

ASTRA: An Action Spotting TRAnsformer for Soccer Videos

Artur Xarles
arturxe@gmail.com

Universitat de Barcelona & Computer Vision Center
Barcelona, Spain

Thomas B. Moeslund
tbm@create.aau.dk
Aalborg University
Aalborg, Denmark

Sergio Escalera
sescalera@ub.edu

Universitat de Barcelona & Computer Vision Center &
Aalborg University
Barcelona, Spain

Albert Clapés
aclapes@ub.edu

Universitat de Barcelona & Computer Vision Center
Barcelona, Spain

ABSTRACT

In this paper, we introduce **ASTRA**, a Transformer-based model designed for the task of Action Spotting in soccer matches. ASTRA addresses several challenges inherent in the task and dataset, including the requirement for precise action localization, the presence of a long-tail data distribution, non-visibility in certain actions, and inherent label noise. To do so, ASTRA incorporates (a) a Transformer encoder-decoder architecture to achieve the desired output temporal resolution and to produce precise predictions, (b) a balanced mixup strategy to handle the long-tail distribution of the data, (c) an uncertainty-aware displacement head to capture the label variability, and (d) input audio signal to enhance detection of non-visible actions. Results demonstrate the effectiveness of ASTRA, achieving a tight Average-mAP of 66.82 on the test set. Moreover, in the SoccerNet 2023 Action Spotting challenge, we secure the 3rd position with an Average-mAP of 70.21 on the challenge set.

1 INTRODUCTION

The field of automatic video analysis has significantly impacted the world of sports in recent years. Various computer vision tasks, such as object detection, tracking, and action localization, have found extensive applications within the sports domain. These applications go beyond analyzing player behavior through detection and tracking, encompassing functionalities like automated data collection or video summarization by identifying crucial actions throughout the footage. It is worth noting that this field has witnessed the emergence of numerous tasks and applications, as extensively reviewed by Thomas et al. [33] and Naik et al. [24].

This paper specifically focuses on the task of Action Spotting, which involves the temporal localization of multiple actions within untrimmed videos. It shares a close relationship with the well-known task of Temporal Action Localization, differing only in the use of a single keyframe to identify each action. While several sports datasets are available to address this task, covering domains such as tennis [43], diving [37], figure skating [17], and gymnastics [28], our primary focus is on soccer. Therefore, to tackle this task, we leverage the SoccerNet-v2 dataset [10], the largest annotated video sports dataset up to date, comprising 550 soccer matches and encompassing 17 distinct actions.

To address the task, we propose **ASTRA** (Action Spotting TRAnsformer), building upon the problem design defined in Soares et al. [31]. This design involves producing time-point detections, each consisting of both a class probability and a temporal displacement over a predefined anchor. ASTRA employs a Transformer encoder-decoder architecture, similar to the one used in DETR [5]. This allows the model to produce outputs with the desired temporal resolution, regardless of its input temporal resolution. Upon analyzing the dataset, we identify three main challenges: a long-tail distribution of the data, where some actions occur infrequently; the non-visibility of certain actions due to replays or camera angles; and noisy labels resulting from the subjective judgment of annotators in determining temporal locations. To address these challenges, we incorporate different techniques into our model. Firstly, we employ a balanced mixup approach to account for the long-tail distribution of the data. Additionally, we integrate audio signals alongside visual signals to improve the detection of non-visible actions. Furthermore, we introduce an uncertainty-aware displacement head that models label uncertainty using a Gaussian distribution. These techniques enhance performance, with ASTRA achieving an Average-mAP of 66.82 on the test split. We further evaluate our model in the SoccerNet 2023 Action Spotting challenge, consisting of 50 matches with hidden ground-truth, where we achieve the 3rd position with a tight A-mAP of 70.21.

The remaining sections of the paper are organized as follows: Section 2 provides a comprehensive review of related work on the task of action spotting. Section 3 introduces ASTRA and outlines its components. Section 4 conducts ablation studies on different aspects of the model, and compares our best solution against state-of-the-art works. Finally, Section 5 concludes the paper, summarizing our key findings and conclusions derived from this research.

2 RELATED WORK

Temporal Action Localization & Action Spotting. Action recognition has undergone significant advancements in recent years, playing a crucial role in video understanding. Initially, methods focused on classifying short trimmed videos [1, 6, 14]. However, with the progress in computer vision, more challenging tasks have emerged. Two prominent tasks in this domain are Temporal Action Localization (TAL) and Action Spotting (AS), which share the objective of temporally locating multiple actions within untrimmed videos. While TAL represents actions as temporal intervals through

the annotation of *begin* and *end* frames, AS represents actions with a *single keyframe*. This distinction offers an advantage for AS in terms of annotation cost, as it requires only one frame per action. Moreover, AS is especially well-suited for capturing actions that are instantaneous or have uncertain start and end times, where a single timestamp can effectively represent them. A concrete example is demonstrated in the SoccerNet-v2 dataset [10], where actions like goals or fouls are typically identifiable at specific temporal points.

Given the inherent similarities between TAL and AS, the methods developed for these tasks often share common components, with their main differences lying in the prediction head. These methods can generally be categorized into two groups: two-stage methods [4, 11, 15, 26, 38, 45] and one-stage methods [9, 18, 20, 23, 29, 31, 41]. In two-stage methods, proposals are first generated and subsequently classified to determine if they correspond to actions or background. These methods tend to be more complex and do not allow for end-to-end training. In contrast, one-stage models directly localize and classify actions in a single step, eliminating the need for proposal generation. These models offer simplicity and often achieve state-of-the-art performance on TAL and AS tasks.

Early one-stage models in temporal action localization utilized anchor windows sampled from sliding windows [3, 22]. For instance, Lin et al. [22] employed a set of anchor windows that were classified into different categories and refined using location offsets and overlap scores. Later, Yang et al. [39] introduced an anchor-free approach that relied on temporal points instead of anchor windows for action localizations. Their work showed the benefits of both anchor-free and anchor-based approaches. Current state-of-the-art methods in temporal action localization are predominantly anchor-free. In particular, ActionFormer [41] and TriDet [29] have achieved remarkable performance in this field. They classify each moment as either background or one of the possible actions. ActionFormer utilizes a transformer encoder architecture with downsizing operations, while TriDet incorporates a Scalable-Granularity Perception (SGP) layer based on CNNs. The SGP layer replaces the self-attention mechanism of ActionFormer to improve both model performance and efficiency. These approaches also utilize temporal regression to refine predictions and obtain more precise results. Another method, TadTR [23], draws inspiration from the DETR model [5] for object detection. TadTR constructs a transformer encoder-decoder architecture with learnable queries representing detection candidates. During training, a bipartite matching problem pairs those candidates with ground-truth actions.

