



BRNO UNIVERSITY OF TECHNOLOGY

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

FACULTY OF INFORMATION TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ

DETECTION AND CLASSIFICATION OF VEHICLES FOR EMBEDDED PLATFORMS

DETEKCE A KLASIFIKACE VOZIDEL PRO VESTAVĚNÉ PLATFORMY

BACHELOR'S THESIS

BAKALÁŘSKÁ PRÁCE

AUTHOR

AUTOR PRÁCE

PATRIK SKALOŠ

SUPERVISOR

VEDOUCÍ PRÁCE

Ing. JAKUB ŠPAÑHEL

BRNO 2023

Bachelor's Thesis Assignment



144778

Institut: Department of Computer Graphics and Multimedia (UPGM)
Student: **Skaloš Patrik**
Programme: Information Technology
Specialization: Information Technology
Title: **Detection and Classification of Vehicles for Embedded Platforms**
Category: Image Processing
Academic year: 2022/23

Assignment:

1. Learn about methods for detecting and classifying objects in an image.
2. Obtain a suitable dataset for detecting different types of vehicles in the image.
3. Find deep learning methods that focus on the object detection problem, primarily methods that support multi-class detection or object detection coupled with subsequent classification.
4. Select appropriate methods for use on embedded platforms and experiment with them.
5. Evaluate the selected methods in an appropriate manner and discuss the results obtained.
6. Create a poster and video presenting your work, its objectives and results.

Literature:

- GE, Zheng, et al. YoloX: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021.
- HUANG, Jonathan, et al. Speed/accuracy trade-offs for modern convolutional object detectors. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017. p. 7310-7311.
- Dle pokynů vedoucího.

Requirements for the semestral defence:

- Completion of the first three points of the assignment
- Considerable work on the fourth point of the assignment

Detailed formal requirements can be found at <https://www.fit.vut.cz/study/theses/>

Supervisor: **Špaňhel Jakub, Ing.**
Head of Department: Černocký Jan, prof. Dr. Ing.
Beginning of work: 1.11.2022
Submission deadline: 10.5.2023
Approval date: 31.10.2022

Abstract

An abstract of the work in English will be written in this paragraph.

Abstrakt

Keywords

Here, individual keywords separated by commas will be written in English.

Klíčová slova

Reference

SKALOŠ, Patrik. *Detection and Classification of Vehicles for Embedded Platforms*. Brno, 2023. Bachelor's thesis. Brno University of Technology, Faculty of Information Technology. Supervisor Ing. Jakub Špaňhel

Detection and Classification of Vehicles for Embedded Platforms

Declaration

I hereby declare that this Bachelor's thesis was prepared as an original work by the author under the supervision of Mr. X The supplementary information was provided by Mr. Y I have listed all the literary sources, publications and other sources, which were used during the preparation of this thesis.

.....
Patrik Skaloš
May 6, 2023

Acknowledgements

Here it is possible to express thanks to the supervisor and to the people which provided professional help (external submitter, consultant, etc.).

Contents

1	Introduction	5
2	Background and Related Work	6
2.1	Common Approaches to Vehicle Detection	6
2.1.1	Vehicle Detection Without a Camera	6
2.1.2	Vehicle Detection Based on Image Processing	7
2.2	Convolutional Neural Networks	7
2.2.1	Architecture of Convolutional Neural Networks	7
2.2.2	Popular Architectures of Convolutional Neural Networks	9
2.3	Object Detection	10
2.3.1	Two-Stage Object Detectors	10
2.3.2	One-Stage Object Detectors	10
2.4	Network Optimization and Compression	11
2.4.1	Network Pruning	12
2.4.2	Knowledge Distillation	12
2.4.3	Weight Quantization	12
2.5	Evaluation Metrics	13
2.6	Vehicle Detection Based on Convolutional Neural Networks	14
2.7	Embedded Platforms for Machine Learning	15
2.7.1	Google Coral	15
2.7.2	Movidius Neural Compute Stick	16
2.7.3	NVIDIA Jetson	16
2.8	Tools and Libraries	17
2.8.1	LabelBox annotation app	17
2.8.2	PyTorch	17
2.8.3	MMDetection Library	17
2.8.4	MMYOLO Library	18
2.8.5	MMDeploy Library	18
2.8.6	ONNX and ONNXRuntime	18
2.8.7	TensorRT	18
3	Datasets	19
3.1	Dataset Criteria	19
3.2	Individual Datasets	20
3.2.1	UA-DETRAC	20
3.2.2	Miovision Traffic Camera Dataset	21
3.2.3	AAU RainSnow Traffic Surveillance Dataset	22
3.2.4	Multi-View Traffic Intersection Dataset	24

