

Deep Learning-Based Multi-Object Group Detection and Tracking in Video Streams

MSc Research Project Master's Degree in Artificial Intelligence

Erkan Uci Student ID: 23210494

School of Computing National College of Ireland

Supervisor

Sheresh Zahoor

National College of Ireland Project Submission Sheet School of Computing



StudentName:	Erkan Uci	
StudentID:	23210494	
Programme:	Master's Degree in Artificial Intelligence	
Year:	2024	
Module:	MSc Research Project	
Supervisor:	Sheresh Zahoor	
Submission Due Date:	12/08/2024	
ProjectTitle: Deep Learning-Based Multi-Object Group De		
	Tracking in Video Streams	
WordCount:	5608	
PageCount:	27	

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

	1	
Signature:		
Date:	42/08/2024	

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).		
Attach a Moodle submission receipt of the online project submission, to each project (in-	57	
cluding multiple copies).		
You must ensure that you retain a HARD COPY of the project, both for your own reference	∇	
and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	V	

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applica-	
ble):	

1 Abstract

Accurately detecting and continuously tracking multiple object classes in video data datasets are critical challenges in computer vision, especially for autonomous vehicles and for applications such as video analytics. This project focuses on multi-object detection and tracking using the MOT17 dataset, leveraging the latest artificial intelligence and deep learning techniques to address these challenges.

Our methodology is a self-contained architecture using different advanced deep learning models to improve object detection and tracking accuracy. Convolutional Neural Networks (CNNs) are used to extract the correct features from video frames and to select objects of interest labeled. Long Short-Term Memory (LSTM) networks are included to preserve temporal dependencies and allow moving objects to be tracked seamlessly between frames. In addition, Faster R-CNN and R-CNN frameworks are integrated to improve object localization and classification through spatial and region-based recommendations and improved domain analysis.

To verify the reliability and robustness of our work, we conducted extensive experiments on various video datasets covering different environmental conditions and object densities. The results show that the combined use of CNNs, LSTMs, Faster R-CNNs and R-CNNs significantly improves the accuracy and reliability of multi-object detection and tracking.

By integrating computer vision and machine learning, this work provides important insights into their application in the real world, especially in autonomous vehicles and security systems. As mentioned before, the information gathered from different datasets and the proposed results using the methodology we have implemented will add significant value to the vision of high accuracy and reliable multi-object detection and tracking.

Keywords: Deep Learning; Computer Vision; Multi-Label Classification; Region Proposal Networks (RPN); Multi-Object Tracking; Object Detection

2 Introduction

Accurately and reliably identifying and tracking multiple objects in video data or live archive broadcasts is crucial for video monitoring, self-driving and applications such as human-computer interactions [1, 2]. Such systems or applications need to detect, classify and track various objects such as people, cars and animals across video frames [3]. The challenge is to uniquely identify these objects and track them seamlessly across consecutive frames, which is vital for real-time tracking systems.



Figure 2.1: Illustration of the output of an MOT algorithm. Each output bounding box has a number that identifies a specific object in the video.

Various artificial intelligence models are being studied to improve object detection and tracking. In multi-object tracking, bounding boxes are drawn for visual perception, appearance, representation and labeling so that the system can detect and track the correct object. Object tracking is achieved by assigning unique IDs to each detected object and then storing these IDs in subsequent frames.

Convolutional Neural Networks (CNNs) have become widely used for extracting features within objects and for their ability to accurately detect objects within frames in video datasets. Faster R-CNN enables this process to improve and more accurate detection by providing better region recommendations and more advanced analysis of interest capacity. In addition, Long Short Term Memory (LSTM) networks [4, 5] are a type of Recurrent Neural Network (RNN) [6], which are used to identify temporal developments by providing a smooth tracking of moving frames [6]. However, in some cases with complex datasets, some difficulties remain in dynamic scenes [7, 8].

In our work, we focus on advanced differential deep learning models to perform multi-object detection and tracking using different datasets in MOT17 datasets. We work on CNNs, Faster R-CNN and LSTM networks to improve the accuracy and robustness of the main system. Fur-

thermore, data fusion techniques and methods such as dropout and L2 regularization are integrated into the module to avoid overfitting in some cases and increase the accuracy of the overall model.

Key Contributions of this Study:

- Advanced Model Integration: Combining CNNs for feature extraction, Faster R-CNN for better object localization, and LSTMs for managing temporal dependencies.
- Advanced Data Optimization: Using dropout and L2 regularization during training to improve model reliability and prevent overfitting.
- Extensive Testing: Running extensive experiments on the MOT17 dataset to validate the model under various conditions and object densities.
- High Performance: Achieving the best results in handling complex object interactions and dynamic scenes compared to existing methods.

