

# Multi Player Tracking in Ice Hockey with Homographic Projections

Harish Prakash, Jia Cheng Shang, Ken M. Nsiempba, Yuhao Chen, David A. Clausi, John S. Zelek

*Vision and Image Processing Group*

*University of Waterloo, Ontario, Canada*

{harish.prakash, jcshang, kmnsiemp, yuhao.chen1, dclausi, jzelek} @ uwaterloo.ca

**Abstract**—Multi Object Tracking (MOT) in ice hockey pursues the combined task of localizing and associating players across a given sequence to maintain their identities. Tracking players from monocular broadcast feeds is an important computer vision problem offering various downstream analytics and enhanced viewership experience. However, existing trackers encounter significant difficulties in dealing with occlusions, blurs, and agile player movements prevalent in telecast feeds. In this work, we propose a novel tracking approach by formulating MOT as a bipartite graph matching problem infused with homography. We disentangle the positional representations of occluded and overlapping players in broadcast view, by mapping their foot keypoints to an overhead rink template, and encode these projected positions into the graph network. This ensures reliable spatial context for consistent player tracking and unfragmented tracklet prediction. Our results show considerable improvements in both the *IDsw* and *IDF1* metrics on the two available broadcast ice hockey datasets.

**Keywords**—Ice hockey; Tracking; Homography; MPNs.

## I. INTRODUCTION

Multi-Object Tracking (MOT) is an important computer vision problem subsuming several tasks such as object detection, localization, re-identification and association across a temporal sequence. It is a highly studied and established problem due to its plethora of applications in robotics, autonomous vehicles, industrial automation, surveillance, and sports. While most existing MOT approaches focus exclusively on tracking crowded pedestrians [1], [2], [3], group dancing [4], and autonomous driving [5], sports tracking is a pivotal task in vision due to its numerous subsequent applications in game analytics & statistics, strategic planning, player evaluation, injury prevention, and crucial game decisions. Player tracking helps save several manual labor hours and human efforts by automating game understanding and player performance assessments. With the advent of deep networks[6], several important strides have been taken to track various sports including soccer [7], handball [8], basketball [9], [10], [11], and volleyball [12], with public MOT datasets [13], [14] to support principled evaluations.

Unlike all the aforementioned sports, ice hockey poses unique challenges to tracking due to its highly *physical* and *fast-paced* nature. Specifically there are three major challenges that exists: (i) the significant occlusion between two or more players at a given instant within the field-of-view (FoV); (ii) the non-linear player dynamics due to

unpredictable player motion, and; (iii) blurs and reduced visibility of players due to frequent camera motion. When faced with these issues, tracking in monocular view increases identity swaps and tracklet fragmentations. Early attempts at tracking ice hockey players were pursued using an ensemble of handcrafted methods. Okuma et al. [15] use Adaboost [16] detection with mixed particle filters [17] for tracking players from television videos. Cai et al. [18] improve upon [15] by utilizing the mean-shift algorithm to stabilize player trajectories, and use rink coordinates (homography) for the particle-filter. However, mixed particle filters are susceptible to identity switches/losses during mutual occlusions, background changes, blurs, and lighting effects, which are often found in hockey. Further, there exists no quantitative evaluation in the above works to show the efficacy of their models.

Recent approaches using deep networks have shown significant improvements in the accuracy of tracking hockey players from broadcast feeds. Vats et al. [19] present a comparison between five different tracking models, fine-tuned on broadcast ice hockey clips and obtain state-of-the-art results using graphs with message passing networks (MPN) [20]. They subsequently use this method to generate tracklets (sequence of player tracks) for player identification & team recognition [21]. As far as we know, this is the only existing benchmark for MOT in ice hockey. However, there exists two major limitations in their approach: first, their model encodes the *bounding box* attributes of players as graphical edge embeddings which leads to identity swaps and tracklet fragmentations. This is because, when there exists heavy occlusion between players as usually observed in broadcast hockey feeds, their bounding boxes overlap significantly in the monocular view (Intersection-over-Union (IoU)  $\uparrow$ ), causing either misassociation of players (identity switch) or a missed connection (lost tracklet) with previous tracks. Second, their node embeddings are encoded solely based on player appearance features, which is ambiguous in hockey due to the fully-covered bulk gear worn by players, similarly colored team jerseys, blur & lighting effects. With these present setbacks, we ask the question: “Given only a monocular broadcast feed, is it possible to declutter occluded players and track their movements with high fidelity?”. Our results show that the answer is **Yes**.