3.2.5	Night and Day Instance Segmented Park Dataset	25
3.2.6	VisDrone Dataset	26
3.3	Dataset processing	26
3.4	Summary of Datasets	27
3.5	Training, Validation and Testing Dataset Split	27
4	Object Detection Models	30
4.1	Model Architectures	30
4.2	Model Configurations	31
4.2.1	Dataset Wrappers	32
4.2.2	Training Augmentation Pipeline	33
5	Experiments	36
5.1	Devices Used in Experiments	36
5.1.1	NVIDIA Jetson Platforms	36
5.1.2	NVIDIA GeForce MX150	36
5.1.3	Intel Core i7-9850H	36
5.1.4	Raspberry PI 4B	37
5.1.5	Software versions	37
5.2	Model Deployment and Optimizations	38
5.3	Experiments and Evaluation	39
6	Results	41
6.1	Precision-Recall Curves of Major Models	41
6.2	Effect of Model's Input Resolution on Inference Speed	42
6.3	Inference Speeds: ONNX Runtime vs. TensorRT	43
6.4	Effect of Model's Shape on Inference Speed	43
6.5	Effect of Batch Size on Inference Speed	44
6.6	Mean Average Precision and Inference Speed Benchmark	45
6.7	Evaluating Models Across Multiple Mean Average Precision Metrics	47
6.8	Inference Speeds on Different Devices	48
6.9	TODO Confusion matrix?	51
6.10	TODO DETRAC vs MIO-TCD?	51
7	Future Work	55
8	Conclusion	56
	Appendices	57
A	Additional Precision-Recall Curves	58
B	Complete Benchmark of All Models on All Devices	59
C	Complete Mean Average Precision Metrics for YOLOv8 Models	62
	Bibliography	64

List of Figures

2.1	Example of computing a 2D convolution on an input vector with an example kernel, from [26].	8
2.2	Visualization of the Intersection over Union calculation.	14
3.1	Example image from the UA-DETRAC dataset.	21
3.2	Example images from the Miovision dataset.	22
3.3	Example image from the AAU RainSnow dataset.	23
3.4	Example image from the drone subset of the Multi-View Traffic Intersection dataset.	24
3.5	Example image from the NDISPark dataset.	25
3.6	Example image from the VisDrone dataset.	26
3.7	Example images from sequences from the DETRAC-UA dataset used for the validation and testing subsets.	29
4.1	Examples of images augmented by the training augmentation pipeline, including the ground truth labels.	33
6.1	Precision-Recall curves are displayed for four major models, including YOLOv8-femto with the smallest input resolution of 352×192 . It is important to note that, in this figure, the YOLOv8-nano model chosen for comparison has an input resolution of 512×288 , although most experiments used the model with an input resolution of 640×384 . This decision was made to simplify the comparison between YOLOv8-nano and YOLOv8 MobileNetV2.	42
6.2	Inference speed (FPS) comparison of YOLOv8-nano, YOLOv8-pico, and YOLOv8-femto models with different input resolutions. The results generally support our hypothesis that the reduction in model's input resolution is directly proportional to the increase in inference speed, although with some deviations. These results were obtained by running tests on NVIDIA Jetson Nano using the TensorRT backend. The tested models were deployed with dynamic shape in the batch size dimension to run the tests with a batch size of 16 to most accurately measure the speed of the smallest models.	44
6.6	Performance comparison of all YOLOv8 models in terms of mAP (mean Average Precision) and inference speed (FPS) on a Raspberry PI 4B with ONNX Runtime inference backend and the batch size of 1. Please note that the X-axis is logarithmic to be able to display all models in a single plot for easier analysis. We also remind the reader that only a single CPU thread was used when testing models on this device.	46

