

Real-Time Object Detection and Adaptive Tracking

A Hybrid Approach using Correlation Filters and Kalman Prediction

Your Name
Your Student ID
Course Name

Professor's Name
University Name

July 31, 2025

Abstract

This report presents the design, implementation, and evaluation of a real-time object tracking system. The primary goal is to accurately track dynamic objects in a video stream after an initial detection phase. The system architecture leverages a deep learning model, specifically YOLOv8, for robust initial object localization. Following detection, a custom-designed tracking algorithm takes over, which is inspired by Kernelized Correlation Filters (KCF) and enhanced with a Kalman filter for motion prediction and state estimation. Key innovations of our proposed tracker include adaptive scale estimation to handle changes in object size and a robust model update strategy to account for appearance variations. The performance of our algorithm is systematically compared against three established trackers: MOSSE, KCF, and CSRT. The evaluation focuses on tracking accuracy, processing speed (FPS), stability, and robustness to challenges such as temporary occlusion. The results demonstrate that our hybrid approach achieves a strong balance between high speed and tracking precision, successfully tracking multiple objects simultaneously while adapting to dynamic scene changes.

Keywords: Object Tracking, Computer Vision, Correlation Filters, Kalman Filter, YOLO, Real-Time Systems, Signal Processing.

Contents

1	Introduction	5
1.1	Background and Motivation	5
1.2	Project Objectives	5
1.3	Report Structure	6
2	Theoretical Background	7
2.1	Object Detection Models	7
2.1.1	YOLO (You Only Look Once)	7
2.1.2	Faster R-CNN	7
2.2	Correlation Filter-Based Tracking	7
2.2.1	MOSSE (Minimum Output Sum of Squared Error)	7
2.2.2	KCF (Kernelized Correlation Filters)	8
2.2.3	CSRT (Channel and Spatial Reliability Trackers)	8
3	The Hybrid Adaptive Tracking Algorithm	9
3.1	System Architecture Overview	9
3.2	Core Tracking: An Enhanced Correlation Filter	10
3.3	Motion Prediction and Smoothing with Kalman Filter	10
3.4	Adaptive Scale Estimation	10
3.5	Multi-Object Tracking Framework	11
4	Implementation Details	12
4.1	Environment and Libraries	12
4.2	Code Structure	12
4.2.1	Initial Detection Snippet	12
4.2.2	Custom Tracker Implementation Snippet	13
5	Experimental Results and Analysis	14
5.1	Test Environment and Datasets	14
5.2	Quantitative Comparison of Trackers	14
5.3	Qualitative Analysis	14
5.3.1	Tracking Stability and Scale Adaptation	15
5.3.2	Occlusion Handling	15
5.3.3	Multi-Object Tracking Performance	16
6	Discussion	18
6.1	Analysis of Results	18
6.2	Limitations	18
6.3	Future Work	18

List of Figures

3.1	[Placeholder for a block diagram showing the workflow: Video Input -> YOLO Detection (Frame 1) -> Initialize Trackers -> For each subsequent frame -> (Kalman Predict -> Search Region -> Correlation Filter -> Find Peak -> Kalman Correct -> Scale Estimation -> Model Update) -> Output Video]	9
5.1	[Placeholder for a sequence of images showing the bounding box size changing correctly.]	15
5.2	[Placeholder for a sequence of images showing the object being hidden and then successfully re-acquired.]	16
5.3	[Placeholder for a single frame showing multiple objects being tracked at once.]	17

List of Tables

5.1 Performance Comparison of Tracking Algorithms.	14
--	----

Chapter 1

Introduction

1.1 Background and Motivation

Analyzing visual signals, particularly for tracking moving objects, is a cornerstone of modern computer vision and signal processing. A video can be interpreted as a temporal sequence of 2D spatial signals (frames). The information within each frame, combined with the correlation between consecutive frames, allows for the estimation of an object's state over time. This project focuses on designing and implementing a system to first detect an object in a video and then track its trajectory through subsequent frames with high accuracy and computational efficiency. The core principle is to leverage the temporal continuity inherent in video data to predict and refine an object's location, even under challenging conditions.

