# Ice hockey player identification via transformers and weakly supervised learning

Kanav Vats    William McNally    Pascale Walters[†]
David A. Clausi    John S. Zelek
Systems Design Engineering, University of Waterloo    [†]Stathletes Inc
{k2vats, wmcnally, dclausi,jzelek}@uwaterloo.ca    pascale.walters@stathletes.com

## Abstract

*Identifying players in video is a foundational step in computer vision-based sports analytics. Obtaining player identities is essential for analyzing the game and is used in downstream tasks such as game event recognition. Transformers are the existing standard in natural language processing (NLP) and are swiftly gaining traction in computer vision. Motivated by the increasing success of transformers in computer vision, we introduce a transformer network for recognizing players through their jersey numbers in broadcast National Hockey League (NHL) videos. The transformer takes temporal sequences of player frames (called player tracklets) as input and outputs the probabilities of jersey numbers present in the frames. The proposed network performs better than the previous benchmark on the same dataset. We implement a weakly-supervised training approach by generating approximate frame-level labels for jersey number presence and use the frame-level labels for faster training. We also utilize player shifts available in the NHL play-by-play data by reading the game time using optical character recognition (OCR) to get the players on the ice rink at a certain game time. Using player-shifts improved the player identification accuracy by 6%.*

## 1. Introduction

Player identification is a problem of fundamental importance in vision-based sports analytics. Identifying players is a key component of player tracking systems [23, 33] that are used by hockey coaches, analysts, and scouts to analyze the game.

Player identification through jersey numbers has been performed using static images [12, 20, 22, 30]. However, inferring jersey number from static images does not take into account the valuable temporal information present in sports videos. To address the issue, Chan *et al.* [6] and Vats *et al.* [33] infer jersey numbers from temporal player sequences called tracklets using an LSTM and

temporal 1D CNN, respectively. Inspired by the increasing success of transformers in computer vision tasks involving both images [5, 9, 21] and videos [1, 10, 14], in this paper, we introduce a transformer for recognizing jersey numbers from player tracklets. The transformer takes as input CNN features of tracklet frames combined with a positional encoding and outputs the probabilities of jersey numbers present in the tracklet. We use the multi-task loss function proposed in Vats *et al.* [30] for training the network. The overall network is illustrated in Fig. 1. The transformer network shows better performance compared to the previous benchmark on the same player identification dataset [33].

One detail common in Chan *et al.* [6] and Vats *et al.* [33] is that all images in a tracklet are annotated with the same label and a tracklet consists of hundreds of frames. As a result, when sampling a fixed number of frames for training, it is possible that the frames may not have a jersey number visible. This leads to inconsistent and slow training. In this paper, we perform weakly-supervised training by generating approximate frame-level labels for tracklet jersey numbers, which leads to faster training.

For further improvement of player identification, we exploit the public NHL play-by-play data that contains information about which players are on the ice at any time of the game. Although the number of players on an NHL team roster is 23, there can be only between 3 and 5 players on the ice for each team at any point in the game (plus one goalie per team). We process this information using an optical character recognition (OCR) system that reads the game time and extracts the players on the ice using a player shift database. We multiply the final jersey number probability vector of a tracklet by a binary vector that encodes which players are present on the ice at a certain time. Using player shift information improves the overall accuracy by 6%. The following items summarize the contributions of this paper:
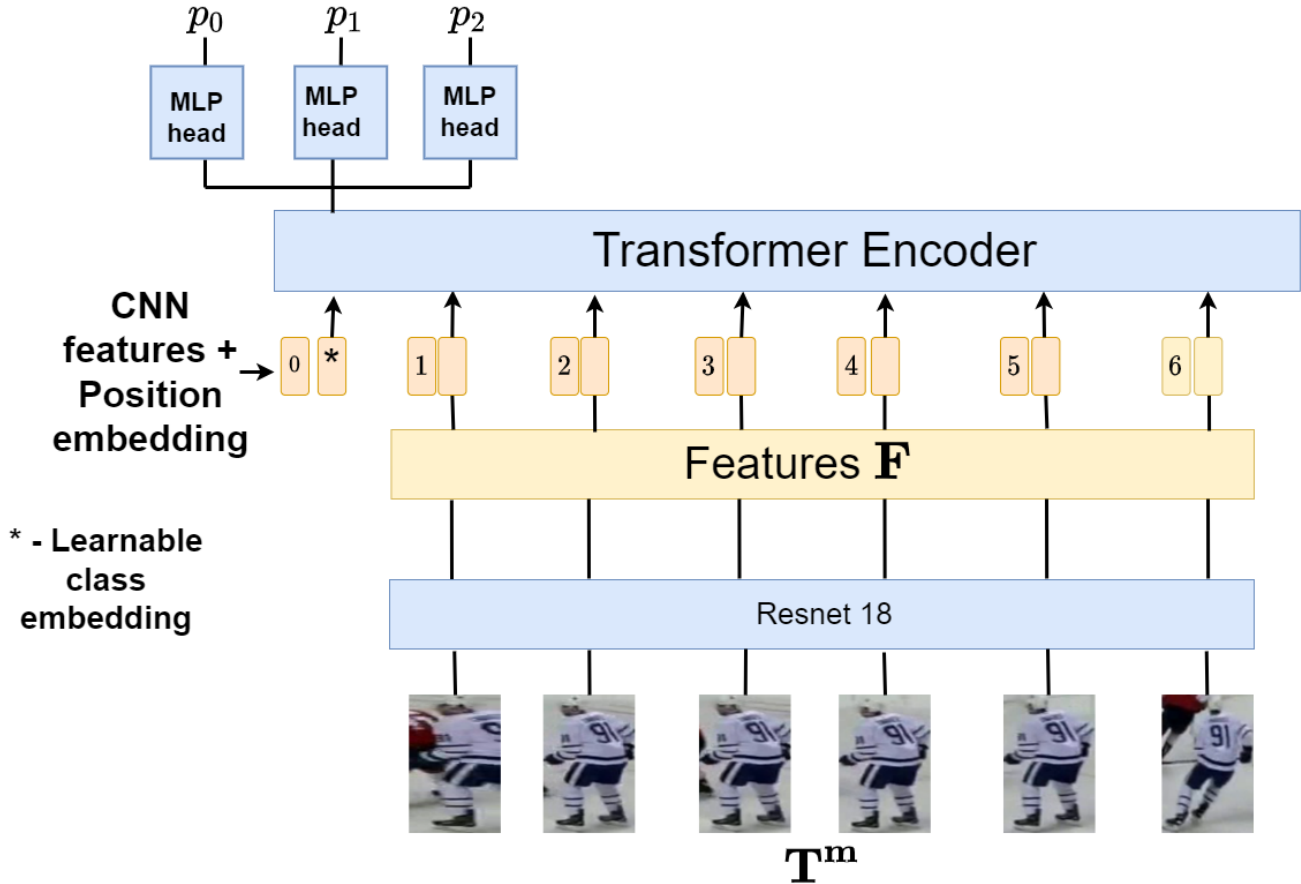
Figure 1. Network architecture for the proposed network. The input to the network is a temporal sequence of $m$ images $\mathbf{T^m}$. Each image in the tracklet is passed through a ResNet18 network to obtain 512 dimensional features $\mathbf{F}$. The features are prepended with the `[class]` token and combined with learnable positional encoding.

