

1 **Resource-Aware Object Detection and Recognition Using** 2 **Edge AI Across the Edge-Fog Computing Continuum***

3 Dragan Stojanović¹, Stefan Sentić¹, Natalija Stojanović¹ and Teodora Stamenković¹

4 Faculty of Electronic Engineering, University of Niš
5 Aleksandra Medvedeva 14, 18000 Niš
6 dragan.stojanovic@elfak.ni.ac.rs, stefan.sentic@elfak.rs, natalija.stojanovic@elfak.ni.ac.rs,
7 teodora.stamenkovic@elfak.rs

8 **Abstract.** Edge computing and edge intelligence have gained significant traction
9 in recent years due to the proliferation of Internet of Things (IoT) devices, the expo-
10 nential growth of data generated at the network edge, and the demand for real-time
11 and context-aware applications. Despite its promising potential, the application of
12 Artificial Intelligence (AI) on the edge faces many challenges, such as edge comput-
13 ing resource constraints, heterogeneity of edge devices, scalability issues, security
14 and privacy concerns, etc. The paper addresses the challenges of deploying deep
15 neural networks (DNNs) for edge intelligence and traffic object detection and recog-
16 nition on a video captured by edge device cameras. The primary aim is to analyze
17 resource consumption and achieve resource-awareness, optimizing computational
18 resources across diverse edge devices within the edge-fog computing continuum
19 while maintaining high object detection and recognition accuracy. To accomplish
20 this goal, a methodology is proposed and implemented that exploits the edge-to-fog
21 paradigm to distribute the inference workload across multiple tiers of the distributed
22 system architecture. The edge-fog related solutions are implemented and evaluated
23 in several use cases on datasets encompassing real-world traffic scenarios and traf-
24 fic objects' recognition problems, revealing the feasibility of deploying DNNs for
25 object recognition on resource-constrained edge devices. The proposed edge-to-fog
26 methodology demonstrates enhancements in recognition accuracy and resource uti-
27 lization, validating the viability of both edge-only and edge-fog based approaches.
28 Furthermore, experimental results demonstrate the system's adaptability to dynamic
29 traffic scenarios, ensuring real-time recognition performance even in challenging
30 environments.

31 **Keywords:** Resource awareness, Traffic Object Recognition, Edge AI, Distributed
32 Neural Networks, Edge-Fog Computing Continuum.

33 **1. Introduction**

34 Edge computing refers to the paradigm of processing data near its source or point of col-
35 lection, rather than relying solely on centralized cloud servers. Edge intelligence involves
36 the integration of artificial intelligence (AI) and machine learning (ML) algorithms into
37 edge devices, enabling them to perform data analytics and decision-making tasks locally.
38 By bringing computation and intelligence closer to the data source, edge computing and
39 edge intelligence offer numerous benefits, including reduced latency, improved bandwidth

* The paper is an extension of the paper presented at RAW 2023 workshop

1 efficiency, enhanced privacy and security, and increased resilience to network failures. AI
2 on edge devices and infrastructure powers real-time, context-aware, and intelligent appli-
3 cations across various fields, such as healthcare, smart cities, industrial automation, trans-
4 portation, and agriculture. However, its potential comes with challenges, including limited
5 resources, diverse device architectures, scalability difficulties, and concerns around secu-
6 rity and privacy. The evolution of edge AI across the edge-fog computing continuum has
7 been extensively explored in recent literature, highlighting the significance of distributed
8 ML and AI in enabling real-time decision-making [20]. It underscores the importance of
9 integrating ML and AI algorithms into edge devices and fog nodes to facilitate intelligent
10 data processing, minimizing latency and network bandwidth usage [8].

11 In recent years, the rapid proliferation of Internet of Things (IoT) devices, includ-
12 ing microcontrollers, single-board computers and smartphones, equipped with built-in or
13 externally connected cameras, has led to their ubiquitous usage across a myriad of appli-
14 cations requiring detection and recognition of objects on the real-time video streams. This
15 widespread adoption has heralded the advent of edge intelligent systems across diverse
16 domains, with a notable emphasis on traffic object detection and recognition. Accurately
17 identifying and classifying traffic objects, such as cars, trucks, motorcycles, bicycles, and
18 pedestrians, is crucial for efficient traffic management, advanced driver assistance sys-
19 tems (ADAS), and autonomous driving. There is an increasing reliance on ML and deep
20 learning (DL) algorithms for video stream analysis for object detection and classification
21 in safety-critical embedded systems and IoT applications, such as autonomous driving
22 systems, surveillance systems and security robots. It emphasizes the need for ML/DL
23 technologies to meet strict timing requirements in real-time systems while maintaining
24 accuracy, given the potentially catastrophic consequences of missed deadlines. Bian et al.
25 in [2] aim to provide a comprehensive exploration of state-of-the-art results in ML/DL-
26 based scheduling techniques, accuracy trade-offs, and security considerations in real-time
27 IoT systems. The potential of Deep Neural Networks (DNNs) to perform efficiently on
28 edge IoT devices is particularly significant, as it harnesses the computational power and
29 availability of these widely used devices. Exploring the intersection of advanced neural
30 network architectures, real-time data processing, and distributed computing paradigms
31 is essential for developing innovative solutions for traffic object detection, classification,
32 and recognition. By combining the proximity and processing capabilities of edge devices,
33 such as smartphones and microcontrollers, with fog servers and cloud infrastructure, re-
34 search efforts aim to enhance the efficiency, accuracy, and responsiveness of traffic object
35 detection systems.

36 This paper addresses the challenges of training and deploying DNNs for traffic object
37 detection and recognition across various edge devices, including Android smartphones,
38 microcontrollers, and single-board computers. The focus is on distributing computational
39 and inference tasks between IoT devices at the far edge, edge servers, and fog servers. We
40 examine different deployment strategies, balancing trade-offs between model size, infer-
41 ence speed, power consumption, and accuracy. Specifically, we evaluate two approaches:
42 (1) quantizing and optimizing the DNN model as TensorFlow Lite for direct deployment
43 on edge devices, and (2) deploying the original DNN model on an edge or fog server.
44 We also explore distributing the object recognition task by dividing it into two inference
45 stages: object detection performed at the edge and object recognition of the detected items
46 carried out on the fog server.

1 Through a series of experiments on traffic object detection and recognition, we eval-
2 uate the performance, accuracy, and resource consumption across different platforms, in-
3 cluding microcontrollers, smartphones, single-board computers, and commodity servers,
4 under various video data parameters and configurations. The results offer valuable insights
5 into the practicality and efficiency of distributing DNNs for traffic object recognition from
6 the edge to the fog. These findings support resource-aware edge intelligence by optimizing
7 computational resource usage while preserving high recognition accuracy. Furthermore,
8 this research lays the groundwork for developing intelligent transportation systems that
9 harness the potential of edge devices, such as Android smartphones and microcontrollers,
10 along with fog computing infrastructure, to improve traffic safety and management.

11 This paper represents the extended version of the paper presented at the RAW 2023
12 Workshop [19]. Besides significant extension of the related work section and detailed ex-
13 planation of the traffic object detection and recognition solutions implemented and pre-
14 sented in the original paper, the main contribution of this paper lies in a new use case for
15 traffic object detection over novel edge devices using a TinyML approach. This involves
16 deploying the traffic object detection solution on a microcontroller (Arduino Nano 33
17 BLE Sense) and a single-board computer (Raspberry Pi 4). Various DNN models suitable
18 for object detection and recognition tasks have been utilized in their original versions and
19 subsequently optimized, quantified and compressed to enable deployment on mentioned
20 edge devices. Extensive experiments have been conducted to evaluate the performance
21 and accuracy of these applications concerning machine learning tasks. The experiments
22 also assessed resource usage, including memory and processing time. The results have
23 been described and analysed in detail. The extended paper provides thorough insights
24 into the implementation and experimental evaluation of various traffic object detection
25 and recognition tasks in different distributed configurations. It also covers the distribution
26 of tasks across the edge-fog computing continuum, from microcontrollers and single-
27 board computers to smartphones and fog servers (PCs).

28 The paper is structured as follows. Section II presents related research work in edge
29 computing and edge intelligence related to object detection and recognition from video
30 streams. Section III presents several strategies and corresponding applications for deploy-
31 ment of DNN for traffic object detection and recognition across various edge devices and
32 a fog server. Section IV presents the experimental evaluation for various traffic object de-
33 tection scenarios and discusses the evaluation results. Section V gives concluding remarks
34 and directions for future research.

35 2. Related Work

36 A growing body of research has focused on harnessing the power of edge devices, such as
37 IoT devices and edge servers, to perform real-time object detection and recognition tasks
38 at the network edge, minimizing latency and bandwidth consumption. The convergence of
39 edge computing and DL methods and techniques has enabled the deployment of resource-
40 efficient object detection and recognition models directly on edge devices, facilitating
41 autonomous decision-making and edge AI applications.

42 Singh and Gill in [18] gives the extensive review of the unique characteristics and
43 advantages of deploying AI algorithms directly on edge devices, enabling real-time infer-
44 ence and decision-making at the network edge. Furthermore, the survey discusses the di-

1 verse applications of Edge AI across domains such as smart cities, healthcare, autonomous
2 vehicles, and industrial automation, highlighting its transformative potential in enhanc-
3 ing efficiency, scalability, and privacy in distributed computing environments. The article
4 delves into the challenges and open research issues associated with Edge AI, such as
5 resource constraints, security concerns, and algorithmic optimizations.

