# CS301 Investigative Studio II

*Vehicle Tracking in Challenging Scenarios using YOLO.*

## Irina Getman

Bachelor of Software Engineering

Yoobee Colleges

**Supervisory team:**

**Main** – Dr. Mohammad Norouzifard

**Co-supervisor** – Aisha Ajmal

Auckland

March 2023

4

# Table of Content

# Abstract

Autonomous driving is a rapidly evolving technology that has the potential to revolutionize the transportation industry. Real-time object is a crucial component for autonomous vehicles, which need to perceive their surroundings and react accordingly.

Various deep learning methods have been proposed to tackle this problem, one of the most popular and efficient algorithms for real-time object detection is YOLO (You Only Look Once), which can process up to 155 frames per second. YOLO uses a single neural network to predict bounding boxes and class probabilities for each object in an image. This makes it faster and more accurate than other methods that use multiple stages or sliding windows. In this paper, we will introduce the basic principles of YOLO, discuss the latest version of the family, YOLOv8, contrast it with the earlier state-of-the-art YOLOv7 and explore its potential applications in autonomous driving scenarios. Finally, we discuss the current research trends and future directions in this field. This paper aims to provide an insight into the current state-of-the-art real-time object detection methods.

**Key words**: Autonomous Driving, Real-time Object Detection, YOLO, vehicle tracking, low visibility.

# Introduction

Autonomous driving technology has been a topic of interest and research for several decades. The goal of this technology is to develop vehicles that can operate without human intervention, using advanced sensors and artificial intelligence algorithms to navigate roads and traffic. Vehicle tracking is a specific application of object detection that involves locating and following vehicles in video sequences. It is an essential part in autonomous driving.

The history of autonomous driving has been a long and fascinating journey that spans centuries and continents. It is believed (Reiser, 2021) that Leonardo da Vinci came up with the first concept of a self-driving vehicle in the 1500s. He designed a cart that could follow a set route using springs and steering devices. However, it was not until the 20th century that more advanced experiments on autonomous cars began and only in our days, autonomous car technology has become more widespread and accessible. In 2017, major car manufacturers introduced the Advanced Driving Assistance System (ADAS) with an aim to aid drivers in avoiding collisions and performing various driving tasks with the assistance of cutting-edge technologies. ADAS has the capability to alert drivers or even take full control of a vehicle in certain situations. According to market research by Canalys (Neowin ·, 2021), 11.2 million cars with level 2 features (hands off) were sold in 2020, marking a growth of 78% compared to the previous year. With the increasing popularity of ADAS, the time when drivers will take their eyes and minds off the road is rapidly approaching. Despite their flaws and the unfortunate fatalities they have caused(Levin & Julia Carrie Wong, 2018), they have the potential to prevent thousands of deaths per year, or about 62% of total traffic deaths. Lane keeping assist accounts for 14,844 of this savings, while pedestrian automatic braking accounts for another 4,106 lives saved (*Advanced Driver Assistance Systems-Data Details - Injury Facts*, 2022).

The ability of autonomous driving systems to perceive and respond to their surroundings in real-time relies heavily on real-time object detection, making it a critical component. Object detection allows autonomous vehicles to identify and track objects, such as other vehicles on the road, enabling them to make informed decisions and take appropriate actions. There are many challenges in real-time object detection for autonomous vehicles, such as varying lighting conditions, occlusions, small objects, fast motion, etc. Therefore, researchers have proposed various algorithms to balance the speed and accuracy of detection.

One popular approach is based on YOLO (You Only Look Once), which is a one-stage detector that processes an image as a whole and outputs bounding boxes and class labels for each object. YOLO has been improved over several versions, with the latest version YOLOv8. Previous version YOLOv7 achieved state-of-the-art performance on various benchmarks. In this paper, we review the

advantages YOLOv8 over its predecessor. We discuss the strengths and limitations of these methods and examine their performance in various driving scenarios.

# Literature Review

## Related works on general Object Detection task

Object detection is the task of locating and classifying objects in an image or video. In recent years, there has been a significant advancement in object detection using deep learning techniques, particularly Convolutional Neural Networks (CNNs). According to Zou et al. (2019),the progress of object detection has gone through two historical periods: traditional object detection period (before 2014) and deep learning-based detection period (after 2014). **In the traditional period**, Viola and Jones pioneered real-time face detection using the sliding window approach, while Histograms of Oriented Gradients (HOG) and Deformable Part-based Model (DPM) were used to detect objects of different sizes.