1. We introduce a weakly-supervised training strategy by obtaining approximate frame-wise jersey labels from a secondary network. The training strategy achieves faster convergence when compared to the naive strategy of not using approximate labels.

2. We introduce a network composed of a transformer encoder for sports jersey number recognition that performs better than the previous benchmark on the dataset [33].

3. We incorporate player shift times into the inference using OCR, allowing the network to focus on the players present at a particular moment in a game. Using player shifts improves player identification accuracy by a further 6%.

## 2. Background

### 2.1. Transformers in computer vision

Following the success of the attention mechanism used in the NLP Transformer [29], many computer vision researchers have opted to incorporate elements of the Transformer into their architectures for image and video recognition. Many of the earlier approaches used CNN feature extractors with a Transformer-based network head. Girdhar *et al*. [14] re-purposed the Transformer architecture for video understanding using a custom multi-head attention unit to process spatio-temporal features that were extracted using an I3D base. Their approach achieved state-of-the-art accuracy on the Atomic Visual Actions dataset [15] using only RGB input. In a similar manner, Gavrilyuk *et al*. [10] used a Transformer encoder to assimilate spatio-temporal features and perform group activity recognition using three different input streams: RGB, optical flow, and 2D pose.

In image classification, Dosovitskiy *et al*. [9] showed

that preliminary feature extraction using CNNs was not necessary. They proposed a pure transformer architecture called the Vision Transformer (ViT) that operated directly on sequences of image patches, or tokens, and found it performed very well on the image classification task. Arnab *et al.* [1] extended the Vision Transformer to video (ViViT) by extracting spatio-temporal tokens from input video. To handle the long sequences of tokens encountered in video, they further proposed factorising the input into spatial and temporal components to improve efficiency. ViViT achieved state-of-the-art accuracy on several action recognition benchmarks.

In other areas of computer vision, Carion *et al.* [5] proposed the Detection Transformer (DETR) for object detection. Using an encoder-decoder Transformer to process CNN-extracted image features, they obtained comparable results to the popular Faster RCNN architecture [26]. Li *et al.* [21] introduced two variants of an encoder-decoder Transformer architecture for single-stage (bottom-up) and two-stage (top-down) human pose estimation. In contrast to previous methods, their Transformer architectures regressed keypoints directly instead of using heatmaps.

## 2.2. Computer vision based sports analytics

Computer vision is currently being applied in many sports analytics problems. Problems such as sports event detection [13, 24, 31], player action recognition [4, 32], sports field registration [17, 28] and sports ball tracking [25, 34] are being solved with the help of computer vision. McNally *et al.* [24] use a hybrid CNN-LSTM network for golf swing sequencing and also introduce a new dataset for the same. Giancola *et al.* [13] introduce a new task of action spotting in soccer for finding anchors of game events in broadcast video. Pidaparthy *et al.* [25] use AlexNet [19] to track the hockey puck in video by minimizing the mean-squared error (MSE) loss between the ground truth and predicted puck coordinates. Sharma *et al.* [28] perform field registration in soccer by computing the transformation between a broadcast image and static field model through nearest neighbour search. Cai *et al.* [4] combine player stick and body pose with optical flow data to perform player level action recognition in ice hockey.