In this work, we formulate MOT as a link prediction

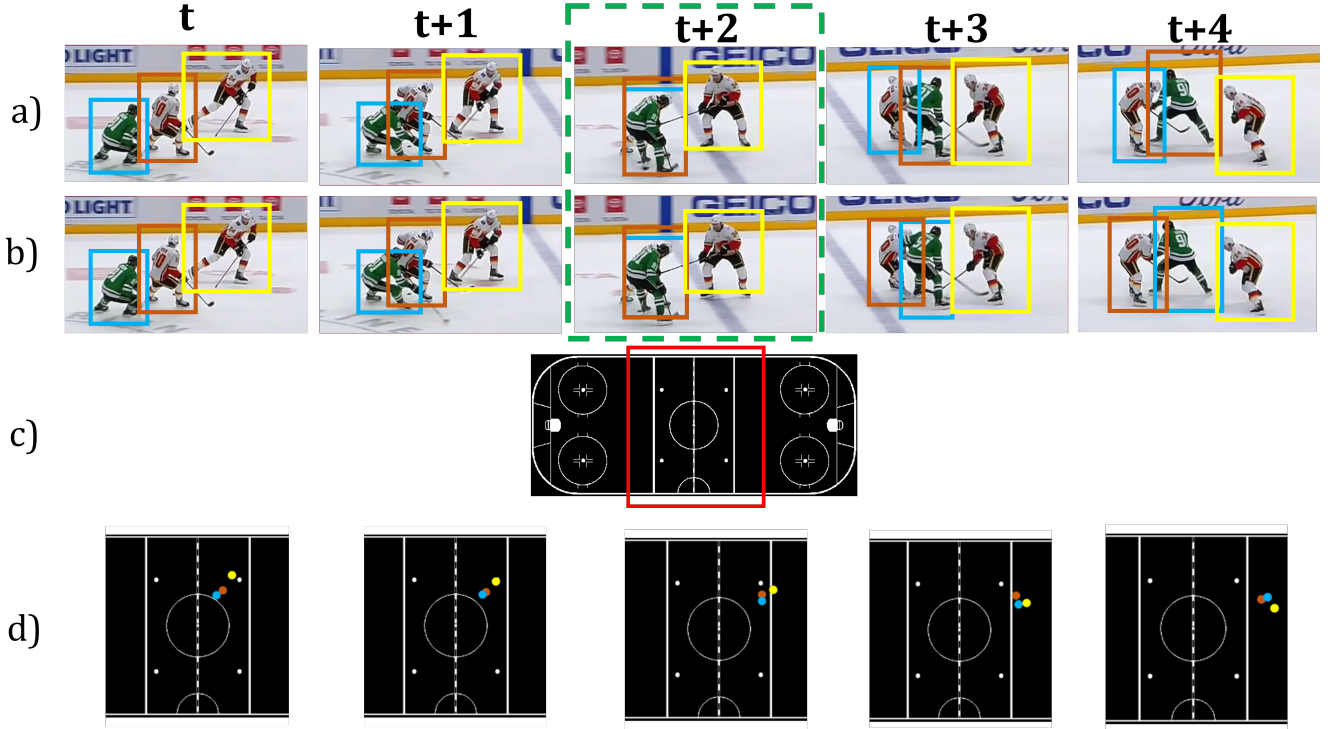


Figure 1. Other trackers vs Our approach. a) During a significant occlusion scenario ( $t + 2$ ), other trackers lead to an ID switch error; b) Our method consistently tracks players before and after occlusion; c) Overhead rink template used for homography projection; d) Player footpoint coordinates mapped to the overhead template. At ( $t + 2$ ), there is a clear distinction between overlapping players from the top view. This information aids our tracker in maintaining player tracklets.

problem in the graph domain since graphical networks offer a natural way of representing players and their relationships. The novelty of our approach lies in coupling this with homography to obtain reliable spatial/positional cues for ice hockey players. Specifically, we map every frame from broadcast view to an overhead view (rink template) using an off-the-shelf homography estimator [22], to obtain the projected footpoints of players in rink coordinates. This yields a *birds-eye* effect, reducing the positional ambiguities due to overlapping players and perspective projection. We encode these projected positions along with player re-identification features as graph embeddings, and propagate information through the graph using the message passing network (MPN) [23], [24] framework. Our model follows the popular *tracking-by-detection* paradigm and we show evaluations with both ground-truth detections and off-the-shelf detector outputs [25], in coherence with the current state-of-the-art (SOTA) benchmark [19]. Through this approach, our model is able to utilize the overhead positional information of players as additional cues to track with consistency during occlusions, blurs, and dynamic player movements.

Our contributions can be summarized as follows:

- We design a simple, yet effective spatio-temporal graphical network and adopt the MPN framework to

track players from broadcast feeds consistently;

- We propose a novel approach based on homography to provide additional positional information to the graph network during occlusions, blurs and non-linear movements, and;
- We show significant improvements in the  $ID_{sw}$  and  $IDF1$  scores on two available broadcast ice hockey datasets.

## II. RELATED WORKS

### A. MOT for Pedestrians

Most existing MOT methods [20], [26], [27], [28], [29], [30], [31], [32], [33] focus exclusively on tracking pedestrians in crowded scenes [1], [2], [3]. Initial approaches combine Kalman filters [34] and Hungarian method [35] for next state motion prediction and object associations respectively. Consequent approaches adopted the two-stage *tracking-by-detection* (TBD) paradigm where: first, all objects present in a sequence are detected; and next, their association features are extracted to link similar objects across frames. Simple Online and Real-time tracking (SORT) [31] establish a TBD baseline for MOT, and argue that SOTA associations can only be obtained with traditional methods [34], [35]. DeepSORT [32] answer SORT's [31] argument, by embedding

appearance features using a ReID network for association, showing lesser identity switches and better tracking results on the MOT16 [2] challenge benchmark. FairMOT [33] combines the detection and ReID steps for Joint Detection and Tracking (JDT) using an anchor-free detector [36]. Tracktor [29] frames tracking as a bounding box regression problem, by converting a detector [25] into a tracker. But, a major limitation with all these methods is their inability to tackle crowded scenes and declutter occluded pedestrians. ByteTrack [28] tries to handle occlusion by associating low-confidence targets too, to retrieve true objects, but fails during significant crowded (overlapping) situations where the object is heavily obstructed. SparseTrack [37] tries to handle occlusion by decomposing dense/crowded pedestrian scenes into sparse subsets using pseudo-depth map estimation, but incur high computational costs and significant processing requirements.

### B. MOT for Sports

Initial methods in sports utilize Kalman filtering [38], [39] and particle filters [15], [18] for tracking, but were unable to preserve identities due to their linear motion model assumptions. Nillius et al. [40] encode player trajectories into track graphs and use a bayesian framework to predict the most likely configuration of player paths. Figueroa et al. [41] track players with multi-cameras by encoding their segmentation blobs as nodes and their relative distances as edges. With YOLOv2 [42] for detection, Acuna et al. [43] track basketball players using SORT [31], while Theagarajan et al. [44] track soccer players using DeepSORT [32]. However, naively extending methods originally designed for pedestrian tracking to sports is a non-trivial task, due to the domain-specific challenges in modeling arbitrary player movements, occlusions, fast camera motions, and perspective projection errors.

### C. Tracking with Graphs

Graph-based formulations offer a flexible way to model target movements, interactions and their relative features. Wang et al. [45] propose a graph-based MOT framework for joint detection and tracking. MOT neural solver [20] exploits the network flow formulation of MOT to define a message passing network for data association. Vats et al. [19] fine-tune the neural solver architecture on the broadcast ice hockey dataset to create the first tracking benchmark for hockey. Luna et al. [46] extend the graph domain to multi-camera tracking, and ReST [47] builds on top with a spatio-temporal network for online-tracking. Both these algorithms use the intrinsic camera parameters to map multiple camera views onto a common ground plane for additional positional cues. But, in our case, we cannot infer camera parameters from broadcast feeds directly. Therefore, we utilize an off-the-shelf homography model [22] trained on a top-view rink template, to map each broadcast frame to the overhead rink

coordinates and obtain homographic footpoint coordinates for the graph MPN.