6 The need to integrate ML techniques into resource-constrained embedded devices,
7 facilitated by advancements in technologies like the IoT and edge computing has given
8 rise to TinyML, an embedded ML technique, a by enabling ML applications on low-cost,
9 resource-constrained devices. However, implementing TinyML comes with challenges
10 such as processing capacity optimization and maintaining model accuracy [9].

11 Shuvo et al. in [17] address the challenges of deploying DNNs on edge devices. It
12 highlights the computational complexity and memory requirements of DNNs, which often
13 necessitate cloud-based processing, leading to latency issues and security concerns. The
14 paper explores optimization techniques at both hardware and software levels to enable
15 efficient DNN deployment on edge devices, focusing on four research directions: novel
16 DL architecture and algorithm design, optimization of existing DL methods, algorithm-
17 hardware co-design, and efficient accelerator design. Through a comprehensive review,
18 the paper provides insights into state-of-the-art tools and techniques for efficient edge in-
19 ference, aiming to facilitate the integration of artificial intelligence capabilities into next-
20 generation edge devices.

21 The research on the distribution of DNNs for resource-aware systems has gained sig-
22 nificant attention in recent years. Several studies have explored different approaches and
23 strategies for optimizing the training and deployment of DNNs in various domains. In the
24 context of object detection and recognition from edge to cloud, several relevant research
25 papers provide valuable insights and inspiration. Bittencourt et al. [3] discuss the inte-
26 gration and challenges of the IoT, fog, and cloud continuum, highlighting the need for
27 efficient resource utilization. Lockhart et al. [12] propose Scission, a performance-driven
28 and context-aware cloud-edge distribution approach for DNNs, emphasizing the impor-
29 tance of considering context and performance in distribution decisions. Cho et al. [4]
30 present a study on DNN model deployment on distributed edges, focusing on distributed
31 inference across edge devices.

32 Lin et al. [11] propose a distributed DNN deployment approach from the edge to the
33 cloud for smart devices, addressing the challenges of efficient utilization of resources in
34 different computing tiers. McNamee et al. [14] advocate for adaptive DNNs in edge com-
35 puting, emphasizing the need for dynamic adaptation to optimize resource usage. Ren et
36 al. [16] provide a survey on collaborative DNN inference for edge intelligence, exploring
37 the collaborative aspects of inference across edge devices. Hanhirova et al. [6] charac-
38 terize the latency and throughput of convolutional neural networks (CNN) for mobile
39 computer vision, providing insights into the perfor-mance aspects of DNNs on resource-
40 constrained devices.

41 Lee et al. [10] propose Transprecise Object Detection (TOD) for maximizing real-
42 time accuracy on the edge, highlighting the importance of accurate object detection for
43 edge scenarios. Parthasarathy et al. [15] introduce DEFER, a distributed edge inference
44 approach for DNNs, focusing on resource-efficient inference in distributed edge environ-
45 ments. Teerapittayanon et al. [21] investigate distributed DNNs over the cloud, edge, and

1 end devices, highlighting the trade-offs between resource utilization and computational
2 capabilities across different components of the system architecture.

3 In line with our research, Dharani et al. in [5] present the utilisation of TinyML and
4 TensorFlow Lite on a mobile phones for image classification, but without more extensive
5 experimental evaluation and discussions. Akhtar et al. in [1] introduce multiple real-time
6 deployable cost-efficient solutions for motorbike detection using state-of-the-art embed-
7 ded edge devices, addressing the critical need for accurate and real-time traffic surveil-
8 lance and road safety. The paper presents an improved baseline accuracy of motorbike
9 detection by developing a custom network based on YOLOv5, offering practical insights
10 for real-world implementation.

11 The aforementioned papers contribute to the understanding of resource aware DNN
12 deployment and optimization techniques in various contexts. Our research aims to provide
13 insights into the efficient training and deployment of DNNs for traffic object detection,
14 classification and recognition, considering the resource utilization from edge to fog in the
15 context of edge devices with various computing capabilities and resources available.

16 **3. Traffic Object Detection and Recognition Across the Edge-Fog** 17 **Continuum**

18 DNNs have shown promising results in various computer vision tasks, including object
19 detection and recognition. In this section we present three case studies related to distri-
20 bution of DNN-based software components in the context of traffic object detection and
21 recognition.

22 **3.1. Traffic Object Recognition**

23 To recognize traffic objects captured by a camera on Android smartphones, we implement
24 two approaches: the first executes entirely on the edge device, while the second divides the
25 process between the edge device and a fog node. The training of the DNN model is based
26 on the TensorFlow Object Detection API and specifically utilizes the SSD MobileNet
27 V2 320x320 coco17 tpu-8 pre-trained model. To enhance the training process, we have
28 utilized video data captured from an Android phone camera, as well as publicly available
29 traffic video datasets such as the Udacity Self Driving Car Dataset [22], INRIA Graz-02
30 (IG02)[13], and the Bike-rider Detector dataset [23]. Manual labeling was applied to these
31 datasets when necessary to ensure accurate annotation of traffic objects, for detection of
32 cars, trucks, motorcycles, bicycles, and pedestrians.

33 The first approach we propose utilizes the computing power and resources available
34 solely on the smartphone, making it an offline approach. This method does not require
35 a network connection for traffic object recognition within the Android application. To
36 accommodate the limited resources available on the smartphone, the TensorFlow model
37 used for recognition is quantized and converted to TensorFlow Lite form. By optimizing
38 the model, we ensure that it can efficiently operate on the smartphone without compro-
39 mising its performance. The trained model is integrated into the Android application,
40 enabling real-time traffic object recognition directly on the smartphone without the need
41 for an internet connection (Figure1). The Figure 2 illustrates the execution flow of object

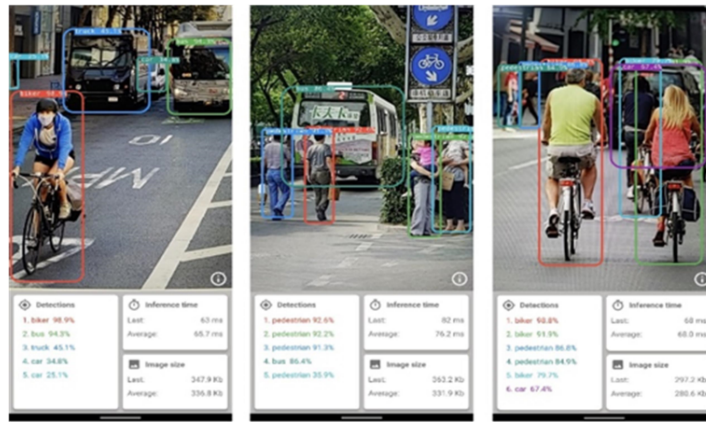


Fig. 1. Android application for traffic object recognition

1 recognition conducted within an Android application. Once the camera image is avail-
 2 able, it is sent to the detector component. In offline mode, the *LocalDetector* component
 3 is employed, utilizing a TensorFlow Lite model deployed on the Android smartphone.
 4 Prior to inputting the image into the model, specific preparations, including scaling and
 5 rotation, must be completed. Furthermore, after detection, it is crucial to parse the results
 6 and generate objects that will be used for GUI creation.

7 The second approach utilizes the advantages of fog computing by offloading part of
 8 the processing to a fog server. In this method, the video captured by the smartphone camera
 9 undergoes preprocessing within the Android application. These preprocessing steps
 10 may include scaling, rotation, and filtering of the captured images to enhance the quality
 11 and clarity of the input data. The preprocessed images are then encoded in Base64 format
 12 and sent to the fog server through a Web socket using the SocketIO library. The fog
 13 server, implemented with Flask/Python, hosts the original TensorFlow model, which con-
 14 ducts object recognition on the received images. The results of this recognition process
 15 are returned to the Android application in JSON format, providing real-time feedback and
 16 visualization of the recognized traffic objects.

17 The execution flow of object recognition in edge-fog scenario is illustrated in Figure
 18 3. The client application captures an image from the camera and processes it before send-
 19 ing it to the server. This processing includes scaling, rotation (considering the device's
 20 sensor), and encoding the image into Base64 format (as a string). The image is trans-
 21 mitted to the server as an emitted event indicating that it is ready for processing. Upon
 22 receipt, the server decodes the Base64 string to reconstruct the image. On the server side,
 23 the TensorFlow library is utilized to run the model and perform object detection within
 24 the frame. The identified objects are then returned to the client in JSON format.