**Deep learning methods** started with the rebirth of CNNs in 2012 by (Krizhevsky et al., 2012), which showed substantially higher image classification accuracy on the ImageNet Large Scale Visual Recognition Challenge (ILSVRC). CNNs can learn features directly from the data and can be retrained and built on pre-existing networks. In the deep learning era, object detection can be grouped into two categories: **two-stage** detection or sparse prediction type and **one-stage** detection or dense prediction.

**Two-stage** detectors use separate networks for region of interest (RoI) extraction and classification and refinement. Regions with CNN features (RCNN) {citation} achieved significant results by combining object proposals and CNNs, but it has notable drawbacks such as multi-stage training, slow detection speed, and requires hundreds of gigabits for storage. SPPNet (Spatial Pyramid Pooling Networks) {citation} introduced a new pooling strategy that eliminated the requirement of fixed-size images for CNN, making it 20 times faster than RCNN. Fast RCNN {citation}fixed the drawbacks of RCNN and SPPnet and achieved higher detection quality with a single-stage training process. Faster RCNN {citation} was the first end-to-end and near-real time deep learning detector, and it introduced Region Proposal Network (RPN) to generate region proposals directly in the network using anchor boxes for object detection. Feature Pyramid Networks (FPN) {citation} used a bottom-up and top-down pathway to construct a pyramid, achieving state-of-the-art single-model results without increasing testing time over a single-scale baseline.

## One-stage detectors. Basic YOLO model

The YOLO method was introduced by Redmon, Divvala, Girshick, and Farhadi (Redmon et al., 2015), who used a single neural network to predict bounding boxes and class probabilities in images, resulting in significant progress in terms of real-time speed and performance. However, YOLO has been found to make more localization errors compared to other state-of-the-art models due to the use of a sum-squared error in the output, which weighs localization errors to be equal with classification errors. This instability could cause training to diverge early on if not addressed properly. YOLO outperformed Fast R-CNN in terms of making fewer background mistakes, but combining these two models produced a significant boost in performance. In contrast, the Single Shot Multibox Detector (SSD), introduced by Liu et al. in 2016, detects objects in images using a single deep neural network that generates scores for the presence of each object category in each default box and produces adjustments to the box to better match the object shape. Although the base YOLO model was capable of processing images at 45 frames per second, the SSD achieves a

faster frame rate of 59 FPS with a slightly lower mAP. However, with a 512x512 image input, the SSD achieves a higher mAP than the combined YOLO and Fast R-CNN models. Another breakthrough model, Retina-Net, was created by Facebook AI Researchers in 2017 to address the foreground-background class imbalance during the training process of one-stage detectors. Retina-Net uses a focal loss approach to reshape the standard cross-entropy loss, which down-weights the loss assigned to well-classified examples, resulting in improved accuracy compared to all existing state-of-the-art two-stage detectors while matching the speed of previous one-stage detectors.

## Evolution of YOLO models

YOLO has undergone an evolution since its introduction. One of its predecessors is the Histogram of Oriented Gradients (HOG) method, first introduced in 1986. This method uses a feature descriptor to detect objects of interest but can be time-consuming. YOLO improves upon this by using a single neural network to directly predict bounding boxes.
There have been several versions of YOLO developed by researchers that became state-of-the-art detectors in various applications. Each version introduces new features and improvements to boost performance and flexibility.
One of the popular most versions is YOLOv5, which comes in four main versions: small (s), medium (m), large (l), and extra-large (x). Each variant offers progressively higher accuracy rates and takes a different amount of time to train. The first official version of YOLOv5 was released by Ultralytics on June 25th, 2020. Since then, Ultralytics (2023) continued to make further improvements to the YOLO architecture introducing cutting-edge object-detection model YOLOv8 in January 2023.
YOLOv7 is a second to last and very successful version of the YOLO object detection system that improves speed and accuracy through several architectural changes. Unlike previous versions that used pre-trained backbones from ImageNet, YOLOv7 is trained entirely on the COCO dataset (*YOLOv7 Paper Explanation: Object Detection and YOLOv7 Pose*, 2022). Some of the major changes introduced in YOLOv7 include the E-ELAN (Extended Efficient Layer Aggregation Network), model scaling for concatenation-based models, trainable BoF (Bag of Freebies) (Wang et al., 2022), planned re-parameterized convolution, and coarse for auxiliary and fine for lead loss. The architecture of YOLOv7 is derived from previous versions such as YOLOv4, Scaled YOLOv4, and YOLO-R. The E-ELAN computational block in the YOLOv7 backbone is designed to improve network efficiency by considering factors such as memory access cost, I/O channel ratio, element-wise operations, activations, and gradient path.

