sigma}(k) \in \mathbb{R}^{\bar{P}}$, and a scale parameter $\boldsymbol{B}(k) \in \mathbb{R}^{\bar{P}}$ (we set $\bar{P} = P/64$ in our experiments). As it will be detailed in Sec. 3.4, high-resolution confidence maps are then generated by fusing the low-resolution confidence maps with the localization vectors and keypoint sizes. Such formulation results in a light network architecture, since the number of parameters for each keypoint-type at the network output is $5\bar{P} = 0.078P$. Moreover, the parameter $\boldsymbol{B}(k)$ is used as weights in the training loss to balance the localization error with respect to each keypoint-type (see below).

The segmentation decoding module then predicts the instance masks of the players by fusing the semantic class scores, offset vectors, and confidence maps of the center-of-the-player keypoint. Furthermore, the pose decoding module leverages the player instance masks to group the predicted body parts into individual skeletons. Finally, the location of the ball is extracted from the confidence map of the ball keypoint.

## 3.2 Network architecture

As illustrated in Fig. 2, the CNN module of DeepSport-Lab has one shared backbone, *i.e.*, a ResNet-50 network [18], for its four head networks. The segmentation head network consists of an Atrous Spatial Pyramid Pooling (ASPP) module involving Atrous convolutions [44] to extract denser feature maps, a modified version of the upsampling module from Panoptic-DeepLab, and two task-specific prediction branches, which outputs the semantic logits and offset vectors, respectively. The upsampling module gradually increases the resolution of the features by leveraging the three skip connections from the backbone. In that module, a bilinear interpolation is applied before each concatenation step. Inspired by the PifPaf framework [25], our three PIF head networks output the PIF parameters for, respectively, the ball, the center of the players, and the other 17 keypoint-types. Each PIF head is equipped with a pixel shuffle [39] operation upsampling the feature maps by a factor of two. The semantic logits and the offset vectors are further upsampled by a bilinear interpolation to reach the input image resolution before being fed to the decoding module.
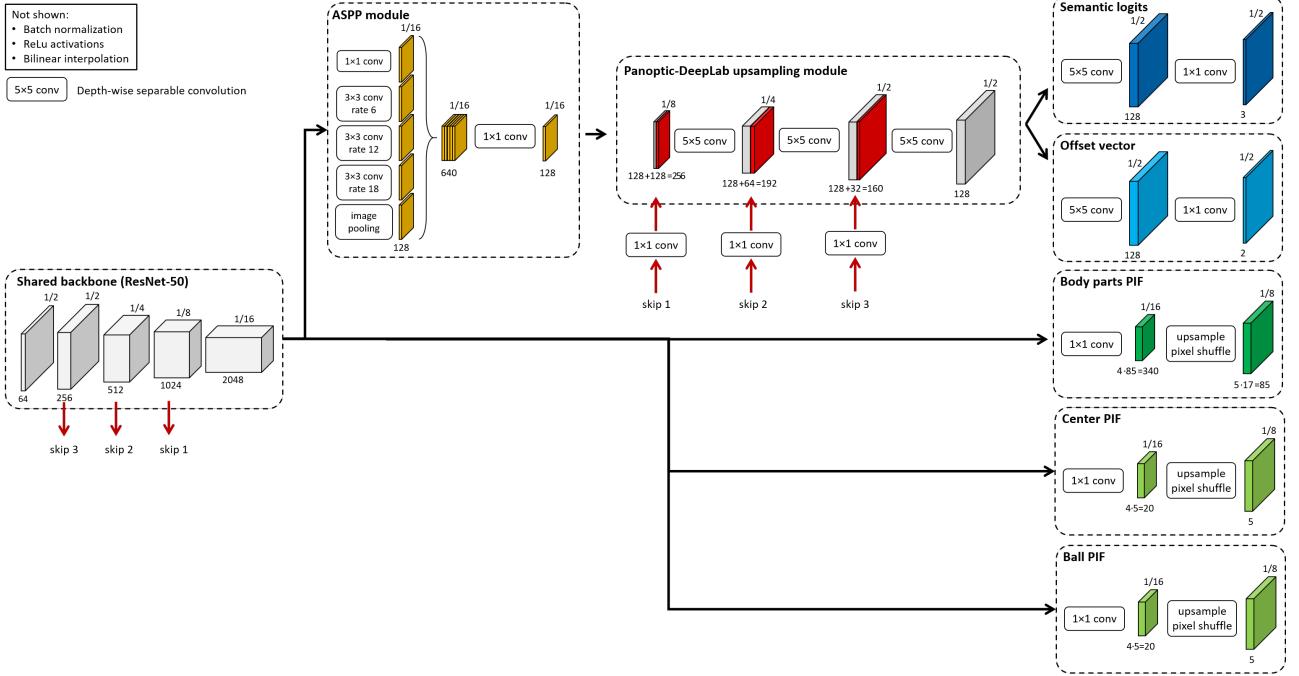
Figure 2: **CNN architecture of DeepSportLab.** Our network adopts an ASPP module, an upsampling module, and a dual-branch for the segmentation head. The upsampling module involves a bilinear interpolation before each concatenation step. The network also adopts a dual-PIF head network outputting the PIFs for the 19 keypoint-types. For lighter visualization, we concatenated the three PIF outputs in Fig. 1.

Note that in Fig. 1 the outputs of the three PIF head networks are concatenated for the sake of simpler visualization. Similarly, the PIF tensors, *e.g.*, $\boldsymbol{S}(k), \boldsymbol{\Psi}(k)$ are defined for all 19 keypoint-types $k$. We recall that our segmentation and PIF head networks are adapted from the Panoptic-DeepLab [6] and PifPaf [25] frameworks, respectively.