25 The second method distributes the computational workload between the smartphone
 26 and the fog server, offloading resource-intensive tasks to a more powerful computing in-
 27 frastructure. This approach enhances the accuracy and robustness of traffic object recog-
 28 nition, particularly in situations where the smartphone's resources are limited. Addition-
 29 ally, leveraging fog computing alleviates the strain on the smartphone's battery and pro-

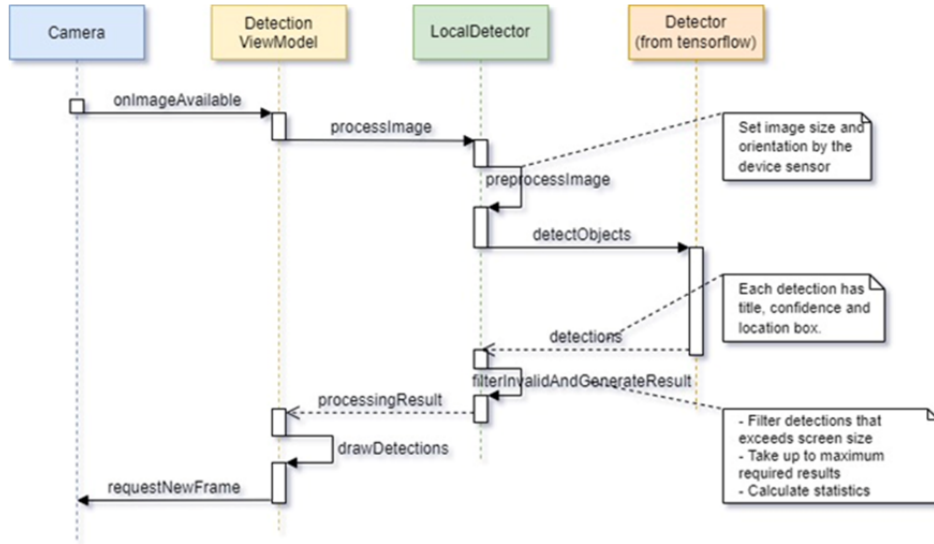


Fig. 2. The sequence diagram during the detection process on Android smartphone.

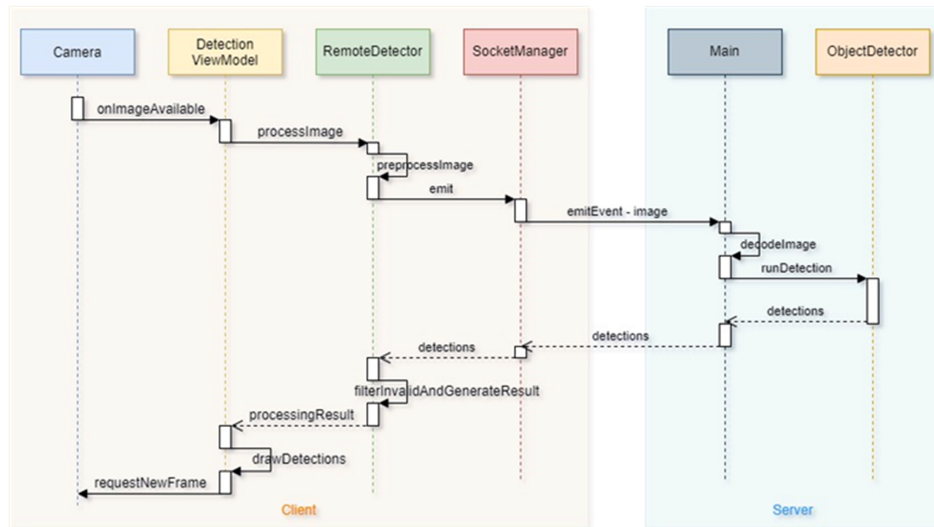


Fig. 3. Sequence diagram for the server-based object detection

1 censing capabilities, resulting in better performance and an improved user experience.
 2 Both methods provide unique advantages regarding resource utilization, real-time per-
 3 formance, and accuracy, addressing various requirements and constraints, which are ex-
 4 perimentally evaluated and discussed in the following section. The source code of An-
 5 droid application implementing traffic object recognition using both method is available
 6 at <https://github.com/drstojanovic/camera>.

7 3.2. Car Model Recognition

8 The second use case is related to car model recognition. To tackle the specific challenge
 9 of car model recognition, we propose a two-stage approach that leverages both edge and
 10 fog computing resources. Our method is built on the pre-trained convolutional neural
 11 network MobileNet V2, known for its outstanding performance across various computer
 12 vision tasks. For training the model specifically for car model recognition, we employ the
 13 Stanford Cars Dataset, which contains over 16,000 images and includes more than 190
 14 different car classes.

15 The execution flow for recognizing car models is illustrated in Figure 4. As depicted,
 16 the process occurs in multiple phases, with some tasks handled on the edge side (client)
 and others on the fog side (server).

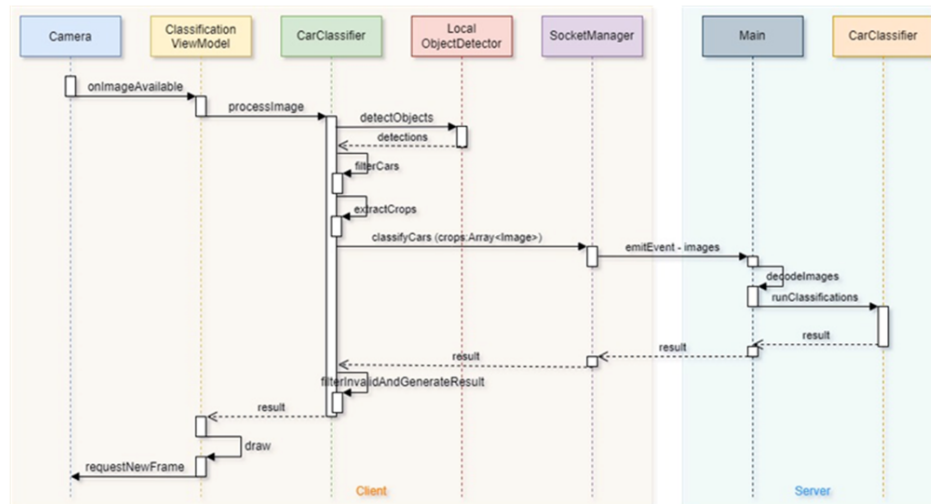


Fig. 4. Execution flow of a car model recognition

17 Since an existing object detector is used, it will return detections for all five classes
 18 of traffic objects: cars, trucks, motorcycles, bikes, and pedestrians. The first step is to
 19 filter out only the class of interest—cars. After capturing the image from the camera, the
 20 following steps are taken:
 21

- 22 1. Object detection is conducted using an edge model, which discards all objects not
 23 belonging to the "car" class.

- 1 2. Based on the detections, the image is cropped, creating a list of car crops.
- 2 3. Each crop is then encoded in Base64 format, and this list of strings is sent to the
- 3 server.
- 4 4. Upon receipt, the server decodes the list back into images. Each crop is resized to fit
- 5 the model's input size (192x192 in this application).
- 6 5. The model is then executed on each crop, providing a list of potential classes for each
- 7 image.
- 8 6. The results are formatted in JSON, resulting in a list of lists whose length corresponds
- 9 to the number of detected cars or crops obtained.
- 10 7. On the edge side, the detection results are merged with the classification results from
- 11 the server, and a bounding box is drawn around each car, displaying the recognized
- 12 class.

The high-level flow of the car model recognition is illustrated in Figure 5. In the first

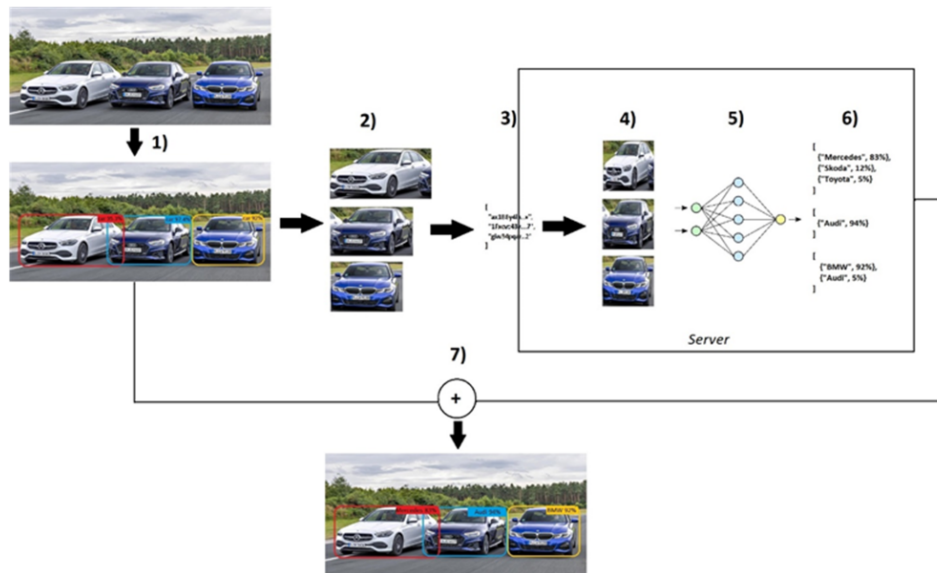


Fig. 5. The steps performed for car model detection.

13 stage of our approach, the Android application takes advantage of edge computing capa-
 14 bilities to detect cars and determine their positions within the video images. We employ
 15 a TensorFlow Lite model, like the one used in the previous implementation, to perform
 16 this initial car detection task on the smartphone. Once the cars are identified, the corre-
 17 sponding regions of interest (ROIs) are extracted from the video frames and preprocessed
 18 to enhance their quality and suitability for subsequent recognition.