## Vehicle tracking

YOLO (You Only Look Once) is a popular object detection algorithm that has been used in several studies for vehicle tracking and real-time object detection.
Real-time object vehicle tracking tasks require high accuracy in object detection and recognition, while also needing to be efficient and fast to operate in real-time, where even small delays can have significant consequences. Several studies conducted the approach to  reduce the computational complexity of the vehicle-tracking models while maintaining high detection accuracy. Bie et al., (2023) propose a novel architecture that combines the YOLOv5n network with a lightweight backbone network. The authors also introduce a multi-scale feature fusion module to improve the detection of small vehicles. Yuan & Xu, (2021) proposed a lightweight vehicle detection algorithm based on an improved YOLOv4 model. The authors introduce a new feature aggregation module to reduce the number of parameters and improve the detection accuracy of small vehicles. They also optimize the convolutional layer and add a new upsampling method to the network. Another approach of reducing the number of convolutional layers and using a smaller input size for images

was introduced by Koay et al., (2021). This study presents a system for real-time vehicle detection in aerial images captured by low-cost unmanned aerial vehicles (UAVs). The system is based on the YOLO (You Only Look Once) object detection algorithm and is designed to run on low-cost edge devices such as Raspberry Pi. All the above proposals contribute to advancing the capabilities of these real-time object detecting systems and making them more practical and accessible in real-world scenarios.

To handle the challenges associated with vehicle detection in real-time scenarios, such as the small size of vehicles, occlusions, the varying lighting and weather conditions, several studies present various solutions. The paper "IDOD-YOLOV7: Image-Dehazing YOLOV7 for Object Detection in Low-Light Foggy Traffic Environments" proposes a new approach to improve object detection in low-light foggy traffic environments. The proposed approach involves using an image-dehazing method to remove the haze and improve the visibility of the images, and then using the YOLOv7 object detection algorithm to detect objects in the dehazed images.(Qiu et al., 2023).

In another study, Dang et al.(2023) presented a new approach to improve real-time traffic sign recognition in bad weather conditions using the YOLOv5 object detection algorithm. The proposed approach involves using a combination of image augmentation techniques and transfer learning to improve the detection accuracy of the YOLOv5 algorithm. The authors first collect a dataset of images of traffic signs in bad weather conditions, including foggy, rainy, and snowy weather. They then apply several image augmentation techniques, including rotation, scaling, and brightness adjustment, to increase the size of the dataset and improve the robustness of the YOLOv5 algorithm. The authors also use transfer learning to fine-tune the YOLOv5 algorithm on the augmented dataset, using a pre-trained model on a large-scale dataset as a starting point. They test the performance of the improved YOLOv5 algorithm on both the augmented dataset and a real-world dataset of traffic signs in bad weather conditions. Even though, this study is not exactly related vehicle tracking tasks, it provides a valuable approach dealing with image recognition in challenging scenarios.