## 3.3 Network supervision

DeepSportLab requires a collection of ground truth data for its supervised training, *i.e.*, binary semantic class scores $\boldsymbol{C}^* \in \{0,1\}^P$, player and ball instance masks, pixel-wise offsets to players' instance centroids $\boldsymbol{\Psi}^* \in \mathbb{R}^{P \times 2}$, and per-keypoint-type (low-resolution) confidence maps $\boldsymbol{S}^*(k) \in \{0,1\}^{\bar{P}}$, localization vectors $\bar{\boldsymbol{\Psi}}^*(k) \in \mathbb{R}^{\bar{P} \times 2}$, and keypoint sizes (or sigma) $\boldsymbol{\Sigma}^*(k) \in \mathbb{R}^{\bar{P}}$. Note that $\boldsymbol{S}^*(k), \bar{\boldsymbol{\Psi}}^*(k)$, and $\boldsymbol{\Sigma}^*(k)$ take non-zero values on only 16 neighboring pixels of each keypoint. The ground truth size parameters are the body part standard deviation values (provided by the COCO dataset) scaled according to the size of the keypoints in the image.

Given the above-mentioned ground truth data, Deep-SportLab is trained by minimizing the following loss function:

$$\begin{aligned} \mathcal{L} = &w_{\text{sem}}\mathcal{L}_{\text{sem}} + w_{\text{off}}\mathcal{L}_{\text{off}} \\ &+ w_{\text{cnf}}\mathcal{L}_{\text{cnf}} + w_{\text{loc}}\mathcal{L}_{\text{loc}} + w_{\text{sig}}\mathcal{L}_{\text{sig}}, \end{aligned} \quad (1)$$

where the first two (and the last three) losses correspond to the segmentation (resp., pose) head network. In (1), $w_{\text{sem}}$, $w_{\text{off}}, w_{\text{cnf}}, w_{\text{loc}}, w_{\text{sig}}$, are the semantic, offset, confidence,

localization, and sigma loss weights, respectively, and

$$\begin{aligned} \mathcal{L}_{\text{sem}} &= \frac{1}{P}\sum_p \text{BCE}(\boldsymbol{C}^*(p), \boldsymbol{C}(p)), \\ \mathcal{L}_{\text{off}} &= \frac{1}{P}\sum_p \|\boldsymbol{\Psi}(p) - \boldsymbol{\Psi}^*(p)\|_2^2, \\ \mathcal{L}_{\text{cnf}} &= \sum_{p,k} \text{BCE}(\boldsymbol{S}^*(k;p), \boldsymbol{S}(k;p)), \\ \mathcal{L}_{\text{sig}} &= \sum_{p,k} |\boldsymbol{\Sigma}(k;p) - \boldsymbol{\Sigma}^*(k;p)|, \\ \mathcal{L}_{\text{loc}} &= \sum_{p,k} \frac{|\bar{\boldsymbol{\Psi}}(k;p) - \bar{\boldsymbol{\Psi}}^*(k;p)|^2}{\boldsymbol{B}^2(k;p)} + \log(\boldsymbol{B}(k;p)), \end{aligned} \quad (2)$$

where BCE denotes the binary cross entropy functions. We consider the mean squared error and mean absolute error for the rest of the sub-losses, summed over the defined variables, *e.g.*, over the keypoints present in the image. For the localization loss we adapt the learnable standard deviation for the regression loss, *i.e.*, inspired by [21]. The purpose of such formalism is to let the network balance the localization error based on the size of the keypoint. Intuitively, while a localization error might be minor for a large keypoint, it can be major for a small keypoint.

In practice, for the sake of numerical stability, we train the network to predict the logarithm of $\boldsymbol{B}$ and $\boldsymbol{\Sigma}$, and use the exponential of the prediction. By doing so, both the localization loss $\mathcal{L}_{\text{loc}}$ and the fusion operator (see Eq. (3) below) avoid any division by zero during the training and inference, respectively.

## 3.4 Inference

DeepSportLab consists in two main inference sub-tasks: *(i)* keypoint decoding, which provides the center of the ball and players as well as the body parts keypoints; *(ii)* player instance segmentation, which also enables pose recognition, as it assigns an instance label to each of the body part keypoints.

Recall that the keypoint confidence maps $S(k)$ are coarse. In order to obtain high-resolution maps, they are fused with the localization vectors $\bar{\Psi}(k)$ and sigma parameters $\Sigma(k)$ through a convolutional operator with unnormalized Gaussian kernel, *i.e.*,

$$\tilde{S}(k;p) := \sum_u S(k;u) \, \exp\big(-\tfrac{\|p-(u+\bar{\Psi}(k;u))\|^2}{2\Sigma^2(k;u)}\big), \quad (3)$$

where $p$ and $u$ are the pixel coordinates in, respectively, high- and low-dimensional pixel grids. The location of the ball is identified by finding the maximum value in the confidence map of the ball keypoint. This amounts to a top-1 detection strategy.

For the next decoding task, each pixel $p$ classified as a player is assigned to only one player instance using a simple center regression. Equivalently, pixels are grouped into individual instances according to their displacements from the center of each player. Concretely, the mask of each player $i$ writes

$$\mathcal{I}_i^{\mathrm{ply}} := \Big\{ p : \tilde{C}(p) \equiv \mathrm{player},$$
$$i = \arg\min_j \left\| o_j^{\mathrm{ply}} - (p + \Psi(p)) \right\|_2 \Big\}. \quad (4)$$

For each of the other 18 keypoint-types $k \in \mathcal{K}\backslash\{\mathrm{ball}\}$, we then compute a set of keypoint instances $\mathcal{O}^k := \{o_1^k, \cdots, o_{N_k}^k\}$, with $N_k$ denoting the number of detected instances for keypoint-type $k$, by finding the maxima within each player's mask in the high-resolution confidence map $\tilde{S}(k)$. We discard the detected keypoint instances with the confidence score less than 0.1. Having identified the mask of each player, the player pose can be simply decoded: the skeleton of a player is formed by collecting body part keypoints whose coordinates belong to the mask of that player.

If multiple keypoint instances of the same type are assigned to one skeleton, we keep only the keypoint instance with the highest confidence score.