6.7	Performance comparison of all YOLOv8 models, including quantized ones, in terms of mAP (mean Average Precision) and inference speed (FPS) on NVIDIA Jetson Xavier NX with TensorRT inference backend and the batch size of 32. Please note that the X-axis is logarithmic to be able to display all models in a single plot for easier analysis.	47
6.8	Benchmark of inference speeds (FPS) on different devices featuring all model architectures. The second plot compares models quantized to FP16 precision on NVIDIA Jetson devices. Please note that in both subplots, the X-axis is logarithmic and that for comparison, YOLOv8-nano model with 512×288 input resolution and YOLOv8-femto model with 352×192 input resolution were selected. A batch size of 8 was selected for tests featuring larger models – YOLOv8-medium, YOLOv8-small and YOLOv8 MobileNetV2, and for the rest of the models, a batch size of 32 was used. The TensorRT inference backend was utilized when testing on NVIDIA Jetson devices, while the ONNX Runtime backend was used on the rest of the devices. Additionally, we remind the reader that for the tests on CPUs, only a single thread was used, utilizing just one of the CPU's cores.	50
6.3	Comparison of inference speeds of all trained YOLOv8 model architectures between the TensorRT and the ONNX Runtime inference backends (of course, both utilizing the devices' GPU). Tests were conducted on NVIDIA Jetson Nano with a batch size of 1 and on NVIDIA Jetson AGX Xavier with a batch size of 32. Please note that for each model architecture, only the model with the largest input resolution was tested and all tested models were deployed as dynamic in the batch size dimension.	52
6.4	Comparison of inference speeds between models with static and dynamic shapes in the batch size dimension. Only the most representative models were selected for comparison, as the rest of the data do not provide additional insights.	53
6.5	Inference speeds achieved by using different batch sizes. Tests were conducted on Raspberry PI 4B, NVIDIA Jetson Nano and NVIDIA Jetson AGX Xavier and only the most representative models were selected for comparison to keep the figure compact. Please note that the horizontal axes are logarithmic to better present the differences in inference speeds. Additionally, inference speed of YOLOv8-medium 640×384 with the batch size of 32 on NVIDIA Jetson Nano is missing because the model did not fit into the memory.	54
A.1	Precision-Recall curves for the rest of the models (which were not included in Figure 6.1).	58

Chapter 1

Introduction

Vehicle detection systems are an important tool for enhancing road safety, and with the rise of convolutional neural networks and the recent advancements in object detection, these systems have become more affordable and accurate. However, high-performance object detectors still require significant computational resources and sacrifices must be made to perform this task in real-time on embedded devices or those with limited computational resources. This paper aims to evaluate the performance trade-offs of state-of-the-art object detectors on low-performance devices and to explore various strategies for improving their efficiency.

A popular family of real-time object detectors is the YOLO (You Only Look Once) series, with the latest, state-of-the-art YOLOv8 model architecture being the primary focus of this paper. To thoroughly evaluate its performance, we benchmark several YOLOv8 models of different sizes, including one trained with the lightweight MobileNetV2 backbone. These models are tested across six different devices, including three NVIDIA Jetson embedded platforms. To accelerate the process of training, deploying and testing different models, this paper utilized PyTorch-based libraries from the OpenMMLab project, including MMYOLO and MMDeploy.

The models were trained on several different datasets that consist of surveillance-type images, most commonly captured by traffic-monitoring cameras mounted on existing infrastructure. The datasets combined contain a total of 1 204 795 object instances in 218 821 images.

In addition to benchmarking the trained models on all devices in terms of their inference speeds (FPS) and mean Average Precision (mAP) metrics, we also aim to explore additional factors that can influence their performance, including weight quantization and the usage of different inference backends, inference batch sizes and input resolutions of the models.

In [chapter 2](#), we begin by explaining essential concepts, tools and libraries relevant to our study. We then proceed with a detailed review of the datasets employed in our research in [chapter 3](#), followed by a discussion on the evaluated object detector models and their training process in [chapter 4](#).