## III. PROPOSED METHOD

Given an ice hockey broadcast feed, the objective of our work is to track multiple players consistently despite prevalent occlusions, blurs, and arbitrary player movements. We leverage Graphical Neural Networks (GNN) for this task (due to their efficiency in modeling relationships) to build a spatio-temporal graph and utilize homography for reliable positional information. Specifically, we project player foot keypoints onto a common overhead rink template and encode these projections along with ReID embeddings [48] as node features and their relative distances as edge features. Since we cannot obtain camera parameters from broadcast feeds directly, we use an off-the-shelf homography estimation model [22] specifically designed for ice hockey, to estimate the homography transformation matrix,  $H \in \mathbb{R}^{3 \times 3}$ . We utilize the message passing network (MPN) framework [23], [24] to propagate player features across the entire graph  $G$ , and update the node and edge embeddings at each message passing step. This helps our model reason globally over the entire sequence for predicting player trajectories. We frame the association between two consecutive frames  $f_i$  and  $f_j$  s.t.  $j > i$ , as a bipartite graph problem, and use a binary classifier with a sigmoid final layer to output association probabilities. During inference, we post-process each graph by pruning to remove low-confidence associations and solve many-to-one violations. Finally, we assign tracklet IDs to nodes wherein, if the node has a prior connection, it inherits the same ID or gets assigned a new ID otherwise. (Ref Fig. II-C)

### A. Problem Formulation

Consider a broadcast hockey sequence (frames)  $S_t = \{F_i \mid i = 1, 2, \dots, n\}$ , where  $n = \frac{t}{\Delta t}$ ;  $t$  = total duration of the sequence, and  $\Delta t$  = duration per frame. For each frame  $F_i$ , assuming that there exists at least one player to track, we have  $P_j^i$  players, where  $j \geq 1$ . For each player  $P$ , the ground truth annotation includes:

$$\{f^{id}, t^{id}, x, y, wd, ht, c, x^{proj}, y^{proj}\}$$

where  $f^{id}$  = frame ID,  $t^{id}$  = player ID (only used during training),  $\{x, y, wd, ht\}$  = bounding box coordinates,  $c$  = annotation confidence score, and  $(x^{proj}, y^{proj})$  = homography coordinates for player footpoints. As per the *tracking-by-detection* paradigm, player annotations are detections from an off-the-shelf detector for the inference stage. In our case, we show both evaluations: first, directly evaluating on the annotated ground truth which gives the best picture of our tracking performance (Ref. Section IV-C), and second, inference and evaluation on the detected output using F-RCNN [25] following the current SOTA [19] benchmark, for a fair comparison.

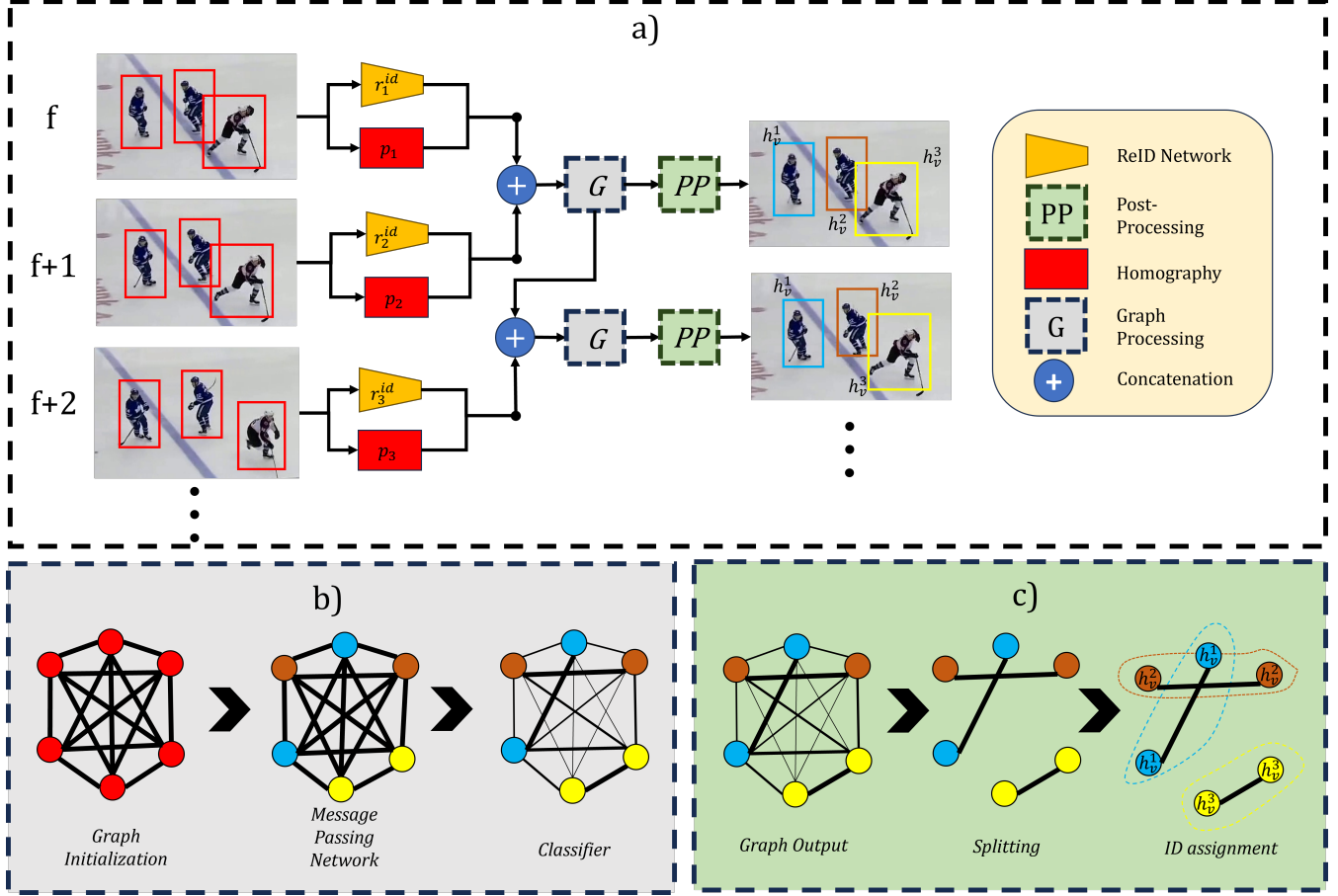


Figure 2. Proposed Approach. a) The general pipeline of our spatio-temporal graph. b)  $G$  denotes the three stages of Graph Initialization, MPN, and Classification. c)  $PP$  denotes the Post-Processing stage where we Prune, solve graph violations, and assign player IDs

We formulate multi-object tracking (MOT) as a bipartite graph matching problem. We deterministically create one graph per frame in the given sequence. Each node in this graph points to one player in that frame and the edges correspond to their relationships with neighboring nodes in other frames. Let us consider the undirected graph  $G_t = (V_t, E_t)$  with bidirectional connections, where  $V_t$  denotes the vertex set and  $E_t$  denotes the edge set at time  $t$ .  $NodeFeatures(\cdot)$  represents the concatenation of appearance features and homographic projections of players at time  $t$ , and  $E_t = V_i \times V_j$  represents a one-one mapping between vertex sets  $v_i$  and  $v_j$  s.t.  $i \neq j$ , meaning that the mapping is not within players in the same frame.