## 2.3. Player identification from static images

Before the advent of deep learning techniques, player identification from images was done with the help of hand crafted features. Although player appearance has been used to identify players in basketball [27], a player's jersey number remains a widely used feature for player identification due to its consistency in the game. Gerke *et al.* [12] was the first to use a CNN for identifying jersey numbers from player images. Vats *et al.* [30] introduce a multi-task loss function for identifying jersey numbers. Li *et al.* [20] use

a spatial transformer network to recognize jersey numbers from player images by warping the jersey number to suitable coordinates. Liu *et al.* [22] augmented the Faster-RCNN [26] network with player pose information for detecting and recognizing jersey numbers from images. Gerke *et al.* [11] also merged their image-based jersey number identification system with player location features on the soccer field.

## 2.4. Player identification from tracklets

Compared to inferring jersey numbers from static images, inferring jersey numbers from player tracklets has been found advantageous [6, 23, 33]. This is because the image sequences provide beneficial temporal information. Lu *et al.* [23] construct a conditional random field (CRF) consisting of feature nodes and identity nodes with appropriate connections and learn the CRF with weakly-supervised learning using a variant of expectation-maximization (EM). Chan *et al.* use a network based on the LRCN network [8] to infer jersey numbers from player tracklets. The final tracklet scores are aggregated using a secondary CNN. Vats *et al.* [33] use 1D temporal convolutions to infer jersey numbers from player tracklets without the use of a secondary CNN.

Our work is related to Lu *et al.* [23] as they also incorporate play-by-play as a prior during CRF training. We incorporate player shift information in a different way through multiplying the jersey number probability vector with binary shift vectors during inference (Section 3.5). We test on a more diverse dataset consisting of 18 teams compared to two in Lu *et al.* and 86 player identities compared to 24 (12 per team) in Lu *et al.*.

## 3. Methodology

### 3.1. Dataset

The player identification tracklet dataset [33] consists of 3510 player tracklets. The dataset is obtained from 84 broadcast NHL videos. The tracklet bounding boxes and identities were annotated manually. The manually annotated tracklets simulate the output of a tracking algorithm. The average length of a player tracklet is 191 frames. Note that the player jersey number is visible in only a subset of tracklet frames. The dataset is divided into 86 jersey number classes including one *null* class representing no jersey number visible. The dataset is heavily imbalanced with the *null* class consisting of 50.4% of tracklet examples.

The training/testing split is done game-wise to avoid any in-game bias. 71 videos are used for training/validation and 13 videos are used for testing. The dataset contains 2829 training tracklets, 176 validation tracklets and 505 test tracklets.

## 3.2. Network architecture

The input to the network is a temporal sequence of $m$ images $\mathbf{T^m} = \{I_i \in \mathbb{R}^{3 \times 300 \times 300}\}_{i=1}^m$ sampled from a player tracklet $\mathbf{T} = \{I_k : I_k \in \mathbb{R}^{300 \times 300 \times 3}\}_{k=1}^n$ of $n$ images. The $m$ images are randomly sampled from the tracklet $T$ serving as a form of data augmentation. The sampling technique is discussed in Section 3.4. The images $\mathbf{T^m}$ are passed through a 2D CNN (Resnet18 [16]) to obtain $m$ features $\mathbf{F} = \{f_i \in \mathbb{R}^{512}\}_{i=1}^m$. The Resnet18 is pretrained on static jersey number images using the image based jersey number dataset introduced by Vats *et al.* [30]. The features $\mathbf{F}$ are input into a transformer encoder consisting of $l$ layers with $h$ multi-headed self-attention heads per layer. Each attention head has a constant dimension of $D_h \in \mathbb{R}^{64}$. Positional encoding $p_i \in \mathbb{R}^{512}$ are added to the features $f_i$. Instead of using fixed positional encoding, the positional encoding is learned. As per the Vision transformer [9], a [class] token similar to BERT [7] is prepended to the CNN features $\mathbf{F}$. The state of the [class] token at the final transformer layer is fed to three multi-layer perceptron (MLP) heads consisting of a layernorm [2] and linear layer. The output of the three MLP heads are three vectors. The first vector $p_0 \in \mathbb{R}^{86}$ denotes the probability distribution of the predicted jersey number considering each jersey number in the dataset as a separate class. The other two vectors $p_1 \in \mathbb{R}^{11}$ and $p_2 \in \mathbb{R}^{11}$ denote the probability distribution of the first and second digit of the predicted jersey number. The one additional class in the 11-dimensional vectors $p_1$ and $p_2$ denotes the absence of a jersey number

We utilize the multi-task loss for jersey number recognition [30] for training the network. Concretely, we let $y_0 \in \mathbb{R}^{86}$ denote the ground truth vector for the holistic jersey number class, and we let $y_1 \in \mathbb{R}^{11}$ and $y_2 \in \mathbb{R}^{11}$ denote the first digit and second digit ground truth vectors respectively. Let

$$\mathcal{L}_0 = -\sum_{i=1}^{86} y_0^i \log p_0^i \qquad (1)$$

be the holistic jersey number component of the loss and

$$\mathcal{L}_1 = -\sum_{j=1}^{11} y_2^j \log p_1^j \qquad (2)$$

and

$$\mathcal{L}_2 = -\sum_{j=1}^{11} y_1^j \log p_2^j \qquad (3)$$

be the digit-wise losses. Instead of using fixed weights for the three losses, the loss weights are learned using the technique introduced in Kendall *et al.* [18], with the overall loss $\mathcal{L}$ given by:

$$\mathcal{L} = \frac{1}{\sigma_1^2}\mathcal{L}_0 + \frac{1}{\sigma_2^2}\mathcal{L}_1 + \frac{1}{\sigma_3^2}\mathcal{L}_2 + \log(\sigma_1) + \log(\sigma_2) + \log(\sigma_3) \qquad (4)$$

where $\{\sigma_i\}_{i=1}^3$ are trainable parameters. The overall network architecture is illustrated in Fig 1.