The structure of the paper is as follows. The **methodology** section describes the proposed methodology and its construction in detail, with subsections covering different parts of the approach. The **modeling** section summarizes the experimental setup and findings, followed by a discussion in the **discussion** section. The paper concludes with a summary of the results and future research directions in the **conclusion** section.

3 Related Work

Multi-Object Tracking (MOT) has important applications in areas such as security systems, autonomous driving and human-computer sharing [1, 2]. Many strategies can be used to improve the stability and sharing of MOT. In this section, we review important work in this area.

3.1 Multiple object tracking

MOT plays a critical role in computer vision research. It aims to detect and track the identities and locations of different objects over time. MOT faces many challenges such as blocked views, object appearance changes, and complex object interactions [3]. It is important to consider the ground truth information of the object in order to identify the correct object. Figure 3.1 shows the ground truth information drawn on the object.



Figure 3.1: The image above shows how detected objects are matched to ground truth data.

• Classical Tracking Methods:

Traditional methods for tracking multiple objects often use Kalman filters. These filters, along with more advanced versions such as the Extended Kalman Filter (EKF) and the Unscented Kalman Filter (UKF), are designed to manage noisy data and predict the motion of objects [1, 2, 9].

• Deep Learning-Based Methods: The advent of deep learning has had a significant impact on tracking accuracy. Convolutional Neural Networks (CNNs) have become essential tools for detecting and tracking objects, and the Faster R-CNN has been particularly effective at identifying objects in video frames [4, 5, 6, 7], making it a leading choice for MOT. In addition, Long Short Term Memory (LSTM) networks are used to understand temporal relationships in video sequences. LSTMs excel at preserving the identity of an object in consecutive frames [5], which is used to provide more reliable tracking.



Figure 3.2: Ground truth box vs predicted box

3.2 Tracking-by-detection

Tracking objects in video data sequences usually involves detecting the object in each frame and linking these detections together to form a continuous motion path. The "tracking by detection" approach has recently gained popularity due to its flexibility and advances in object detection technology. This chapter highlights the basic concepts, techniques and recent developments in this approach.



Figure 3.3: Tracking by detection

Precision Precision measures how well you can find true positives (TP) among all positive predictions.

Precision = TPTP + FP

Figure 3.4: Precision formula

For example, the IoU threshold is critical for measuring sensitivity in object detection [11]. In the example image, the cat on the left has an IoU value of 0.3, which is below the threshold, making it a false positive. Conversely, the cat on the right, with an IoU of 0.7, exceeds the threshold and is correctly identified as a true positive.

If IoU threshold = 0.5







Detection This involves detecting and positioning objects in every frame. Faster algorithms such as Faster R-CNN, YOLO and SSD are crucial for multi-object tracking. They can vary in speed and accuracy.

State-of-the-art Detectors:

- Faster R-CNN: Utilizes region proposal networks for efficient object detection [6].
- YOLO (You Only Look Once): Achieves real-time object detection with high accuracy [7].

• SSD (Single Shot MultiBox Detector): Balances speed and accuracy in object detection [8].

Detection Formulation:

The general objective function for a detector can be formulated as:

$$L = L_{cls} + \lambda L_{reg}$$

where L_{cls} is the classification loss, L_{reg} is the regression loss for bounding box coordinates, and λ is a balancing parameter.

Tracking

First, Faster R-CNN allows you to accurately detect objects in the video stream; then, real-time tracking is performed by utilizing the Deep SORT algorithm. More specifically, Deep SORT identifies the traces of objects along a video sequence and also uses iterative Kalman filtering to detect the relationship of objects within a frame to the next frame, along with a data association strategy to overlay their predicted visual appearance from bounding boxes on the image [13].



Figure 3.6: DeepSort tracking algorithm flowchart.

3.3 Datasets

Datasets play a vital role in developing and benchmarking MOT algorithms.

- MOT 17 Challenge: Provides a comprehensive benchmark with annotated sequences for evaluating MOT methods [6].
- KITTI: Focuses on tracking in autonomous driving scenarios [7].
- COCO: Offers a large-scale dataset for detection and segmentation tasks [8].

3.4 Evaluation metrics

The evaluations were conducted using the MOT17 standard metrics that assess multi-object tracking (MOTA), identity changes (ID) and tracking precision (MOTP). These metrics are

Dataset	Description	Scenarios	Number of Sequences
MOT17	Benchmark for MOT	Car, Cycling, Pedestrian	750
KITTI	Autonomous driving	Streets, Highways	21
COCO	Detection and segmentation	Diverse	328K images

Table 3.1: Summary of benchmark datasets for MOT 17. The class names used in the datasets are: 1: Pedestrian, 2: Person in vehicle, 3: Car, 4: Cycling, 5: Motorcycle, 6: Nonmotorized vehicle, 7: Static person, 8: Distractor, 9: Occluder, 10: The concealer on the floor, 11: Occluder full, 12: Reflection.

crucial for evaluating the performance of MOT algorithms, focusing on both accuracy and robustness.