**Node Formulation** Each node  $v_i \in V_t$  represents one unique player  $i$  found in  $F_t$  and contains: frame ID (timestamp)  $f_i^{id} \in \mathbb{R}^1$ , player ID  $t_i^{id} \in \mathbb{R}^1$  (ground-truth, only used for training), bounding box coordinates  $b_i = \{x_i, y_i, wd_i, ht_i\} \in \mathbb{R}^4$ , ReID appearance features  $r_i^{id} \in \mathbb{R}^{512}$  and homographic projection coordinates  $p_i = \{x_i^{proj}, y_i^{proj}\} \in \mathbb{R}^2$ .

The ReID features are generated for each node  $i$  via an off-the-shelf Re-Identification network [48], described by:

$$r_i^{id} = ReID(b_i |_{crop}) \quad (1)$$

where  $b_i |_{crop}$  denotes the cropped bounding box for the  $i^{th}$  player. The homographic projections are estimated by projecting the bottom-mid point ( $\sim$ foot keypoint) of the bounding box from the monocular broadcast view. The left footpoint  $f_l = (x_i + \frac{wd_i}{2})$  and the right footpoint  $f_r = (y_i + ht_i)$  for player  $i$  are projected as:

$$p_i = H_i(f_l, f_r) \quad (2)$$

with  $H_i$  being the  $3 \times 3$  Homography matrix (Ref. Eq. 5)

**Edge Formulation** Each edge,  $e_{ij} \in E_t$  is represented as the interconnection between two players from two distinct frames. It is encoded as the concatenation of relative appearance  $\Delta r_{ij}^{id}$  and positional  $\Delta p_{ij}$  features between the pair of nodes  $v_i$  and  $v_j$ , where  $e_{ij} = v_i \times v_j, i \neq j$ . This can be represented as:



$$\Delta r_{ij}^{id} = [\|r_i^{id} - r_j^{id}\|_1, \text{cosine\_similarity}(r_i^{id}, r_j^{id})] \quad (3)$$

$$\Delta p_{ij} = [\|p_i - p_j\|_1, \|p_i - p_j\|_2] \quad (4)$$

$[\cdot, \cdot]$  denotes the concatenation of Euclidian distance & Cosine Similarity for  $\Delta r_{ij}^{id}$ , and Euclidian & Manhattan distance for  $\Delta p_{ij}$ . This consideration is inspired from [46] to obtain higher-dimensional distinctive features.

### B. Homographic Projection

Due to the monocular nature of broadcast ice hockey sequences, there exists high levels of player occlusion and dynamic camera movement effects (blurs, pans, tilts, zooms). This limits the scope of tracking players consistently when they're completely obscured or remain hidden, even if for very short intervals. To provide the tracker with reliable positional cues in such scenarios, we propose an approach using homography to warp player positions from the broadcast video feed onto an common overhead rink template. This helps reduce the variance in frames present across the sequence due to camera motion, and provides a pseudo *top-view* tracking effect to disentangle overlapping players.

At any given frame  $F_t$ , a player's  $P_i$  foot keypoint coordinates represents their exact point of contact with the ice, which when projected to the overhead rink plane provides additional positional cues for uncluttered tracking. For the player  $P$  with footpoints  $(P_{fx}, P_{fy})$  in broadcast view:

$$p_i = s \begin{bmatrix} P_{x'} \\ P_{y'} \\ 1 \end{bmatrix} = H \begin{bmatrix} P_{fx} \\ P_{fy} \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & 1 \end{bmatrix} \begin{bmatrix} P_{fx} \\ P_{fy} \\ 1 \end{bmatrix} \quad (5)$$

where  $p_i$  = homographic projection of the  $i^{th}$  node,  $(P_{x'}, P_{y'})$  = projected homogenous player footpoints in the overhead view,  $s$  = scale factor, and  $H = 3 \times 3$  homography matrix.

But this is a non-trivial task, since the camera parameters for broadcast feeds are unknown and thereby, we do not know the values for  $H$ . Therefore, we tackle this issue using an off-the-shelf homography-estimator [22] pre-trained on an ice hockey top-view rink template, to map each broadcast frame onto the overhead view and obtains its respective homographic projection.

### C. Temporal Graph Association

To facilitate association of players across frames, we design a simple temporal graph network, to correlate player features. Inspired by [47], we iteratively correlate the learned graph  $G_{t-1}^T$  at time  $t-1$  with the next graph  $G_t$  at time  $t$  to form a new graph  $G_t^T$ . Assuming this as the  $n^{th}$  iteration, each node  $v_{t-1} \in G_{t-1}^T$  contains aggregated embeddings from  $n-1$  iterations and is connected with the new nodes

$v_t \in G_t$  to form the new temporal graph  $G_t^T$  at time  $t$  (Ref algorithm 1). The edge set thus created can be denoted as:

$$e_{t-1,t} = \begin{cases} 1, & \text{if } v_i \times v_j \forall v_i \in G_{t-1}, v_j \in G_t \\ 0, & \text{otherwise} \end{cases}$$

This aggregation of uniterated graphs with learned graphs helps propagate learned features and assign consistent identities for the same player (Ref Fig. II-C)