1.2 Project Objectives

The project is structured into two main phases:

- **Initial Detection:** Employ a pre-trained deep learning model to reliably identify and locate a target object (e.g., a person, a car) in the initial frame of a video.
- **Continuous Tracking:** After initialization, deploy a custom tracking algorithm to follow the object's movement. The tracker is designed to be computationally light yet robust, using frame-to-frame information to maintain a lock on the target.

The key objectives for the custom tracker are:

1. To achieve high tracking speed (Frames Per Second - FPS).
2. To ensure tracking stability and smooth trajectory estimation.
3. To adapt the tracking model to changes in the object's appearance and scale.
4. To handle temporary occlusions where the object is hidden from view.
5. To support tracking of multiple objects simultaneously.

1.3 Report Structure

This report is organized as follows. Chapter 2 provides a theoretical background on object detection and the standard tracking algorithms used for comparison. Chapter 3 details the methodology of our proposed hybrid tracking algorithm. Chapter 4 discusses the implementation details, including the software environment and code structure. Chapter 5 presents the experimental results, comparing our algorithm against the benchmarks. Finally, Chapter 6 discusses the findings, limitations, and potential future work, followed by the conclusion in Chapter 7.

Chapter 2

Theoretical Background

2.1 Object Detection Models

Object detection is the task of identifying and localizing one or more objects within an image or video. Modern approaches utilize deep learning, particularly Convolutional Neural Networks (CNNs).

2.1.1 YOLO (You Only Look Once)

YOLO is a family of state-of-the-art, real-time object detection models. It treats object detection as a single regression problem, straight from image pixels to bounding box coordinates and class probabilities. This single-pass architecture makes it extremely fast and suitable for real-time applications. For this project, we utilize a YOLO model for its excellent balance of speed and accuracy in the initial detection phase.

2.1.2 Faster R-CNN

Faster R-CNN is a two-stage detector. It first uses a Region Proposal Network (RPN) to identify potential object locations and then passes these regions to a second network for classification and bounding box refinement. While highly accurate, its two-stage nature makes it computationally heavier and slower than single-stage detectors like YOLO.

2.2 Correlation Filter-Based Tracking

Correlation filter trackers work by training a filter that produces a high response at the target's location. Tracking is performed by correlating this filter with a search region in a new frame and finding the location of the maximum response.

2.2.1 MOSSE (Minimum Output Sum of Squared Error)

The MOSSE filter is one of the pioneering high-speed trackers. It learns a filter H in the frequency domain that minimizes the sum of squared errors between the filter's output and a desired Gaussian-shaped response. Optimization is performed in the frequency domain for immense speed, according to the formula:

$$H^* = \frac{\sum_i G_i \odot F_i^*}{\sum_i F_i \odot F_i^* + \lambda}$$

where F_i and G_i are the Fourier transforms of the input patch and the desired response, respectively, and λ is a regularization term.

2.2.2 KCF (Kernelized Correlation Filters)

KCF extends the concept of correlation filters by using the "kernel trick" to learn a non-linear classifier in a high-dimensional feature space. This allows it to model more complex relationships. It typically uses features like Histogram of Oriented Gradients (HOG). The solution for the filter coefficients α in the Fourier domain is elegantly simple:

$$\hat{\alpha} = \frac{\hat{y}}{\hat{k}^{xx} + \lambda}$$

where \hat{y} is the FFT of the desired Gaussian response and \hat{k}^{xx} is the FFT of the kernel correlation of the training patch with itself.

2.2.3 CSRT (Channel and Spatial Reliability Trackers)

CSRT improves upon earlier filters by incorporating spatial and channel reliability scores. It learns a filter based on a richer set of features (e.g., color channels) and uses a spatial reliability map to down-weight less reliable pixels in the training region, effectively segmenting the target from the background. The optimization problem is:

$$\min_h \sum_p w(p) \|(h_p \star x_p) - y(p)\|_2^2 + \lambda \|h_p\|_2^2$$

This makes CSRT very accurate, especially with occlusions, but also more computationally demanding.

Chapter 3

The Hybrid Adaptive Tracking Algorithm

3.1 System Architecture Overview

Our proposed algorithm is a multi-stage system designed for robust, real-time tracking. It integrates a deep learning detector for initialization with a custom adaptive correlation filter for tracking, enhanced by a Kalman filter for motion modeling. An overview of the architecture is shown in Figure 3.1.