Similar techniques have also demonstrated state-of-the-art performance in the task of action spotting on SoccerNet, as further discussed in Giancola et al. [13]. For instance, E2E-Spot [18] proposes a 2D CNN backbone with Gate Shift Modules [32], which incorporate temporal context and produce per-frame predictions using a Gated Recurrent Unit [8] layer. This model operates directly on the raw video frames, providing increased flexibility compared to using pre-extracted features. However, that introduces additional complexity and computational cost during training. Soares et al. [31] achieve SOTA results by defining a set of dense anchors (i.e. one anchor per input token), similar to ActionFormer or TriDet, to represent

temporal positions. These anchors are then classified into different action classes and temporally refined. The model uses pre-extracted features obtained from various pre-trained video backbones and utilizes a U-Net-like architecture for the model’s trunk.

Our approach takes inspiration from the Transformer encoder-decoder architecture of TadTR and DETR. We employ a similar architecture, with learnable queries in the decoder, where each generated query represents a specific temporal position, akin to the dense anchors proposed in Soares et al.’s work. Furthermore, we also leverage a set of strong pre-extracted features to train our model, providing a solid foundation for accurate action localization, and avoiding the added complexity when using raw frames.

Uncertainty estimation. Uncertainty estimation techniques have demonstrated their potential to enhance the performance of regression models by providing reliability estimates and accounting for potential errors in predicted values [7, 36, 40]. This becomes particularly valuable when dealing with inherently uncertain ground-truth data, characterized by measurement errors, noise, or label ambiguity. For instance, Tang et al. [40] approached Action Quality Assessment (AQA) task by modelling the quality score as a Gaussian distribution, maximizing the log-likelihood function to estimate both mean and variance. Similarly, Xie et al. [36] and Chen et al. [7] also employed Gaussian distributions for temporal regression in TAL. However, they used the Kullback-Leibler divergence to fit their models. In our work, we adopt a similar Gaussian distribution for modelling temporal displacements, and like in [40], we maximize the log-likelihood function for fitting purposes.

Multimodal approaches. In addition to the visual modality, certain methods for action classification, TAL or AS incorporate additional modalities. These modalities can include optical flow [21, 35, 44] or audio [19, 25, 27, 34] among others, and they differ in how they fuse these modalities. Specifically, for AS in SoccerNet, an approach that combines different modalities is the one in Vanderplaetse and Dupont [34]. They extract features from the log-mel spectrogram of the audio using a VGG-inspired model and explore various fusion techniques. In our work, we adopt a similar approach, leveraging a VGG-inspired model to extract audio features from the log-mel spectrogram. However, we further fine-tune the backbone model during training. We perform an early fusion of different features, merging them at the input of the Transformer encoder.

3 METHODS

Problem definition. Action spotting involves the identification and precise localization of actions within an untrimmed video X . Given the video input or a representation of it, the objective is to identify and locate all the actions occurring in the video, represented as $A = \{a_1, \dots, a_N\}$. The number of actions, denoted as N , may vary across different videos. Each action instance a_i comprises an action class c_i and its corresponding temporal position t_i , forming a pair $a_i = (c_i, t_i)$. Here, $c_i \in \{1, \dots, C\}$ represents the action class index, with C being the total number of distinct action classes.

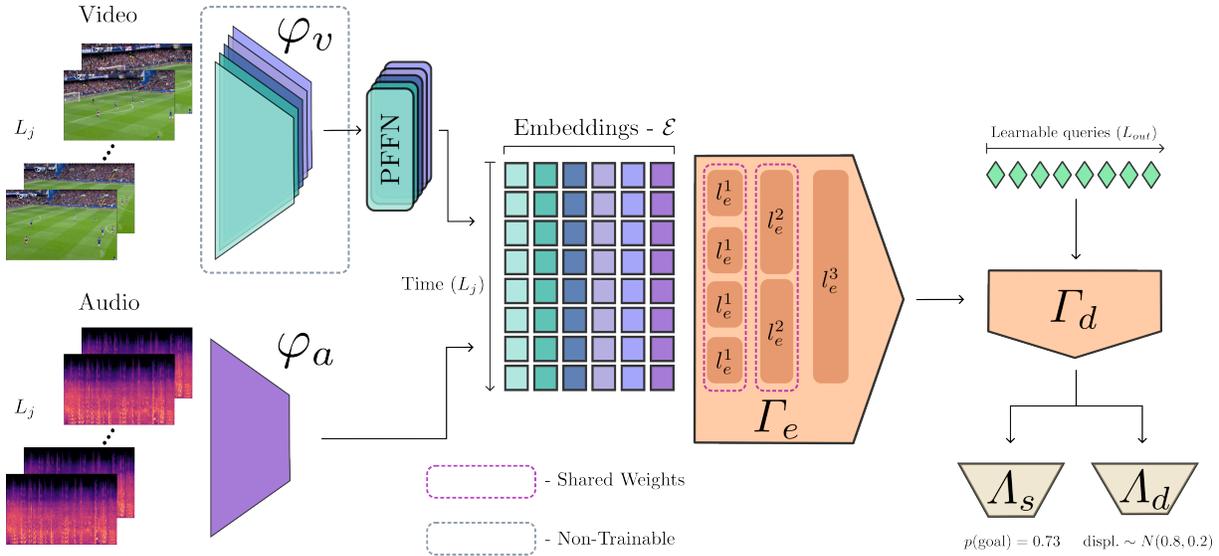


Figure 1: ASTRA (Action Spotting TRAnsformer) architecture: Visual and audio embeddings from different backbones (φ_v and φ_a) are combined and processed through a Transformer encoder-decoder (Γ_e and Γ_d). The resulting embeddings are then utilized by a classification head (Λ_s) for temporal location classification, and a displacement head (Λ_d) to further refine predictions.