20 The preprocessed and cropped car images are then sent from the Android application
 21 to the fog server. The fog server application handles the second stage of the car model
 22 recognition process. It employs a trained TensorFlow model, specifically designed for the

1 Stanford Cars Dataset, to identify the model of each car in the received images and sends
 2 the results back to the Android application (Figure6). The fog server's superior compu-
 3 tational resources and processing power allow it to perform more intensive tasks, such
 as detailed car type classification. To enable communication between the Android appli-

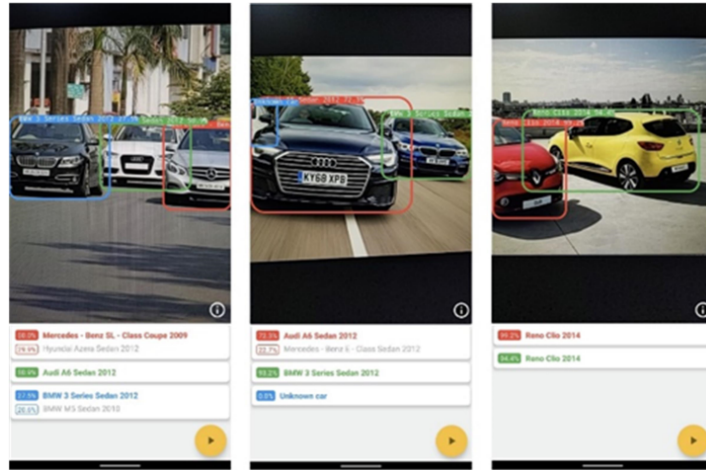


Fig. 6. Car model detection in the Android application

4
 5 cation and the fog server, we use WebSockets. This approach facilitates real-time, bidi-
 6 rectional communication, allowing for the efficient transfer of preprocessed car images
 7 from the smartphone to the server and the return of recognition results to the mobile ap-
 8 plication. By utilizing WebSockets, we ensure a seamless and responsive user experience
 9 throughout the car model recognition process.

10 By leveraging both edge and fog computing resources, our two-stage approach op-
 11 timizes the distribution of computational tasks. The edge computing performed on the
 12 smartphone efficiently detects cars and extracts relevant regions of interest (ROIs), which
 13 reduces the volume of data sent to the fog server. This strategy minimizes bandwidth us-
 14 age and latency. The more resource-intensive task of car model recognition is offloaded
 15 to the fog server, taking advantage of its greater computational capabilities and trained
 16 model. This approach maximizes the utilization of computing resources and enhances the
 17 overall performance and accuracy of car model recognition in our system. The source code
 18 for traffic object recognition server is available at GitHub link <https://github.com/drstojanovic/trafficAssistantServer>
 19 while its docker image can be pulled from https://hub.docker.com/r/stefan2708/ta_server.

20 3.3. Vehicle detection on microcontrollers and single-board computers

21 The third use case focuses on vehicle detection on low-resource edge devices, such as
 22 microcontrollers and single-board computers. The goal is to detect and recognize cars in
 23 video streams captured by integrated cameras. This involves implementing various object
 24 detection model architectures and optimizing them for execution on devices like the Ar-
 25 duino Nano 33 BLE Sense and Raspberry Pi 4. Part of the model training and evaluation

1 was conducted using the Edge Impulse platform [7], an online MLOps platform that sup-
2 ports the training, testing, and deployment of ML/DL models across a wide range of edge
3 devices. Additionally, transfer learning was applied using TensorFlow Lite Model Maker.
4 The training dataset, consisting of 499 images with a total of 4,281 vehicles, was sourced
5 from Kaggle. In addition to collecting data directly from the edge devices, datasets were
6 also uploaded to the Edge Impulse platform in YOLO txt format. In this format, each im-
7 age has an associated .txt file listing the detected objects, where each line in the file repre-
8 sents a single object, containing its class and the normalized coordinates of its bounding
9 box.

10 The initial implementation of vehicle detection was performed on the Arduino Nano
11 33 BLE Sense, a resource-constrained edge device. Given the limited capabilities of this
12 development board, the trained models operate on smaller image dimensions. However,
13 increasing image size leads to longer inference times and larger model sizes. Unfortu-
14 nately, due to the dataset containing small objects, models trained on low-resolution im-
15 ages yielded poor results. For the Arduino Nano application, two models, FOMO Mo-
16 bileNetV2 0.1 and FOMO MobileNetV2 0.35 were trained. The values 0.1 and 0.35 rep-
17 resent the alpha parameters in the MobileNetV2 architecture, indicating the network's
18 width by scaling the number of channels in each layer. These hyperparameters help bal-
19 ance model size, computational efficiency, and accuracy. FOMO models were trained for
20 various numbers of epochs using both RGB and Grayscale images, with various hyperpa-
21 rameter configurations, to find the optimal trade-off between performance and accuracy.
22 The Edge Impulse platform was used to generate binary files for model inference on the
23 Arduino device. Inference can be initiated with the command *edge-impulse-run-impulse*.
24 Object detection results from the OV7675 camera attached to the Arduino Nano are dis-
played in the browser, as shown in Figure 7. The Raspberry Pi 4 Model B (RPi) offers

This shows a live preview of the camera on your device. Classification is running on the device, not on this computer.



Time per inference: 316 ms.

Fig. 7. Object detection from Arduino camera shown in browser

25 greater computing resources, enabling the training of a broader range of models. Beyond
26 the FOMO algorithms, the following models were also trained using the Edge Impulse
27 platform:
28

- 1 1. MobileNetV2 SSD FPN-Lite: This model was pre-trained on the COCO 2017 dataset
2 with images of size 320x320. It consists of three parts:
 - 3 – Basic network (MobileNetV2): Provides high-level features for classification or
4 detection. By removing the fully connected and softmax layers and adding a de-
5 tection network, the model can determine object locations in the image.
 - 6 – Detection network (Single Shot Detector, SSD): Detects multiple objects in an
7 image using a single convolutional network. SSD models are faster and more
8 efficient as they simultaneously predict object classes and regions containing ob-
9 jects.
 - 10 – Feature Pyramid Network (FPN): Utilizes an input image of a single size to ge-
11 nerate feature maps for different sizes, facilitating the detection of objects of vari-
12 ous sizes.
- 13 2. YOLOv5: This model, part of the You Only Look Once (YOLO) family, performs
14 object detection in a single pass. Introduced in 2020, YOLOv5 incorporates the Ef-
15 ficientDet architecture, built on EfficientNet, to optimize both resource usage and
16 accuracy. Unlike its predecessor, YOLOv5 abandons anchor-based detection, instead
17 relying on a convolutional layer to predict bounding box coordinates directly.

18 YOLOv5 is used in the transfer learning process, leveraging prior training on a larger
19 dataset. This enables the model to learn from a broader data set and improve its generaliza-
20 tion capabilities. During the training of the YOLOv5 model on Edge Impulse, the model
21 size can be selected: XS, S, M, and L. Due to resource limitations, the smallest model with
22 1.9M parameters, sized at 3.78 MB, was chosen. Models from Edge Impulse are down-
23 loaded in the EIM (Edge Impulse Model) format. EIM files are binary files for Linux and
24 macOS that encapsulate the complete impulse built on the Edge Impulse platform, includ-
25 ing signal processing, model, and inference blocks. EIM files are architecture-specific and
26 allow direct inference execution on the RPi device.

27 Inference on the RPi device is performed on images generated using a camera con-
28 nected to the Arduino development board. Frames are captured on the Arduino and trans-
29 mitted to the RPi device using serial communication. The OV7675 camera records images
30 sized at 320x240 in RGB565 format. Within the Python script on the RPi device, bytes are
31 read from the serial port, and conversion from RGB565 to RGB888 format is carried out.
32 The converted image is passed to the `run_inference` function, which performs pre-
33 processing and inference (Figure 8). TensorFlow Lite Model Maker¹ is a library designed
34 for training TensorFlow Lite models with custom datasets. It leverages transfer learning to
35 minimize the amount of training data required and reduce training time. For object detec-
36 tion, the library offers five versions of EfficientDet-Lite models, each differing in memory
37 usage, detection latency, and mean Average Precision (mAP). These models are deployed
38 on the RPi for real-time vehicle detection. In this setup, frames transmitted from the Ar-
39 duino Nano's camera are used for detection. However, the inference process differs here,
40 as it utilizes the TensorFlow Lite Interpreter, which requires a tensor as input. This neces-
41 sitates preprocessing the image and generating a tensor with the correct shape and data
42 type. Finally, the function completes the process by drawing bounding boxes and sav-
43 ing the annotated image. The source code of the Arduino Nano and RPi applications are
44 available at <https://github.com/drstojanovic/object-detection-rpi>.

¹ https://www.tensorflow.org/lite/models/modify/model_maker



Fig. 8. Detection of vehicles using DL models deployed on RPi

1 **4. Resource-Aware Experimental Evaluation**

2 This section presents an experimental evaluation of previously described use cases and the
3 corresponding methods for distributing DNNs in traffic object detection and recognition.
4 In the first use case, we assess two solutions: one executed entirely within the Android
5 application using a TensorFlow Lite model, and the other performed on the fog server. In
6 the third use case, we evaluate solutions implemented on Arduino Nano and Raspberry Pi
7 (Rpi) devices. To measure performance, accuracy, and resource consumption (CPU, mem-
8 ory, and energy), we conduct experiments using various video data parameters and con-
9 figurations. The objective is to explore the trade-offs between different approaches across
10 the edge-fog computing continuum and analyze their behavior under varying conditions.
11 For evaluation, we use representative datasets that include a variety of traffic scenarios
12 and object types, ensuring coverage of diverse lighting conditions, weather patterns, and
13 traffic densities.

14 **4.1. Smartphone and Server implementation**

15 To evaluate the Android application solution, we measure its performance and accuracy
16 directly on the smartphone. The computational requirements, including CPU usage, mem-
17 ory consumption, and energy consumption, are analyzed using the Android Profiler tool.
18 Furthermore, we assess the accuracy of traffic object recognition by comparing the appli-
19 cation's outputs with ground truth annotations from the datasets.