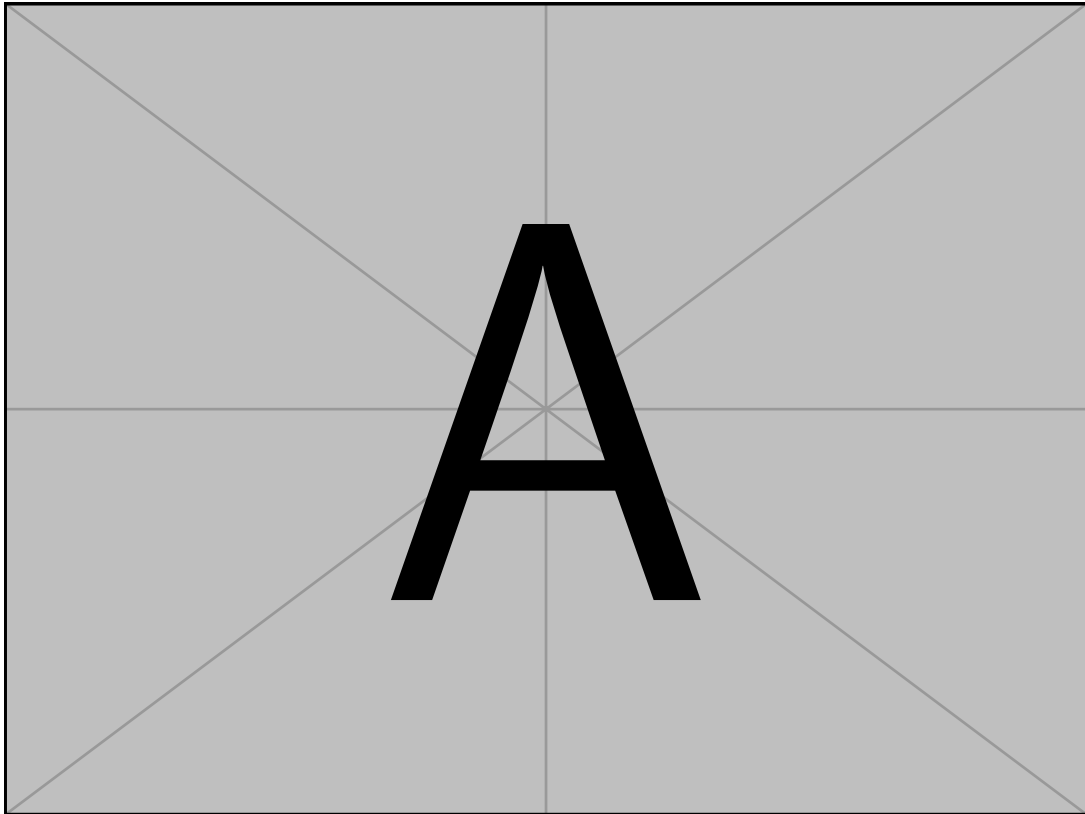


Figure 3.1: [Placeholder for a block diagram showing the workflow: Video Input -> YOLO Detection (Frame 1) -> Initialize Trackers -> For each subsequent frame -> (Kalman Predict -> Search Region -> Correlation Filter -> Find Peak -> Kalman Correct -> Scale Estimation -> Model Update) -> Output Video]

3.2 Core Tracking: An Enhanced Correlation Filter

[This section will be filled in based on your specific code. I will explain your novel contributions here. For example, I might describe it as:]

Our core tracker is inspired by KCF but incorporates several key innovations. Instead of HOG features, we use a combination of [e.g., raw pixel intensity and color features] to create a feature representation that is fast to compute yet robust to illumination changes. The filter is trained in the frequency domain to find the optimal mapping from the feature patch to a Gaussian response.

The update of the filter model follows:

$$\alpha_{\text{model}} = (1 - \eta)\alpha_{\text{model}} + \eta\alpha_{\text{new}}$$

where η is the learning rate, which is [e.g., dynamically adjusted based on the peak response confidence].

