# Football Player Tracking and Re-Identification
# in a Single Video Feed

## 1. Objective

The primary objective of this computer vision project was to develop a robust and efficient solution for identifying and consistently tracking individual football players within a given 15-second high-definition video feed (15sec_input_720p.mp4). A critical requirement was to implement re-identification capabilities, ensuring that players maintain their unique assigned identity even if they temporarily leave the camera's view or are occluded and subsequently re-enter the frame. This work simulates a fundamental component for real-time sports analytics systems.

## 2. Approach and Methodology (Final Solution)

The most successful and robust approach adopted for this project leverages the integrated tracking capabilities provided by the Ultralytics YOLOv11 model.

Key Components:

*Object Detection (Ultralytics YOLOv11):*

A custom-trained and fine-tuned YOLOv11 model (best.pt) serves as the backbone for player detection. YOLOv11 is a single-stage object detector known for its excellent balance of speed and accuracy, making it ideal for real-time applications.

The model processes each video frame to identify instances of "players" (and potentially the "ball" as per the training, though only players are actively tracked and visualized in the final output). It provides precise bounding box coordinates, confidence scores, and class labels for each detected object.

*Integrated Tracking (model.track() / ByteTrack):*

Instead of a custom tracking implementation, the project utilizes the highly optimized model.track() method provided directly by the Ultralytics library.

This function internally incorporates advanced tracking algorithms, prominently ByteTrack, which is leveraged when essential dependencies like supervision are installed. ByteTrack is a "tracking-by-detection" algorithm that excels at associating low-confidence detections and handling short-term occlusions, significantly improving track consistency.

Mechanism:

Detection-to-Track Association: ByteTrack associates newly detected bounding boxes with existing tracks. It employs a two-stage association strategy: first matching high-confidence detections, then using the remaining unassociated tracks and low-confidence detections to recover objects that might have been temporarily missed or occluded.

Motion Prediction: A Kalman filter is integrated to predict the future position of existing tracks, aiding in the association process and providing smooth trajectories even when detections are sporadic.

Track Management: The tracker intelligently manages the lifecycle of each track. The persist=True parameter in model.track() is crucial, as it maintains track history across frames, enabling the re-

identification of players who might have exited and re-entered the field of view. This ensures that the same player consistently retains their unique track_id.

## Video Processing and Visualization (OpenCV):

The OpenCV library (cv2) is used for all video input/output operations. It facilitates reading frames from the 15sec_input_720p.mp4 file and writing the processed frames to the output_tracked.avi video.

OpenCV is also instrumental in visualizing the tracking results, specifically for drawing accurate bounding boxes around detected players and rendering their unique track_id on each frame.

## Workflow:

The entire tracking process is orchestrated within a single Python script (tracking.py):

*Initialization:* The Ultralytics YOLOv11 model is loaded, and video capture (cv2.VideoCapture) and video writer (cv2.VideoWriter) objects are set up.

*Stream Processing:* The model.track() method is invoked with the input video stream (source=VIDEO_PATH, stream=True). Key parameters like conf=0.5 (detection confidence threshold) and persist=True (for re-identification) are configured.

*Frame-by-Frame Annotation:* The model.track() method yields results for each frame. For every frame:

> It retrieves the original image (result.orig_img).

> It accesses the result.boxes attribute to obtain bounding box coordinates (xyxy), confidence scores (conf), class IDs (cls), and crucially, the assigned track_id for each detected object.

> The code iterates through these tracked objects, filtering specifically for the 'player' class.

For each player, a green bounding box (cv2.rectangle) and their unique ID: {track_id} label (cv2.putText) are drawn onto the frame.

*Output Generation:* The annotated frame is then written to the output_tracked.avi file by the cv2.VideoWriter.

*Cleanup:* Upon completion or interruption, video resources are gracefully released.