**Remark 1** *Notice that our pose decoding algorithm described above is very fast and simple. Compared with Pif-Paf [25], DeepSportLab does not require to learn the part affinity fields with dimension $\bar{P}\times 19\times 7$ for pose decoding. Moreover, the greedy decoding algorithm in PifPaf [25], which is similar to the PersonLab's [33], starts from a body part and finds the next connected body part (among other candidate parts, given the part affinity fields) by computing a pixel-wise association score.*

## 4 Experiments

Our model is trained using a mix of the COCO [28] and DeepSport [1] datasets and is compared against Panoptic DeepLab [6] for player instance segmentation, OpenPif-Paf [26] for player pose estimation, and BallSeg [42] for ball detection.

**DeepSport dataset.** It consists of 672 images of professional basketball games captured from 29 arenas involving a large variety of game configurations and various lighting conditions [1]. Each image captures half of the basketball court with a resolution between 2Mpx and 5Mpx. The resulting images have a definition varying between 65px/m (furthest point on court in the arena with the lowest resolution cameras) and 265px/m (closest point on court in the arena with the highest resolution cameras). We extract 100 (out of 672) images for the validation and 64 images for the testing such that the arenas in the test set are neither present in the training nor the validation sets. The dataset contains approximately 380 annotations of ball masks and 5500 annotations of player masks and poses. The poses (composed of only 4 keypoint-types) are only used for evaluation. The images are scaled to keep humans with similar height compared to those found in the COCO dataset [28].

**COCO dataset.** In order to learn the pose of players, we consider a subset of the COCO dataset [28] consisting of images containing only humans or balls. We also filter out the images containing at least one person with the area of larger than 10% of the whole image. This filtering results in 42271 and 2356 images for the training and validation, respectively.

**Mismatch between DeepSport and COCO pose annotations.** The COCO dataset, used to train the body part keypoints, contains pose annotations with 17 body parts; whereas, in DeepSport dataset used at testing, they are identified by four body parts: head, hip, foot 1, and foot 2, *i.e.*, agnostic about the skeletons facing forward or backward. Hence, the quality computation phase during testing requires a delicate treatment. First, the skeletons predicted with the COCO convention are converted into the DeepSport skeleton format as follows. The locations of the head and hip for metric computations are defined, respectively, as the middle of the left- and right-ear and the middle of the left- and right-hip. Since foot 1 and foot 2 labels are interchangeable in DeepSport dataset, the keypoint metric (sensitive to inversion) is computed for the two assignments of the feet, *i.e.*, assigning foot 1/foot 2 to either left-ankle/right-ankle or right-ankle/left-ankle, and the assignment with a higher resulting value is considered. We note that this conversion of skeleton introduces error in quality computation of predicted poses, but since that computation is the same for different models, performed comparisons remain fair.

**Training setup.** Each batch of data during training involves equal share of images from DeepSport and COCO

| Method | bDQ | pSQ | pDQ | pEQ | | | ms | MB |
|---|---|---|---|---|---|---|---|---|
| | | | | AP | AR | $F_1$ | | |
| DeepSportLab | 52.07 | 80.3 | 90.1 | 87.5 | 82.1 | 42.4 | $436 \pm 105$ | 1757 |
| Pan.-DeepLab [6] | – | 82.2 | 91.4 | – | – | – | $69 \pm 7$ | 1809 |
| OpenPifPaf [26] | – | – | – | 88.5 | 79.6 | 41.9 | $155 \pm 18$ | 1623 |
| BallSeg [42] | 46.16 | – | – | – | – | – | $14 \pm 6$ | 1239 |

Table 1: **Performance on DeepSport's *test***: Comparison of the proposed multi-task DeepSportLab with SoA models addressing each individual tasks. Panoptic-Deeplab and BallSeg are trained on the DeepSport dataset, while PifPaf is trained on COCO. DeepSportLab uses both DeepSport and COCO datasets to train the player and ball instances segmentation and only uses COCO to train the player pose. DeepSportLab compares favorably against its counterparts in term of quality measures, while reducing by a factor 3 the required memory.

| player keypoints | player masks | Decoder | pSQ | pDQ | pEQ | | |
|---|---|---|---|---|---|---|---|
| | | | | | AP | AR | $F_1$ |
| DeepSportLab | DeepSportLab | | 80.3 | 90.1 | 87.5 | 82.1 | 42.4 |
| DeepSportLab | Oracle | DeepSportLab | 100 | 100 | 87.7 | 83.8 | 42.9 |
| OpenPifPaf [26] | Pan.DeepLab [6] | | 82.2 | 91.4 | 87.3 | 82.1 | 42.3 |

Table 2: **DeepSportLab Decoder study.** Oracle data or off-the-shelf models are used to analyse the sensitivity of DeepSportLab's decoder to the masks and keypoints predictions. The use of instance ground-truth masks show that the DSL instance segmentation is good enough to decode the player pose. DSL also compares favorably to a decoding using a combination of OpenPifPaf [26] and Pan.DeepLab [6].

datasets. Therefore, during every epoch, the learning algorithm works over a different subset of the COCO dataset. We trained our models with SGD optimizer with a learning rate set to $3 \cdot 10^{-4}$, momentum to $0.95$, and weight decay to $10^{-5}$. We perform random horizontal flip, scaling, and rotation followed by a random $641 \times 641$ cropping during training with batch size of $4$. We set the weights of the training sub-losses as $w_{sem} = 10$, $w_{off} = 0.1$, $w_{cnf} = 20$, $w_{loc} = 10$, and $w_{sig} = 10$. We initialized the network's backbone and player's keypoints head using the pre-trained weights taken from Pif-Paf framework.

**Quality measures.** This paper tackles a multi-objective task, hence, requires multiple quality measures for expressing the performance of individual sub-tasks. A detailed explanation of the quality measures is provided in the Supplementary Document. We define the ball Detection Quality (bDQ) as in the BallSeg framework [42], *i.e.* the Area under the Curve (AUC) associated with the Receiver Operating Characteristic (ROC) curve: the True Positive (TP) rate vs. the False Positive (FP) rate as a function of the detection score threshold.

We compute player Segmentation Quality (pSQ) and player Detection Quality (pDQ) identical to two components of the Panoptic Quality (PQ) measure introduced in [22], but restricted to only player (or person) class.