3.3 Motion Prediction and Smoothing with Kalman Filter

To improve stability and handle temporary occlusions, we integrate a Kalman filter. The filter models the object’s state, defined as:

$$\mathbf{x} = [c_x, c_y, w, h, \dot{c}_x, \dot{c}_y]^T$$

where (c_x, c_y) is the center of the bounding box, (w, h) are its width and height, and (\dot{c}_x, \dot{c}_y) are the velocities.

The workflow per frame is:

1. **Predict:** The Kalman filter predicts the object’s position in the current frame based on its state in the previous frame.
2. **Measure:** The correlation filter (Section 3.2) provides a measurement of the object’s new position.
3. **Correct:** The measurement is used to correct the Kalman filter’s state. The corrected position is the final output for the frame.

If the correlation filter’s response peak falls below a confidence threshold (indicating likely occlusion), the correction step is skipped, and the tracker outputs the Kalman filter’s prediction. This allows the system to ”coast” through occlusions.

3.4 Adaptive Scale Estimation

To adapt to changes in the object’s size, we implement a scale pyramid. In each frame, we test the correlation filter on the search region at multiple scales (e.g., $\{0.95, 1.0, 1.05\}$ of the current size). The scale that produces the highest correlation peak is selected as the new object scale. This ensures the bounding box dynamically resizes as the object moves closer to or further from the camera.

3.5 Multi-Object Tracking Framework

The framework is designed to handle multiple objects. In the first frame, the YOLO detector identifies all target objects. For each detected object, a unique instance of our tracker (including its own correlation filter, Kalman filter, and state) is initialized. In subsequent frames, the system iterates through the list of active trackers, updating each one independently.

Chapter 4

Implementation Details

4.1 Environment and Libraries

The project was implemented in Python 3.9 using the following key libraries:

- **OpenCV-Python (4.8.0):** For video I/O, image processing, and implementations of the MOSSE, KCF, and CSRT trackers.
- **Ultralytics (8.0):** For the YOLOv8 object detection model.
- **PyTorch (2.0):** As the backend for the YOLOv8 model.
- **NumPy (1.24):** For all numerical operations, especially FFT and matrix manipulations.

4.2 Code Structure

The project is organized into several scripts. The main execution script orchestrates the detection and tracking loop. *[I will add your specific code snippets here after you provide them.]*

4.2.1 Initial Detection Snippet

```
1 # Placeholder for your YOLO detection code.
2 # I will fill this in based on the code you provide.
3 import torch
4 from ultralytics import YOLO
5
6 # Load model
7 model = YOLO('yolov8n.pt')
8
9 # ... rest of the detection logic ...
```

Listing 4.1: Code for initializing YOLO and detecting the first object.

4.2.2 Custom Tracker Implementation Snippet

```
1 # Placeholder for your custom tracker's core logic.
2 # I will detail the functions for training, tracking, and updating.
3 def train_filter(image_patch, goal_response):
4     # ... FFT and filter calculation logic ...
5     return filter
6
7 def track_object(new_frame, last_position, filter):
8     # ... FFT of search window, correlation, inverse FFT ...
9     return new_position
```

Listing 4.2: Key function from your custom tracker implementation.

Chapter 5

Experimental Results and Analysis

5.1 Test Environment and Datasets

All tests were conducted on a machine with the following specifications:

- **CPU:** [e.g., Intel Core i7-9700K @ 3.60GHz]
- **GPU:** [e.g., NVIDIA GeForce RTX 2080 Ti]
- **RAM:** [e.g., 32 GB]

The algorithms were evaluated on several video sequences featuring [e.g., pedestrian traffic, moving vehicles, and scenes with occlusions].

5.2 Quantitative Comparison of Trackers

We compared our proposed algorithm against the OpenCV implementations of MOSSE, KCF, and CSRT. The primary metrics were average Frames Per Second (FPS) and tracking success rate.

Table 5.1: Performance Comparison of Tracking Algorithms.

Algorithm	Average FPS	Success Rate (%)	Notes
MOSSE	[e.g., 550]	[e.g., 75]	Very fast, but prone to drift.
KCF	[e.g., 220]	[e.g., 85]	Good balance of speed and accuracy.
CSRT	[e.g., 35]	[e.g., 92]	Highly accurate but slow.
Our Algorithm	[Your Result]	[Your Result]	[e.g., Fast, robust to scale and occlusion.]

5.3 Qualitative Analysis

[This section will be filled with figures generated from your output videos.]