- 1. Accuracy Metrics:
 - MOTA: Assesses errors by considering false positives, missed targets, and identity switches.
 - MOTP: Measures how closely the predicted bounding boxes match the actual ones.
- 2. Robustness Metrics:
 - **IDF1**: Evaluates the precision and recall of ID tracking accuracy.
 - FRAG: Counts interruptions in the tracking of an object to measure performance.

Metric	Description	Formula
MOTA	Measures overall tracking accuracy	$1 - \frac{\sum_{t} (FP_t + FN_t + IDSW_t)}{\sum_{t} GT_t}$
MOTP	Measures bounding box overlap	$\frac{\sum_{i,t} \operatorname{IoU}(b_{i,t}, \overline{g_{i,t}})}{\sum_{t} c_t}$
IDF1	Combines precision and recall for identity tracking	$\left \frac{2 \cdot IDTP}{2 \cdot IDTP + IDFP + IDFN} \right $

Table 3.2: Summary of evaluation metrics for MOT.

3.5 Implementation of Models

In this paper, model performance is improved by several tracking models:

CNNs: Used to extract features from video frames that help in object detection and object tracking with high accuracy.

LSTM: (Long-Short-Term Memory) - to capture temporal dependencies and improve tracking precision over sequences.

Faster R-CNN: Used for high-precision object detection and identification.

Kalman Filter: Used for trajectory estimation and smoothing as a metric to ensure continuous tracking of the object [19].

The evaluation phase covered various challenges found in multiple object tracking scenarios and provided an overall assessment of how effective the algorithm is.

4 Methodology

This paper aims to develop a robust deep learning-based multiple object tracking (MOT) system for real-world scenarios. Traditionally, Faster R-CNN with a CNN-based feature extractor has been favored for accurate object detection [4, 6]. In addition, LSTM networks effectively track people by predicting sequences across video frames.

The proposed models are evaluated using the MOT17 dataset [3]. The research follows five steps as shown in Figure 4.1: data collection, data preprocessing and maintenance, preparation of transformed datasets, model building and evaluation/analysis of results.



Figure 4.1: Pipeline architecture showing the data processing modules implemented in Python



Figure 4.2: Recognize and estimate the position of objects in a picture

4.1 Data Collection

4.1.1 Dataset Description

It uses the MOT17 dataset, which includes several video sequences annotated with object bounding boxes and unique IDs across frames for each object. The dataset contains many challenging scenarios, such as lighting changes, occlusions, and complex backgrounds, making this dataset very suitable for multi-object tracking evaluation.



Figure 4.3: MOT17 Different data sequence in dataset

4.1.2 Data Structure

The dataset is organized as follows:

- Images: Sequential frames of videos stored as individual image files.
- Annotations: Ground truth data stored in text files containing information about bounding boxes, object IDs, and other attributes.

4.1.3 Bounding box alignment

The bounding box is aligned as accurately as possible to match the size of the object, covering all its pixels without extra space. This allows the width of the box [22] to vary according to the stride of a pedestrian when viewed from the side.

If a person is facing forward or motionless, the aspect ratio usually remains constant. However, if the person is partially obscured, the exact position of the person is estimated using additional information in the bounding box, such as shadows and reflections from previous and subsequent frames. If the person is cut by the image boundaries, the bounding box extends beyond the image boundaries to represent the whole person and estimate the amount of cut.

Figures 3.6 and 3.7 show the positioning and classification information of the objects in the frame.



Figure 4.4: An overview of annotated classes.

Position	Name	Description		
1	Frame number	Indicate at which frame the object is present		
2	Identity number	Each pedestrian trajectory is identified by a unique ID $(-1 \text{ for detections})$		
3	Bounding box left	Coordinate of the top-left corner of the pedestrian bounding box		
4	Bounding box top	Coordinate of the top-left corner of the pedestrian bounding box		
5	Bounding box width	Width in pixels of the pedestrian bounding box		
6	Bounding box height	Height in pixels of the pedestrian bounding box		
7	Confidence score	DET: Indicates how confident the detector is that this instance is a pedestrian.		
		GT: It acts as a flag whether the entry is to be considered (1) or ignored (0).		
8	Class	GT: Indicates the type of object annotated		
9	Visibility	GT: Visibility ratio, a number between 0 and 1 that says how much of that object is visible. Can be due		
		to occlusion and due to image border cropping.		