The pose Estimation Quality (pEQ) is measured based on the Object Keypoint Similarity (OKS) criteria defined in [37]. The components of pEQ measure are the variants of Average Precision (AP) and Average Recall (AR) thresholded at different OKS values.

**Computational resources.** Time and memory performance were measured at inference on batches of one single $641 \times 641$ pixels image, on machine configured with an NVidia V100 and a 32 cores Intel Xeon Gold 5217 running at 3GHz. Inference time is given in milliseconds (ms) and GPU memory in megabytes (MB).
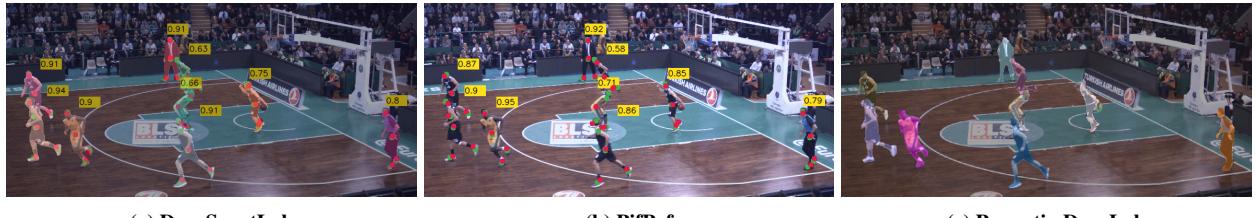
**Evaluation.** We evaluate the performance of all three tasks on DeepSport's test set. Predicted and annotated human instances outside the basketball court are discarded during the computation of quality criteria.

Table 1 shows that DeepSportLab only lacks 1% in terms AP but improves by 2.5% in terms of AR compared to OpenPifPaf [26], which leads to an improvement of 0.5% in terms of $F_1$ in total. In terms of player segmentation quality, DeepSportLab lacks by only 1.9% and 1.3% in terms of pSQ and pDQ, respectively, compared to Panoptic-DeepLab [6]. Since DeepSportLab outputs three different tasks at the same time, obtaining such a good segmentation quality reveals the excellence of the multi-task training.

A visual analysis of our results compared to OpenPif-Paf [26] and Panoptic-DeepLab [6] is presented in Fig. 3. It reveals that body part keypoints are generally assigned to their correct instance based on our computationally straightforward assignment strategy (see Sec. 3.4).

Regarding ball detection, our multi-task DeepSport-Lab framework is compared to the state-of-the-art BallSeg [42], which formulates the ball detection problem as a segmentation task. BallSeg uses an ICNet network [47] trained on balls from the original DeepSport dataset [1] (i.e. ball size ranges between 15 and 45 pixels). The performances provided in Table 1 result from a training where images were scaled to keep humans with similar height compared to those found in the COCO dataset (i.e. ball size ranges between 7 and 18 pixels). At that scale, DeepSportLab significantly outperforms BallSeg.

The computational comparison of the different models demonstrates the benefit of the unified framework. DeepSportLab can be deployed on devices with less than 2GB of memory, while the combination of its three counterparts would require almost 5GB. Inference time is only informative due to the different levels of optimization of each implementation. Actually, when combining different task-specific models, inference times do not add up when running computations in parallel.

**(a) DeepSportLab**      **(b) PifPaf**      **(c) Panoptic-DeepLab**

Figure 3: **Pose recognition and mask segmentation samples.** Body parts are shown with red and green colors as prediction and ground truth, respectively. The numbers highlighted in yellow are OKS values for their corresponding prediction-annotation matching. In addition, colors are used to show the segmented player masks.

| Dataset used at training | bDQ | pSQ | pDQ | pEQ | | |
|---|---|---|---|---|---|---|
| | | | | AP | AR | $F_1$ |
| COCO & DeepSport | 52.07 | 80.3 | 90.1 | 87.5 | 82.1 | 42.4 |
| COCO | 23.19 | 75.2 | 82.7 | 87.6 | 81.6 | 42.3 |

Table 3: **Sport-specific dataset study.** DeepSportLab evaluated on DeepSport dataset [1] shows limited performances when trained exclusively on COCO dataset [28]. This demonstrate the need for a sport specific dataset. Note: the ball detection quality (bDQ) reported when training only on COCO has been achieved by filtering out training images exempt of balls from the official COCO dataset. Without that filtering, bDQ evaluated on Deep-Sport dataset is far worse.

## 5 Ablation studies

**DeepSportLab Decoder.** Our proposed pose decoder is original in that it uses the segmentation mask to assign the keypoints to each instance. In Table 2, this strategy is applied on oracle data as well as on off-the-shelf predictions form OpenPifPaf [26] and Panoptic DeepLab [6]. The use of oracle masks instead of the intermediate predictions increases the $F_1$ Score by only 0.5%. This shows that, on the DeepSport dataset, the instance segmentation is good enough to associate the player keypoints, meaning that the pose estimation task is only limited by the keypoints detection. Using off-the-shelf models to produce the intermediate mask and the keypoints predictions does not improve the pose quality, revealing that our keypoint prediction is competitive. Additional investigations considering oracle data and error breakdown are provided in supplementary material.

**Sport Specific Dataset.** While containing a fair amount of balls, the COCO dataset is not rich enough to address the objective of DeepSportLab. This is illustrated in Table 3, where DeepSportLab was trained exclusively on COCO and evaluated on the DeepSport dataset [1]. The ball detection quality (bDQ) value reported for COCO is significantly smaller than the one obtained when training on DeepSport. This demonstrates the need of having a task specific dataset to reach good performance when working on sport data.

**DeepSportLab for individual tasks.** Table 4 evaluates the penalty or gain induced by the multi-task objective compared to an individual training of each task on the DeepSportLab backbone. The instance segmentation task benefits from the joined training, while the ball detection task does not. We understand this observation by the fact that the players mask are tightly coupled to their pose while the ball is not.