Subsequently, in [chapter 5](#), we delve into the specifics of our experiments, outlining the methodologies employed, providing an overview of the devices used for testing and discussing the important aspects to consider when interpreting the results. These results are then presented in detail in [chapter 6](#), where we share our findings and insights gained from the experiments. Finally, we suggest potential areas for future work that were beyond the scope of this paper.

Chapter 2

Background and Related Work

2.1 Common Approaches to Vehicle Detection

Camera-based object detectors based on convolutional neural networks are evaluated in this paper. However, to provide further context to the problem, other commonly used solutions for the problem of vehicle detection are briefly explained in this section – non-camera-based object detectors and camera-based image processing techniques.

2.1.1 Vehicle Detection Without a Camera

Although camera-based vehicle detection systems have gained significant popularity due to advancements in computer vision and machine learning, there are alternative techniques that do not rely on cameras for vehicle detection. These methods offer different advantages, mainly reduced computational requirements or improved performance in certain environmental conditions. Since each detector has its own limitations, a vehicle detection system typically consists of different types of detectors to overcome these limitations. Following is a summary of non-camera-based techniques widely used for vehicle detection tasks.

Magnetic Sensors

One of the simplest solutions to detect vehicles is to use an induction loop.^[4] The sensor, composed of a single wire, can be buried under concrete and detect the presence of metal objects passing by. With a controller, induction loops typically provide data about vehicle presence, but using more advanced algorithms, speed, approximate classification and much more can be determined from this simple sensor. Although the design is very simple, installation is not and induction loops can even get damaged over time, while repairs call for temporary shutdowns of roads. It is worth noting that there are many other types of magnetic sensors, which can provide more detailed and accurate data with simpler sensor installation, but the idea stays the same.

Ultrasonic Sensors

Another type of inexpensive and simple detector is an ultrasonic sensor. ^[32] These sensors can operate in a variety of conditions and are typically mounted on existing infrastructure. These distance measurement sensors are usually used to detect the presence and distance of nearby vehicles and their speed, but they can be utilized in many different ways to even provide a simple shape of a vehicle to classify it.

Radars and Lidars

Radars (Radio Detection and Ranging), which can also be mounted on existing infrastructure above ground, work similarly to ultrasonic sensors, but a single radar can oversee a much wider area of a road. [11] They are usually used to detect the presence, speed, heading and shape of a vehicle. The shape can then be used to predict a class of the vehicle. Lidars (Light Detection and Ranging) can be used in a similar way, but are far more accurate, which makes them better at detecting shapes and locations of objects.

2.1.2 Vehicle Detection Based on Image Processing

Image processing techniques take a camera feed as input and in addition to detecting the presence, speed, heading and shape of a vehicle, they can also provide its color, license plate and countless other characteristics, limited mostly by the software, not by the detector. These camera-based systems, which are relatively inexpensive, can be mounted on existing infrastructure, do not emit any energy and are highly versatile, while providing a large amount of data.

An algorithm applies a series of pre-defined filters and transformations to an image to extract patterns and features that resemble vehicles. Although modern convolutional neural networks (explained in the following [section 2.2](#)) are generally more effective at recognizing more complex patterns, image processing can still be extremely helpful for simpler tasks or as a component of a larger vehicle detection system.

A huge amount of research has been conducted on using image processing to solve the problem of vehicle detection [37]. Many of the techniques proposed can operate in real-time and are very reliable in typical environmental conditions. However, they can be sensitive to changes in illumination or have their performance affected by rain, snow, shadows, occlusions, or noise. While there are methods for addressing some of these challenges, such as shadow removal techniques [24], they add to the computational requirements of the system. Overall, our research suggests that the image processing approach is well-suited for simpler tasks or systems with lower accuracy requirements, but it is often difficult to implement and may lack versatility in more complex or demanding scenarios.