20 Similarly, for the fog server solution, we assess its performance, accuracy, and re-
21 source consumption. The server's computational requirements, such as CPU usage, mem-
22 ory utilization, and energy consumption, are analyzed. Additionally, the recognition re-
23 sults from the fog server are compared with ground truth annotations to evaluate the so-
24 lution's accuracy. A timeline diagram illustrating CPU and memory consumption, along
25 with energy usage for traffic object recognition performed on the Android smartphone
26 and the fog server, is presented in Figure 9. The maximum CPU utilization during local
27 detection on the Android smartphone reached 33%, while server-based detection peaked
28 at 14%. The maximum RAM usage during local detection was 316 MB, compared to 240
29 MB for server-based detection. The TensorFlow model size is 11.2 MB; however, with
30 model quantization, it can be reduced to 3.2 MB.

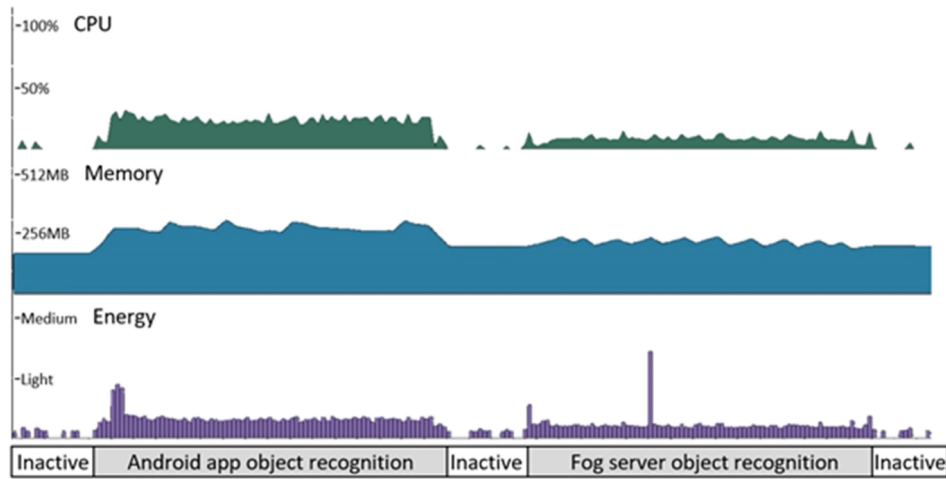


Fig. 9. CPU, memory, and energy usage for edge (smartphone) and edge server solutions

1 During the experimental evaluation, we varied the video data parameters and con-
 2 figurations to analyze the performance of both solutions under different conditions. This
 3 involved adjusting factors such as video resolution, frame rate, lighting conditions, and
 4 traffic densities. By conducting experiments across a range of scenarios, we aim to pro-
 5 vide a thorough assessment of each solution’s performance and resource utilization. The
 6 experimental evaluation is carried out using appropriate benchmarking tools and metrics
 7 to ensure reliable and meaningful results. We measure the execution time, resource uti-
 8 lization, and accuracy of both solutions across various datasets and configurations. The
 9 collected data is analyzed and compared to identify the strengths and weaknesses of each
 10 approach. For all experiments, a Google Pixel 4 phone was used, while the server appli-
 11 cation ran on a Lenovo Legion laptop (CPU: i5-9300H, RAM: 16.0 GB, GPU: NVIDIA
 12 GeForce GTX 1650). The server is configured on a local network to facilitate access by
 13 the mobile application.

14 The first experiment aimed to assess the accuracy of recognition with varying image
 15 resolutions and object sizes. The phone’s camera was directed at a computer monitor dis-
 16 playing an image of a car that consistently shrank in size. Testing was conducted for both
 17 operational modes at each of the four available resolutions in the application settings. The
 18 image quality was set to the maximum (100 %, with no compression). The reliability of
 19 detection, expressed as percentages, is presented in Table 1. This experiment showed that
 20 the reliability values for detection, and thus the overall detection quality, are quite compa-
 21 rable for both edge detection and server-based detection. A slight advantage was noted for
 22 edge detection at lower image resolutions, likely because of the extra image processing
 23 involved in sending the image to the server (such as encoding and decoding in Base64
 24 format). One potential solution to improve performance is to use a more advanced image
 25 transport method, like implementing one of the transfer protocols specifically designed for
 26 image handling. We also assessed how performance (execution speed) depends on image
 27 resolution and image quality, with values expressed in milliseconds (ms). This experiment
 28 aimed to evaluate the impact of image compression on detection speed and image size.

Table 1. Dependency of detection accuracy on image resolution and object size

Object size	Resolution							
	640 x 640		512 x 512		300 x 300		160 x 160	
	Edge	Fog	Edge	Fog	Edge	Fog	Edge	Fog
100%	98.9	97.8	98.7	98.7	98.9	98.3	90	84.3
80%	95.4	95	95.1	95.4	94.7	93.1	82.4	82.1
50%	89.2	90	87.4	88.7	89.9	87.9	83.4	80.5
30%	47.1	42.3	40.1	32.8	31.7	19.7	23.1	17.4
15%	12.3	12.7	10.7	12.1	10.2	2.1	7.4	1.8

1 The phone's camera was directed at a computer monitor displaying a consistent image of
 2 a car. Testing was conducted for both operating modes across all four available resolutions
 3 and image quality settings of 100 %, 70 %, 50 %, 30 %, and 20 % achieved through JPEG
 4 image compression.

5 Table 2 illustrates how detection speed varies with changes in image characteristics.
 6 The rows represent a decrease in image quality, while the columns show a reduction in
 7 resolution. The values are presented in milliseconds, and it is evident that local detection
 8 operates at a significantly higher speed compared to server-based detection.

Table 2. Dependency of performance (speed of execution) on image resolution and image quality (in ms)

Quality	Resolution							
	640 x 640		512 x 512		300 x 300		160 x 160	
	Edge	Fog	Edge	Fog	Edge	Fog	Edge	Fog
100%	61.4	153.2	58	137	55.9	103.2	45.8	83.9
70%	56.5	112.2	50.4	99.8	43.9	87.7	43	76.9
50%	55.5	108.2	49.8	98.9	43.6	86.7	42.9	76.4
30%	55.3	103.8	49.1	93.1	43.4	84.7	43.4	75
20%	54.4	101.8	48.9	90.2	43.9	84.1	42.8	74.9

9 Regarding local detection, the difference in detection speed between settings A (640x640
 10 resolution and 100 % quality) and settings B (160x160 resolution and 20 % quality)
 11 amounts to 18.6 milliseconds per frame. Although the difference is difficult to measure
 12 precisely, it is evident that the detection quality for smaller objects is noticeably lower
 13 with settings B.

14 When it comes to server-based detection, internet connection speed becomes a sig-
 15 nificant factor. Since a GPU was utilized, the model execution itself was short (45ms),
 16 with most of the time being spent on image transportation. During the testing, an internet
 17 speed of 55.8 Mbps for download and 7.7 Mbps for upload was used. In the case of edge
 18 detection, the byte size of the image does not have as much influence as it does for server-
 19 based detection, as each image is transported to the server, and any reduction in image
 20 size contributes to increased detection speed. The difference in detection speed between
 21 settings A (640x640 resolution and 100 % quality) and settings B (160x160 resolution

1 and 20 % quality) amounts to 78.3 milliseconds per frame. A significant degradation in
 2 detection quality on the server compared to edge detection was observed when decreasing
 3 the image resolution.

4 By assessing the performance, accuracy, and resource consumption of both the An-
 5 droid application and the fog server solutions, we aim to shed light on the trade-offs asso-
 6 ciated with each approach. This experimental evaluation will enhance our understanding
 7 of how these distribution methods perform in traffic object recognition tasks. Ultimately,
 8 this analysis will aid in choosing the most suitable solution based on specific require-
 9 ments, including resource availability, real-time performance, and accuracy.

10 4.2. Arduino Nano and RPi Solutions for Object Detection

11 Testing and evaluating vehicle detection models involved comparing key parameters es-
 12 sential for implementing AI at the edge. Resource utilization metrics, such as RAM and
 13 flash memory usage, are particularly important for edge devices. Additionally, the time re-
 14 quired for object detection is critical for real-time decision-making scenarios. A compari-
 15 son of model accuracy was also conducted. The objective is to compare models developed
 16 using the Edge Impulse platform and TensorFlow Lite tools, exploring various combina-
 17 tions of preprocessing and hyperparameters. Furthermore, this comparison encompasses
 18 devices with differing resources, specifically the Arduino Nano and RPi.

19 The comparison between the FOMO 0.1 and FOMO 0.35 models on the Arduino
 20 Nano has been conducted. All results are presented for the quantized versions of the mod-
 21 els, utilizing integer 8-bit values. The EON model format was chosen due to its lower re-
 22 source overhead compared to TFLite models. Figure 10 illustrates the comparison of these
 23 models across different image sizes, focusing on model accuracy and inference time. The
 24 diagrams indicate that the accuracy of the FOMO 0.1 and FOMO 0.35 models is compa-
 25 rable across all image sizes. However, FOMO 0.1 shows significantly better performance
 26 for smaller dimensions, such as 64x64, whereas the advantage shifts toward the FOMO
 27 0.35 algorithm for larger dimensions. Inference time increases with image size, reach-
 28 ing 3 seconds for dimensions of 160x160. Furthermore, the inference time of the FOMO
 29 0.35 model diverges considerably from that of the FOMO 0.1 model as the dimensions
 30 increase.