Method overview. Our solution, ASTRA, leverages embeddings \mathcal{E} from multiple modalities to achieve its goals. Specifically, ASTRA is built upon $|\mathcal{E}| - 1$ pre-computed visual embeddings, complemented by an additional audio embedding derived from the log-mel spectrogram of the audio using a VGG-inspired backbone. The network responsible for generating this audio embedding is jointly trained with the ASTRA model. The embeddings are input to the model in clips spanning a duration of T seconds. These features from each backbone are processed in parallel streams, where Pointwise Feed-Forward Networks (PFFN) project them to a common dimension d . The projected embeddings are then combined in the subsequent Transformer encoder-decoder module, with learnable queries in the decoder. Inspired by the architecture proposed in DETR, this module enables ASTRA to handle different input and output temporal dimensions (L_{in} and L_{out} , respectively) and facilitates a straightforward fusion of multiple embeddings. To enhance ASTRA’s ability to capture fine-grained details, we introduce a temporally hierarchical architecture for the Transformer encoder. This architecture enables the encoder to attend to more local information in the initial layers and reduces the computational cost. Finally, ASTRA employs two prediction heads to generate classification and displacement predictions for the anchors introduced by Soares et al. [31]. These anchors correspond to specific temporal positions and class actions, as described in their work. Additionally, we adopt their suggestion of employing a radius for both classification and displacement (r_c and r_d , respectively) to define the temporal range around a ground-truth action within which it can be detected.

Furthermore, to account for label uncertainty, ASTRA adapts the prediction head responsible for displacement by modeling them as Gaussian distributions instead of deterministic temporal positions. This allows ASTRA to capture temporal location uncertainty and provide a more comprehensive representation of the actions.

Additionally, ASTRA incorporates a balanced mixup technique to improve model generalization and accommodate the long-tail distribution of the data.

We illustrate the ASTRA architecture in Figure 1, and further details are discussed in the subsequent sections.

3.1 Input embeddings

As previously mentioned, ASTRA is built upon $|\mathcal{E}|$ embeddings proceeding from multiple modalities, $|\mathcal{E}| - 1$ pre-extracted visual embeddings and an additional audio embedding. Let $E_j \in \mathbb{R}^{L_j \times d_j}$ denote the sequence of features associated to the j -th embedding, where $j \in \{1, \dots, |\mathcal{E}|\}$. Here, L_j represents the temporal dimension (i.e. L_{in} for each embedding), and d_j represents the feature dimension specific to that embedding. It is important to note that different embeddings may have varying temporal or feature dimensions.

Visual embeddings. The $|\mathcal{E}| - 1$ visual embeddings used as inputs to our model are obtained from the Baidu Research repository.¹ They are extracted using distinct backbones (φ_v) fine-tuned on the SoccerNet dataset for action classification. With a receptive field of 5 seconds, each embedding is computed with a stride of 1. To ensure a consistent feature dimension d across all embeddings, PFFNs are employed. These PFFNs project the embeddings through two linear layers with ReLU activation, applying dropout with probability p .

Audio embedding. For the additional audio embedding, we employ a VGG-inspired backbone [16] (φ_a), which is pre-trained on the AudioSet dataset [12]. The backbone takes the log-mel spectrogram of the audio as input and is fine-tuned during the training of the ASTRA model. We further replace the last linear layer of the backbone to produce the desired feature dimension of d . In line with

¹<https://github.com/baidu-research/vidpress-sports>

common practice, we feed the backbone with log-mel spectrogram segments, each spanning a duration of T_a seconds. Consequently, we obtain the audio embedding E_a with $L_a = \lfloor T/T_a \rfloor$ and $d_a = d$.

3.2 Transformer encoder-decoder

After aligning the feature dimension of all embeddings in \mathcal{E} , they are passed into the Transformer encoder-decoder module. Prior to that, a learnable encoding specifying temporal position and backbone source is added. Then, the enriched tokens (i.e., feature vectors corresponding to specific embeddings and temporal positions) undergo a Hierarchical Transformer encoder, where they progressively interact with tokens that are further apart in terms of temporal distance. This hierarchical structure enables the model to attend to fine-grained local details in the early layers while gradually incorporating broader context in the subsequent layers. In the Transformer decoder, a set of L_{out} learnable queries, representing temporal positions, is introduced. These queries evolve and capture relevant information from the Transformer encoder output tokens during the self-attention and cross-attention mechanisms in the decoder.

Hierarchical Transformer encoder (Γ_e). The Hierarchical Transformer encoder is composed of n_e vanilla Transformer encoder layers. Each layer applies the standard multi-head self-attention with h heads, followed by a two-layered PFFN with a widening factor of α , dropout, layer normalization, and residual connections. To incorporate the temporal hierarchy, in each layer l_e , where $i \in \{1, \dots, n_e\}$, the input clip of T seconds is divided into 2^{n_e-i} segments. Tokens within the same temporal segment are processed together within the layer. Importantly, all segments within a layer share the same transformer encoder layer, ensuring weight sharing and parameter efficiency.

Transformer decoder (Γ_d). The transformer decoder is composed of n_d vanilla Transformer decoder layers. Each layer applies the standard multi-head self-attention and multi-head cross-attention, each with h heads, followed by a two-layered PFFN with a widening factor of α , dropout, and residual connections. Unlike the hierarchical structure in the encoder, in the decoder, all tokens within the same clip interact with each other.

The Transformer encoder-decoder module in ASTRA provides two main advantages over other methods in TAL or AS:

- (1) Flexible handling of input and output temporal dimensions, L_{in} and L_{out} . While the input temporal dimension is typically fixed, ASTRA allows for a different output temporal dimension, providing the ability to customize the temporal resolution. This flexibility is particularly beneficial in our AS task, as highlighted in Section 4.4.
- (2) Seamless integration of multiple embeddings with varying temporal dimensions. Unlike other methods that require embeddings to have the same temporal dimension and concatenate them along the feature dimension, ASTRA can accommodate individual embeddings as separate tokens, allowing for diverse temporal resolutions.

3.3 Prediction heads

The evolved queries, representing L_{out} temporal positions uniformly distributed over the T seconds, are input to two prediction heads: the classification head and the uncertainty-aware displacement head. Figure 2 provides a visual representation of the predictions produced by these prediction heads.

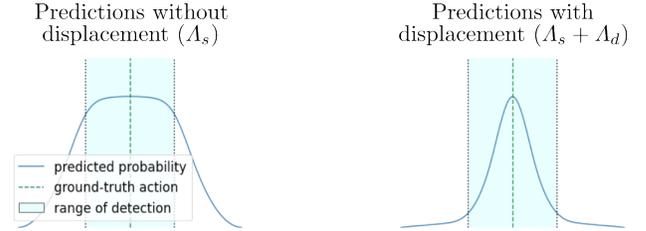


Figure 2: Example ground-truth action prediction with and without displacements. The utilization of only the classification head results in predictions spanning the entire range of detection of the ground-truth action (left), whereas incorporating displacements refines the predictions, aligning them closer to the actual temporal position of the ground-truth action (right).

Classification head (Λ_s). The classification head consists of two linear layers that project the evolved queries to the desired output feature dimension $C + 1$, representing the C different actions plus an additional background class. We incorporate dropout of p and use the ReLU activation function as an intermediate activation. Finally, a sigmoid activation function is applied so the output for each pair of temporal position and action class represents the probability of a ground truth action of that class occurring within the range of detection of the corresponding temporal position.