## 3.3. Training details

For handling the severe class imbalance in the dataset, the *null* class tracklets are sampled with a probability of $p_s = 0.1$ [33]. The network is trained with Adam optimizer with an initial learning rate of $0.0001$ and a batch size of 16. The learning rate is reduced by a factor of $\frac{1}{5}$ after 2500 iterations and again after 5000 iterations. Several data augmentation techniques such as random rotation by $\pm 10$ degrees, randomly cropping $300 \times 300$ pixel patches from the tracklet images and color jittering are used while training. Each augmentation technique is used on a per-tracklet basis instead of a per-frame basis. The experiments are preformed on two NVIDIA P-100 GPUs.

## 3.4. Training through approximate labels

The tracklets present in the training set can contain hundreds of frames such that the jersey number is only visible in a small subset of frames. Previous approaches in the literature [6, 33] sample a fixed number of frames randomly from a tracklet without any information of where the jersey number is actually visible. Therefore certain sampled tracklets with a non-null jersey number class may not have a jersey number visible. A toy example depicting such a scenario is shown in Fig. 2. This leads to inconsistent training signals which results in slow/unstable training as we demonstrate in experiments. To address this issue, we create frame-level labels indicating the frames in the tracklet where the jersey number is visible.

To generate these frame level labels, let $\mathcal{M}$ be a model trained to predict a jersey number in static images and let $\mathbf{T} = \{I_k : I_k \in \mathbb{R}^{300 \times 300 \times 3}\}_{k=1}^n$ be a training tracklet consisting of $n$ images $I_k$. The model $\mathcal{M}$ is run on every image $I_k$ to obtain the probability $p_k$ of whether a jersey number is visible in the image $I_k$. This gives $n$ probability scores $\{p_k \in [0,1]\}_{k=1}^n$. The $n$ probability scores are thresholded with a binary threshold $\phi$ to obtain $n$ binary values $\mathbf{B} = \{b_k \in \{0,1\}\}_{k=1}^n$. The value of $b_k$ denotes the presence of jersey number in a tracklet frame.

$$b_k = 1 \text{ if jersey number present in frame} \qquad (5)$$
$$b_k = 0 \text{ otherwise} \qquad (6)$$

The algorithm to obtain approximate labels in summarized in Algorithm 1. The model $\mathcal{M}$ is a ResNet18 [16] pretrained on a a jersey number dataset consisting of static images [30].
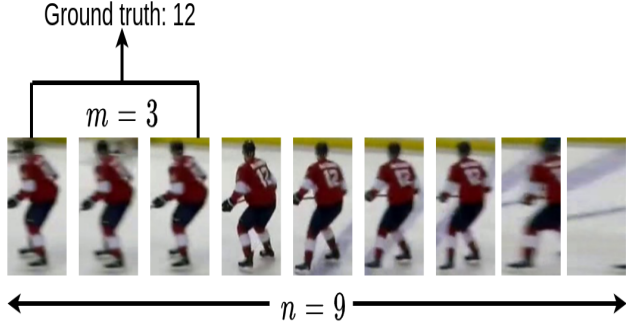
Figure 2. Toy exampling for a tracklet (of length $n = 9$ frames) where sampling $m = 3$ consecutive frames from the start leads to a sequence with ground truth 12 with no jersey number visible.

After precomputing $\mathbf{B}$, let $\mathbf{T^m} = \{I_i \in \mathbb{R}^{300 \times 300 \times 3}\}_{i=l}^{l+m}$ where $l >= 1$ and $l + m <= n$ be the $m$ images randomly sampled from a tracklet $\mathbf{T}$ for training. The corresponding $\mathbf{B^m} = \{b_i \in \{0, 1\}\}_{i=l}^{l+m}$ where $l > 1$ and $l + m <= n$ has at least one $b_i = 1$. This ensures that at least one image with a visible jersey number is present in the sampled tracklet.

For implementation, we let $\mathbf{I}$ denote the indices in the vector $B$ for which $b_k = 1$. We randomly sample an index $start\_idx$ from $\mathbf{I}$ and then sample $m$ frames from the tracklet $\mathbf{T}$ starting from index $start\_idx$ to $start\_idx + m$. A random offset $o \in [0, m)$ is subtracted from $start\_idx$ to ensure that the sampled tracklet $\mathbf{T^m}$ may have a non-zero jersey number label at any sampled frame (and not necessarily always at the beginning). The algorithm is provided in Algorithm 2.