31 Figure 11 illustrates the memory utilization for the same models. Although flash mem-
 32 ory usage remains consistent across various image dimensions, maximum RAM utiliza-
 33 tion rises sharply as the dimensions increase. The diagram also shows the available RAM
 34 on the Arduino development board; models that exceed this limit cannot be executed on
 35 the device.

36 The evaluation of the models executed on the Raspberry Pi includes those generated
 37 on Edge Impulse as well as models created using TensorFlow Model Maker. Figure 12
 38 illustrates the changes in accuracy and inference time for different models with varying
 39 image sizes. The YOLOv5 model demonstrates superior accuracy in all scenarios, except
 40 for the smallest image size. However, in terms of object detection speed, the YOLOv5
 41 algorithm requires up to seven times longer to make decisions compared to the FOMO
 42 algorithms.

43 The resource utilization diagram (Figure 13) reveals that the maximum RAM usage
 44 with FOMO algorithms exceeds that of the YOLO algorithm at higher dimensions. Flash

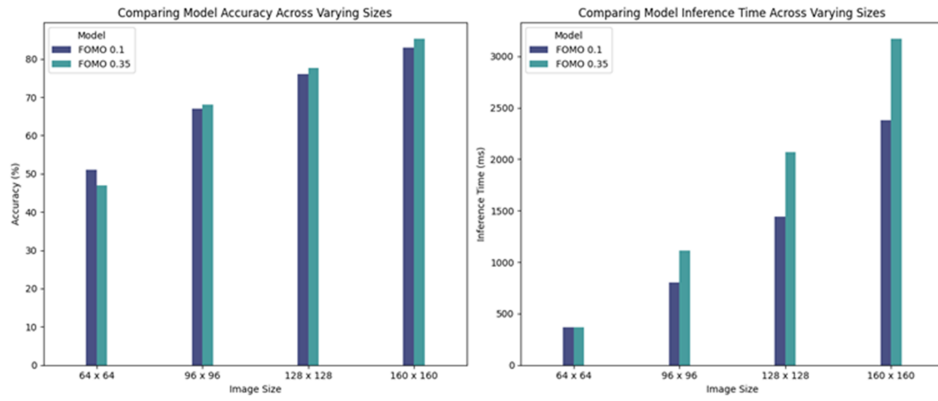


Fig. 10. Comparison of NN models based on accuracy and inference time.

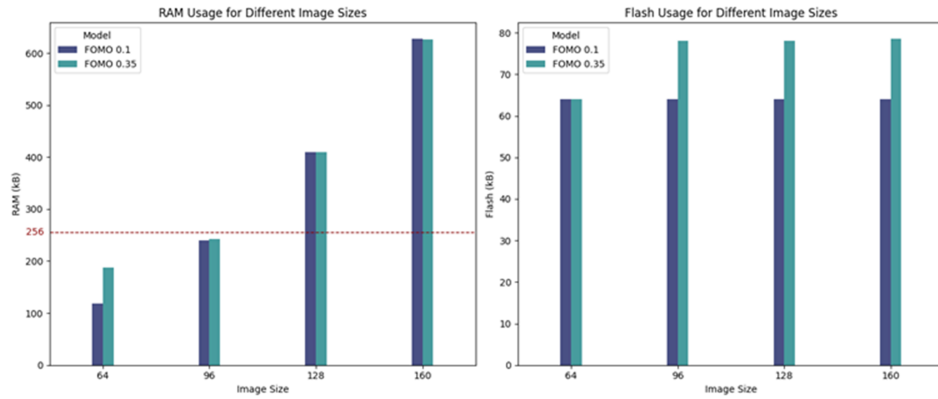


Fig. 11. Comparison of NN models based on memory utilization

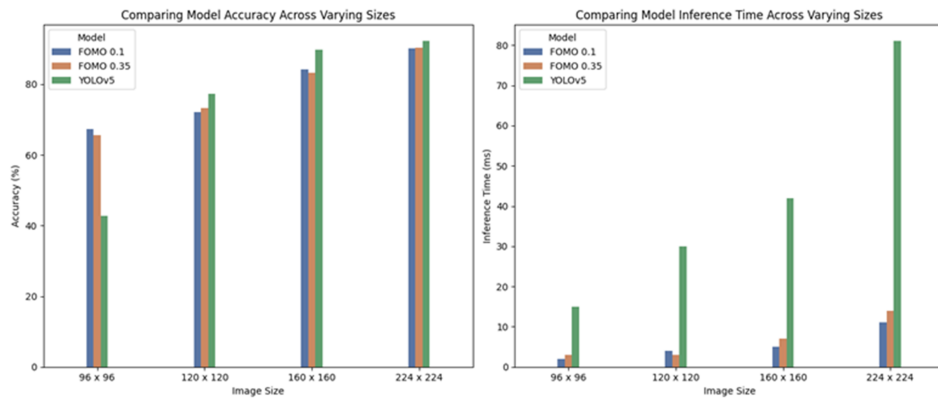


Fig. 12. Comparison of NN models based on accuracy and inference time on RPi

1 memory utilization remains consistent across all dimensions. Notably, the FOMO algo-
 2 rithms require approximately 70 KB, while the YOLOv5 model occupies 1.8 MB. On

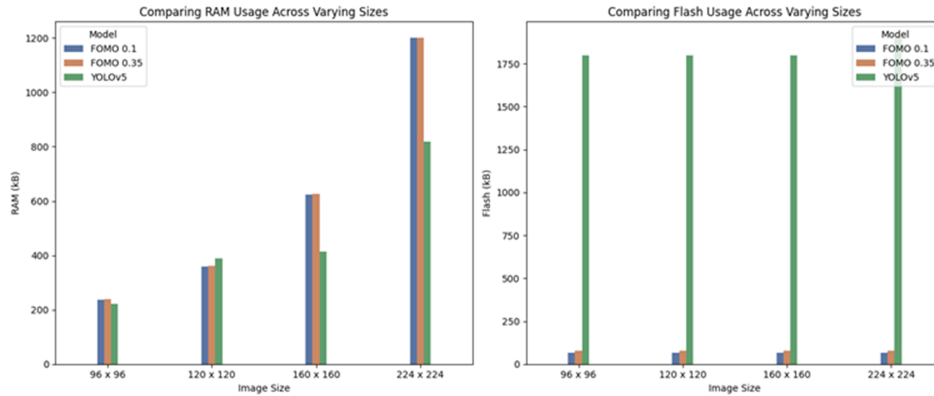


Fig. 13. Comparison of NN models based on memory utilization on RPi

2
 3 Edge Impulse, several model optimizations are available. Figure 14 illustrates the differ-
 4 ences between EON models with float32 values and their quantized versions using int8
 5 values. Notably, quantization has a substantial effect on the size of the MobileNetV2 SSD
 and YOLOv5 models, as well as on inference time. The numerical comparison of the

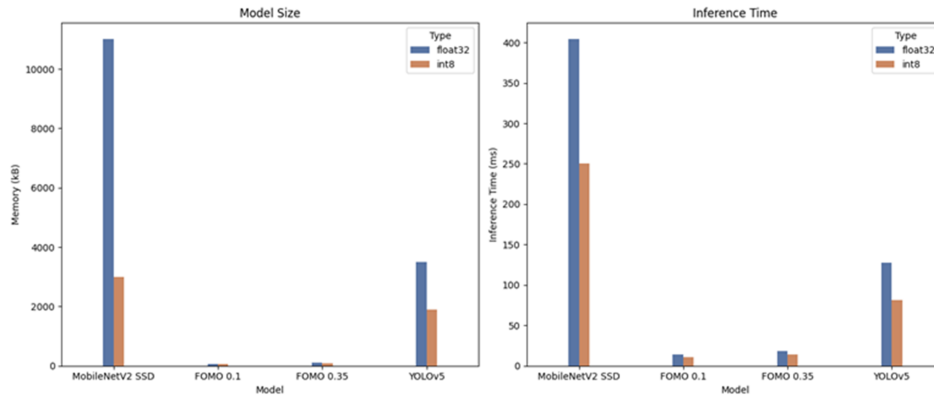


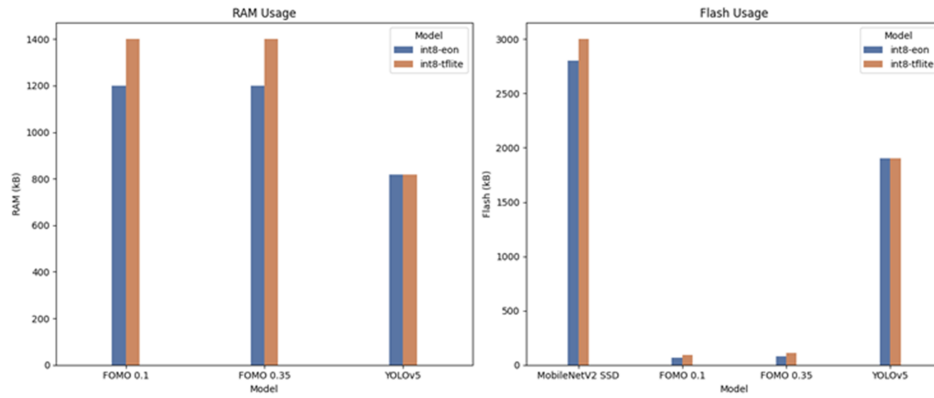
Fig. 14. Comparison between float32 and int8 NN model versions

6
 7 models' size and inference time values is given in Table 3. The comparison of models
 8 compiled using EON and TensorFlow Lite compilers is shown in Figure 15, and pre-
 9 sented in Table 4). The objective of the EON compiler is to reduce resource overhead,
 10 while maintaining unchanged accuracy and inference time. The values are presented for

Table 3. The NN model size and inference time values

	Model size (float32)	Model size (int8)	Inference time (float32)	Inference time (int8)
MobileNetV2 SSD	11 MB	3 MB	404 ms	249 ms
FOMO 0.1	66.9 KB	64.4 KB	14 ms	11 ms
FOMO 0.35	102.7 KB	78.5 KB	18 ms	14 ms
YOLOv5	3.5 MB	1.9 MB	128 ms	81ms

1 the int8 versions of the models. Resource optimization using the EON compiler is not supported for YOLOv5.