Uncertainty-aware displacement head (Λ_d). The uncertainty-aware displacement head is composed of two linear layers with a dropout rate of p , using the ReLU activation function. Two additional linear layers are constructed in parallel, taking the previous output as input. Both layers project the evolved queries to a feature dimension of $C + 1$ and apply dropout. The first layer utilizes a linear activation function and outputs the estimated mean displacement for each query and action class. The second layer employs an exponential activation function to generate positive values representing the estimated variance. This allows us to model the displacement as a Gaussian distribution, capturing the uncertainty associated with the displacement estimation for each temporal position and action class. These estimated displacements are then used to refine the predictions produced by the classification head.

In summary, the prediction heads produce predictions $\hat{y}_i^{t,c} = (\hat{s}_i^{t,c}, \hat{d}_i^{t,c} \sim N(\{\hat{\mu}_i^{t,c}, \hat{\sigma}_i^{t,c}\}))$ for each temporal position $t \in \{1, \dots, L_{out}\}$ and action class $c \in \{1, \dots, C + 1\}$ in a given clip sample i . The first element corresponds to the classification score, while the second represents the estimated Gaussian distribution for the displacement, indicated by the predicted mean and variance.

3.4 Data augmentation

ASTRA is strengthened by integrating a diverse set of data augmentation techniques applied to the input features. These techniques are designed to improve the model’s generalization capability and, for some of them, to also account for the long-tail distribution of the data. These techniques are as follows:

Balanced mixup. Similar to traditional mixup [42], virtual training examples are generated by creating linear combinations of pairs of examples using a parameter λ sampled from a beta distribution $\lambda \sim \text{Beta}(\alpha_m, \beta_m)$. However, our approach has a distinction. Instead of sampling both samples from the same original distribution, the second element of the pair is sampled from a balanced distribution created using a queue. This queue contains two samples from each action class and is updated at each batch iteration.

Temporal dropout and temporal switch. Treating a temporal sequence as the set of tokens corresponding to the same temporal position, we introduce two techniques. Firstly, in temporal dropout, we randomly drop entire temporal sequences with probability p_{td} . In the dropped positions, we substitute them with learnable tokens. Secondly, in temporal switch, we randomly swap the positions of consecutive pairs of temporal sequences with probability p_{ts} .

3.5 Training details

The model is trained using a combination of a classification loss (\mathcal{L}_c) and a displacement loss (\mathcal{L}_d). For classification, we employ a binary cross-entropy focal loss for all actions, temporal positions, and data samples, as formulated in Equation 1, where ground-truth labels are denoted as $y_i^{t,c} = (s_i^{t,c}, d_i^{t,c})$. The hyperparameter γ adjusts the rate at which easy examples are down-weighted.

$$\mathcal{L}_c = -\frac{1}{N \cdot L_{out} \cdot (C+1)} \left(\sum_{i=1}^N \sum_{t=1}^{L_{out}} \sum_{c=1}^{C+1} |s_i^{t,c} - \hat{s}_i^{t,c}|^\gamma \cdot (s_i^{t,c} \cdot \ln(\hat{s}_i^{t,c}) + (1 - s_i^{t,c}) \cdot \ln(1 - \hat{s}_i^{t,c})) \right) \quad (1)$$

For the displacement loss (\mathcal{L}_d), it is only applied within the r_d seconds radius of ground-truth actions and is based on the negative log-likelihood function of the target Gaussian distribution. We formulate \mathcal{L}_d , similar to Zhang et al. [40], as shown in Equation 2.

$$\mathcal{L}_d = -\frac{1}{N_{dis}} \cdot \sum_{i=1}^N \sum_{t=1}^{L_{out}} \sum_{c=1}^{C+1} \mathbb{1}_{\{\exists d_i^{t,c}\}} \cdot \left(\frac{\alpha_L}{(\hat{\sigma}_i^{t,c})^2} \cdot |d_i^{t,c} - \hat{\mu}_i^{t,c}|^2 + (1 - \alpha_L) \cdot \ln[(\hat{\sigma}_i^{t,c})^2] \right) \quad (2)$$

In the above equation, N_{dis} is the total number of ground-truth displacements (i.e. inside the range of detection of a ground-truth action). Additionally, α_L is a weight that balances the attention paid to uncertain information, with larger values of α_L focusing more on the uncertainty, while smaller values tend to result in a more typical single point estimation of the displacement.

To effectively merge both losses, we introduce a weight w_c on the classification loss. This weight ensures that both losses are scaled to the same range of values.

3.6 Inference

At inference time, the data augmentation techniques are disabled. Moreover, the temporal position classifications are refined by incorporating the displacement estimations, represented by the mean of the estimated Gaussian distribution. To reduce the number of candidate actions, Soft Non-Maximum Suppression [2] is applied with a 1D adaptation as proposed in the work by Soares et al. [31].

4 RESULTS

In this section, we provide an overview of the dataset used in our study, highlighting its key characteristics and challenges. We also discuss the implementation details, the evaluation metric and protocols employed for assessing the proposed models, and present a comprehensive analysis of all ablation experiments conducted. Lastly, we provide a detailed evaluation of our best-performing model, including its performance on the 2023 SoccerNet challenge.

4.1 Dataset

SoccerNet-v2 is a comprehensive dataset comprising 550 soccer matches from major European competitions. Among these matches, 500 have publicly available annotations with keyframes depicting 17 different actions. Table 1 provides a breakdown of these actions and their frequencies in the annotated matches. The remaining 50 matches serve as a hidden ground-truth evaluation set, accessible only to the organizers for assessing the submitted predictions.

While solving the task of action spotting for this dataset, we encounter three main difficulties:

Long-tail data. Like many real-world datasets, SoccerNet exhibits a highly unbalanced distribution. As shown in Table 1, certain actions occur much more frequently than others. This imbalance poses a challenge, as a model that disregards the long-tail distribution may perform well on the predominant head classes but struggle to adequately address the less frequent tail classes. Overfitting to these tail classes is also a concern. This issue is further aggravated when considering the task of classifying every temporal position L_{out} and introducing an additional background class. To mitigate this problem, our approach incorporates balanced mixup. This technique forces the model to iterate more times through clips (or mixtures of clips) containing actions, particularly those from the tail end of the distribution. By leveraging this approach, our aim is to enhance the model’s ability to handle long-tail actions, while simultaneously improving its generalization capabilities across all classes, as typically observed in mixup approaches.