---

#### Algorithm 1 Model Inference

---

- 1: Given g.t  $y_{e_{ij}}$ , players  $P_k \in G_{t-1} \oplus G_t$ ,  $i = 0$  &  $w = 2$ ;
  - 2: **for**  $t = 0$  to  $(T-1)$  **do**
  - 3:   Initialize  $h_{v_i}^0, h_{v_j}^0$  and  $h_{e_{ij}}^0$  for  $G_{t-1} \oplus G_t$
  - 4:   **if**  $i > 0$  **then**
  - 5:     Replace  $G_{t-1}$  with  $G_t^{learned}$
  - 6:   **end if**
  - 7:    $h_{v_i}^0, h_{v_j}^0 = f_v^{FE}(h_{v_i}^0, h_{v_j}^0)$
  - 8:    $h_{e_{ij}}^0 = f_e^{FE}(h_{e_{ij}}^0)$
  - 9:   **while**  $(l \geq L)$  **do**
  - 10:      $h_{v_i}^L, h_{v_j}^L, h_{e_{ij}}^L = \text{MPN}(G_{t-1} \oplus G_t, h_{v_i}^{l-1}, h_{v_j}^{l-1})$
  - 11:      $\hat{y}_{e_{ij}} = f_{cls}^L(h_{e_{ij}}^L)$
  - 12:   **end while**
  - 13:  $G_t^{learned} = \text{Post-Processing}(G_{t-1} \oplus G_t, \hat{y}_{e_{ij}}^L)$
  - 14:  $i++$ ;
  - 15: **end for**
- 

### D. Message Passing Network

We adopt the Message Passing Network (MPN) structure as introduced by [24] to propagate the node and edge information across the graph  $G$ . Message passing intuitively helps the graphs learn their neighboring features; each edge learns about the projection and appearance feature of its neighboring nodes and each node learns about the geometric features of its neighboring edges. To begin with, we initialize the node embeddings and edge embeddings as:

$$h_{v_i}^{(0)} = f_v^{FE}([r_{v_i}^{id}, p_{v_i}]) \quad (6)$$

$$h_{e_{ij}}^{(0)} = f_e^{FE}([\Delta r_{e_{ij}}^{id}, \Delta p_{e_{ij}}]) \quad (7)$$

Note that we add  $p$  and  $\Delta p$  to both the node and edge features respectively to propagate homographic (positional) information throughout the graph. Given these initial embeddings, as standardized by well-established methods [19], [20], [47], [23], [46], we perform  $L$  iterations of edge updates and node updates as two separate steps.

Network	Layer	Input	Output
$f_v^{FE}(\cdot)$	FC+GELU	515	128
		4	32
$f_e^{FE}(\cdot)$	FC+GELU	4	8
		8	6
$f_v^{ME}(\cdot)$	FC+GELU	38	64
		64	32
$f_e^{ME}(\cdot)$	FC+GELU	70	32
		32	6
$f^{cls}(\cdot)$	FC+GELU	6	4
	FC+Sigmoid	4	1

Table I  
DETAILS OF EACH MLP ENCODER.  
FC - FULLY CONNECTED, GELU ACTIVATION [49]

**Edge Update** We utilize a learnable multi-layer perceptron (MLP) to perform edge encoding for  $l = \{1, \dots, L\}$  steps using the source and destination nodes connected by the edge:

$$(v \rightarrow e) \quad h_{e_{ij}}^{(l)} = f_e^{ME}([h_{v_i}^{(l-1)}, h_{v_j}^{(l-1)}, h_{e_{ij}}^{(l-1)}]) \quad (8)$$

where,  $f_e^{ME}$  is the edge encoder. This leads to the sharing of appearance and projection embeddings from the neighboring nodes  $h_{v_i}, h_{v_j}$  to its connecting edge  $h_{e_{ij}}$ .

**Node Update** Similar to Eq. (8), we utilize a learnable MLP to perform node update for  $l = \{1, \dots, L\}$  steps, using the aggregated messages coming from its neighboring nodes:

$$(e \rightarrow v) \quad h_{v_i}^{(l)} = \sum_{j \in \mathcal{N}(v_i)} f_v^{ME}([h_{v_j}^{(l-1)}, h_{e_{ij}}^{(l-1)}]) \quad (9)$$

where,  $f_v^{ME}$  is the node encoder and  $\mathcal{N}(v_i)$  denotes the neighboring nodes of  $v_i$ . Note that  $f_v^{ME}$  and  $f_e^{ME}$  are two separate networks with different dimensions, but share the same MLP architecture (Ref. Table I)

In both the updates, the MLP encodes all the information into a higher-dimensional feature space. The message passing step  $L$  is akin to the receptive field found in convolutional neural networks [6] and a higher value of  $L$  corresponds to farther propagation of information in the graph, but at the cost of computation.

**Classification** We propose to learn the association between nodes as a *link prediction task* by framing player tracking as a graph partition problem. That is, after  $L$  iterations, we perform binary classification to predict the edge probabilities  $\hat{y}_{e_{ij}}$  connecting nodes  $v_i$  and  $v_j$ , as:

$$\hat{y}_{e_{ij}} = f^{cls}(h_{e_{ij}}) \quad (10)$$

where  $f^{cls}$  is a learnable MLP with a sigmoid final layer to output probabilities. During inference, the edges with

low confidence scores (weak connections) are removed. During training, the binary ground truth labels  $y_{e_{ij}}$  and their corresponding predictions  $\hat{y}_{e_{ij}}$  are compared to find the sigmoid focal loss [50].

#### E. Post-Processing

We adopt a post-processing step during inference to prune and resolve violations in our final graphs, and assign consistent tracklet IDs.

**Pruning.** This is the first step in refining the predicted edge confidence scores  $\hat{y}_{e_{ij}}$  by our classifier. We define a confidence threshold hyperparameter  $\xi$ , where:

$$\text{Pruned}(\hat{y}_{e_{ij}}) = \begin{cases} \hat{y}_{e_{ij}}, & \text{if } \hat{y}_{e_{ij}} > \xi \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

This helps eliminate the low-confidence edges and retain only the most strongest correspondences. For our experiments, we've used  $\xi = 0.9$  to only retain the most strongest edges. Alongside, similar to [47], we remove many-to-one edge violations, wherein, there can only exist at most one connection between two nodes within the connected graphs. This ensures that we represent each player identity through a unique connected component, since no two different players can share the same identity at the same time.

**Assigning IDs.** The final graph contains unique connected components where each component represents one unique player tracklet. As we iterate temporally, we assign a new identity when there exists no previous connected components for a node or propagate the same identity otherwise.

$$\text{ID}(h_{v_j}) = \begin{cases} \text{SameID}, & \text{if } h_{v_j} = h_{v_i}; \forall j \in G_t, i \in G_{t-1} \\ \text{NewID}, & \text{otherwise} \end{cases} \quad (12)$$

## IV. IMPLEMENTATION DETAILS

This section contains details about datasets used, training scheme, evaluation metrics and the final results.

#### A. Datasets

We experiment our methodology on the two available ice hockey tracking datasets - first, similar to the current SOTA benchmark [19], we train and test on the broadcast hockey dataset for a fair comparison; second, we evaluate on the public VIP Hockey Tracking Dataset (VIP-HTD) [51] to demonstrate the generalization of our method. Both datasets contain side-of-the-rink broadcast videos with occlusions, blurs, and challenging player movements.