**Fig. 15.** Comparison of EON and TensorFlow Lite models

2

Table 4. Review of memory utilization for EON and TensorFlow Lite models

	RAM (EON)	RAM (TensorFlowLite)	Flash (EON)	RAM (TensorFlowLite)
MobileNetV2 SSD	/	/	2.8 MB	3 MB
FOMO 0.1	1.2 MB	1.4 MB	64.4 KB	94.4 KB
FOMO 0.35	1.2 MB	1.4 MB	110.5 KB	78.5 KB
YOLOv5	817.5 KB	817.5 KB	1.9 MB	1.9 MB

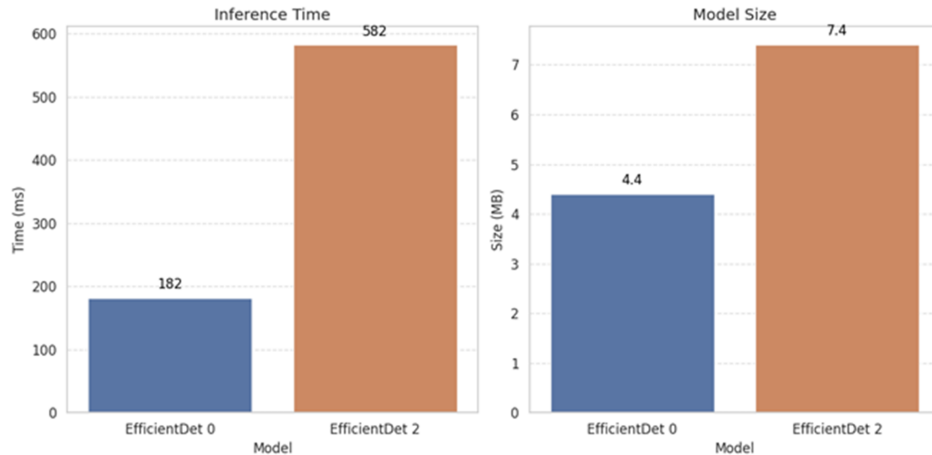
3 TensorFlow Model Maker enables training of five versions of EfficientDet models (1-
4 5). The results of the trained EfficientDet2 model are presented in Table 5. Due to the
5 potential loss of model accuracy resulting from optimization, an evaluation of the Tensor-
6 Flow Lite model is provided using mAP (mean Average Precision) metric for evaluation.
7 The analysis of the results concludes that the differences in accuracy between TensorFlow
8 and TensorFlow Lite models are nearly indistinguishable. The TensorFlow model's pre-
9 cision is 45.35 %. At a 50 % overlap threshold, the model's precision is 81.3 %, while
10 for a 75 % overlap threshold, the value is nearly halved. Detection precision for differ-

Table 5. Evaluation of the TensorFlow Lite model

	TensorFlow model	TensorFlow Lite model
AP	0.4535	0.4426
AP50	0.81	0.8
AP75	0.48	0.46
APs	0.003	0.0027
APm	0.35	0.3
APl	0.66	0.65
ARmax1	0.08	0.08
ARmax10	0.5	0.5
ARmax100	0.58	0.54
ARs	0.06	0.045
ARm	0.55	0.49
ARl	0.75	0.71
AP_car/	0.45	0.44

ent sizes (APs – small, APm – medium, APl – large) indicates that the model performs well with larger objects. Additionally, the model exhibits better accuracy when there are multiple objects in the image (ARmax10, ARmax100). There is room for improvement in detecting small objects and under very strict overlap criteria.

EfficientDet 0 and 2 versions were deployed on the RPi device, and a comparison of inference time and model size is given in Figure 16. Memory analysis during the ob-

**Fig. 16.** Comparison of EfficientDet 0 and 2

ject detection process was performed and the results are given in Table 6. Both models require a similar amount of memory, with only minor differences in usage. EfficientDet 0 consumes slightly fewer resources than EfficientDet 2; however, both models function effectively within the available system resources. Through comprehensive testing

Table 6. Memory utilization of EfficientDet models (in KB)

	total	used	free	shared	buff/cache	available
EfficientDet 0	3794	599	1937	75	1256	3043
EfficientDet 2	3794	613	1915	80	1265	3025

1 and evaluation of NN models created using various tools, significant advancements in
 2 technologies for deploying models on edge devices have become apparent. By leveraging
 3 advanced optimization techniques, it is feasible to execute complex model architectures
 4 on resource-constrained devices. For example, models trained on Edge Impulse occupy
 5 considerably less space and consume fewer resources compared to TensorFlow Lite mod-
 6 els. Conversely, the strength of the TensorFlow Lite framework lies in its broader selection
 7 of model architectures and the ability to utilize established models for object detection.

8 5. Conclusions

9 The proposed approach and experimental evaluations provide valuable insights into the
 10 challenges and opportunities of deploying DNNs for traffic object detection and recog-
 11 nition across the edge-fog computing continuum. One of the key takeaways is that dis-
 12 tributing tasks between edge devices and fog servers offers a balanced trade-off between
 13 performance, resource consumption, and accuracy. By leveraging edge devices like smart-
 14 phones for object detection, we can reduce data transmission and latency. We found that
 15 current smartphones perform well in both accuracy and speed, with only 33% CPU uti-
 16 lization during inference. This raises the question of whether offloading detection tasks to
 17 the fog is necessary, given that the communication overhead may outweigh any potential
 18 performance benefits. The results suggest that, in many cases, performing both detection
 19 and recognition on the mobile device itself is a more efficient strategy, especially for real-
 20 time applications.

21 Meanwhile, the fog server, with its higher computational capacity, is well-suited well-
 22 suited for more complex tasks such as car model recognition, where detailed feature ex-
 23 traction and processing are required. This delegation not only reduces the computational
 24 burden on edge devices like smartphones or embedded platforms but also enables the
 25 deployment of heavier models that might otherwise exceed the capabilities of smaller
 26 devices. For instance, running computationally intensive architectures such as YOLOv5
 27 or EfficientDet on the fog server ensures that these models can operate without com-
 28 promising performance or requiring extensive optimization, as would be necessary on
 29 edge devices. This division of labor improves the system’s overall efficiency by allowing
 30 lightweight tasks, such as object detection, to occur locally on edge devices while offload-
 31 ing more demanding tasks to the fog server. Consequently, the combination of fast local
 32 inference with sophisticated fog-based recognition creates a robust pipeline that balances
 33 speed and accuracy across different layers of the edge-fog continuum.

34 While the smartphone proved capable, the resource constrained embedded devices,
 35 such as RPi or Arduino Nano, could still benefit from fog or cloud offloading, especially
 36 for complex models. However, our focus in this study was on evaluating lightweight archi-
 37 tectures directly on these devices. Future research could explore hybrid approaches, where
 38 RPi or Arduino devices collaborate with fog or cloud systems for more demanding tasks,

1 striking a balance between local processing and offloading. We found that the quantized
 2 and minimized model deployed on Arduino Nano, RPI device, and Android smartphones
 3 achieved reasonable performance with efficient resource usage. The use of quantization
 4 and model optimization techniques, especially through TensorFlow Lite and Edge Im-
 5 pulse, proved crucial for deploying DNNs on resource-constrained devices like the Ar-
 6 duino Nano and RPi. Our results show that while Edge Impulse models consume fewer
 7 resources, TensorFlow Lite offers greater flexibility through access to a wider variety of
 8 model architectures. The experiments also highlighted the trade-offs between model size,
 9 inference time, and accuracy, emphasizing the importance of fine-tuning hyperparameters
 10 based on the specific hardware and use case.

11 The ability to achieve near real-time detection using optimized models demonstrates
 12 that modern AI technologies can effectively run on low-power edge devices. Our findings
 13 suggest that combining edge and fog computing provides a scalable and efficient way to
 14 implement more demanding AI solutions, such as specific object recognition on video
 15 stream. The experiments also revealed that model size and computational overhead must
 16 be carefully managed, particularly on microcontrollers, where even slight increases in
 17 image resolution or model complexity can lead to significant resource constraints.

18 However, there are still several avenues for future research in this area. Some potential
 19 directions include:

- 20 – Exploration of federated learning techniques tailored for object detection and recog-
 21 nition tasks at the edge for aggregating model updates from distributed edge nodes
 22 efficiently while accounting for resource constraints.
- 23 – Investigation of online continual learning techniques for adaptive object detection and
 24 recognition at the edge that enable edge devices to incrementally learn from streaming
 25 data while retaining knowledge learned from previous tasks.
- 26 – Research resource-efficient model adaptation techniques in edge-fog environments
 27 that dynamically adjust model complexity and capacity based on available computa-
 28 tional resources and streaming data characteristics.
- 29 –

30 By further exploring these research directions, we can continue to advance the field
 31 of DNN model distribution and slicing across the edge-fog-cloud computing continuum,
 32 enabling resource-aware and efficient object detection, classification and recognition sys-
 33 tems for various real-world applications.