---

**Algorithm 1:** Algorithm for creating approximate frame-wise jersey number labels.

---

**1 Input**: Player tracklet $\mathbf{T}$, Image-wise jersey number model $\mathcal{M}$ , Threshold $\phi$
**2 Output**: Frame for labels $\mathbf{B}$
**3 Initialize**: $\mathbf{B} = null$
**4 for** $I_k \in \mathbf{T}$ **do**
**5**     $p_k = \mathcal{M}(I_k)$
**6**     **if** $p_k > \phi$ **then**
**7**         $\mathbf{B}$.append(1)
**8**     **else**
**9**         $\mathbf{B}$.append(0)
**10**    **end**
**11 end**

---

### 3.5. Incorporating player shifts

To incorporate player shifts for improving player identification performance, the game time in the video needs to

---

be synced with the player shifts database, denoted by $\mathcal{S}$. $\mathcal{S}$ contains player shifts according to game time along with the corresponding jersey number and team affiliations. To read game time from broadcast video clips, the EasyOCR[1] library was used. Let $t_s$ denote the starting game time and $t_e$ denote the ending game time of a short video clip obtained using OCR. The player shifts $S'$ that are present in the game time between $t_s$ and $t_e$ are extracted from the player shift database $\mathcal{S}$. The set $S'$ can be expressed as a union $S' = S_h \cup S_a$ where $S_h$ and $S_a$ are the subsets of home and away shifts present in the set $S'$. Let the sets $\mathcal{H}$ and $\mathcal{A}$ denote the jersey numbers corresponding to $S_h$ and $S_a$ respectively.

Given a test video, player tracking and team identification are performed to obtain player tracklets [33]. We assign a single jersey number probability vector $p_{jn}$ and team affiliation (home, away or referee) to each tracklet using the inference algorithm discussed in Vats *et al.* [33]. We then construct *shift vectors* $v_h \in \mathbb{R}^{86}$ and $v_a \in \mathbb{R}^{86}$ that encode the jersey numbers present in the home and away teams. Let $null$ denote the no-jersey number class and $j$ denote the index associated with jersey number $n_j$ in $p_{jn}$ vector.

$$v_h[j] = 1, \text{if } n_j \in \mathcal{H} \cup \{null\} \quad (7)$$
$$v_h[j] = 0, \; otherwise \quad (8)$$

similarly,

$$v_a[j] = 1, \text{if } n_j \in \mathcal{A} \cup \{null\} \quad (9)$$
$$v_a[j] = 0, \; otherwise \quad (10)$$

Based on whether the player tracklet belongs to the home or the away team, the final player identity $Id$ is computed as

$$Id = argmax(p_{jn} \odot v_h) \quad (11)$$

---

[1]Found online at: https://github.com/JaidedAI/EasyOCR

(where $\odot$ denotes element-wise multiplication) if the tracklet belongs to the home team, otherwise,

$$Id = argmax(p_{jn} \odot v_a) \qquad (12)$$

if the player belongs to the away team.

## 4. Results

We compare the performance of the proposed network with Vats *et al.* [33], which is the current state-of-the art on the dataset. The network performs better than Vats *et al.*, demonstrating the effectiveness of the proposed approach. The results are shown in Table 5.

We also re-implement Chan *et al.* [6] from scratch due to unavailability of publicly-available code and dataset. The proposed approach obtains $10.1\%$ more accuracy than Chan *et al.*. The reasons for better accuracy of the proposed approach compared to Chan *et al.* are: (1) Chan *et al.* use a temporal receptive field of only 16 frames whereas the proposed approach has a more than double receptive field of $40$ frames. (2) lack of data augmentation such as random rotation, color jittering in Chan *et al.* (3) the dataset used in our work is half the size and much more skewed ($50.4\%$ *null* class) compared to Chan et al. due to which their late fusion network overfits on our dataset. (4) Chan *et al.* does not incorporate techniques to handle dataset class imbalance.

We also compare the proposed weakly-supervised training scheme making use of approximate labels to sampling frames randomly from any point in the tracklet (not using approximate frame labels) [6, 33]. The proposed scheme of training with the help of approximate labels improves the training convergence as illustrated in Fig. 4. The validation accuracy curves are shown in Fig. 5. The reason for improved convergence with the proposed training scheme is that all the tracklet mini-batches sampled using approximate labels have the jersey number visible which results in a consistent training signal.

### 4.1. Ablation studies

The number of transformer layers $l$, the number of attention heads $h$ and length of sequence for training/evaluation $m$ are important parameters affecting the overall performance. Hence, an ablation study is performed to determine the best value for each parameter.

#### 4.1.1 Attention heads

We perform an ablation study to determine to best value of the number of attention heads per transformer layer $h$. The values of $h \in \{2, 4, 6, 8, 10\}$ were tested while keeping the number of transformer layers $l$ and sequence length for training/evaluation $m$ constant ($l = 2, m = 30$). The value of $h = 8$ showed the best performance with an accuracy of

$83.6\%$ and a weighted F1 score of $84.2\%$. Table 2 shows the accuracy and F1 score values at the different values of $h$ tested. Using more than 8 attention heads resulted in a performance decrease due to overfitting.

#### 4.1.2 Transformer layers

We determine to best value of the number of transformer layers $l$ by testing $l \in \{2, 4, 6, 8\}$ while keeping the number of attention heads per layer $h$ and the sequence length used for training/evaluation $m$ constant ($h = 8, m = 30$). From Table 3, the best accuracy value of $83.37\%$ and F1 score of $83.85\%$ was obtained with $l = 2$. The performance of the network declines after increasing the transformer layers from $l = 2$ to $l = 8$. This is because of overfitting since the number of parameters in the model increases around four times from $\sim 3.2$ million when $l = 2$ to $\sim 12.6$ million when $l = 8$ with no significant improvement in accuracy.