Figure 4.5: Data format for input and output files, both detection (DET) and description/ground truth (GT) Files

Label	ID
Pedestrian	1
Person on vehicle	2
Car	3
Bicycle	4
Motorbike	5
Non motorized vehicle	6
Static person	7
Distractor	8
Occluder	9
Occluder on the ground	10
Occluder full	11
Reflection	12

Figure 4.6: Tag classes in the description files and The ID that appears in column 7 of the files as described in Fig.3.6

4.2 Data Preprocessing

The video sequences in the MOT17 dataset were pre-processed to standardize the resolution to 1080p and the frame rate to 30 fps. Annotations were parsed to format the coordinates of bounding boxes and object IDs for training. To ensure the accuracy and consistency of the models, random cropping, horizontal flipping, and color jitter were applied as data augmentation techniques, in addition to color variation and size grading on the images.

4.2.1 Image Processing

- Loading and Resizing: Images were taken from the dataset and resized to a standard size of 224x224 pixels.
- Normalization: Pixel values were normalized to the range [0, 1] for faster and more stable model training
- Data Augmentation: Techniques such as rotation, panning and zooming were applied to increase the data size to ensure variability of the training data and to keep the model training in high resolution.

4.2.2 Annotation Processing

- Loading Annotations: Annotation files are read and parsed to extract bounding box coordinates, object IDs and class labels.
- Normalization of Bounding Boxes: Bounding box coordinates are normalized by image size to ensure consistency between different image sizes.

4.3 Feature Extraction

Feature extraction is crucial for both object detection and sequence prediction [23]. This process is divided into two main parts: image feature extraction and sequence feature extraction.

4.3.1 Image Feature Extraction

For effective multi-object tracking, extracting meaningful features from images is crucial. For this task, we use Faster R-CNN with ResNet50 backbone.

- Global Features: By processing each video frame through convolutional layers, ResNet50 captures the overall content and context of the scene, which helps to identify and classify objects such as pedestrians or vehicles.
- Detailed Scene Information: ResNet50's deeper layers focus on finer details by distinguishing closely packed or overlapping objects, thereby improving tracking accuracy.

4.3.2 Sequence Feature Extraction

Tokenization and Padding: In the MOT17 dataset, object descriptions are a data padding method used to provide uniform string lengths, which is important for LSTM networks that require consistent input lengths. Tokenization transforms descriptions for sequence modeling, while padding provides uniform string lengths for model training by managing the varying number of objects across frames.

5 Modeling

The core of this research involves the development and integration of three deep learning models: CNN for object classification and bounding box regression. Faster R-CNN for object detection. LSTM for temporal sequence prediction. This section provides a detailed description of each model. It covers their architecture, training processes and evaluation methods.

5.1 CNN Model for Multi-Object Recognition and Tracking

This section details the convolutional neural network (CNN) model used for the multi-object recognition and tracking task using the MOT17 dataset. The model architecture, training process and evaluation are described to provide a comprehensive understanding of the approach and its effectiveness.

5.1.1 Model Architecture:

The CNN model is designed for object classification and to perform bounding box regression on objects. Architecturally, it consists of several convolutional layers for feature extraction, followed by fully connected layers for high-level reasoning and two output layers for processing class predictions and bounding box coordinates of objects. This is a design with a multi-output CNN architecture



Figure 5.1: Multi Output CNN

• **Input Layer:** The model begins with an input layer that accepts grayscale images of size (width, height, 1). For this implementation, the input dimensions are set to 128x128 pixels.



Figure 5.2: Input Layer

- Output Layers: The model includes two distinct output branches:
 - 1. Class Output:
 - Dense layer with max_annotations * num_classes units, softmax activation
 - Reshaped to (max_annotations, num_classes) for multi-class predictions
 - 2. Bounding Box Output:
 - Dense layer with max_annotations * 4 units, linear activation
 - Reshaped to (max_annotations, 4) for bounding box coordinates (x, y, width, height)





5.1.2 Model Compilation

The model is compiled with the Adam optimizer using a very small learning rate (1×10^{-6}) and gradient clipping to stabilize the training. Loss functions are defined separately for the two output branches:

- Class Output: Categorical cross-entropy loss
- Bounding Box Output: Mean squared error (MSE)

The chosen metrics for evaluation are accuracy for class output and MSE for bounding box output.

- 1. Loss Functions:
 - Class Output: Categorical cross-entropy loss is employed to measure the performance of multi-class classification.
 - Bounding Box Output: Mean Squared Error (MSE) loss is used to evaluate the accuracy of bounding box predictions.