2.2 Convolutional Neural Networks

[[Lepšie ich popísať? Sú základom projektu, no]]

Convolutional neural networks (CNNs) are a class of deep learning models that have gained widespread popularity in recent years, mainly thanks to their ability to learn hierarchical features from raw image data. CNNs are particularly useful in the field of computer vision for various tasks, including image classification, object detection and segmentation. In the context of vehicle detection, these models have demonstrated a superior combination of accuracy, efficiency and flexibility compared to traditional detection methods explained earlier. For example, compared to object detection based on image processing, the main advantage of CNNs is that they are able to learn automatically from raw image data instead of relying on pre-defined filters and transformations. [20]

2.2.1 Architecture of Convolutional Neural Networks

A typical CNN consists of several key components, including convolutional layers, pooling layers and fully connected layers. In this subsection, details of these key components and

their roles in the context of object detection are explained. Please note that a modern convolutional neural network consists of more building blocks which will not be discussed here to keep this introduction short.

Convolutional Layers

Convolutional layers form the backbone of a CNN and perform 2D convolution operations on the input data using trainable kernels to detect specific patterns or features. A kernel, also known as a filter, is responsible for detecting a particular feature, such as an edge or a texture. It is practically a matrix, usually small in spatial dimensionality and its weights are trainable parameters adjusted during the training phase. [26]

A two-dimensional convolution is a mathematical operation that involves computing element-wise multiplications of the convolution kernel and the corresponding sub-region of the input image. This process is typically repeated throughout the entire image in a sliding manner, resulting in a new image (matrix) as an output called a feature map. It can be better explained by Figure 2.1.

Many filters (often of different types) are used in a typical CNN and together produce a set of feature maps that capture countless different aspects of the input image.

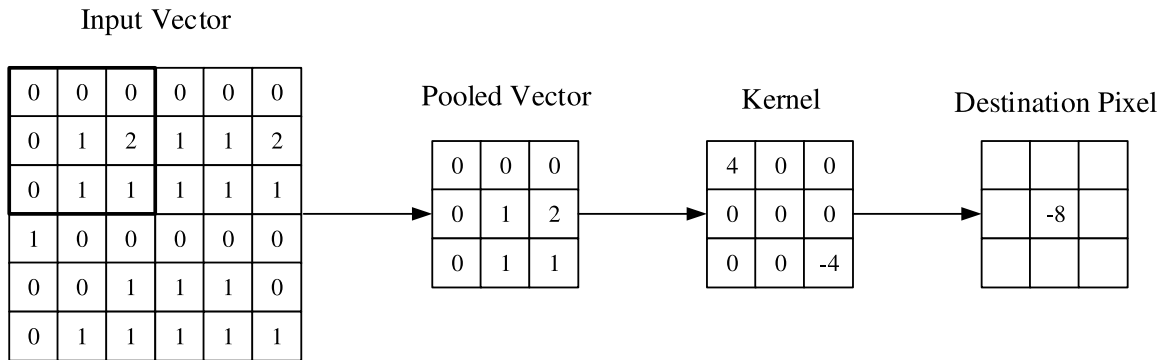


Figure 2.1: Example of computing a 2D convolution on an input vector with an example kernel, from [26].

Pooling Layers

The pooling layers serve to gradually reduce the spatial dimensions of feature maps output from the convolutional layers while retaining as much information as possible. These layers improve the computational efficiency of the CNN by reducing the number of its parameters. Generally, the pooling operation takes a window or a filter and moves it over the input feature map in strides, most commonly taking the maximum value of the inputs within the window. This operation is called max-pooling. [26]

Fully-connected Layers

A fully-connected layer, also known as a dense layer is a type in which each neuron in the next layer is connected to each neuron in the previous one. In a CNN, it is typically used at the end of the network to classify the input data (features extracted by the previous convolutional layer). [26]

2.2.2 Popular Architectures of Convolutional Neural Networks

[[Treba vôbec toto? Nejak to neviem naviazať na nasledujúce subsections]]

Over the years, various CNN architectures have been developed to address different challenges and requirements. In this subsection, we review several influential architectures which achieved state-of-the-art¹ performance in a wide range of tasks, including image classification, object detection or semantic segmentation.