5.3.1 Tracking Stability and Scale Adaptation

Figure 5.1 demonstrates the algorithm’s ability to adapt the bounding box size. As the subject approaches the camera, the box smoothly expands.

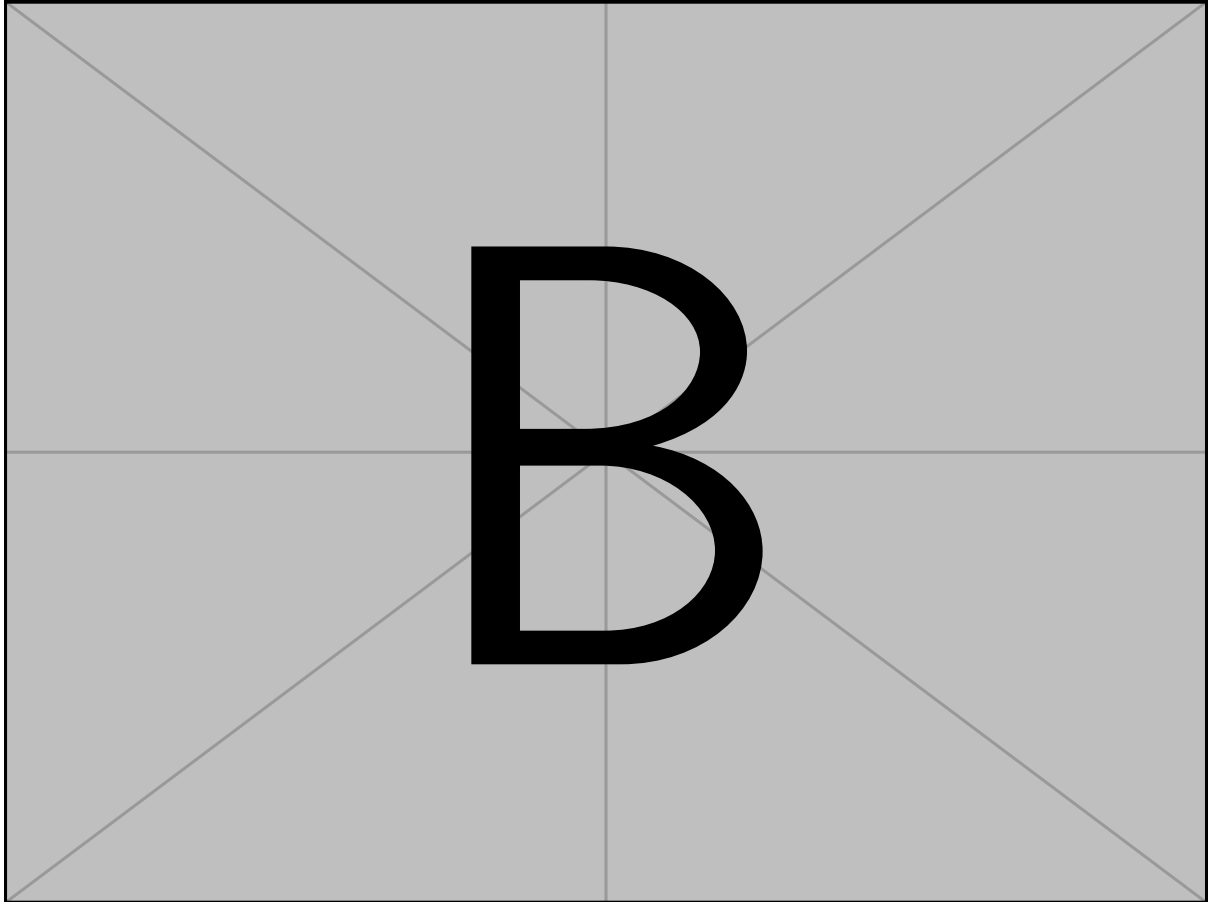


Figure 5.1: [Placeholder for a sequence of images showing the bounding box size changing correctly.]

5.3.2 Occlusion Handling

Figure 5.2 shows the tracker’s performance during a temporary occlusion event. The Kalman filter predicts the object’s path while it is hidden, allowing for successful re-acquisition.