**Broadcast Tracking Dataset [19]** This dataset contains 84 broadcast clips sampled from 25 NHL games, with an average duration of  $\sim 36$  seconds per clip. The dataset has a  $1280 \times 720p$  resolution (Standard Definition) at a frame rate of 30fps, with a train:validation:test ratio of 58:13:13. We follow the same training and testing scheme as [19] for equal

comparison, and show superior results with our method (Ref Table II)

**VIP-HTD [51]** This *public* dataset contains 22 broadcast hockey clips sampled from 8 NHL games, with both 30 & 60Hz frame rates, recorded at  $1280 \times 720p$  resolution. We perform cross-dataset validation (trained on broadcast dataset; tested on VIP-HTD) on all the 7 test clips in this dataset to showcase the generalizability of our method for any given broadcast hockey feed (Ref Table III)

### B. Training Details.

Given the player bounding boxes and tracklet IDs, we find the footpoint projection coordinates using an off-the-shelf homography model [22] trained specifically for NHL ice-hockey rinks. This model is currently the SOTA for hockey and helps predict highly accurate overhead projections. Next, we exploit the OSNet architecture [48] as our ReID network for player feature extraction (Eq. (1)), pre-trained on the ImageNet dataset [52]. We encode the 512-D ReID feature vectors along with the 3-D (homogenous coordinates) homography features into 32-D node embeddings (Eq. (6)), and the 4-D edge features (Eq. (7)) into 6-D edge embeddings. We run the MPN for  $L = 6$  iterations and pass the final graph output into our binary classifier for predicting edge probabilities.

---

#### Algorithm 2 Model Training

---

```

1: Given g.t  $y_{e_{ij}}$ , players  $P_k \in G_{t-1} \oplus G_t$  &  $w = 2$ ;
2: for  $t = 1$  to  $(T - 1)$  do
3:   Initialize  $h_{v_i}^0, h_{v_j}^0$  and  $h_{e_{ij}}^0$  for  $G_{t-1} \oplus G_t$ 

4:    $(h_{v_i}^0, h_{v_j}^0) = f_v^{FE}(h_{v_i}^0, h_{v_j}^0)$ 

5:    $h_{e_{ij}}^0 = f_e^{FE}(h_{e_{ij}}^0)$ 

6:   while  $l \leq L$  do:
        $h_{v_i}^L, h_{v_j}^L, h_{e_{ij}}^L = \text{MPN}(G_{t-1} \oplus G_t, h_{v_i}^{l-1}, h_{v_j}^{l-1})$ 
        $\hat{y}_{e_{ij}} = f^{cls}(h_{e_{ij}}^L)$ 
7:   end while
8:    $\text{graph\_loss} = \text{sigmoid\_focal\_loss}(\hat{y}_{e_{ij}}, y_{e_{ij}})$ 
9:    $\text{graph\_loss.backward}()$   $\triangleright$  Backpropagation
10:   $\text{optimizer.step}()$   $\triangleright$  Update parameters
11: end for

```

---

During training, we utilize the ground-truth player annotations and calculate the prediction losses using Focal Loss [50]

$$\text{Focal Loss} = \sum_l \sum_{e_{ij} \in G_{t-1} \cup G_t} \mathcal{L}(\hat{y}_{e_{ij}}^l, y_{e_{ij}}) \quad (13)$$

where,  $\hat{y}_{e_{ij}}^l$  is the edge prediction at iteration  $l$  and  $y_{e_{ij}}$  is the ground-truth indicator function:

$$\mathbb{1}(y_{e_{ij}}) = \begin{cases} 1, & \text{if } v_i = v_j \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

that is, for every player match,  $y_{e_{ij}} = 1$  else 0. We use Adam Optimizer [53] without weight decays to update our model parameters. The learning rate (LR) is initialized at 0.01, with a gradual warmup for the first 10 epochs and a cosine annealing schedule (min. LR = 0.001) thereafter. We trained our model for 30 epochs with a batch size of 16 on a single NVIDIA GeForce RTX 4090 GPU with 24 GB RAM, 2.5GHz clock speed, and performed validation after every 2 training epochs.

### C. Evaluation Metrics

Most common evaluation metrics used in popular SOTA tracking methods are the Multi-Object Tracking Accuracy (MOTA)[54] and IDF1 score [55]. With respect to our problem context, they can be defined as:

**MOT Accuracy:** It is calculated as the complement of three distinct errors -

- False Positives (FP): No. of false players detected;
- False Negatives (FN): No. of true players missed;
- Identity Switches (IDsw): No. of identity swaps/re-initializations made for players within the field-of-view.

$$\text{MOTA} = 1 - \frac{\sum_t \text{FN}_t + \text{FP}_t + \text{IDsw}_t}{\sum_t \text{GT}_t} \quad (15)$$

where,  $\text{GT}_t$  denotes the ground-truth annotations.

**IDF1 Score:** It is defined as the ratio of correctly identified players over the average number of ground-truth and computed identities:

$$\text{IDF1} = 2 \times \frac{\text{TP}_{id}}{(2 \times \text{TP}_{id}) + \text{FP}_{id} + \text{FN}_{id}} \quad (16)$$

where,  $\text{TP}_{id}, \text{FP}_{id}, \text{FN}_{id}$  are True Positive, False Positive and False Negative player identities respectively. Alternatively, IDF1 score can also be defined as the harmonic mean of ID Precision and ID Recall.

The FPs and FNs used to calculate MOTA relies solely on the detector's quality. Even if the tracker consistently associates players, the MOTA will be skewed if there exists high FP and FN detections, as they have twice the weightage of IDsw in MOTA. Since this doesn't give a clear picture of the tracker's association capabilities, we focus only on the IDsw score and the IDF1 score as the *key* metrics in player tracking. These metrics are especially relevant in ice hockey, since they measure how consistently a player is tracked with respect to his original identity. Therefore, our primary objective is to have  $\downarrow$  IDsw and  $\uparrow$  IDF1 score for consistent player tracking. Our preferred evaluations are directly based on ground-truth annotations; but, to be consistent with [19] we use similar detection outputs to show our results.