Few topic-related studies mentioned using **Kalman** filter(Zhao et al., 2020),(Lin et al., 2021) A Kalman filter is a mathematical algorithm that uses a series of measurements over time to estimate the state of a system and predict future states. It is commonly used in control and signal processing applications for estimating and filtering noisy sensor data. In the context of vehicle tracking, the Kalman filter is often used to estimate the position, velocity, and acceleration of a vehicle based on measurements from sensors such as cameras or radar. By combining these measurements with a model of the vehicle's motion, the Kalman filter can track the vehicle's position and velocity over time, even in the presence of noisy or missing measurements.
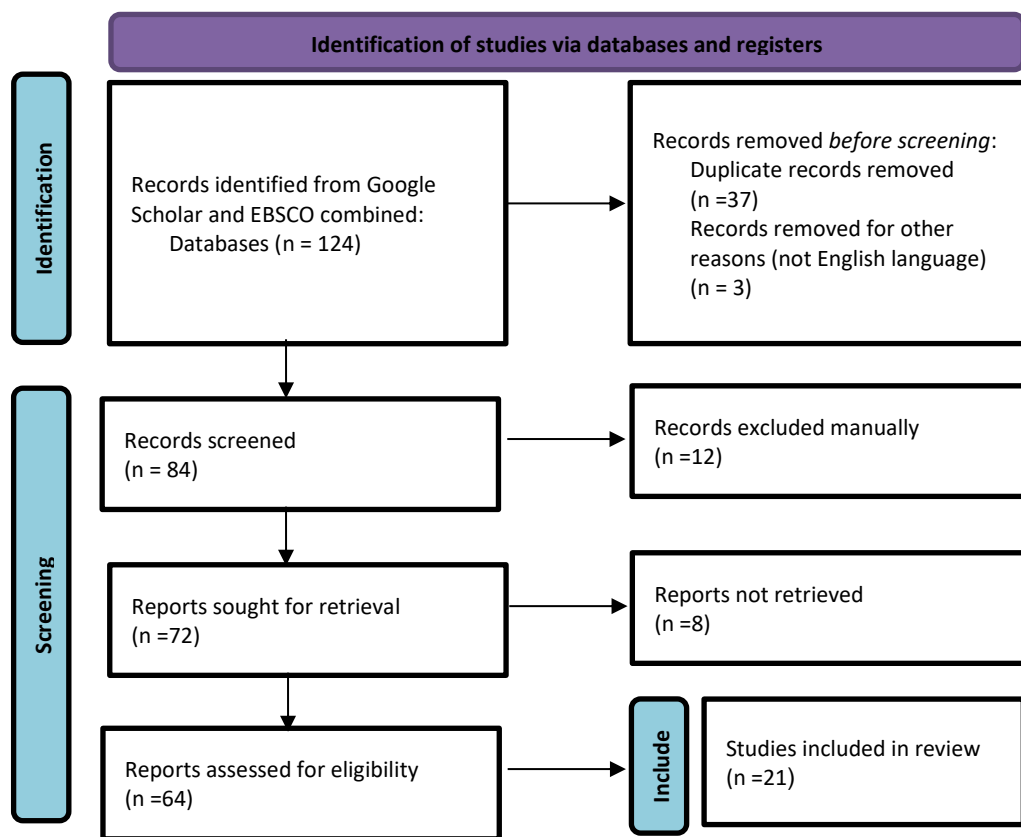
Zhao et al. (2020) presented a novel algorithm for tracking vehicles and switching their identification (ID) in driving recording sensors. The algorithm uses a combination of a Kalman filter and a Hungarian algorithm to track vehicles across multiple frames and associate them with their corresponding IDs. The authors also introduce a new switching mechanism based on the number of frames a vehicle is tracked without detection. The proposed algorithm is evaluated on a dataset of real-world driving scenarios and achieves high tracking accuracy and low switching error rates compared to other tracking algorithms. The authors suggest that the proposed algorithm can be applied to various driving recording sensor systems, such as black box recorders and dash cameras.

Lin et al. (2021) introduced a system that can count vehicles, estimate their speed, and classify them based on their type in real-time. The system is based on a combination of virtual detection zones and the YOLO (You Only Look Once) object detection algorithm. The authors first define virtual detection zones along a road, which are used to track vehicles as they pass through. They then apply the YOLO object detection algorithm to the images captured by cameras placed along the road, to detect and classify vehicles passing through the detection zones. The system uses a Kalman filter to estimate the speed of each vehicle based on its position in the virtual detection zones over time. The system can also classify vehicles into different types, such as cars, trucks, and motorcycles, based on their size and shape.
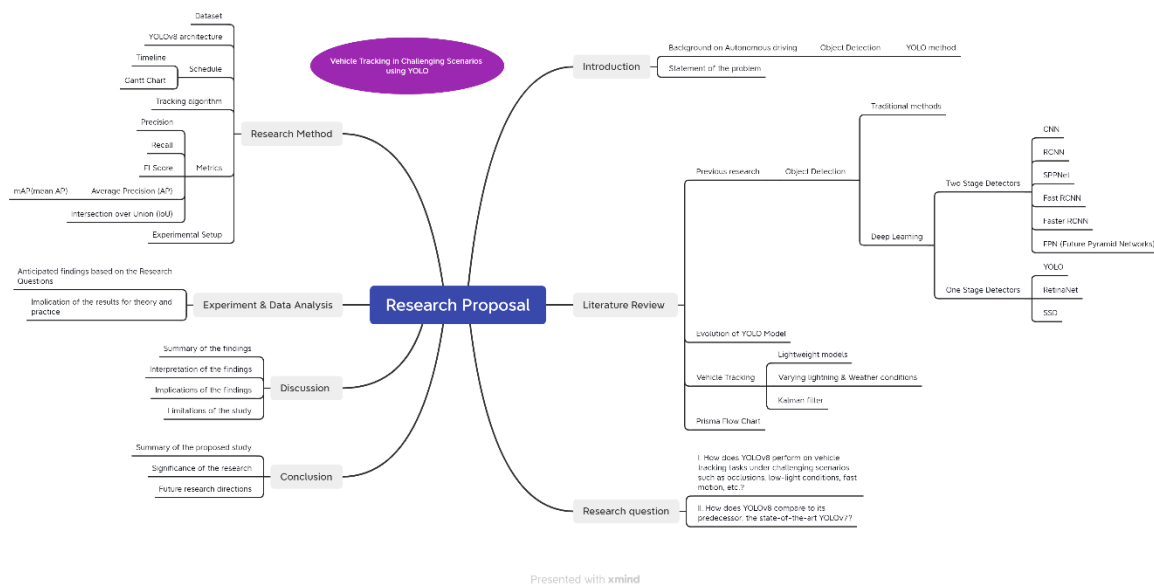
Overall, the literature suggests that YOLO-based approaches are effective for real-time vehicle tracking and object detection. These approaches have shown high accuracy and real-time performance in complex traffic scenarios and aerial videos. However, further research is needed to explore the potential of the newest YOLO-based approaches versions 8 and 7 for the autonomous driving applications performing in real-time and under various challenging scenarios.