#### 4.1.3 Sequence length

We determine the best value of the training and evaluation sequence length $m$ by keeping the transformer layers and number of attention heads per layer constant. The values of $m \in \{10, 20, 30, 40, 50\}$. From Table 4, the lowest performance was shown by $m = 10$ with an accuracy of $81.58\%$. Increasing $m$ to 20 improved the accuracy and F1 score due to increase in receptive field of the network. However, the accuracy between $m = 20$ to $m = 50$ remained the same. The best performance was obtained by $m = 40$ with an accuracy of $83.37\%$ and F1 score of $84.14\%$. Further increasing sequence length $m$ beyond 40 did not improve performance.

## 5. Result of incorporating player shifts

We evaluate the network on the player tracklets obtained by running a tracking algorithm [3, 33] on the 13 test videos. This evaluation is different from the evaluation done in Section 4 since the player tracklets are now obtained from the player tracking algorithm (rather than being manually annotated). The accuracy obtained by incorporating player shifts using OCR into player identification is compared to two baselines: (1) not incorporating any kind of roster/shift information, and (2) using player rosters available at the start of the game instead of player shifts [33].

From Table 1, not using any shifts/roster data obtains a mean accuracy of $82.02\%$, that is $4.12\%$ greater than Vats *et al.* [33] . Incorporating player shifts obtains the best mean accuracy of $87.97\%$, which is $\sim 6\%$ more than not using any shift or roster data. In fact, every video except the first video in the test set obtains equal or more accuracy when using the player shift data. This is because using
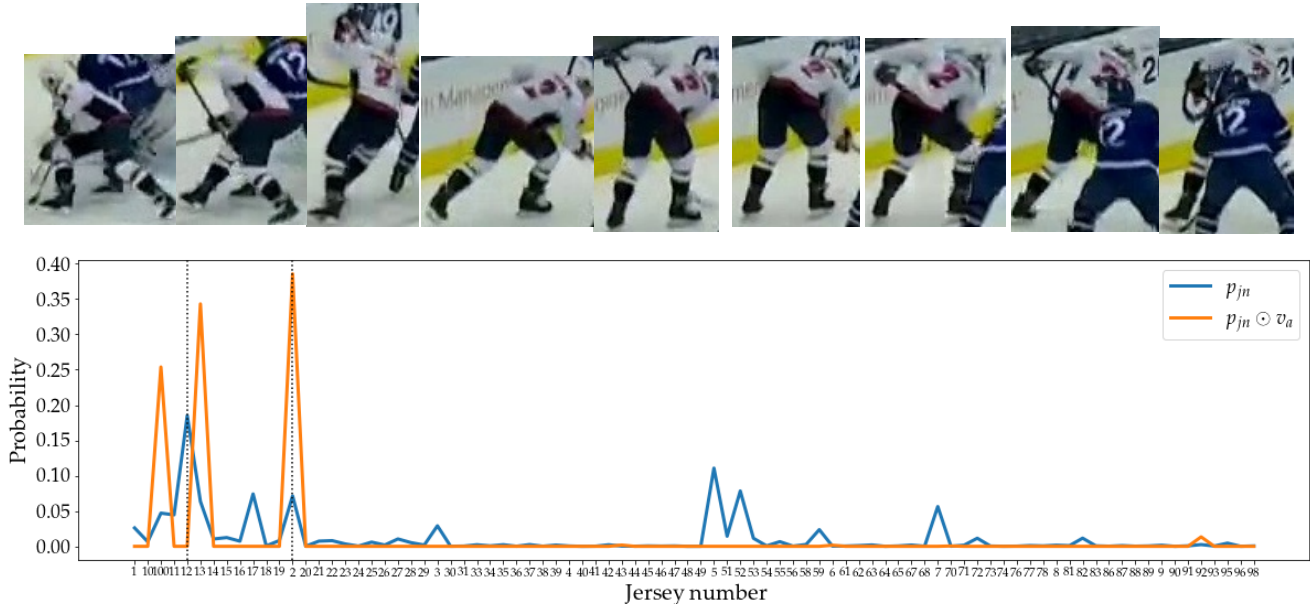
Figure 3. **Top row:** Example of a 'hard' tracklet where the ground truth jersey number 2 is tilted. There is also a heavy occlusion with the opposition player with jersey number 12. Note that the original tracklet contains 89 frames, however, only a subset of frames is shown here due to space constraints. **Bottom row:** For the tracklet shown in the top row, $p_{jn}$ is the probability of jersey number present in the tracklet (blue color). Orange color line is the normalized probability $p_{jn} \odot v_a$, i.e, the probability of jersey number multiplied by the shift vector $v_a$. For $p_{jn}$ the highest confidence value exists for jersey number 12 (first vertical line from left), which is incorrect. Multiplying with the shift vector $v_a$ corrects the mistake by making the system focus only on the jersey number present in the away team during the game shift, after which the probability of the correct jersey number 2 (second vertical line from left) becomes the greatest.