Method	MOTA % $\uparrow$	FP $\downarrow$	FN $\downarrow$	IDsw $\downarrow$	IDF1 $\uparrow$ %
SORT [31]	92.4	2403	5826	673	53.7
DeepSORT [32]	94.2	1881	4334	528	59.3
FairMOT [33]	91.9	<b>1179</b>	7568	768	61.5
Tracker [29]	94.4	1706	<b>4216</b>	687	56.5
Hockey MOT [19]	94.5	1653	4394	<b>431</b>	62.9
Hockey MOT <sup>†</sup> [19]	-	-	-	1056	71.8
<b>Ours</b>	<b>95.4</b>	1924	4323	453	<b>71.3</b>
<b>Ours<sup>†</sup></b>	-	-	-	<b>151</b>	<b>95.1</b>

Table II  
EVALUATION RESULTS ON THE BROADCAST ICE HOCKEY DATASET. <sup>†</sup> INDICATES GROUND-TRUTH ANNOTATIONS USED FOR EVALUATION

NHL teams (clips)	No.of frames	FPS (Hz)	IDsw $\downarrow$	IDF1 $\uparrow$
CAR vs. BOS	2606	60	115	85.4
CAR vs. NYR	2137	60	71	81.1
CGY vs. DAL	1330	60	41	77.4
CHI vs. TOR	1618	30	152	60.3
CNTRL vs. PAC	1671	30	136	71.7
PIT vs. SJ	1391	30	164	70.1
STL vs. SJ	2258	60	108	86.2
Hockey MOT			787	80.2
CAR vs. BOS	2606	60	11	89.3
CAR vs. NYR	2137	60	7	93.1
CGY vs. DAL	1330	60	6	92.4
CHI vs. TOR	1618	30	6	90.9
CNTRL vs. PAC	1671	30	13	90.5
PIT vs. SJ	1391	30	7	97.4
STL vs. SJ	2258	60	7	96.3
<b>Ours</b>	13,011	-	<b>60</b>	<b>92.84</b>

Table III  
CROSS-DATASET VALIDATION ON THE PUBLIC VIP-HOCKEY TRACKING DATASET

#### D. Results

We report the results of our model’s performance on the test-sets of the broadcast ice hockey dataset [19] and the public VIP-HTDataset [51]. It is to be noted that we reproduce the results of our benchmark [19], under similar hardware and testing conditions as our own method’s evaluations, to avoid any discrepancies. From Table II during ground-truth evaluations, our model outperforms the SOTA model by a large 23.3%  $\uparrow$  in *IDF1* score, and  $10\times \downarrow$  in *IDsw*. This is due to the ability of our model to handle heavy occlusions and blurs prevalent in these videos. We see similar trends with the F-RCNN [25] detection inputs, where our model surpasses all methods by a 8.4% $\uparrow$  in *IDF1*. Our MOT Accuracy is still higher than all other methods, despite incurring more FPs and FNs due to the ability of our tracker to recover robustly from mistakes made by the detector. We show qualitative results for all 7 videos below the reference section. In Table III, we cross-validate our model on the VIP-HTDataset [51] showing a clear superiority in both *IDsw* and *IDF1* metrics, compared to [19]. This asserts two important things: (i) our model generalizes well to unseen hockey feeds, despite of varying environmental conditions; (ii) Our model’s performance isn’t affected by the  $\uparrow$  in frame rate.

#### V. CONCLUSION

We present a novel approach based on the combination of graphical neural networks and homography to effectively track ice hockey players in broadcast feeds. We project player footpoints to an overhead rink template to maintain consistent positional cues, especially during occlusions and blurry situations. This provides a pseudo ‘top-view’ effect to disentangle overlapping players and maintain their trajectories. Message passing network (MPN) is used to aggregate player features and model their temporal relationships, followed by a classifier to predict player association probabilities. We achieve as significant  $\uparrow$  in *IDF1* and  $\downarrow$  in *IDsw*, when compared to both the current SOTA benchmark and a public tracking dataset. We believe that our work can also benefit various other sports in the future.

#### ACKNOWLEDGMENT

This work was supported by Stathletes through the Mitacs Accelerate Program and the Natural Sciences and Engineering Research Council of Canada (NSERC).

## REFERENCES

- [1] L. Leal-Taixé, A. Milan, I. Reid, S. Roth, and K. Schindler, "MOTChallenge 2015: Towards a benchmark for multi-target tracking," 2015.
- [2] A. Milan, L. Leal-Taixé, I. Reid, S. Roth, and K. Schindler, "MOT16: A benchmark for multi-object tracking," 2016.
- [3] P. Dendorfer, H. Rezatofighi, A. Milan, J. Shi, D. Cremers, I. Reid, S. Roth, K. Schindler, and L. Leal-Taixé, "MOT20: A benchmark for multi object tracking in crowded scenes," 2020.
- [4] P. Sun, J. Cao, Y. Jiang, Z. Yuan, S. Bai, K. Kitani, and P. Luo, "Dancetrack: Multi-object tracking in uniform appearance and diverse motion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [5] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [6] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional neural networks," *Neural Information Processing Systems*, vol. 25, 01 2012.
- [7] M. Manafifard, H. Ebadi, and H. A. Moghaddam, "A survey on player tracking in soccer videos," *Computer Vision and Image Understanding*, vol. 159, pp. 19–46, 2017.
- [8] M. Buric, M. Ivasic-Kos, and M. Pobar, "Player tracking in sports videos," in *2019 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, 2019, pp. 334–340.
- [9] M.-C. Hu, M.-H. Chang, J.-L. Wu, and L. Chi, "Robust camera calibration and player tracking in broadcast basketball video," *IEEE Transactions on Multimedia*, vol. 13, no. 2, pp. 266–279, 2011.
- [10] W.-L. Lu, J.-A. Ting, J. J. Little, and K. P. Murphy, "Learning to track and identify players from broadcast sports videos," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 7, pp. 1704–1716, 2013.
- [11] H.-T. Chen, C.-L. Chou, T.-S. Fu, S.-Y. Lee, and B.-S. P. Lin, "Recognizing tactic patterns in broadcast basketball video using player trajectory," *Journal of Visual Communication and Image Representation*, vol. 23, no. 6, pp. 932–947, 2012.
- [12] M. Takahashi, K. Ikeya, M. Kano, H. Ookubo, and T. Mishina, "Robust volleyball tracking system using multi-view cameras," in *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 2016, pp. 2740–2745.
- [13] Y. Cui, C. Zeng, X. Zhao, Y. Yang, G. Wu, and L. Wang, "Sportsmot: A large multi-object tracking dataset in multiple sports scenes," 2023.
- [14] A. Cioppa, S. Giancola, A. Deliege, L. Kang, X. Zhou, Z. Cheng, B. Ghanem, and M. Van Droogenbroeck, "Soccernet-tracking: Multiple object tracking dataset and benchmark in soccer videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3491–3502.
- [15] K. Okuma, A. Taleghani, D. Freitas, J. Little, and D. Lowe, "A boosted particle filter: Multitarget detection and tracking," vol. 3021, 05 2004.
- [16] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, vol. 1, 2001, pp. I–I.
- [17] Vermaak, Doucet, and Perez, "Maintaining multimodality through mixture tracking," in *Proceedings Ninth IEEE International Conference on Computer Vision*, 2003, pp. 1110–1116 vol.2.
- [18] Y. Cai, N. Freitas, and J. Little, "Robust visual tracking for multiple targets," 05 2006, pp. 107–118.
- [19] K. Vats, M. Fani, D. A. Clausi, and J. S. Zelek, "Evaluating deep tracking models for player tracking in broadcast ice hockey video," 2022.
- [20] G. Brasó and L. Leal-Taixé, "Learning a neural solver for multiple object tracking," 2020.
- [21] K. Vats, P. Walters, M. Fani, D. A. Clausi, and J. Zelek, "Player tracking and identification in ice hockey," 2021.
- [22] J. C. Shang, Y. Chen, M. J. Shafiee, and D. A. Clausi, "Rink-agnostic hockey rink registration," 2023.
- [23] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu, "Relational inductive biases, deep learning, and graph networks," 2018.
- [24] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," 2017.
- [25] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *CoRR*, vol. abs/1506.01497, 2015. [Online]. Available: <http://arxiv.org/abs/1506.01497>
- [26] Y.-H. Wang, J.-W. Hsieh, P.-Y. Chen, M.-C. Chang, H. H. So, and X. Li, "Smiletrack: Similarity learning for occlusion-aware multiple object tracking," 2024.
- [27] K. Yi, K. Luo, X. Luo, J. Huang, H. Wu, R. Hu, and W. Hao, "Ucmtrack: Multi-object tracking with uniform camera motion compensation," 2024.
- [28] Y. Zhang, P. Sun, Y. Jiang, D. Yu, F. Weng, Z. Yuan, P. Luo, W. Liu, and X. Wang, "Bytetrack: Multi-object tracking by associating every detection box," 2022.
- [29] P. Bergmann, T. Meinhardt, and L. Leal-Taixé, "Tracking without bells and whistles," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, oct 2019. [Online]. Available: <https://doi.org/10.1109/2Ficcv.2019.00103>