## PRISMA flow chart

To conduct a systematic literature review, I have selected two databases: Google Scholar and EBSCO. The keywords: vehicle tracking, YOLO, low visibility, autonomous driving were included to construct search string. The publication dates were set to include the published articles between 2019 and 2023.

## Mind Map



Presented with xmind

# Research question

I.    How does YOLOv8 perform on vehicle tracking tasks under challenging scenarios such as occlusions, low-light conditions, fast motion, etc.?

II.   How does YOLOv8 compare to its predecessor, the state-of-the-art YOLOv7?

# Research Methods

Here is an outline of the methodology that will be conducted in the next phases.

I.    The dataset: We will describe the COCO dataset used for training and testing the YOLOv8 model, including its size, composition, and any pre-processing steps taken.

II.   The YOLOv8 model: we will provide details on the YOLOv8 model used in the research, including its architecture, hyperparameters, and training procedure.

III.  The tracking algorithm: we will describe the algorithm used for tracking vehicles in challenging scenarios, including any modifications made to adapt it to the specific problem.

IV.   Evaluation metrics:
   a.  Precision: Precision measures the proportion of true positive detections among all positive detections. It is calculated as the ratio of true positives to the sum of true positives and false positives.
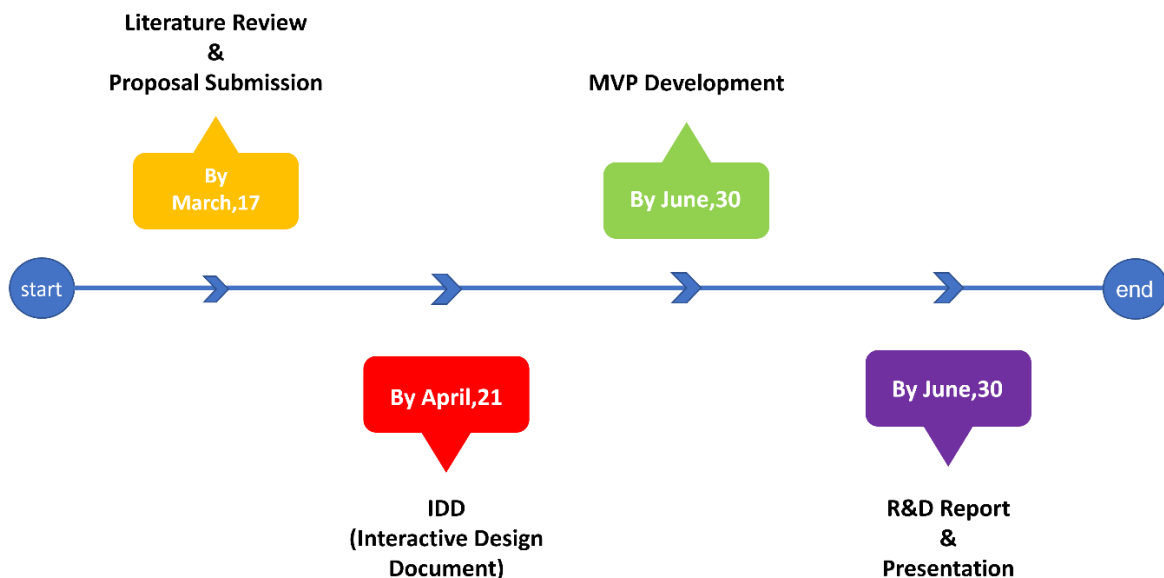
b. Recall: Recall measures the proportion of true positive detections among all actual objects. It is calculated as the ratio of true positives to the sum of true positives and false negatives.