Non-visible actions. In the original videos from SoccerNet, not all annotated actions are directly visible due to replays or camera transitions that may occlude the actions. This is accentuated in some actions as kick-offs, clearances or indirect free-kicks, as illustrated in Figure 3. Consequently, the model needs to rely on contextual information and extrapolation to predict these actions. To address this challenge, we incorporate audio into the model, assuming that the broadcast commentary or the audience reactions can assist in identifying some of the actions that are not visually observable in the videos.

Action	Ball out of play	Throw-in	Foul	Indirect FK	Clearance	Shot on target	Shot off target	Corner	Substitution
Absolute frequency	31810	18918	11674	10521	7896	5820	5256	4836	2839
Action	Kick-off	Direct FK	Offside	Yellow Card (YC)	Goal	Penalty	Red Card (RC)	YC -> RC	
Absolute frequency	2566	2200	2098	2047	1703	173	55	46	

Table 1: Frequency of occurrence of the different actions in the 500 games with publicly available annotations.

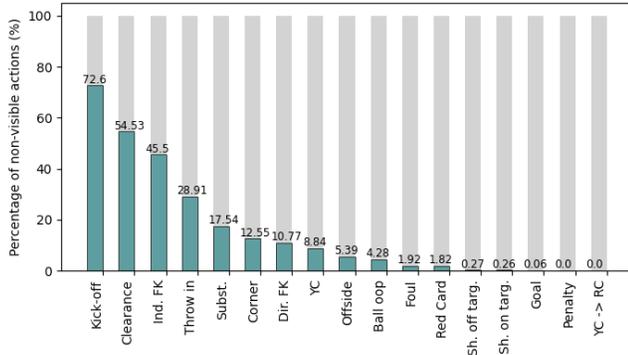


Figure 3: Percentage of non-visible ground-truth actions for each action class.

Noisy labels. Annotating actions can be subjective, leading to varying degrees of clarity in indicating the exact temporal positions of actions. While some actions have distinct temporal indications, others are more complex and rely on the annotator’s judgement. This subjectivity introduces noise in the temporal annotations. Furthermore, the presence of non-visible actions exacerbates the annotation challenge, as the temporal selection of these actions relies solely on the annotator’s subjective judgment. To address these issues, we introduce uncertainty estimation techniques. We incorporate an uncertainty-aware displacement head that models displacements as Gaussian distributions rather than deterministic values. By capturing the inherent uncertainty in the ground-truth data, we aim to mitigate the impact of noisy labels.

4.2 Evaluation metric & Evaluation protocols

Evaluation metric. The performance of the model’s spotted actions is assessed using the Average-mAP. This metric quantifies the Area Under the Curve (AUC) of the mean Average Precision (mAP) at different tolerances δ . The mAP is computed by averaging the Average Precision (AP) values across different action classes. The AP summarizes the precision-recall curve into a single value, representing the average precision across all recalls. It can be computed as follows:

$$AP = \sum_{s=0}^{S-1} (Recalls(s) - Recalls(s+1)) \cdot Precisions(s) \quad (3)$$

Here S denotes the total number of thresholds considered, while $Recalls(s)$ and $Precisions(s)$ refer to the recall and precision, respectively, at threshold s . We can further denote Average-AP as the per-class Average Precision averaged across the different tolerances.

In SoccerNet, two versions of this metric are commonly employed. A loose metric with tolerances ranging from 5 to 60 seconds, and a tight metric with tolerances between 1 and 5 seconds. We present results for both metric versions, but we primarily focus on

the tight metric to guide decisions regarding our model’s components because it aligns with the SoccerNet 2023 challenge.

Evaluation protocols. We employ two different evaluation protocols to train and evaluate our models:

- **Protocol 1 (P1):** The dataset of 500 annotated matches is divided into three sets: train, validation, and test. The train split is used for model training, the validation split is utilized to determine the optimal stopping point during training, and the test split is employed to quantitatively compare different models in our ablation experiments, evaluating their performance using the previously introduced metric. Each model is trained five times with different seeds, and the average of these runs is reported.
- **Protocol 2 (P2):** The model selected based on P1 is trained using the combined data from train, validation and test sets. Subsequently, the trained model is used to generate predictions on the challenge split and obtain the final performance.

4.3 Implementation details

We employed PyTorch for model implementation, using Adam optimizer with base LR of $5e^{-5}$. Optimization encompassed Learning Rate Warmup via Cosine Decay over 50 epochs, with 3 epochs for initial warmup. The model was fed with clips of $T = 50$ seconds, using an embedding dimension of $d = 512$. We utilized $|\mathcal{E}| = 6$ embeddings, 5 corresponding to visual data, and the remaining one to audio data. The temporal dimension of the input visual features was set to $L_j = T$, $\forall j \in \{1, \dots, |\mathcal{E}| - 1\}$, while for the audio $L_a = \lfloor T/T_a \rfloor$ with $T_a = 0.96$ seconds. Furthermore, we set the temporal output dimension as $L_{out} = 2 \times T$. In our model, we employed $n_e = 3$, $n_d = 3$, $h = 8$ and $\alpha = 4$ along with $C = 17$. Additionally, we set $p = 0.4$ for dropout. The radii were set experimentally as $r_c = 2$ and $r_d = 3$ seconds, and the loss function parameters were set as $\gamma = 1$, $\alpha_L = 0.3$, and $w_c = 100$. We also applied balanced mixup with parameters $\alpha_m = 1$ and $\beta_m = 0.6$, along with $p_{td} = 0.5$ and $p_{ts} = 0.3$. The window size for Soft Non-Maximum Suppression was experimentally determined for each action, with values ranging from 5 to 14 seconds. Code is available at <https://github.com/arturxe2/ASTRA>.

4.4 Ablations

The results in terms of tight and loose Average-mAP are summarized in Table 2. These results are obtained on the Test split following the evaluation protocol P1.

The base model (M0) is the simplest and differs significantly from our solution ASTRA. It only utilizes visual embeddings and replaces the Hierarchical Transformer encoder with a vanilla Transformer encoder. Moreover, it does not employ any data augmentation techniques. Additionally, it lacks the focal loss term in the classification loss, and the uncertainty-aware displacement head is substituted with a typical regression head that utilizes mean squared error as

Model	Added feature	loose A-mAP	tight A-mAP
Base models			
M0	-	75.21	62.38
M1	M0 + Hierarchical TE	75.42	62.32
M2	M1 + $r_c = 2, r_d = 3$	74.65	63.97
Data Augmentations			
M3	M2 + Mixup (0.6, 0.6)	75.61	64.97
M4	M2 + Balanced mixup (1, 0.6)	76.55	65.82
M5	M4 + Other augmentations	77.49	66.07
Output dimension			
M6	M5 + $L_{out} = T$	77.72	64.24
Additional improvements			
M7	M5 + Focal loss	78.02	66.09
M8	M7 + Uncertainty	78.14	66.63
M9 (ASTRA)	M8 + Audio modality	78.09	66.82

Table 2: Ablation experiments with Average-mAP results on the Test split under P1 evaluation.

the displacement loss. Furthermore, it utilizes the optimal values in Soares et al. [31] for the radii, $r_c = 3$ and $r_d = 6$. This model achieves a tight Average-mAP of 62.38.