34 **6. Acknowledgements**

35 Research presented in this paper is supported by the Ministry of Science and Techno-
 36 logical Development, Republic of Serbia, 451-03-65/2024-03/200102, Erasmus+ Project
 37 FAAI: The Future is in Applied AI (WP4 - Artificial Intelligence framework for training
 38 in HE) - 2022-1-PL01-KA220-HED-000088359 and COST Action CERCIRAS - Con-
 39 necting Education and Research Communities for an Innovative Resource Aware Society
 40 - CA19135.

1 References

- 2 1. Akhtar, A., Ahmed, R., Yousaf, M.H., Velastin, S.A.: Real-time motorbike detection: Ai on the
3 edge perspective. *Mathematics* 12(7) (2024), [https://www.mdpi.com/2227-7390/](https://www.mdpi.com/2227-7390/12/7/1103)
4 [12/7/1103](https://www.mdpi.com/2227-7390/12/7/1103)
- 5 2. Bian, J., Arafat, A.A., Xiong, H., Li, J., Li, L., Chen, H., Wang, J., Dou, D., Guo, Z.: Machine
6 learning in real-time internet of things (iot) systems: A survey. *IEEE Internet of Things Journal*
7 9(11), 8364–8386 (2022)
- 8 3. Bittencourt, L., Immich, R., Sakellariou, R., Fonseca, N., Madeira, E., Curado, M., Villas, L.,
9 DaSilva, L., Lee, C., Rana, O.: The internet of things, fog and cloud continuum: Integration
10 and challenges. *Internet of Things* 3-4, 134–155 (2018)
- 11 4. Cho, E., Yoon, J., Baek, D., Lee, D., Bae, D.: Dnn model deployment on distributed edges. In:
12 Bakaev, M., Ko, I.Y., Mrissa, M., Pautasso, C., Srivastava, A. (eds.) *ICWE 2021 Workshops.*
13 *Communications in Computer and Information Science*, vol. 1508. Springer, Cham (2021)
- 14 5. Dharani, A., Kumar, S.A., Patil, P.N.: Object detection at edge using tinyml mod-
15 els. *SN Computer Science* 5(1), 11 (11 2023), [https://doi.org/10.1007/](https://doi.org/10.1007/s42979-023-02304-z)
16 [s42979-023-02304-z](https://doi.org/10.1007/s42979-023-02304-z)
- 17 6. Hanhirova, J., Kämäräinen, T., Seppälä, S., Siekkinen, M., Hirvisalo, V., Ylä-Jääski, A.: La-
18 tency and throughput characterization of convolutional neural networks for mobile computer
19 vision. In: *Proceedings of the 9th ACM Multimedia Systems Conference (MMSys '18)*. pp.
20 204–215. New York, NY, USA (2018)
- 21 7. Hymel, S., Banbury, C.R., Situnayake, D., Elium, A., Ward, C., Kelcey, M., Baaijens, M.,
22 Majchrzycki, M., Plunkett, J., Tischler, D., Grande, A., Moreau, L., Maslov, D., Beavis,
23 A., Jongboom, J., Reddi, V.J.: Edge impulse: An mlops platform for tiny machine learning.
24 *ArXiv abs/2212.03332* (2022), [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:254366602)
25 [254366602](https://api.semanticscholar.org/CorpusID:254366602)
- 26 8. Iftikhar, S., Gill, S.S., Song, C., Xu, M., Aslanpour, M.S., Toosi, A.N., Du, J., Wu, H.,
27 Ghosh, S., Chowdhury, D., Golec, M., Kumar, M., Abdelmoniem, A.M., Cuadrado, F., Vargh-
28 ese, B., Rana, O., Dustdar, S., Uhlig, S.: Ai-based fog and edge computing: A system-
29 atic review, taxonomy and future directions. *Internet of Things* 21, 100674 (2023), [https:](https://www.sciencedirect.com/science/article/pii/S254266052200155X)
30 [//www.sciencedirect.com/science/article/pii/S254266052200155X](https://www.sciencedirect.com/science/article/pii/S254266052200155X)
- 31 9. Kallimani, R., Pai, K., Raghuvanshi, P., Iyer, S., López, O.L.A.: Tinyml: Tools, appli-
32 cations, challenges, and future research directions. *Multimedia Tools and Applications*
33 83(10), 29015–29045 (2024), <https://doi.org/10.1007/s11042-023-16740-9>,
34 [dOI: 10.1007/s11042-023-16740-9](https://doi.org/10.1007/s11042-023-16740-9)
- 35 10. Lee, J.K., Varghese, B., Woods, R., Vandierendonck, H.: Tod: Transprecise object detection
36 to maximize real-time accuracy on the edge. In: *Proceeding of the 5th EEE 5th International*
37 *Conference on Fog and Edge Computing*. pp. 53–60 (2021)
- 38 11. Lin, C.Y., Wang, T.C., Chen, K.C., Lee, B.Y., Kuo, J.J.: Distributed deep neural network de-
39 ployment for smart devices from the edge to the cloud. In: *Proceedings of the ACM Mobi-*
40 *Hoc Workshop on Pervasive Systems in the IoT Era*. p. 43–48. *PERSIST-IoT '19*, Association
41 for Computing Machinery, New York, NY, USA (2019), [https://doi.org/10.1145/](https://doi.org/10.1145/3331052.3332477)
42 [3331052.3332477](https://doi.org/10.1145/3331052.3332477)
- 43 12. Lockhart, L., Harvey, P., Imai, P., Willis, P., Varghese, B.: Scission: Performance-driven
44 and context-aware cloud-edge distribution of deep neural networks. In: *Proceedings of the*
45 *IEEE/ACM 13th International Conference on Utility and Cloud Computing (UCC)*. pp. 257–
46 268. Leicester, United Kingdom (2020)
- 47 13. Marszałek, M., Schmid, C.: Inria annotations for graz-02 dataset (2007), [https://lear.](https://lear.inrialpes.fr/people/marszalek/data/ig02/)
48 [inrialpes.fr/people/marszalek/data/ig02/](https://lear.inrialpes.fr/people/marszalek/data/ig02/), accessed: 2024-10-1
- 49 14. McNamee, F., Dustdar, S., Kilpatrick, P., Shi, W., Spence, I., Varghese, B.: The case for adap-
50 tive deep neural networks in edge computing. In: *Proceedings of the IEEE 14th International*
51 *Conference on Cloud Computing (CLOUD)*. pp. 43–52 (2021)

- 1 15. Parthasarathy, A., Krishnamachari, B.: Defer: Distributed edge inference for deep neural net-
2 works. In: Proceedings of the 14th International Conference on COMMunication Systems &
3 NETworkS (COMSNETS). pp. 749–753 (2022)
- 4 16. Ren, W., Qu, Y., Dong, C., Jing, Y., Sun, H., Wu, Q., Guo, S.: A survey on collaborative dnn
5 inference for edge intelligence. *Machine Intelligence Research* 20, 370–395 (2023)
- 6 17. Shuvo, M.M.H., Islam, S.K., Cheng, J., Morshed, B.I.: Efficient acceleration of deep learning
7 inference on resource-constrained edge devices: A review. *Proceedings of the IEEE* 111(1),
8 42–91 (2023)
- 9 18. Singh, R., Gill, S.S.: Edge ai: A survey. *Internet of Things and Cyber-Physical Systems*
10 3, 71–92 (2023), [https://www.sciencedirect.com/science/article/pii/
11 S2667345223000196](https://www.sciencedirect.com/science/article/pii/S2667345223000196)
- 12 19. Stojanovic, D., Sentic, S., Stojanovic, N.: Towards resource-efficient dnn deployment for traf-
13 fic object recognition: From edge to fog. In: Zeinalipour, D., Blanco Heras, D., Pallis, G.,
14 Herodotou, H., Trihinas, D., Balouek, D., Diehl, P., Cojean, T., Furlinger, K., Kirkeby, M.H.,
15 Nardelli, M., Di Sanzo, P. (eds.) Euro-Par 2023: Parallel Processing Workshops. pp. 30–39.
16 Springer Nature Switzerland, Cham (2024)
- 17 20. Taheri, J., Dustdar, S., Zomaya, A., Deng, S.: *Edge Intelligence: From Theory to Prac-*
18 *tice*. Springer International Publishing (2023), [https://books.google.es/books?
19 id=cCV5zwEACAAJ](https://books.google.es/books?id=cCV5zwEACAAJ)
- 20 21. Teerapittayanon, S., McDanel, B., Kung, H.T.: Distributed deep neural networks over the cloud,
21 the edge and end devices. In: Proceedings of the IEEE 37th International Conference on Dis-
22 tributed Computing Systems (ICDCS). pp. 328–339. Atlanta, GA, USA (2017)
- 23 22. Udacity: Udacity self-driving car dataset (2024), [https://public.roboflow.com/
24 object-detection/self-driving-car](https://public.roboflow.com/object-detection/self-driving-car), accessed: 2024-10-1
- 25 23. Yong-Hah, R.: Bike-rider detector (2019), [https://github.com/yonghah/
26 bikerider-detector](https://github.com/yonghah/bikerider-detector), accessed: 2024-10-1