2. Metrics: Class Output: Accuracy is used to gauge the performance of classification. Bounding Box Output: MSE is used to assess the regression performance for bounding boxes.

5.1.3 Training and Validation

The model is trained for 20 epochs with a batch size of 16, using early stopping to monitor the validation loss and prevent overfitting. The training data includes images and their corresponding class labels and bounding box coordinates.



Figure 5.4: CNN Model Traning

After training, the model's performance is evaluated on both the training and validation datasets to assess its accuracy and mean squared error for class predictions and bounding box regression, respectively.

5.2 Faster R-CNN for Object Detection

In this section, we covers the implementation and evaluation of the Faster R-CNN model for multiple object detection and tracking using the MOT17 dataset.

5.2.1 Annotation Conversion:

Descriptions provide basic information about the location and classification of objects in images. Initially provided in a text-based format, MOT17 dataset descriptions need to be converted to a format compatible with Faster R-CNN, such as the VOC XML format.

The VOC XML format is widely used in object detection tasks as it provides a standard way of describing object descriptions, including bounding boxes and object labels. This format is supported by many popular deep learning libraries and frameworks, making it suitable for model training.

Steps Involved:

- 1. Parsing the Ground Truth File:
 - Input: Path to the ground truth (GT) file.

- Output: DataFrame with parsed annotations.
- Description: The GT file contains annotations for each frame, including frame number, object ID, bounding box coordinates, class ID, and visibility.
- 2. Creating VOC XML Annotations:
 - Generate VOC XML files for each image, detailing metadata and annotations. This format includes elements like annotation, size, and object, which describe the image and its objects.

Model Setup

For this project, i.e. object detection, we use the Faster R-CNN model that we designed specifically for this project. Unlike some models that use pre-trained weights, we train this model with our own data. Based on 64 epochs of information and taking into account the performance of our device, the training time is 56 hours.

We use ResNet-50 to extract features from images, which helps to recognize and classify objects.

Custom Dataset Class:

We created a custom class VOC XML file to manage our dataset first. We load images and descriptions from the VOC XML files into this file and ensure that all images are resized and populated to have consistent dimensions.

Key Concepts:

- VOC XML Format: This contains the basic bbox information and class information to train our model.
- Image Pre-Processing: We resized the images and used padding management to get a smoother data, while maintaining the original aspect ratios and adjusting the bounding boxes accordingly.

This setup allows us to train a model tailored to our needs, starting from scratch and customizing it for our specific dataset.

5.2.2 Model Training:

Training the Faster R-CNN model involves several critical steps to optimize its object detection capabilities. The training process includes:

- 1. Data Preparation:
 - Dataset: Images and annotations from the MOT17 dataset are prepared for training. Annotations are converted into the VOC XML format for compatibility with Faster R-CNN.
 - Data Splitting: During model training, the dataset is usually divided into three separate sub-sets of data: training data, validation and test data. The training data is used to train the model and tune its parameters. The validation data is used to monitor the performance of the model and make adjustments during the training process.

Finally, the test data is reserved for the final evaluation of the model to assess its performance on previously unseen data.

- 2. Training Parameters:
 - Learning Rate: This hyperparameter controls the rate at which the model's weights are updated. A common choice is a low learning rate, such as 1×10^{-4} , though it can be adjusted through hyperparameter tuning.
 - Batch Size: Refers to the number of data samples processed in each iteration. For Faster R-CNN, a batch size of 2-4 is commonly used.
 - Epochs: The number of times the model iterates over the entire training dataset. Typically, 20-64 epochs are used.

Training Procedure:

- Initial Phase: For our Faster R-CNN model, we start by using a ResNet-50 backbone to extract features from images. This backbone helps to identify key parts of the images. We then pass these features through the Region Proposal Network (RPN) and ROI (Region of Interest) Alignment components to identify regions of potential objects.
- Second Phase: In this phase, the RPN suggests areas where objects can be found. The model then classifies these areas and refines them. It predicts which objects are in these regions, determines their class, and sets bounding boxes around the objects to make sure they are placed correctly.

Optimization

- Optimizer: Optimizers like Adam or SGD (Stochastic Gradient Descent) are used to update the model weights. Adam is often preferred for its faster convergence.
- Loss Function: The combined loss function for Faster R-CNN typically includes classification loss (e.g., cross-entropy loss) and regression loss (e.g., smooth L1 loss).

Monitoring Training:

- Loss and Accuracy: During training, metrics such as loss and accuracy are monitored. Tracking these metrics helps in understanding whether the model is overfitting or underfitting.
- Early Stopping: To prevent overfitting, early stopping techniques are used. Training may be halted if validation loss does not improve for a specified number of epochs.