In M1, we introduce the Hierarchical Transformer encoder. As shown in Table 2, the results are comparable to M0. There is a slight increase in performance on the loose Average-mAP by +0.21, while a minor reduction is observed in the tight metric by -0.06. Despite these similar results, we opt to use the Hierarchical Transformer encoder as it offers computational cost reduction compared to the vanilla Transformer encoder. Moreover, through experimental modifications, we adjust the radii of action detection to $r_c = 2$ and $r_d = 3$ resulting in an improvement of +1.65 on the tight metric.

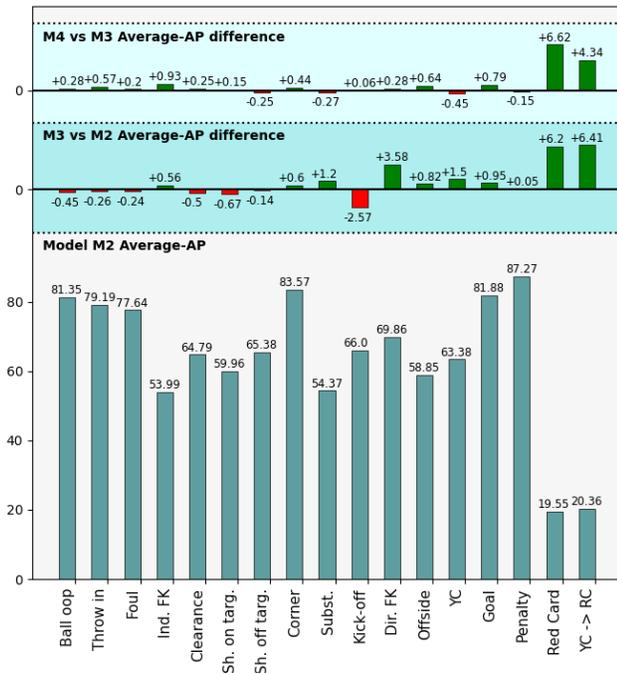


Figure 4: Per-class results comparison of models M2, M3, and M5. The figure displays Average-AP scores for each action in M2 (bottom), the differences between M2 and M3 (middle), and the differences between M3 and M4 (top).

As anticipated, the introduction of data augmentation techniques, such as typical mixup, enhances the model’s generalization capabilities. Specifically, when using the best set of tried parameters $(\alpha_m, \beta_m) = (0.6, 0.6)$ in M3, we observe an additional improvement of +1.00 in the tight Average-mAP. Figure 4 illustrates that these improvements are most pronounced in tail actions, such as red cards or second yellow cards. This can be attributed to the fact that tail actions are more prone to overfitting, thus benefiting greatly from improved generalization. By adapting the typical mixup approach to our proposed balanced mixup, and utilizing the best of the tried parameter combinations $(\alpha_m, \beta_m) = (1, 0.6)$, we achieve an additional improvement of +0.85. As depicted in Figure 4, these improvements are observed across most of the action classes, although the changes are relatively smaller in head classes. Once again, the impact on tail classes is particularly notable. These results demonstrate the effectiveness of our proposed balanced mixup technique in handling long-tail data. Furthermore, the introduction of additional augmentations such as temporal dropout and temporal switch leads to a further performance boost of +0.25.

Model M6 serves as a demonstration of the importance of an adequate output temporal dimension. As seen in Table 2, when employing a smaller output temporal dimension (i.e. $T_{out} = T$) there is a noticeable decrease in performance with respect to M5 by -1.84. This finding empowers the use of the Transformer encoder-decoder module that allows the output dimension to not be restricted by the input dimension. Other modifications to M5, such as introducing a focal loss term in the classification loss (M7), also led to a slight improvement in performance, particularly in the loose metric.

Furthermore, the inclusion of the uncertainty-aware displacement head in M8 resulted in a notable enhancement of +0.54, demonstrating the effectiveness of this module. Figure 5 presents a visualization of the average predicted variability associated with each action prediction. Notably, actions with higher variability are primarily those with high non-visibility (e.g., kick-off, clearance, indirect free-kick, throw-in) or actions that require the annotator’s judgment for precise temporal annotation, such as offsides. It is in these actions with high variability that the module seems to show the most improvement, supporting our hypothesis that the uncertainty-aware displacement head performs better in handling noisy labels compared to a typical regression head.

Finally, the inclusion of audio in M9 further enhances the model, contributing an additional +0.19 improvement and resulting in a 66.82 tight Average-mAP for the ASTRA model. In Figure 6 (bottom), we can observe the diverse scores for each individual action.

Ensemble of ASTRAs. To further enhance the results for the SoccerNet Action Spotting Challenge 2023, we explore the use of an ensemble comprising modifications of ASTRA models. As shown in Figure 6, the removal of different aspects of ASTRA leads to models that maintain strong overall performance while exhibiting diverse predictions. Each of the models demonstrates improved performance for specific actions. The diversity among the models within the ensemble is crucial for achieving effective ensembling. With this in mind, we propose an ensemble that combines our final ASTRA model with the models depicted in Figure 6. These additional models remove audio, focal loss, and uncertainty components, respectively. For each temporal position, we average the

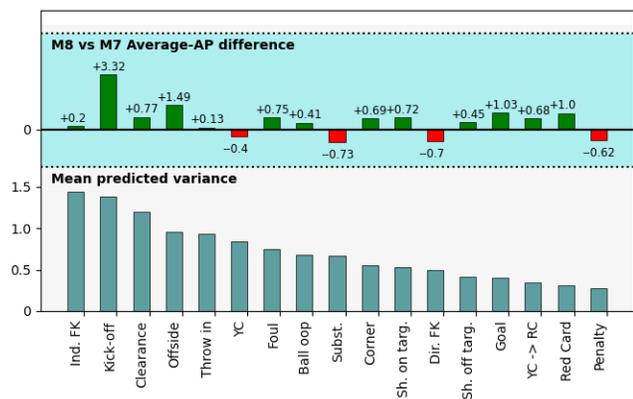


Figure 5: Analysis of M8 model: mean predicted displacement variance in temporal locations with classification probability greater than 0.5 (bottom) and difference in Average-AP between M7 and M8 (top).

predictions of all models in the ensemble. By employing this ensemble approach, we achieve a tight Average-mAP of 67.60 (+0.78). This result emphasizes the ability of appropriately diverse models in an ensemble to provide a slight improvement over the individual performance of ASTRA.