## 3. Techniques Tried and Their Outcomes

The project's evolution involved a systematic exploration of different tracking methodologies, driven by the goal of achieving robust re-identification.

Initial Attempt: *DeepSORT Integration with ResNet-based Feature Extractor*

Methodology: The first approach involved an integration of YOLOv11 for detection with the DeepSORT (Deep Simple Online Realtime Tracking) algorithm. DeepSORT requires appearance features for re-identification. To fulfill this, a custom resnet_extractor.py was developed to generate fixed-size feature vectors from cropped player images using a pre-trained ResNet model. These extracted features, along with bounding box coordinates and confidence scores from YOLOv11, were then fed into the DeepSORT tracker.

Outcome: While this setup enabled basic multi-object tracking, it suffered from significant ID switching issues. Players frequently received new IDs after even brief occlusions or when they were in close proximity to other players. The generic nature of ResNet features, not specifically trained for person re-identification in sports, proved insufficient for consistently distinguishing individuals, leading to unreliable re-identification accuracy.

Second Attempt: *DeepSORT Integration with OSNet Feature Extractor (via torchreid)*

Methodology: Recognizing the limitations of generic ResNet features for person re-identification, the methodology was refined to incorporate OSNet (Omni-Scale Network). OSNet is a state-of-the-art neural network specifically designed for person re-identification, known for its strong discriminative capabilities across various scales. It was integrated into the pipeline via the torchreid library. The objective was to leverage OSNet's superior appearance embeddings to improve track continuity.

Outcome: This iteration showed notable improvement in re-identification accuracy compared to the previous ResNet approach, as OSNet's features were indeed more distinct. However, ID switching persisted in more challenging scenarios, particularly during rapid, complex movements, prolonged occlusions, or when players were tightly clustered.

Final Solution: *Ultralytics model.track() (Leveraging ByteTrack)*

Methodology: Driven by the persistent challenges and the desire for a more robust and streamlined solution, the project pivoted to utilizing the built-in model.track() function provided by the Ultralytics framework. This function encapsulates advanced, highly optimized tracking algorithms (such as ByteTrack, which is seamlessly integrated when supervision is available) that are specifically designed to work in conjunction with YOLO detections.

Outcome: This approach dramatically improved the accuracy and reliability of player re-identification. ByteTrack's ability to consider both high and low confidence detections, combined with its efficient motion model, resulted in highly consistent IDs for players throughout the video, even when they moved out of frame and reappeared.

## 4. Challenges Encountered

The journey to a robust tracking system presented several key challenges:

*Consistent ID Assignment:* The most significant hurdle was consistently assigning and maintaining a unique ID for each player across consecutive frames, especially when they were involved in fast movements, passed behind other players, or briefly left the field of view. Initial methods struggled with track fragmentation and ID flips.

*Feature Robustness for Re-identification:* Developing or integrating appearance features that were sufficiently robust and discriminative to differentiate between similarly uniformed players in dynamic conditions was crucial. Generic features proved inadequate, and even specialized Re-ID networks required careful tuning within the tracking framework.

*Integration Complexity and Efficiency:* Manually stitching together separate detection, feature extraction, and tracking components (as in the DeepSORT attempts) introduced considerable code complexity, debugging overhead, and potential performance bottlenecks. Ensuring efficient data flow between these modules was a continuous challenge.

## 5. Future Improvements and Considerations

Given more time and resources, the project could be further enhanced by:

*Robustness in Extreme Conditions:*

Investigating StrongSORT, an extension of DeepSORT that integrates advanced re-identification networks (like OSNet or BoT-SORT) and better motion models to potentially achieve even higher accuracy in extremely challenging scenarios, such as very dense player clusters or prolonged, complex occlusions.

*Scalability for Large-Scale Deployment:*

Testing and optimizing the solution for scalability on longer video clips or across multiple camera feeds simultaneously. This would involve exploring distributed processing frameworks or more advanced stream management techniques to ensure consistent performance in broader applications.