## 6 Conclusions

DeepSportLab is a framework handling pose estimation, instances segmentation and ball detection tasks, central to team sports analysis. Its architecture, based on a shared backbone, makes it more practical to deploy compared to existing off-the-shelf solutions, since it dramatically reduces the memory requirements, without affecting the prediction accuracy. It proposes a new ball detection approach and a novel pose decoding algorithm based on the instance masks and showing interesting performances. Two main lessons are drawn in terms of image interpretation problem formulation: *(i)* adopting part intensity fields to locate the ball appears to be as effective than formulating this problem as a high-resolution image segmentation problem; and *(ii)* assigning pose keypoints to their respective instances based on the spatial offsets predicted for instance segmentation sounds like a promising solution to reduce the decoding cost of PifPaf's pose recognition. In terms of perspectives, we observe that the bottleneck for more accurate player pose estimation lies in the prediction of keypoints, and not in their assignment to instances. Fundamental research is also desired to extend our multi-task framework to more generic datasets.

| Method | bDQ | pSQ | pDQ | pEQ AP | pEQ AR | pEQ $F_1$ | ms | MB |
|---|---|---|---|---|---|---|---|---|
| DeepSportLab | 52.07 | 80.3 | 90.1 | 87.5 | 82.1 | 42.4 | $436 \pm 106$ | 1757 |
| DeepSportLab - Ball only | 54.72 | – | – | – | – | – | $73 \pm 12$ | 1688 |
| DeepSportLab - Instances masks only | – | 73.4 | 86.6 | – | – | – | $164 \pm 37$ | 1753 |
| DeepSportLab - Poses (masks + keypoint) only | – | 79.2 | 90.2 | 87.8 | 82.2 | 42.5 | $496 \pm 107$ | 1755 |

Table 4: **DeepSportLab on individual tasks study.** DeepSportLab was trained on the individual tasks and compared with the result of the combined training. Only the ball doesn't benefit from being trained jointly with another task.

# References

[1] DeepSport data set. https://sites.uclouvain.be/ispgroup/Softwares/DeepSport. retrieved on Feb. 24th, 2021.

[2] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *IEEE transactions on pattern analysis and machine intelligence*, 43(1):172–186, 2019.

[3] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017.

[4] Fan Chen and Christophe De Vleeschouwer. Formulating team-sport video summarization as a resource allocation problem. *IEEE Transactions on Circuits and Systems for Video Technology*, 21(2):193–205, 2011.

[5] Fan Chen and Christophe Vleeschouwer. Automatic production of personalized basketball video summaries from multi-sensored data. pages 565 – 568, 10 2010.

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

[7] Anthony Cioppa, Adrien Deliege, Noor Ul Huda, Rikke Gade, Marc Van Droogenbroeck, and Thomas B Moeslund. Multimodal and multiview distillation for real-time player detection on a football field. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 880–881, 2020.

[8] Anthony Cioppa, Adrien Deliege, and Marc Van Droogenbroeck. A bottom-up approach based on semantics for the interpretation of the main camera stream in soccer games. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1765–1774, 2018.

[9] Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative loss function. *arXiv preprint arXiv:1708.02551*, 2017.

[10] Moritz Einfalt, Charles Dampeyrou, Dan Zecha, and Rainer Lienhart. Frame-level event detection in athletics videos with pose-based convolutional sequence networks. In *Proceedings Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports*, pages 42–50, 2019.

[11] Mykyta Fastovets, Jean-Yves Guillemaut, and Adrian Hilton. Athlete pose estimation from monocular tv sports footage. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1048–1054, 2013.

[12] Alireza Fathi, Zbigniew Wojna, Vivek Rathod, Peng Wang, Hyun Oh Song, Sergio Guadarrama, and Kevin P. Murphy. Semantic instance segmentation via deep metric learning. *arXiv preprint arXiv:1703.10277*, 2017.

[13] Iván Fernández, Fan Chen, Fabien Lavigne, Xavier Desurmont, and Christophe Vleeschouwer. Browsing sport content through an interactive h.264 streaming session. 06 2010.

[14] Sebastian Gerke, Antje Linnemann, and Karsten Müller. Soccer player recognition using spatial constellation features and jersey number recognition. *Computer Vision and Image Understanding*, 159:105–115, 2017.

[15] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.

[16] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.

[17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[19] Jihye Hwang, Sungheon Park, and Nojun Kwak. Athlete pose estimation by a global-local network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 58–65, 2017.

[20] Maxime Istasse, Julien Moreau, and Christophe De Vleeschouwer. Associative embedding for team discrimination. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.

[21] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018.

[22] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9404–9413, 2019.

[23] Muhammed Kocabas, Salih Karagoz, and Emre Akbas. Multiposenet: Fast multi-person pose estimation using pose residual network. In *Proceedings of the European conference on computer vision (ECCV)*, pages 417–433, 2018.

[24] Shu Kong and Charless C. Fowlkes. Recurrent pixel embedding for instance grouping. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9018–9028, 2018.

[25] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11977–11986, 2019.

[26] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. OpenPifPaf: Composite Fields for Semantic Keypoint Detection and Spatio-Temporal Association. *arXiv preprint arXiv:2103.02440*, March 2021.

[27] Xiaodan Liang, Liang Lin, Yunchao Wei, Xiaohui Shen, Jianchao Yang, and Shuicheng Yan. Proposal-free network for instance-level object segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2978–2991, 2017.

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

[29] Keyu Lu, Jianhui Chen, James J Little, and Hangen He. Lightweight convolutional neural networks for player detection and classification. *Computer Vision and Image Understanding*, 172:77–87, 2018.

[30] Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8837–8845, 2019.

[31] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *Advances in neural information processing systems*, pages 2277–2287, 2017.

[32] David Novotny, Samuel Albanie, Diane Larlus, and Andrea Vedaldi. Semi-convolutional operators for instance segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 86–102, 2018.

[33] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–286, 2018.

[34] George Papandreou, Tyler Zhu, Nori Kanazawa, Alexander Toshev, Jonathan Tompson, Chris Bregler, and Kevin Murphy. Towards accurate multi-person pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4903–4911, 2017.