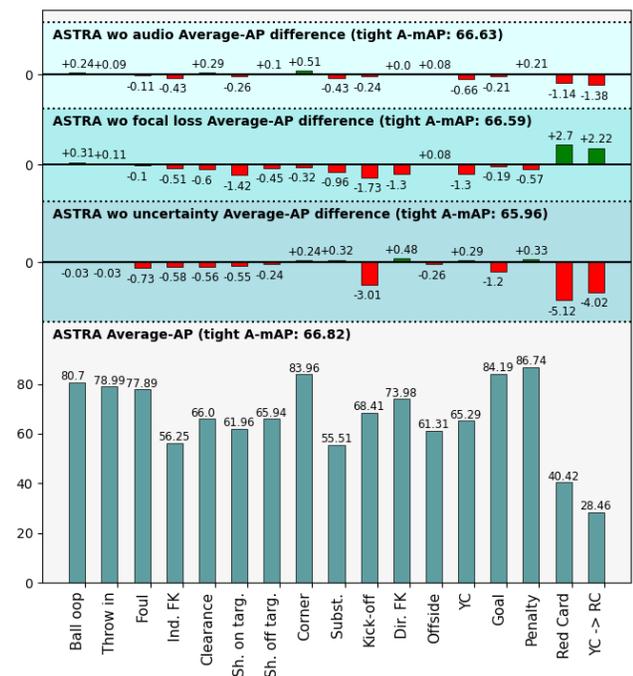


Figure 6: Per-class results of ASTRA and ensemble models. The figure illustrates Average-AP scores for each action in ASTRA (bottom), the differences when removing uncertainty (middle-down), differences when removing focal loss (middle-top), and differences when removing audio (top).

4.5 Results on challenge split

For the evaluation of ASTRA models in the challenge split, we followed the evaluation protocol P2. The results, presented in Table 3, showcase ASTRA’s performance in comparison to the top models in the SoccerNet 2023 Action Spotting challenge. Notably, ASTRA achieves a tight Average-mAP score of 69.43. With the implementation of the ensemble approach, we observe a further improvement, reaching an Average-mAP of 70.21. This result secures the 3rd position in the challenge, surpassing the previous baseline by a +1.88 margin. It is worth noting that ASTRA’s performance stands close to the winning solutions, with a difference gap of 1.10 points from the current SOTA. Additionally, our method achieves the best results on the loose metric and on non-visible actions. The incorporation of label uncertainty modeling and the inclusion of audio input likely contribute to these results, especially in scenarios where label noise is pronounced for non-visible actions.

Model	Tight Average-mAP			Loose Average-mAP		
	All	Vis.	Non vis.	All	Vis.	Non vis.
1- SDU_VSISLAB	71.31	76.29	54.09	78.56	81.67	69.13
2- mt_player	71.10	77.22	58.5	78.79	82.02	77.62
3a- ASTRA (ensemble)*	70.21	75.08	62.34	79.27	81.85	79.39
3b- ASTRA	69.43	74.40	61.10	79.02	81.70	79.47
4- team_ws_action	69.17	75.18	59.12	76.95	80.39	75.92
5- CEA_LVA	68.38	74.79	47.68	73.98	78.57	61.75
Baseline- Yahoo [30]	68.33	73.22	60.88	78.06	80.58	78.32

Table 3: Comparison of ASTRA with other state-of-the-art models in terms of Average-mAP in the SoccerNet 2023 challenge. Metrics include loose and tight evaluations for all actions (All), visible actions (Vis.), and non-visible actions (Non vis.). Results marked with an asterisk (*) indicate slight differences from the official SoccerNet challenge results due to minor modifications in the ASTRA architecture.

5 CONCLUSION

This work presented ASTRA, a model designed to address the task of action spotting in soccer matches. Ablation studies demonstrate the effectiveness of different modules within the model in tackling the challenges inherent to the task and the dataset, such as the need for precise spots, the long-tail distribution of the data, the non-visibility in some actions, and the issue of noisy labels. Additionally, ASTRA achieves good results in the SoccerNet 2023 Action Spotting challenge. It surpasses the previous SOTA performance by +1.88, and its performance is in close proximity to that of the challenge winners.

Acknowledgements. This work has been partially supported by the Spanish project PID2022-136436NB-I00 and by ICREA under the ICREA Academia programme.