c. F1-score: The F1-score is the harmonic mean of precision and recall, and provides a balanced measure of the performance of an object detection algorithm.

d. Average Precision (AP): Average Precision is commonly used to evaluate the performance of object detection algorithms on benchmark datasets such as COCO. It is calculated by plotting a precision-recall curve for the algorithm and computing the area under the curve.

   i. mAP, or mean Average Precision, is a commonly used evaluation metric for object detection algorithms. It is calculated by taking the mean of the Average Precision (AP) values for each class in the dataset.

   ii. mAP provides an overall measure of the algorithm's performance across all classes in the dataset. A higher mAP value indicates better performance of the object detection algorithm.

e. Intersection over Union (IoU): IoU measures the overlap between the predicted bounding box and the ground truth bounding box. It is calculated as the ratio of the intersection area to the union area of the two bounding boxes.

   These metrics provide different perspectives on the performance of an object detection algorithm, and it's common to report multiple metrics when evaluating an algorithm.
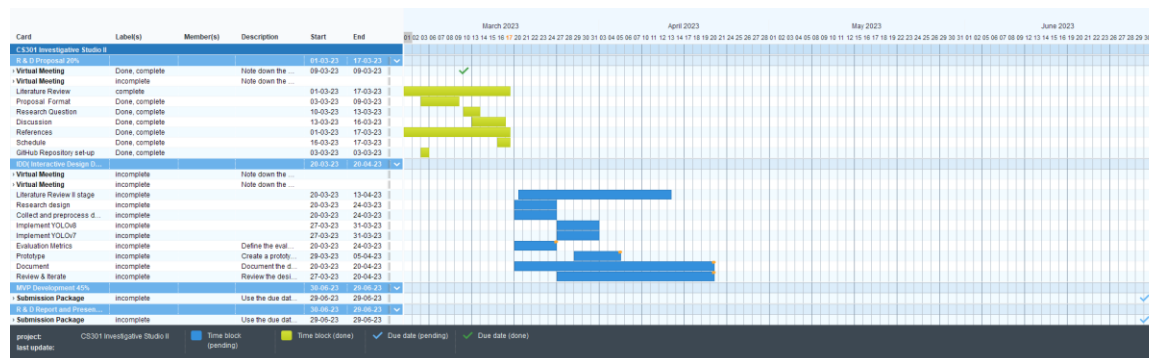
V. Experimental setup: details on the experimental setup, including hardware and software used, and any other relevant information.

# Schedule for experiments/implementation

## Timeline

## Gantt Chart



# Conclusion

To summarize, autonomous driving is a rapidly evolving technology that has the potential to transform the transportation industry. Real-time object detection is a crucial aspect of this technology, allowing autonomous vehicles to perceive their surroundings and react accordingly. Numerous deep learning methods have been proposed to address this challenge, among which YOLO (You Only Look Once) stands out as one of the most popular and efficient algorithms for real-time object detection, capable of processing up to 155 frames per second. In this paper, we conducted a literature review of object detection with a specific focus on vehicle detection studies. While YOLO is a well-documented and widely used algorithm in various applications, YOLOv8, the latest version of the family, has not been extensively researched yet. Hence, our findings may fill a potential research gap and contribute to the academic community. Additionally, we propose a comparison between YOLOv8 and the earlier state-of-the-art YOLOv7.

Our proposed research methodology entails utilizing the COCO dataset to train and test the YOLOv8 model and implementing an algorithm that addresses vehicle tracking in challenging scenarios. We plan to evaluate the performance of the YOLOv8 model using several metrics, including precision, recall, F1-score, average precision (AP), mAP, and intersection over union (IoU). Additionally, we will provide details on the experimental setup, including information on the hardware and software used.