[35] Miran Pobar and Marina Ivasic-Kos. Mask r-cnn and optical flow based method for detection and marking of handball actions. In *2018 11th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*, pages 1–6. IEEE, 2018.

[36] Vito Reno, Nicola Mosca, Roberto Marani, Massimiliano Nitti, Tiziana D'Orazio, and Ettore Stella. Convolutional neural networks based ball detection in tennis games. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1758–1764, 2018.

[37] Matteo Ruggero Ronchi and Pietro Perona. Benchmarking and error diagnosis in multi-instance pose estimation. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[38] Melike Şah and Cem Direkoğlu. Evaluation of image representations for player detection in field sports using convolutional neural networks. In *International Conference on Theory and Applications of Fuzzy Systems and Soft Computing*, pages 107–115. Springer, 2018.

[39] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874–1883, 2016.

[40] Daniel Speck, Pablo Barros, Cornelius Weber, and Stefan Wermter. Ball localization for robocup soccer using convolutional neural networks. In *Robot World Cup*, pages 19–30. Springer, 2016.

[41] Jonas Uhrig, Eike Rehder, Björn Fröhlich, Uwe Franke, and Thomas Brox. Box2pix: Single-shot instance segmentation by assigning pixels to object boxes. In *IEEE Intelligent Vehicles Symposium (IV)*, 2018.

[42] Gabriel Van Zandycke and Christophe De Vleeschouwer. Real-time cnn-based segmentation architecture for ball detection in a single view setup. In *Proceedings of the 2nd International Workshop on Multimedia Content Analysis in Sports*, pages 51–58, 2019.

[43] Brandon Victor, Zhen He, Stuart Morgan, and Dino Miniutti. Continuous video to simple signals for

swimming stroke detection with convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 66–75, 2017.

[44] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.

[45] Dan Zecha, Moritz Einfalt, Christian Eggert, and Rainer Lienhart. Kinematic pose rectification for performance analysis and retrieval in sports. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1791–1799, 2018.

[46] Dan Zecha, Moritz Einfalt, and Rainer Lienhart. Refining joint locations for human pose tracking in sports videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.

[47] Hengshuang Zhao, Xiaojuan Qi, Xiaoyong Shen, Jianping Shi, and Jiaya Jia. Icnet for real-time semantic segmentation on high-resolution images. In *Proceedings of the European conference on computer vision (ECCV)*, pages 405–420, 2018.

# A  Quality measures

As mentioned in the main document, we tackle a multi-objective task, hence, requires multiple quality measures for expressing the performance of individual sub-tasks. The details of the quality measures used in our work is provided in the following.

**Ball Detection Quality (bDQ).**  As in the BallSeg framework [42], given a threshold $\tau$ in the dynamic range of the confidence scores, a predicted ball keypoint is identified as a True Positive (TP) (or False Positive (FP)) detection , if its location lies inside (respectively, outside) the ground truth mask and its predicted confidence is greater than $\tau$: for the predicted ball $o^{\text{ball}}$,

$$o^{\text{ball}} \in \text{TP}^{\text{ball}}(\tau) \Leftrightarrow o^{\text{ball}} \in \mathcal{I}^{\text{ball*}}, \; \tilde{S}(\text{ball}, o^{\text{ball}}) \geq \tau$$
$$o^{\text{ball}} \in \text{FP}^{\text{ball}}(\tau) \Leftrightarrow o^{\text{ball}} \notin \mathcal{I}^{\text{ball*}}, \; \tilde{S}(\text{ball}, o^{\text{ball}}) \geq \tau$$
$$(5)$$

where $\mathcal{I}^{\text{ball*}}$ denotes the semantic mask of the ball. By repeating this procedure for all images, we obtain the TP and FP sets associated with the full set of images. The TP rate (TPr) and FP rate (FPr) ratios are defined as

$$\text{TPr}(\tau) := \frac{|\text{TP}^{\text{ball}}(\tau)|}{|\{\text{images with annotated ball}\}|},$$
$$\text{FPr}(\tau) := \frac{|\text{FP}^{\text{ball}}(\tau)|}{|\{\text{all images}\}|}.$$

The bDQ is then computed as the area under the ROC curve associated with the TPr and FPr.

**Player Segmentation Quality (pSQ).**  Compared with the bDQ, a predicted player mask is identified as a TP mask (otherwise, an FP), if its IoU with one of the ground truth player masks is higher than the threshold of 0.5:

$$(\mathcal{I}_i^{\text{ply}}, \mathcal{I}_j^{\text{ply*}}) \in \text{TP}^{\text{ply}} \Leftrightarrow \exists j, \; \text{IoU}(\mathcal{I}_i^{\text{ply}}, \mathcal{I}_j^{\text{ply*}}) \geq 0.5. \quad (6)$$

The pSQ is then defined as the averaged IoU over the TP pairs:

$$\text{pSQ} := \frac{1}{|\text{TP}^{\text{ply}}|} \sum_{(u,v) \in \text{TP}^{\text{ply}}} \text{IoU}(u, v). \quad (7)$$

**Player Detection Quality (pDQ).**  Having identified the TP set as in (6), the pDQ is then defined as the $F_1$-score:

$$\text{pDQ} := \frac{2|\text{TP}^{\text{ply}}|}{N_{\text{ply}} + N_{\text{ply}}^*}. \quad (8)$$

Note that the pSQ (7) and pDQ (8) criteria are the segmentation and recognition quality components of the Panoptic Quality (PQ) measure introduced in [22], *i.e.*, concretely, the PQ for player segmentation reads $\text{PQ} := \text{pSQ} \cdot \text{pDQ}$.