REFERENCES

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. 2016. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675* (2016).
- [2] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S Davis. 2017. Soft-NMS—improving object detection with one line of code. In *Proceedings of the IEEE international conference on computer vision*. 5561–5569.
- [3] Shyamal Buch, Victor Escorcia, Bernard Ghanem, Li Fei-Fei, and Juan Carlos Niebles. 2019. End-to-end, single-stream temporal action detection in untrimmed videos. (2019).
- [4] Shyamal Buch, Victor Escorcia, Chuanqi Shen, Bernard Ghanem, and Juan Carlos Niebles. 2017. Sst: Single-stream temporal action proposals. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 2911–2920.
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 213–229.
- [6] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6299–6308.
- [7] Yunze Chen, Mengjuan Chen, Rui Wu, Jiagang Zhu, Zheng Zhu, Qingyi Gu, and Horizon Robotics. 2020. Refinement of Boundary Regression Using Uncertainty in Temporal Action Localization. In *BMVC*.
- [8] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* (2014).
- [9] Anthony Cioppa, Adrien Deliege, Silvio Giancola, Bernard Ghanem, Marc Van Droogenbroeck, Rikke Gade, and Thomas B Moeslund. 2020. A context-aware loss function for action spotting in soccer videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13126–13136.
- [10] Adrien Deliege, Anthony Cioppa, Silvio Giancola, Meisam J Seikavandi, Jacob V Dueholm, Kamal Nasrollahi, Bernard Ghanem, Thomas B Moeslund, and Marc Van Droogenbroeck. 2021. Soccernet-v2: A dataset and benchmarks for holistic understanding of broadcast soccer videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4508–4519.
- [11] Victor Escorcia, Fabian Caba Heilbron, Juan Carlos Niebles, and Bernard Ghanem. 2016. Daps: Deep action proposals for action understanding. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14*. Springer, 768–784.
- [12] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 776–780.
- [13] Silvio Giancola, Anthony Cioppa, Adrien Deliege, Floriane Magera, Vladimir Somers, Le Kang, Xin Zhou, Olivier Barnich, Christophe De Vleeschouwer, Alexandre Alahi, et al. 2022. SoccerNet 2022 challenges results. In *Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports*. 75–86.
- [14] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Freund, Peter Yianilos, Moritz Mueller-Freitag, et al. 2017. The "something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*. 5842–5850.
- [15] Fabian Caba Heilbron, Juan Carlos Niebles, and Bernard Ghanem. 2016. Fast temporal activity proposals for efficient detection of human actions in untrimmed videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1914–1923.
- [16] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. 2017. CNN architectures for large-scale audio classification. In *2017 IEEE international conference on acoustics, speech and signal processing (icassp)*. IEEE, 131–135.
- [17] James Hong, Matthew Fisher, Michaël Gharbi, and Kayvon Fatahalian. 2021. Video pose distillation for few-shot, fine-grained sports action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9254–9263.
- [18] James Hong, Haotian Zhang, Michaël Gharbi, Matthew Fisher, and Kayvon Fatahalian. 2022. Spotting Temporally Precise, Fine-Grained Events in Video. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*. Springer, 33–51.
- [19] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. 2019. Epic-fusion: Audio-visual temporal binding for egocentric action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5492–5501.
- [20] Chuming Lin, Chengming Xu, Donghao Luo, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. 2021. Learning salient boundary feature for anchor-free temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3320–3329.
- [21] Ji Lin, Chuang Gan, and Song Han. 2019. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*. 7083–7093.
- [22] Tianwei Lin, Xu Zhao, and Zheng Shou. 2017. Single shot temporal action detection. In *Proceedings of the 25th ACM international conference on Multimedia*. 988–996.
- [23] Xiaolong Liu, Qimeng Wang, Yao Hu, Xu Tang, Shiwei Zhang, Song Bai, and Xiang Bai. 2022. End-to-end temporal action detection with transformer. *IEEE Transactions on Image Processing* 31 (2022), 5427–5441.
- [24] Banoth Thulasya Naik, Mohammad Farukh Hashmi, and Neeraj Dhanraj Bokde. 2022. A comprehensive review of computer vision in sports: Open issues, future trends and research directions. *Applied Sciences* 12, 9 (2022), 4429.
- [25] Alessandro Pieropan, Giampiero Salvi, Karl Pauwels, and Hedvig Kjellström. 2014. Audio-visual classification and detection of human manipulation actions. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 3045–3052.
- [26] Zhiwu Qing, Haisheng Su, Weihao Gan, Dongliang Wang, Wei Wu, Xiang Wang, Yu Qiao, Junjie Yan, Changxin Gao, and Nong Sang. 2021. Temporal context aggregation network for temporal action proposal refinement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 485–494.
- [27] Muhammad Bilal Shaikh, Douglas Chai, Syed Mohammed Shamsul Islam, and Naveed Akhtar. 2022. MAiVAR: Multimodal Audio-Image and Video Action Recognizer. In *2022 IEEE International Conference on Visual Communications and Image Processing (VCIP)*. IEEE, 1–5.
- [28] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. 2020. Finegym: A hierarchical video dataset for fine-grained action understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2616–2625.
- [29] Dingfeng Shi, Yujie Zhong, Qiong Cao, Lin Ma, Jia Li, and Dacheng Tao. 2023. Tridet: Temporal action detection with relative boundary modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18857–18866.
- [30] Joao VB Soares and Avijit Shah. 2022. Action spotting using dense detection anchors revisited: Submission to the SoccerNet Challenge 2022. *arXiv preprint arXiv:2206.07846* (2022).
- [31] João VB Soares, Avijit Shah, and Topojoy Biswas. 2022. Temporally Precise Action Spotting in Soccer Videos Using Dense Detection Anchors. In *2022 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2796–2800.
- [32] Swathikiran Sudhakaran, Sergio Escalera, and Oswald Lanz. 2020. Gate-shift networks for video action recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 1102–1111.
- [33] Graham Thomas, Rikke Gade, Thomas B Moeslund, Peter Carr, and Adrian Hilton. 2017. Computer vision for sports: Current applications and research topics. *Computer Vision and Image Understanding* 159 (2017), 3–18.
- [34] Bastien Vanderplaetse and Stephane Dupont. 2020. Improved soccer action spotting using both audio and video streams. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 896–897.
- [35] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. 2016. Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision*. Springer, 20–36.
- [36] Ting-Ting Xie, Christos Tzelepis, and Ioannis Patras. 2020. Boundary uncertainty in a single-stage temporal action localization network. *arXiv preprint arXiv:2008.11170* (2020).
- [37] Jinglin Xu, Yongming Rao, Xumin Yu, Guangyi Chen, Jie Zhou, and Jiwen Lu. 2022. Finediving: A fine-grained dataset for procedure-aware action quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2949–2958.
- [38] Mengmeng Xu, Chen Zhao, David S Rojas, Ali Thabet, and Bernard Ghanem. 2020. G-tad: Sub-graph localization for temporal action detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10156–10165.
- [39] Le Yang, Houwen Peng, Dingwen Zhang, Jianlong Fu, and Junwei Han. 2020. Revisiting anchor mechanisms for temporal action localization. *IEEE Transactions on Image Processing* 29 (2020), 8535–8548.
- [40] Boyu Zhang, Jiayuan Chen, Yinfei Xu, Hui Zhang, Xu Yang, and Xin Geng. 2021. Auto-Encoding Score Distribution Regression for Action Quality Assessment. *arXiv preprint arXiv:2111.11029* (2021).
- [41] Chen-Lin Zhang, Jianxin Wu, and Yin Li. 2022. Actionformer: Localizing moments of actions with transformers. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IV*. Springer, 492–510.
- [42] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412* (2017).

- [43] Haotian Zhang, Cristobal Sciotto, Maneesh Agrawala, and Kayvon Fatahalian. 2021. Vid2player: Controllable video sprites that behave and appear like professional tennis players. *ACM Transactions on Graphics (TOG)* 40, 3 (2021), 1–16.
- [44] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. 2018. Temporal relational reasoning in videos. In *Proceedings of the European conference on computer vision (ECCV)*. 803–818.
- [45] Xin Zhou, Le Kang, Zhiyu Cheng, Bo He, and Jingyu Xin. 2021. Feature combination meets attention: Baidu soccer embeddings and transformer based temporal detection. *arXiv preprint arXiv:2106.14447* (2021).