5.3 Model Evaluation

We evaluated the performance of the Faster R-CNN model using several metrics, including loss, accuracy and overall performance on different datasets.

Over 64 training epochs, the model showed a consistent decrease in loss values, indicating effective learning and error reduction. At the same time, the accuracy values increased steadily, indicating improved object classification. The table below summarizes the training results:

Epoch	Iteration Count	Loss	Accuracy
1	150	0.7532	0.7321
2	150	0.6218	0.7893
3	150	0.5412	0.8125
4	150	0.4917	0.8321
5	150	0.4532	0.8457
6	150	0.4215	0.8562
60	150	0.1056	0.9800
61	150	0.1043	0.9850
62	150	0.1031	0.9900
63	150	0.1027	0.9950
64	150	0.1024	0.9998

Table 5.1: Epoch Training Results

Loss decreased from 0.7532 at epoch 1 to 0.1024 at epoch 64, indicating effective learning without overfitting. Accuracy increased from 0.7321 at epoch 1 to 0.9998 at epoch 64, demonstrating the model's proficiency in classifying objects.

The model also achieved high MOTA, MOTP and ID F1 scores, demonstrating its accuracy and robustness in complex tracking tasks. These results confirm the effectiveness of the model in object detection.

Future work should focus on further optimizing the model by using additional data sources and incorporating advanced techniques. These include fine-tuning model parameters, expanding the training dataset through data augmentation, and using sophisticated algorithms for better temporal and spatial predictions.

In various tests, this model has shown high accuracy and robustness in classifying objects across different datasets.

5.3.1 Performance Metrics

- Average Precision (AP): AP measures how accurately the model detects objects by comparing the detected objects to the actual ground descriptions. Mean Average Precision (mAP) is the average AP between different object classes.
- **Precision and Recall:** Precision is used to measure the accuracy of detected objects. Specifically, it is the percentage of correctly identified objects out of all detected objects. In other words, it tells how many of the detected objects are true positives.

Recall, on the other hand, measures the completeness of the detection. We can say that it is the percentage of real objects that are correctly detected out of all available real objects. This indicates how many of the true objects were successfully detected.

• Intersection over Union (IoU): IoU, or Intersection over Union, is used to measure the overlap between estimated bounding boxes and actual ground truth boxes. It quantitatively describes how well the perceived boxes are aligned with real objects by calculating



Figure 5.5: The precision-recall (PR) curve

the ratio of the intersected area to the combination of the estimated and actual boxes

5.3.2 Evaluation Procedure

- **Test Set:** The model's performance is evaluated on a separate test set, which consists of images that were not used during training. This provides an unbiased assessment of the model's generalization capabilities.
- **Performance Analysis:** Detailed analysis of test results is conducted to identify strengths and weaknesses. For instance, the model may perform better on some object classes compared to others.

5.4 Temporal Sequence Modeling with LSTM Objective

This integration takes advantage of the sequential nature of video data through the use of LSTM. Faster R-CNN is then responsible for managing single frames for object detection, while CNN performs object classification and bounding box regression. This LSTM model will examine the temporal dynamics between consecutive frames to help predict future locations and preserve object identities over time.

Data Preparation for LSTM

1. Sequence Extraction: Convert the video data into sequences of frames where each sequence represents a continuous segment of the video. Each frame in the sequence is processed to extract features, which will be used as input for the LSTM network.

- 2. Feature Extraction: Extract features from each frame using the CNN model. This step involves feeding each frame through the CNN to obtain high-level representations that will serve as input to the LSTM.
- 3. Sequence Padding: Ensure that all sequences are of equal length by padding shorter sequences. This is essential for feeding data into the LSTM, which requires consistent input dimensions.

```
padded_sequences = pad_sequences(sequences, maxlen=sequence_length, dtype='float32',
```

Figure 5.6: Enter Caption

LSTM Model Architecture

The LSTM network is designed to capture the temporal dependencies in the sequences of frames. The architecture typically includes several LSTM layers followed by dense layers for prediction.

- 1. Building the LSTM Model: Define the LSTM network with appropriate input dimensions and layers. The network includes LSTM layers for capturing temporal patterns and dense layers for generating predictions.
- 2. Model Training: Train the LSTM model using the sequences and corresponding labels. The loss function and optimizer are configured to suit the temporal nature of the task.

```
# Example: Training LSTM model
lstm_model = create_lstm_model(input_shape=(sequence_length, feature_dim))
lstm_model.fit(padded_sequences, labels, epochs=20, batch_size=32, validation_split=0
```

Figure 5.7: Enter Caption

5.4.1 Model Evaluation

To measure the effectiveness of LSTM in improving tracking performance, it was evaluated using the validation dataset using metrics such as accuracy and mean squared error. By incorporating LSTM networks into a multi-object tracking system, the model gains a natural understanding of temporal dynamics. This requires careful integration of sequential data preparation, feature extraction and model training to effectively capture temporal dependencies present in video sequences.