Table 1. Overall player identification accuracy for 13 test videos. The mean accuracy for identification increases by $5.95\%$ after including the player shift data.

| Video number | Ours w/ shift data | Ours w/ roster data | Ours w/o shift/roster data | Vats *et al.* [33] w/o shift/roster data |
|:---:|:---:|:---:|:---:|:---:|
| 1 | 90.70% | 95.35% | 90.60% | 90.60% |
| 2 | **91.43%** | 85.71% | 74.29% | 57.1% |
| 3 | **87.72%** | 87.72% | 84.2% | 84.2% |
| 4 | **80.00%** | 76.0% | 72.00% | 74.0% |
| 5 | **83.33%** | 83.33% | 81.48% | 79.6% |
| 6 | **90.00%** | 90.0% | 90.00% | 88.0% |
| 7 | **85.07%** | 80.60% | 73.13% | 68.6% |
| 8 | **93.75%** | 93.75% | 91.6% | 91.6% |
| 9 | **94.45%** | 93.18% | 88.6% | 88.6% |
| 10 | **93.02%** | 88.37% | 83.72% | 86.04% |
| 11 | **82.22%** | 80.00% | 71.11% | 44.44% |
| 12 | **84.85%** | 84.85% | 84.85% | 84.85% |
| 13 | **86.11%** | 83.33% | 80.56% | 75.0% |
| Mean | **87.97%** | 86.32% | 82.02% | 77.9% |

player shifts helps the algorithm focus on a smaller subset of possible players present at a particular time. The lower accuracy of the first test video is due to inaccuracies in the shifts database. Using the player roster obtains an accuracy $86.32\%$, which is just $1.65\%$ lower than the accuracy obtained when using player shifts, which demonstrates that even if player shifts are not available, using the available roster can provide performance comparable to using

player shift data. Fig. 3 shows an example of a tracklet where incorporating player shifts corrects the prediction of the model that does not use any shift or roster information.

## 6. Conclusion

In this paper, we introduced and implemented a transformer network for identifying players from player tracklets. We introduce a novel, weakly-supervised training tech-
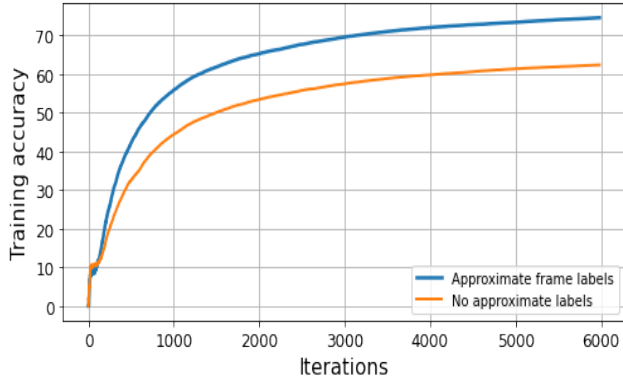
Figure 4. Training curves corresponding to a network with transformer layers $l = 2$, attention heads per layers $h = 8$ and training sequence length $m = 40$ Training with approximate labels makes the network converge faster while training.
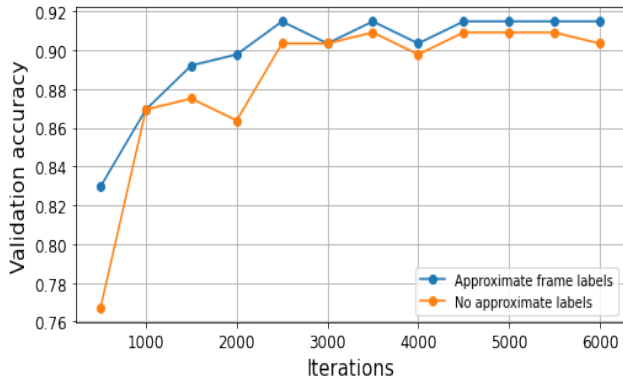


Figure 5. Validation accuracy curves corresponding to a network with transformer layers $l = 2$, attention heads per layers $h = 8$ and training sequence length $m = 40$. The initial accuracy at iteration number 500 is 6.2% higher when training with approximate labels (blue color curve). The network also converges faster and obtains a higher accuracy value with approximate label based training.

Table 2. Ablation study to determine the best value of attention heads per layer $h$ keeping number of layers $l$ and sequence length $m$ constant ($l = 2, m = 30$).

| $h$ | Accuracy | F1 score |
|---|---|---|
| 2 | 82.97% | 83.65% |
| 4 | 82.97% | 83.32% |
| 6 | 83.17% | 83.74% |
| 8 | **83.37 %** | **83.85%** |
| 10 | 82.38% | 82.90% |

nique with the help of approximate labels to significantly speed up training. We also use a player shift database to significantly improve player identification accuracy on test

Table 3. Ablation study to determine the best value layers $l$ keeping number of attention heads $h$ and sequence length $m$ constant ($h = 8, m = 30$).

| $l$ | Accuracy | F1 score |
|---|---|---|
| 2 | **83.37 %** | **83.85%** |
| 4 | 81.98% | 82.74% |
| 6 | 81.58% | 82.17% |
| 8 | 82.77% | 83.17% |

Table 4. Ablation study to determine the best value of training and evaluation sequence length $m$ keeping number of attention heads $h$ and number of layers $l$ constant ($h = 8, l = 2$).

| $m$ | Accuracy | F1 score |
|---|---|---|
| 10 | 81.58% | 81.75% |
| 20 | 83.37% | 83.76% |
| 30 | 83.37% | 83.85% |
| 40 | **83.37%** | **84.14%** |
| 50 | 83.37% | 84.07% |

Table 5. The result of the best performing model ($h = 8, l = 2, m = 30$) compared with the previous state-of-the-art on the dataset.

| Model | Accuracy | F1 score |
|---|---|---|
| Proposed | **83.37 %** | **84.14 %** |
| Vats *et al.* [33] | 83.17% | 83.19% |

videos. However, player identification is even more challenging when the jersey number of the player is not visible. Considering the fact that players in team sports such as ice hockey don't move randomly by follow roles such as defender, forward etc,future work will focus on improving player identification by incorporating a prior based on player positional data (e.g., left wing, center, right wing, defense, *etc.*).