[[skontrolovať či som niečo nesplietol a či je všetko správne]]

LeNet-5

LeNet-5 [19], introduced by Yann LeCun and his team in 1998 is considered one of the first successful applications of convolutional neural networks. Designed for handwritten digit recognition, LeNet-5 was the source of inspiration for modern CNNs with its combination of convolutional, pooling and fully connected layers.

AlexNet

Developed by Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton in 2012, AlexNet [17] marked a significant breakthrough in the field of deep learning and won the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) by a considerable margin. It popularized the use of deep CNNs for image classification and featured the use of Rectified Linear Units (ReLU) as activation functions, dropout technique for regularization and data augmentations for training.

VGGNet

VGGNet [31], proposed by Karen Simonyan and Andrew Zisserman in 2014, is best known for its uniform architecture, consisting of a series of stacked convolutional layers with 3×3 filters followed by max-pooling. VGGNet demonstrated that deeper networks generally achieve better performance, achieving top results in the ILSVRC with its 16-layer and 19-layer variants (VGG-16 and VGG-19). However, its success lies in being very computationally expensive.

ResNet

In 2015, ResNet (Residual Network) [13] was introduced by Kaiming He and his team, addressing the degradation problem that occurs in very deep CNNs. ResNet incorporates skip connections, enabling the network to effectively “skip” layers during the training process. This innovation allowed ResNet to scale up to hundreds of layers while improving performance, which was thought impossible. ResNet achieved state-of-the-art results in various computer vision tasks and has inspired many subsequent architectures.

MobileNet

MobileNet [14], developed by Andrew G. Howard and his team at Google in 2017, is a popular light-weight convolutional neural network architecture that has been widely used on mobile and embedded devices with limited computational resources. One of its key

¹The term “state-of-the-art” refers to methods that have achieved superior results compared to previous best methods in a specific task or application.

features is its use of depthwise separable convolutions. Traditional convolutional layers perform a full convolution on the input, but depthwise separable ones perform a depthwise convolution followed by a pointwise convolution, which reduces the number of computations required while maintaining accuracy. Another advantage is the model’s small size and low computational requirements compared to previously discussed networks allowing for real-time use, of course by sacrificing some accuracy.

MobileNetV2

MobileNetV2 [29], an evolution of the original MobileNet proposed by Mark Sandler and his team at Google in 2018. The architecture introduces the concept of *Inverted Residuals* and *Linear Bottlenecks*, aiming to enhance the network’s representational capacity while maintaining low computational complexity. *Inverted Residual* blocks leverage the idea of residual connections to mitigate the vanishing gradient problem and improve the training phase. *Linear Bottlenecks*, on the other hand, reduce the computational cost of convolutional layers by decreasing the number of channels before applying a convolutional operation and restoring the original number of channels afterward. These advancements enable MobileNetV2 to achieve superior accuracy and efficiency compared to its predecessor.

2.3 Object Detection

While the CNNs discussed in section 2.2 can be used for image classification², the object detection task also involves localization – marking all objects in the input image by a bounding box. However, these CNNs can and often are used as backbones³ to extract features from the input images. Modern object detectors can be divided into two categories – one-stage detectors and two-stage detectors.

2.3.1 Two-Stage Object Detectors

In the two-stage object detection task, the first stage selects region proposals (selecting regions that are likely to contain an object) and passes them to the second stage for classification. In 2013, the R-CNN framework [12] (Regions with CNN features) was designed by Girshick *et al.*, replacing the old and inefficient sliding window detection technique, marking a breakthrough in object detection. [20]

Although accurate, object detectors based on the R-CNN architecture (including fast R-CNN and faster R-CNN) are generally computationally expensive compared to one-stage object detectors and won’t be discussed here in detail.

2.3.2 One-Stage Object Detectors

When building a real-time object detector, most commonly a one-stage detector is selected. Although their accuracy is generally lower than that of two-stage detectors, their speed allows for real-time object detection even on resource-constrained devices.

²Image classification is the task of classifying an image into a class category – for example recognizing whether an image contains a cat or a dog.

³The backbone of an object detector is responsible for extracting features from the input image.

Single Shot MultiBox Detector (SSD)

The Single Shot MultiBox Detector (SSD) is an object detection algorithm that aims to provide a balance between accuracy and computational efficiency. Proposed by Wei Liu *et al.* in their 2015 paper [21], SSD addresses the challenge of detecting objects in images with varying scales and aspect ratios in real-time.