In summary, this paper has conducted a thorough review of the relevant literature and presented a timeline for future research, including advancements in real-time object detection methods such as YOLOv8.

# References

*Advanced Driver Assistance Systems-Data Details—Injury Facts*. (2022). Injury Facts.

    https://injuryfacts.nsc.org/motor-vehicle/occupant-protection/advanced-driver-assistance-

    systems/data-details/#:~:text=ADAS%20technologies%20have%20the%20potential

Bie, M., Liu, Y., Li, G., Hong, J., & Li, J. (2023). Real-time vehicle detection algorithm based on a

    lightweight You-Only-Look-Once (YOLOv5n-L) approach. *Expert Systems with Applications*,

    *213*, N.PAG-N.PAG. Academic Search Complete.

Dang, T. P., Tran, N. T., To, V. H., & Tran Thi, M. K. (2023). Improved YOLOv5 for real-time traffic signs

    recognition in bad weather conditions. *The Journal of Supercomputing*, 1–19.

Koay, H. V., Chuah, J. H., Chow, C.-O., Chang, Y.-L., & Yong, K. K. (2021). YOLO-RTUAV: Towards Real-

    Time Vehicle Detection through Aerial Images with Low-Cost Edge Devices. *Remote Sensing*,

    *13*(21), 4196. Academic Search Complete.

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional

    Neural Networks. *Advances in Neural Information Processing Systems*, *25*.

    https://proceedings.neurips.cc/paper/2012/hash/c399862d3b9d6b76c8436e924a68c45b-

    Abstract.html

Levin, S. & Julia Carrie Wong. (2018). *Self-driving Uber kills Arizona woman in first fatal crash*

    *involving pedestrian*. The Guardian; The Guardian.

    https://www.theguardian.com/technology/2018/mar/19/uber-self-driving-car-kills-woman-

    arizona-tempe

Lin, C.-J., Jeng, S.-Y., & Lioa, H.-W. (2021). A Real-Time Vehicle Counting, Speed Estimation, and

    Classification System Based on Virtual Detection Zone and YOLO. *Mathematical Problems in*

    *Engineering*, 1–10. Academic Search Complete.

Qiu, Y., Lu, Y., Wang, Y., & Jiang, H. (2023). IDOD-YOLOV7: Image-Dehazing YOLOV7 for Object

Detection in Low-Light Foggy Traffic Environments. *Sensors*, *23*(3), 1347.

Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2015). You Only Look Once: Unified, Real-Time

Object Detection. *ArXiv.Org*. https://doi.org/10.48550/arXiv.1506.02640

Reiser, A. (2021, August 9). History of Autonomous Cars. *TOMORROW'S WORLD TODAY®*.

https://www.tomorrowsworldtoday.com/2021/08/09/history-of-autonomous-cars/

*Ultralytics | Revolutionizing the World of Vision AI*. (2023). Ultralytics. https://ultralytics.com/yolov8

Wang, C.-Y., Bochkovskiy, A., & Liao, H.-Y. M. (2022). *YOLOv7: Trainable bag-of-freebies sets new*

*state-of-the-art for real-time object detectors* (arXiv:2207.02696). arXiv.

https://doi.org/10.48550/arXiv.2207.02696

*YOLOv7 Paper Explanation: Object Detection and YOLOv7 Pose*. (2022). LearnOpenCV – Learn

OpenCV, PyTorch, Keras, Tensorflow with Examples and Tutorials.

https://learnopencv.com/yolov7-object-detection-paper-explanation-and-

inference/#YOLOv7-Object-Detection-Inference

Yuan, D. L., & Xu, Y. (2021). Lightweight Vehicle Detection Algorithm Based on Improved YOLOv4.

*Engineering Letters*, *29*(4), 1544–1551. Academic Search Complete.

Zhao, Y., Zhou, X., Xu, X., Jiang, Z., Cheng, F., Tang, J., & Shen, Y. (2020). A Novel Vehicle Tracking ID

Switches Algorithm for Driving Recording Sensors. *Sensors (14248220)*, *20*(13), 3638–3638.

Academic Search Complete.

Zou, Z., Shi, Z., & Guo, Y. (2019). *Object Detection in 20 Years: A Survey*.

https://arxiv.org/pdf/1905.05055.pdf

## Extras

Github link

*Mohammad Norouzifard*

Trello board