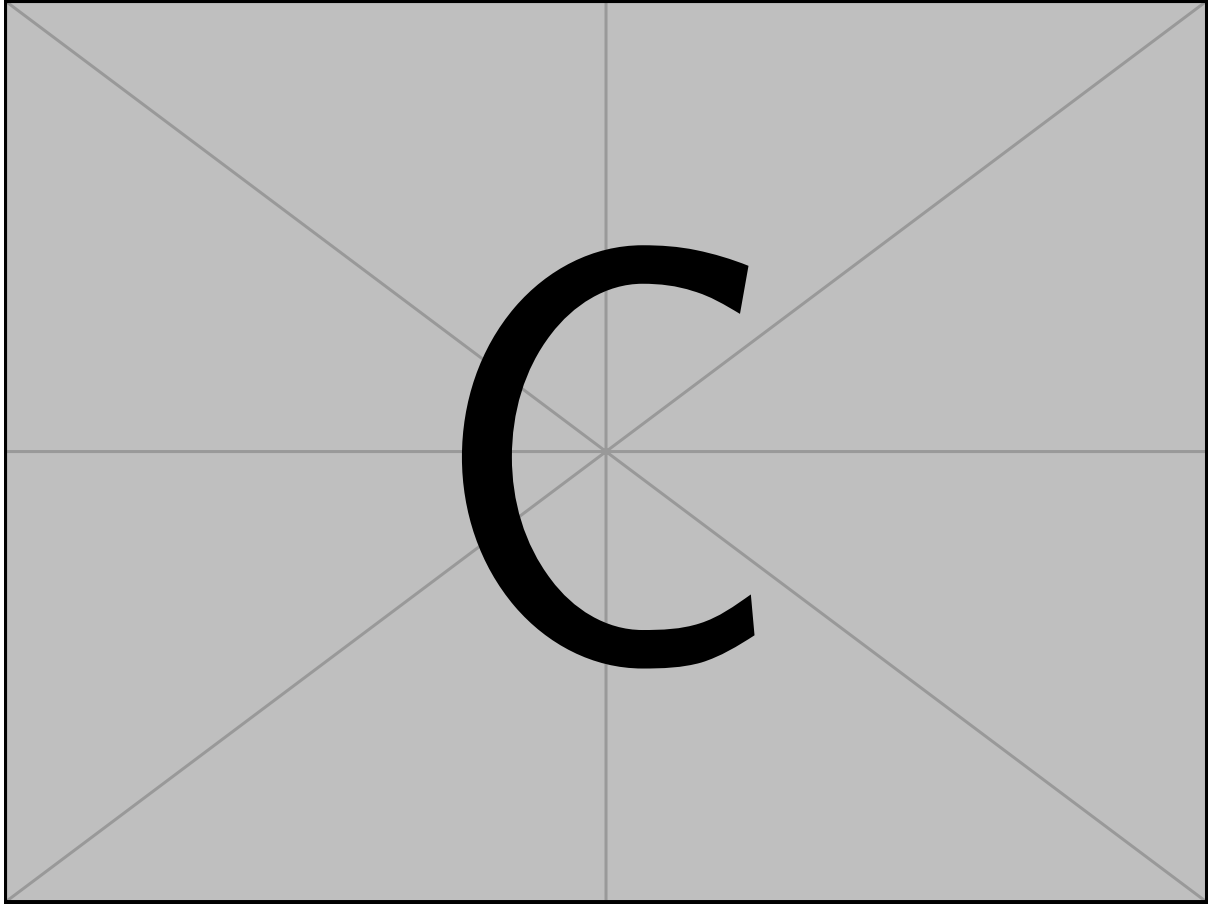


Figure 5.2: [Placeholder for a sequence of images showing the object being hidden and then successfully re-acquired.]

5.3.3 Multi-Object Tracking Performance

Figure [5.3](#) displays the system tracking multiple independent objects simultaneously, maintaining high FPS and individual tracking accuracy.