## 7. Acknowledgment

## References

[1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6836–6846, October 2021. 1, 3

[2] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv*, abs/1607.06450, 2016. 4

[3] Guillem Braso and Laura Leal-Taixe. Learning a neural solver for multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 6

[4] Zixi Cai, Helmut Neher, Kanav Vats, David A. Clausi, and

John Zelek. Temporal hockey action recognition via pose and optical flows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019. 3

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 213–229, Cham, 2020. Springer International Publishing. 1, 3

[6] Alvin Chan, Martin D. Levine, and Mehrsan Javan. Player identification in hockey broadcast videos. *Expert Systems with Applications*, 165:113891, 2021. 1, 3, 4, 6

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. 4

[8] Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):677–691, 2017. 3

[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 1, 2, 4

[10] Kirill Gavrilyuk, Ryan Sanford, Mehrsan Javan, and Cees G. M. Snoek. Actor-transformers for group activity recognition. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 836–845, 2020. 1, 2

[11] 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. Computer Vision in Sports. 3

[12] S. Gerke, K. Müller, and R. Schäfer. Soccer jersey number recognition using convolutional neural networks. In *2015 IEEE International Conference on Computer Vision Workshop (ICCVW)*, pages 734–741, 2015. 1, 3

[13] Silvio Giancola, Mohieddine Amine, Tarek Dghaily, and Bernard Ghanem. Soccernet: A scalable dataset for action spotting in soccer videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018. 3

[14] R. Girdhar, J. Joao Carreira, C. Doersch, and A. Zisserman. Video action transformer network. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 244–253, Los Alamitos, CA, USA, jun 2019. IEEE Computer Society. 1, 2

[15] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al.

Ava: A video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018. 2

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4

[17] Namdar Homayounfar, Sanja Fidler, and Raquel Urtasun. Sports field localization via deep structured models. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4012–4020, 2017. 3

[18] 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 (CVPR)*, June 2018. 4

[19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. 3

[20] G. Li, S. Xu, X. Liu, L. Li, and C. Wang. Jersey number recognition with semi-supervised spatial transformer network. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1864–18647, 2018. 1, 3

[21] Ke Li, Shijie Wang, Xiang Zhang, Yifan Xu, Weijian Xu, and Zhuowen Tu. Pose recognition with cascade transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1944–1953, June 2021. 1, 3

[22] H. Liu and B. Bhanu. Pose-guided R-CNN for jersey number recognition in sports. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2457–2466, 2019. 1, 3

[23] Wei-Lwun Lu, J. Ting, J. J. Little, and K. P. Murphy. Learning to track and identify players from broadcast sports videos. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 35(07):1704–1716, jul 2013. 1, 3

[24] William McNally, Kanav Vats, Tyler Pinto, Chris Dulhanty, John McPhee, and Alexander Wong. Golfdb: A video database for golf swing sequencing. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2553–2562, 2019. 3

[25] Hemanth Pidaparthy and James H. Elder. Keep your eye on the puck: Automatic hockey videography. *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1636–1644, 2019. 3

[26] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. 3

[27] Arda Senocak, Tae-Hyun Oh, Junsik Kim, and In So Kweon. Part-based player identification using deep convolutional representation and multi-scale pooling. In *Proceedings of the*

*IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018. 3

[28] Rahul Anand Sharma, Bharath Bhat, Vineet Gandhi, and C. V. Jawahar. Automated top view registration of broadcast football videos. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 305–313, 2018. 3

[29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 2

[30] Kanav Vats, Mehrnaz Fani, David A. Clausi, and John Zelek. Multi-task learning for jersey number recognition in ice hockey. In *Proceedings of the 4th International Workshop on Multimedia Content Analysis in Sports*, MMSports'21, page 11–15, New York, NY, USA, 2021. Association for Computing Machinery. 1, 3, 4

[31] Kanav Vats, Mehrnaz Fani, Pascale Walters, David A. Clausi, and John Zelek. Event detection in coarsely annotated sports videos via parallel multi-receptive field 1d convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020. 3

[32] Kanav Vats, Helmut Neher, David A. Clausi, and John Zelek. Two-stream action recognition in ice hockey using player pose sequences and optical flows. In *2019 16th Conference on Computer and Robot Vision (CRV)*, pages 181–188, 2019. 3

[33] Kanav Vats, Pascale Walters, Mehrnaz Fani, David A Clausi, and John S. Zelek. Player tracking and identification in ice hockey. *ArXiv*, abs/2110.03090, 2021. 1, 2, 3, 4, 5, 6, 7, 8

[34] X. Zhang, T. Zhang, Y. Yang, Z. Wang, and G. Wang. Real-time golf ball detection and tracking based on convolutional neural networks. In *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 2808–2813, 2020. 3