- [30] P. Chu, J. Wang, Q. You, H. Ling, and Z. Liu, "Transmot: Spatial-temporal graph transformer for multiple object tracking," *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 4859–4869, 2021.
- [31] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," *CoRR*, vol. abs/1602.00763, 2016. [Online]. Available: <http://arxiv.org/abs/1602.00763>
- [32] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," *CoRR*, vol. abs/1703.07402, 2017. [Online]. Available: <http://arxiv.org/abs/1703.07402>
- [33] Y. Zhang, C. Wang, X. Wang, W. Zeng, and W. Liu, "Fairmot: On the fairness of detection and re-identification in multiple object tracking," *International Journal of Computer Vision*, vol. 129, no. 11, 2021. [Online]. Available: <http://dx.doi.org/10.1007/s11263-021-01513-4>
- [34] R. E. Kalman, "A new approach to linear filtering and prediction problems," 1960.
- [35] H. Kuhn, "The hungarian method for the assignment problem," *Naval Research Logistic Quarterly*, vol. 2, 05 2012.
- [36] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *CoRR*, vol. abs/1904.07850, 2019. [Online]. Available: <http://arxiv.org/abs/1904.07850>
- [37] Z. Liu, X. Wang, C. Wang, W. Liu, and X. Bai, "Spasetrack: Multi-object tracking by performing scene decomposition based on pseudo-depth," 2023.
- [38] S. Iwase and H. Saito, "Parallel tracking of all soccer players by integrating detected positions in multiple view images," in *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, vol. 4, 2004, pp. 751–754 Vol.4.
- [39] M. Xu, J. Orwell, and G. Jones, "Tracking football players with multiple cameras," in *2004 International Conference on Image Processing, 2004. ICIP '04.*, vol. 5, 2004, pp. 2909–2912 Vol. 5.
- [40] P. Nillius, J. Sullivan, and S. Carlsson, "Multi-target tracking - linking identities using bayesian network inference," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2, 2006, pp. 2187–2194.
- [41] P. Figueroa, N. Leite, R. Barros, I. Cohen, and G. Medioni, "Tracking soccer players using the graph representation," in *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, vol. 4, 2004, pp. 787–790 Vol.4.
- [42] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 779–788, 2015.
- [43] D. Acuna, "Towards real-time detection and tracking of basketball players using deep neural networks." [Online]. Available: <https://api.semanticscholar.org/CorpusID:31248790>
- [44] R. Theagarajan and B. Bhanu, "An automated system for generating tactical performance statistics for individual soccer players from videos," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 2, pp. 632–646, 2021.
- [45] Y. Wang, K. Kitani, and X. Weng, "Joint object detection and multi-object tracking with graph neural networks," 2021.
- [46] E. Luna, J. C. SanMiguel, J. M. Martínez, and P. Carballeira, "Graph neural networks for cross-camera data association," 2022.
- [47] C.-C. Cheng, M.-X. Qiu, C.-K. Chiang, and S.-H. Lai, "Rest: A reconfigurable spatial-temporal graph model for multi-camera multi-object tracking," 2023.
- [48] K. Zhou, Y. Yang, A. Cavallaro, and T. Xiang, "Omni-scale feature learning for person re-identification," in *ICCV*, 2019.
- [49] D. Hendrycks and K. Gimpel, "Gaussian error linear units (gelus)," 2023.
- [50] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," 2018.
- [51] H. Prakash, Y. Chen, S. Rambhatla, D. Clausi, and J. Zelek, "Vip-htd: A public benchmark for multi-player tracking in ice hockey," in *Computer Vision and Intelligent Systems*, 2024.
- [52] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [53] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2017.
- [54] K. Bernardin and R. Stiefelwagen, "Evaluating multiple object tracking performance: The clear mot metrics," *EURASIP Journal on Image and Video Processing*, vol. 2008, 01 2008.
- [55] E. Ristani, F. Solera, R. S. Zou, R. Cucchiara, and C. Tomasi, "Performance measures and a data set for multi-target, multi-camera tracking," 2016.