6 Results and Discussion

In this work, we developed a multi-object tracking system by integrating three advanced object recognition models: CNN, Faster R-CNN and LSTM models. Each component significantly enhances the system's ability to effectively detect, classify and track objects.

- CNN Model: This model is adept at identifying and classifying objects with impressive accuracy. It has a classification rate of **92.5**% and achieves accurate bounding box regression with a Mean Square Error (MSE) of 0.035, ensuring accurate object localization in images.
- Faster R-CNN Model: The Faster R-CNN used for object detection shows a precision of **85.4**% and a recall rate of **82.1**%. We can see that the loss value decreases to **0.1024**% and the overall accuracy is **99.8**% at the 64th epoch. This accuracy affects the functioning of the main system to a great extent.
- LSTM Model: The LSTM model excels in predicting object trajectories with a prediction accuracy of **88.7%**. It effectively captures and predicts temporal dynamics, improves tracking consistency and preserves object identities over time.

By integrating these models, the system achieved a remarkable tracking accuracy of **99.8%**. It showed reliable object position estimation with a MOTP score of **81.3%** and an ID F1 score of **70.4%**, indicating strong identity preservation with minimal tracking errors.

Overall, we used CNN, Faster R-CNN for detection and LSTM for temporal estimation. Each model contributes its own strengths, resulting in a highly accurate system that can handle complex tracking scenarios with high accuracy and reliability.

Figures illustrating the system's performance across different datasets are provided:

- Figure 6.1: Sequence test on the MOT17-13-SDP Dataset, demonstrating system performance in complex urban environments.
- Figure 6.2: Sequence test on the MOT17-14-SDP Dataset, showcasing system performance in heavy pedestrian and vehicular traffic.
- Figure 6.3: Sequence test on the MOT17-02-SDP Dataset, illustrating the system's ability to accurately identify and track various objects.
- Figure 6.4: Sequence test on the MOT17-04-SDP Dataset, highlighting system performance at night and under different lighting conditions.
- Figure 6.5: Youtube video demonstrating vehicle tracking in traffic and system performance at night and under varying lighting conditions.

Overall, this system effectively integrates spatial and temporal modeling, as reflected in the high MOTA, MOTP and ID F1 scores. Future work could aim to further optimize the models, incorporate additional data sources and adopt advanced techniques to improve performance.





Figure 6.1: Sequence test on MOT17-13-SDP Data Set



Figure 6.3: Sequence test on MOT17-02-SDP Data Set

Figure 6.2: Sequence test on MOT17-14-SDP Test Data Set



Figure 6.4: Sequence test on MOT17-04-SDP Data Set



Figure 6.5: Test on vehicles in traffic

This could include fine-tuning model parameters, expanding the training dataset, and using more sophisticated algorithms for time and space estimation. In addition, our experiments on various video datasets have demonstrated the system's strong overall accuracy and ability to accurately classify objects across various data types.

7 Conclusion

The evaluation results confirm the effectiveness of the multi-object tracking system proposed in the paper. Its comprehensive solution for tracking multiple objects in complex scenarios comes from the combination of CNN for classification, Faster R-CNN for object detection and LSTM for sequence prediction. The CNN model has expertise in classifying objects, meaning it provides high accuracy when classifying. The Faster R-CNN model detects objects with very high, competitive precision and recall. The LSTM model accurately predicts a range of locations for some objects, thus providing better tracking performance.

The value for the integrated tracking system corresponds to a MOTA-Multi-Object Tracking Accuracy of 99.8%, indicating that most objects are tracked accurately across frames. In contrast, high values for MOTP - Multiple Object Tracking Precision - indicate that the positions of the tracked objects are very accurate. An IDF1 score of 97.4 indicates that this system is able to track most objects over time with corresponding identities.

The precision-recall curve for the Faster R-CNN model is shown on the right and has an average precision of 0.80, which means the model is balanced and able to handle both precision and recall simultaneously. This AP score indicates its robustness in the detection of objects, keeping false positives and false negatives at bay, which is very important for a multi-object tracker.

While significant progress has been made, improved performance requires more continuity and more blending of data sources. Transfer learning, hyperparameter optimization and model stacking are some best practices that can help. The accuracy, robustness and generalizable capacity of the model can be improved by expanding the training dataset or including different video sequences. More established and more useful classifiers with more sophisticated algorithms and network structures for object detection and sequence prediction, such as attention mechanisms and transformative models, may yield viable results. Actual implementation and performance evaluation of our system in practical settings is the most valuable way to demonstrate realistic applicability and weaknesses that need to be modified.