Unlike the two-stage object detectors which consist of a region proposal stage followed by a classification stage, SSD directly predicts both the object locations and their corresponding class labels in a single forward pass of the network.

Built upon a base network such as VGG-16, SSD adds a series of convolutional layers to capture information at different scales, allowing detection of various object sizes. It uses anchor boxes placed uniformly across feature maps to handle varying aspect ratios.

You Only Look Once (YOLO)

YOLO (You Only Look Once) is a family of one-stage object detectors that prioritize real-time performance. Introduced by Joseph Redmon *et al.* in 2015 [28], YOLO divides the input image into a grid, where each grid cell is responsible for predicting bounding boxes and class probabilities for objects located within that cell. YOLOv1 has undergone many iterations each improving upon the previous version’s performance and efficiency—from YOLOv2 to YOLOv8 and including YOLOX, YOLOR and several more.

In this work, we focus on the latest, YOLOv8 state-of-the-art object detector [16] developed by Ultralytics, which, at this time, doesn’t have a paper released. However, information about the detector can be drawn from a post by OpenMMLab [25].

Built upon the YOLOv5, the YOLOv8 detector updated the head⁴ module to be a decoupled one, separating the classification and detection heads, while the new backbone and neck⁵ modules are based on the YOLOv7 ELAN concept. In terms of training strategy, YOLOv8 extends the training epochs from 300 to 500 and the data augmentation process is modified during the final 10 epochs, with a reduction in the intensity of augmentations. Additionally, the loss function undergoes a revision, incorporating the TaskAlignedAssigner from TOOD and introducing Distribution Focal Loss for regression loss, which further enhances the detector’s accuracy and efficiency.

2.4 Network Optimization and Compression

As CNNs have grown deeper and more complex to improve accuracy, the computational and memory requirements have increased significantly. This makes it difficult to deploy CNNs on embedded devices with limited resources. To address this challenge, researchers have developed various model optimization and compression techniques that aim to reduce the computational cost and memory footprint while trying to maintain the model’s accuracy. A brief overview of the most significant model optimization and compression techniques is given in this section. However, in this paper, only one of these techniques—weight quantization will be used.

⁴The head of an object detector predicts the object classes and bounding box coordinates using features input by the neck.

⁵The neck module of an object detector refines features input from the backbone.

2.4.1 Network Pruning

Typically, there are many parameters in a CNN which were not utilized during the training phase and do not contribute to the network's performance. The pruning model optimization technique aims to simply remove these redundant parameters from the model while maintaining accuracy. Of course, more aggressive pruning can be performed, removing even more parameters, including useful but still less important ones, although sacrificing some accuracy.

Various pruning techniques exist, including weight pruning (removal of a single weight), neuron pruning (removal of an entire neuron and its connections), filter pruning and layer pruning.

Although the pruning technique was developed to reduce the model's storage requirements, it is also used to reduce its computational requirements.

[[cite <https://link.springer.com/article/10.1007/s10462-020-09816-7>]]

2.4.2 Knowledge Distillation

Knowledge distillation is a model compression technique in which a smaller, more compact student network is trained to mimic the behavior of a larger, high-performing teacher network (or an ensemble of them) to learn the teacher model's generalization capability. The student network learns from outputs from the teacher network instead of the ground truth labels. For more effective knowledge distillation methods, these outputs often include intermediate feature maps of the teacher network. Generally, the student network cannot achieve accuracy as high as the teacher network, but when performed correctly, it typically achieves higher accuracy than if trained the conventional way.

[[cite 2022 structural KD for object detection]]

[[vymazať alebo rephrase:]] Because this method practically involves training two models, with one being much larger than the other, it requires more training time than training just the student model would. It is however very useful when training and comparing several student networks.

Although this technique was initially considered highly advantageous for addressing the problem proposed in this paper, upon careful analysis, we assumed that the associated complexity would pose significant obstacles, primarily in terms of increased development time. Therefore, this technique was not employed in the current study; however, its exploration is highly recommended for future work.