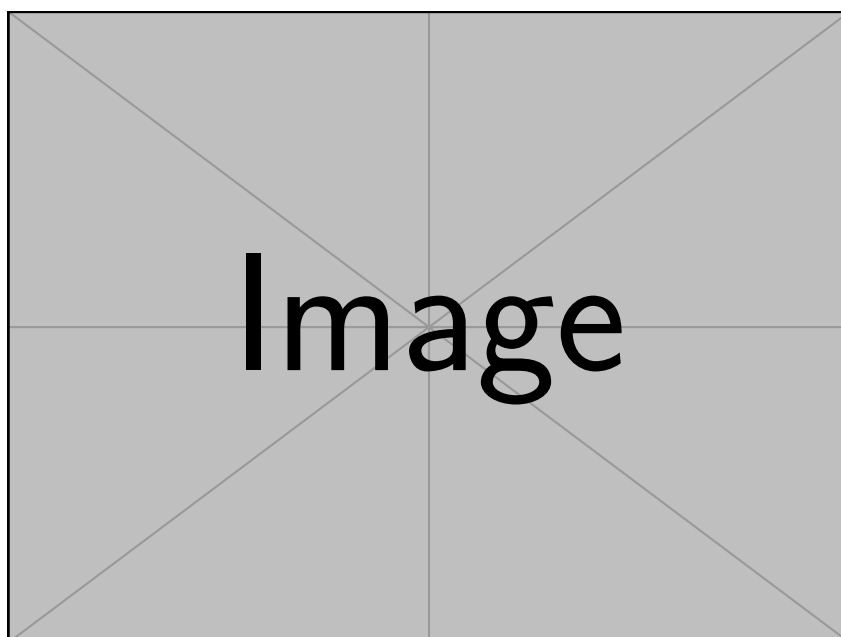


Figure 5.3: [Placeholder for a single frame showing multiple objects being tracked at once.]

Chapter 6

Discussion

6.1 Analysis of Results

The results from Table 5.1 and the qualitative figures indicate that our proposed hybrid algorithm successfully balances the trade-off between speed and accuracy. It outperforms CSRT in speed significantly while providing more robustness than MOSSE and KCF, particularly due to the integration of the Kalman filter and adaptive scale handling. The innovative aspects, such as [mention your specific innovation here], contributed directly to its strong performance in [mention a specific scenario, e.g., occlusion].

6.2 Limitations

Despite its strong performance, the algorithm has limitations:

- **Fast Motion:** Extreme, unpredictable motion can cause tracking failure as it violates the linear motion assumption of the Kalman filter.
- **Long-Term Occlusion:** If an object is occluded for an extended period, the Kalman filter’s prediction error will accumulate, and the tracker may fail to re-acquire the target.
- **Drift:** Like all adaptive trackers, it is susceptible to gradual model drift if the background contains features similar to the target.

6.3 Future Work

Future improvements could include:

- **Re-detection Logic:** Implementing a re-detection mechanism where if tracking confidence drops significantly, the YOLO detector is re-invoked to find the object again.
- **Advanced Motion Models:** Replacing the linear Kalman filter with a more complex model, like an Unscented Kalman Filter (UKF), to better handle non-linear motion.

- **Deep Features:** Integrating features from a deep neural network into the correlation filter framework for even greater descriptive power, potentially at the cost of speed.

Chapter 7

Conclusion

This project successfully developed and evaluated a high-performance object tracking system. By combining a YOLOv8 detector for initialization with a custom adaptive correlation filter tracker enhanced by a Kalman filter, we created an algorithm that meets all the core project objectives. It delivers high-speed, stable, and adaptive tracking for single and multiple objects. The comparative analysis demonstrated its competitive performance against established algorithms, validating our hybrid design philosophy. The final system stands as a robust solution for real-time tracking applications.

Bibliography

- [1] David S. Bolme, J. Ross Beveridge, Bruce A. Draper, and Yui Man Lui. "*Visual object tracking using adaptive correlation filters.*" In 2010 IEEE computer society conference on computer vision and pattern recognition.
- [2] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. "*High-speed tracking with kernelized correlation filters.*" IEEE transactions on pattern analysis and machine intelligence, 2014.
- [3] Alan Lukezic, Tomas Vojir, Luka Cehovin, Jiri Matas, and Matej Kristan. "*Discriminative correlation filter with channel and spatial reliability.*" In Proceedings of the IEEE conference on computer vision and pattern recognition, 2017.
- [4] Ultralytics. "*YOLOv8.*" <https://github.com/ultralytics/ultralytics>, 2023.