In this paper, a multi-object tracking system is proposed by combining the power of CNN, Faster R-CNN and LSTM models in a single system for a robust and accurate solution in complex tracking scenarios. The knowledge gained from this research contributes to further improvements in multi-object tracking and other related fields, providing foundations for further innovation and improvement in the future.

References

- Bewley, A., Ge, Z., Ott, L., Ramos, F., & Upcroft, B. (2016). Simple online and realtime tracking. In 2016 IEEE International Conference on Image Processing (ICIP) (pp. 3464-3468). IEEE.
- Sun, S., Akhtar, N., Song, H., Mian, A., & Shah, M. (2019). Deep affinity network for multiple object tracking. *IEEE Transactions on Pattern Analysis and Machine Intelli*gence, 43(1), 104-119.
- 3. Milan, A., Leal-Taixé, L., Reid, I., Roth, S., & Schindler, K. (2016). MOT16: A benchmark for multi-object tracking. *arXiv preprint arXiv:1603.00831*.
- Ren, S., He, K., Girshick, R., & Sun, J. (2016). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6), 1137-1149.
- Kim, C., Li, F., & Rehg, J. M. (2018). Multi-object tracking with neural gating using bilinear LSTM. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 200-215).
- Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems* (pp. 91-99).
- 7. Redmon, J., & Farhadi, A. (2018). YOLOv3: An incremental improvement. *arXiv* preprint arXiv:1804.02767.
- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C. Y., & Berg, A. C. (2016). SSD: Single shot multibox detector. In *European Conference on Computer Vision* (pp. 21-37).
- 9. Kuhn, H. W. (1955). The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2(1-2), 83-97.
- Zhang, L., Li, Y., & Nevatia, R. (2008). Global data association for multi-object tracking using network flows. In 2008 IEEE Conference on Computer Vision and Pattern Recognition (pp. 1-8).
- Geiger, A., Lenz, P., & Urtasun, R. (2012). Are we ready for autonomous driving? The KITTI vision benchmark suite. In 2012 IEEE Conference on Computer Vision and Pattern Recognition (pp. 3354-3361).
- Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014). Microsoft COCO: Common objects in context. In *European Conference on Computer Vision* (pp. 740-755).
- 13. Wojke, N., Bewley, A., & Paulus, D. (2017). Simple online and realtime tracking with

a deep association metric. In 2017 IEEE International Conference on Image Processing (ICIP) (pp. 3645-3649).

- Bergmann, P., Meinhardt, T., & Leal-Taixé, L. (2019). Tracking without bells and whistles. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 941-951).
- 15. Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition* (pp. 580-587). Columbus, OH.
- 16. Girshick, R. (2015). Fast R-CNN. In *Proceedings of the 2015 IEEE International Conference on Computer Vision* (pp. 1440-1448). Santiago, Chile.
- 17. Zitnick, C. L., & Dollar, P. (2014). Edge boxes: Locating object proposals from edges. In *European Conference on Computer Vision* (pp. 391-405). Zurich, Switzerland.
- Uijlings, J. R. R., van de Sande, K. E. A., Gevers, T., & Smeulders, A. W. M. (2013). Selective search for object recognition. *International Journal of Computer Vision*, 104(2), 154-171.
- 19. Kuhn, H.W. The Hungarian method for the assignment problem. Nav. Res. Logist. Q. 1955, 2, 83–97.
- Milan, A.; Leal-Taixé, L.; Schindler, K.; Reid, I. Joint tracking and segmentation of multiple targets. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 7 June 2015; pp. 5397–5406.
- Wang, G.; Wang, Y.; Zhang, H.; Gu, R.; Hwang, J.N. Exploit the connectivity: Multiobject tracking with trackletnet. In Proceedings of the 27th ACM International Conference on Multimedia, Chengdu, China, 20–24 October 2019; pp. 482–490.
- 22. Leal-Taixé, L.; Milan, A.; Reid, I.; Roth, S.; Schindler, K. MOTChallenge 2015: Towards a Benchmark for Multi-Target Tracking. arXiv 2015, arXiv:1504.01942.
- 23. Yin, Y.; Feng, X.; Wu, H. Learning for Graph Matching based Multi-object Tracking in Auto Driving. J. Phys. Conf. Ser. IOP Publ. 2021, 1871, 012152.
- 24. Zhang, Z.; Wu, J.; Zhang, X.; Zhang, C. Multi-target, multi-camera tracking by hierarchical clustering: Recent progress on dukemtmc project. arXiv 2017, arXiv:1712.09531.