**Pose Estimation Quality (pEQ).**  For pose estimation task we use the OKS criteria, *i.e.*, for every pair of the predicted pose $\Upsilon_i$ and ground truth pose $\Upsilon_j^*$, it is defined as

$$\text{OKS}(\Upsilon_i, \Upsilon_j^*) := \text{mean}_k \exp\left(-\frac{\|o_i^k - o_j^{k*}\|_2^2}{2s_j^2 \kappa_k^2}\right), \quad (9)$$

where the mean is taken over the annotated body part keypoints, $s$ denotes the square root of the area of the bounding-box tightly containing all the body parts, and $\kappa_k$ is the per-keypoint-type scale constant controlling falloff. The predicted skeletons are then sorted according to their confidence scores defined as the average over the body part confidence scores: from the pixel-wise confidence map in (Eq. 2 in the main document)

$$\Upsilon_i^{\text{conf}} = \frac{1}{17} \sum_{k \in \mathcal{K}_{\text{part}}} \tilde{S}(k; o_i^k). \quad (10)$$

Next, the ordered predictions are assigned to the ground truths, with which they have the highest OKS value. Once the matching is complete, the set of TP skeletons $\text{TP}^{\text{skl}}(\tau)$ with respect to the OKS threshold $\tau$ is determined. Concretely, for a fixed OKS threshold $\tau$ (ranging from 0.5 to 0.95), a pair of predicted and ground truth skeletons is identified as a TP, if their OKS is higher than $\tau$. The Precision (Pr) and Recall (Re) values are then computed as

$$\text{Pr}(\tau) := \frac{|\text{TP}^{\text{skl}}(\tau)|}{N_{\text{ply}}}, \quad \text{Re}(\tau) := \frac{|\text{TP}^{\text{skl}}(\tau)|}{N_{\text{ply}}^*}. \quad (11)$$

Finally, the Average Precision (AP) and Average Recall (AR) values read, respectively,

$$\text{AP} := \text{mean}_\tau \text{Pr}(\tau), \quad \text{AR} := \text{mean}_\tau \text{Re}(\tau). \quad (12)$$

**Remark 2** *The quality metrics above are defined per-image; however, in practice, we compute the bDQ, pSQ (7), pDQ (8), and AP and AR (12) over all the images in the validation or test sets.*

**Remark 3** *As required for computing pEQ, the values of $\kappa_k$ associated with the body parts are set according to the convention of DeepSport dataset,* i.e., $\kappa_{\text{head}} = 0.15$, $\kappa_{\text{hip}} = 0.2$, *and* $\kappa_{\text{foot 1}} = \kappa_{\text{foot 2}} = 0.2$.

# B  Ablation studies

## B.1  Decoding with Oracle Data

The importance of the accuracy of each output for the decoding process can be obtained by using the oracle data instead of their corresponding network outputs. This study is helpful to find out whether the error propagates from one block to the other. Table 5 compares the metrics when different permutations of oracle data were used on Deep-Sport dataset. The first message drawn from this study

| player centroid | offset vectors | semantic masks | PQ | pSQ | pDQ | pEQ | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | AP | AR | $F_1$ |
| – | – | – | 72.3 | 80.3 | 90.1 | 87.5 | 82.1 | 42.4 |
| ✓ | – | – | 71.4 | 79.7 | 89.6 | 87.2 | 82.7 | 42.5 |
| – | ✓ | – | 75.7 | 82.9 | 91.3 | 87 | 81.8 | 42.2 |
| – | – | ✓ | 86.7 | 94.5 | 91.7 | 87.2 | 82.5 | 42.4 |
| ✓ | ✓ | – | 77 | 83 | 92.8 | 86.2 | 82.3 | 42.1 |
| ✓ | – | ✓ | 87 | 94.5 | 92.1 | 87.7 | 83.7 | 42.8 |
| – | ✓ | ✓ | 93.5 | 98.3 | 95.1 | 88 | 82.7 | 42.6 |
| ✓ | ✓ | ✓ | 100 | 100 | 100 | 87.7 | 83.8 | 42.9 |

Table 5: **DeepSportLab Decoding with oracle data.** Note that $PQ := pSQ \cdot pDQ$.

| Method | PQ | pSQ | pDQ | pEQ | | |
|---|---|---|---|---|---|---|
| | | | | AP | AR | $F_1$ |
| DeepSportLab | 34.3 | 75.3 | 45.5 | 43.7 | 44.2 | 22 |
| DeepSportLab - oracle center | 52.1 | 78.1 | 66.7 | 57.1 | 60.7 | 29.4 |
| DeepSportLab - oracle segmentation | 100 | 100 | 100 | 63.5 | 67.2 | 32.6 |
| OpenPifPaf [26] | – | – | – | 66.9 | 70.9 | 45.4 |
| Pan.-DeepLab [6] | 48.4 | 78.6 | 61.5 | – | – | – |

Table 6: **Comparison of different methods evaluated on COCO's validation set.** Three different cases are considered for DeepSportLab: (1) Decoding with network's outputs, (2) Decoding with oracle centers, and (3) Decoding with oracle segmentation masks.

is that our player segmentation is good enough to associate the PIF keypoints because the increase of $F_1$ Score is only by 0.5% when using all of the oracle data compared to when using none. Thus, for further improve the pose estimation task, the PIF keypoints should be trained better.

## B.2   Keypoints Error Breakdown

Further improvements can be achieved once we know the source of error. Ronchini and Perona [37] break the estimated body keypoints in 5 different categories based on their calculated KS, *i.e.* keypoint similarity between the keypoint $o$ of a detection $\Upsilon$ and $o^*$ of an annotation $\Upsilon^*$. KS is calculated using (9) without the mean over all keypoints. If KS of $o$ and $o^*$ is higher than 0.85, this prediction is considered Good. Jitter happens when KS drops between 0.5 and 0.85. In case that KS is less than 0.5, $o$ can be either a Miss, Swap, or Inversion. In our case, since we switch the right and left feet in case of wrong detection, the Inversion will never occur. This is because foot 1 and foot 2 labels are interchangeable in DeepSport dataset (See Section 4 in the main document). Next, Swap happens when $o$ is wrongly associated to another skeleton. Miss happens when $o$ is predicted, but not in the right location, and it was not a Swap. Finally, FN KP happens when the keypoint is not detected at all. Fig. 4 shows the examples for each of these categories. Fig. 5 depicts the error breakdown based on the error category and type of keypoint.

## B.3   Evaluation on COCO dataset

As stated in the main text, our main contribution is to come up with a multi-task framework specific to sports scenes. However, as an ablation study, the model was also evaluated on COCO dataset [28] which is very much diverse in terms of both the scenery and the size of people in images. Table 6 shows metrics evaluated on COCO's validation set. Note that in this experiment, only the keypoints visible in the image are considered for the evaluation.

DeepSportLab decoder is studied in three different cases: (1) When using the network outputs, (2) When using oracle centroid of humans, and (3) When using the full human mask oracle. In the first case, due to the diversity of people size in COCO images, the segmentation task falls short in terms of pDQ which leads to error in pEQ. When adding the oracle center, the PQ increases significantly, suggesting that the center needs more training. When using the oracle masks, (*i.e.* PQ = 100), pEQ increases by 45.3% and 34.2% in terms of AP and AR, respectively. This shows the importance of the segmentation masks on big and challenging datasets such as COCO. It is worth mentioning that training in this case needs a lot of hyper-parameter tuning and optimization. Our computational resources certainly did not allow to fully explore the parameter space. Obtaining more competitive results on COCO dataset is seen as a future work for this framework.
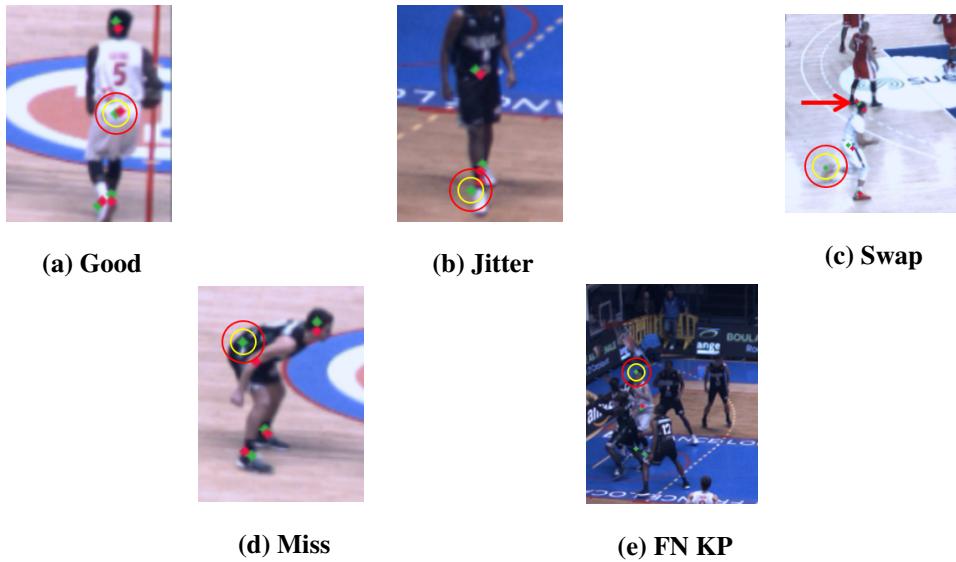
Figure 4: **Error samples.** Green and red dots show the annotated and predicted keypoints, respectively. Yellow and Red circles resemble the borders from which the KS will be less than 0.85 and 0.5, respectively. In (c), the red arrow points toward the wrongly predicted foot where the swap occurs.
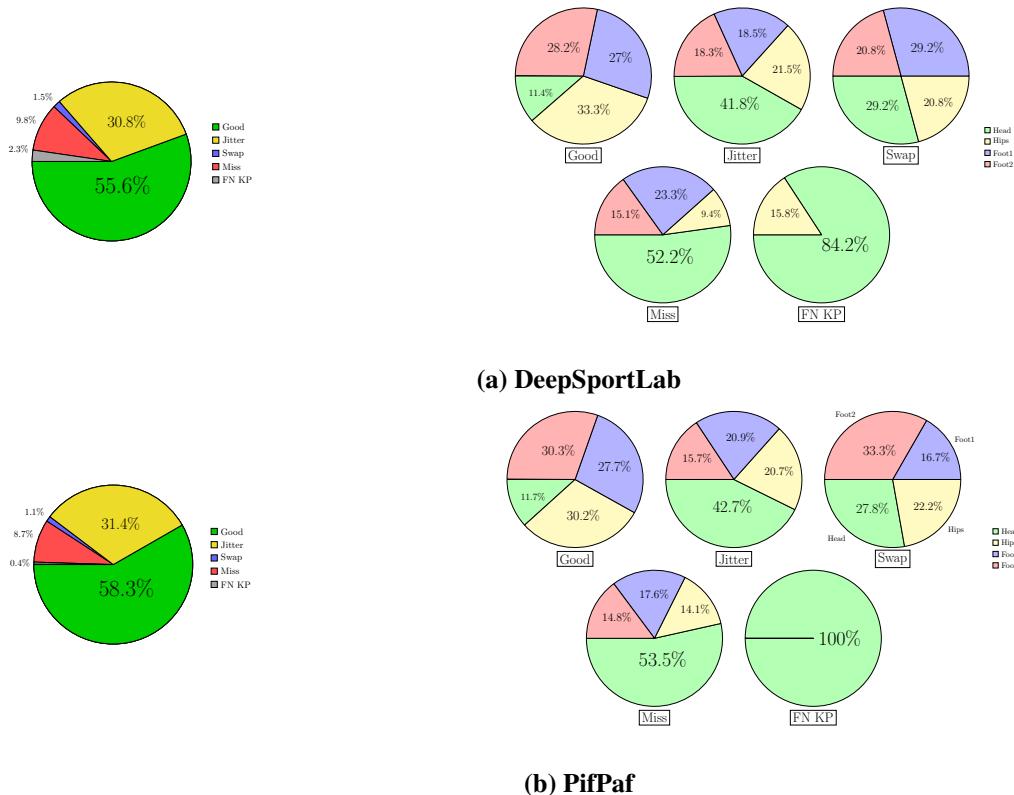


Figure 5: **Error breakdown.** Pie charts on the left show the distribution of keypoints in 5 categories based on their KS. Pie charts on the right show the distribution of each type of error based on the keypoints type.