2.4.3 Weight Quantization

[[todo [6] nasledujúce tri odstavce]]

In convolutional neural networks, weights and biases are stored as 32-bit floating-point numbers (FP32), which provide precision often unnecessary for the CNNs to be accurate. Quantization is the process of reducing the number of bits used to represent these parameters and generally decreases the storage and computational requirements of the network.

Although reducing the precision of the model's parameters decreases its accuracy, the drop in accuracy is typically insignificant when compared to the substantial benefits in storage and computational efficiency gained, which are essential factors when developing a high-performance object detector.

Most commonly, models parameters are quantized to either a 16-bit floating-point (FP16) representation or an 8-bit integer (INT8) representation. Quantizing to FP16 usually

doesn't require any post-quantization steps. However, quantizing to an integer representation, such as INT8, is a different process. Because of the limitations of the integer representation and the distinction from floating-point representations, weight calibration⁶ during the quantization process or even model fine-tuning⁷ after it is recommended to maintain the highest possible prediction accuracy. This, of course, only applies to post-training quantization (PTQ) while several other quantization techniques are available—mainly training the model with weights in the desired representation from the beginning (quantization-aware training - QAT). Furthermore, a more advanced technique called the mixed-precision quantization (MPQ) can be used to quantize parameters in a more nuanced manner by assigning different precisions to individual parameters or parameter groups depending on their sensitivity to numerical errors [33].

While the weight quantization technique can significantly decrease a model's computational requirements, it is crucial to note that not all devices are compatible with every parameter representation. Therefore, it's important to verify whether a specific device supports computations with the desired number representations.

Most deep learning libraries which can be used for model deployment or inference on a target device support weight quantization and offer tutorials on how it's done. These libraries include TensorRT, ONNX Runtime, OpenVINO and many others. Ones relevant to this project will be discussed later in this chapter.

2.5 Evaluation Metrics

Evaluation metrics are crucial for assessing the performance of different object detection models, enabling comparison between different architectures and tracking improvements during training. In this section, we provide a brief explanation and background of the mean Average Precision (mAP) metric, which will be used to evaluate the detectors trained in this project. The information in this section is derived from a blog post on the V7Labs website [30].

Precision and Recall

Precision is a measure of the proportion of true positive detections out of all detections (true positives and false positives), while recall measures the proportion of true positive detections out of all ground truth⁸ objects (true positives and false negatives) in the dataset. Therefore, precision represents how many predictions of the model are correct and recall measures how many of the ground truth objects were predicted by the model.

If the precision and recall values are computed at different confidence score⁹ thresholds, they can be plotted one against the other to obtain a precision-recall curve, which can be used to visually evaluate the overall accuracy of the model.

⁶Weight calibration during weight quantization refers to the process of adjusting the quantized weights to minimize the loss of accuracy that occurs due to the reduction in numerical precision.

⁷Model fine-tuning refers to further training of a model with lower learning rate to refine the model's parameters.

⁸Ground truth refers to the actual, true data used as a reference for comparison with the model's predictions.

⁹Confidence score is a number output from a detector for each detection stating the model's confidence that the detection is correct.

Average Precision

The Average Precision (AP) is a metric widely used to represent the performance of an object detector. It is simply calculated as the area under the precision-recall curve and ranges from 0 to 1, where AP of 1 indicates perfect precision and recall at all thresholds.

Mean Average Precision

To evaluate multi-class object detectors, mean Average Precision (mAP) is used instead of the previously explained Average Precision. Computing the mAP involves finding the AP for each class and calculating the average over the number of classes.

However, for object detection tasks, precision is typically calculated with different thresholds of the Intersection over Union (IoU) metric. The IoU metric measures the overlap between two bounding boxes—between the model’s prediction and the ground truth bounding box, and represents the quality of the alignment between the two boxes. Calculating the IoU simply means dividing the intersection area of the two bounding boxes over their union. When calculating precision, the IoU metric is used to determine whether a predicted bounding box should be considered a true positive (IoU is higher than the defined IoU threshold) or a false positive (IoU is lower than the threshold). For visualization, see Figure 2.2.

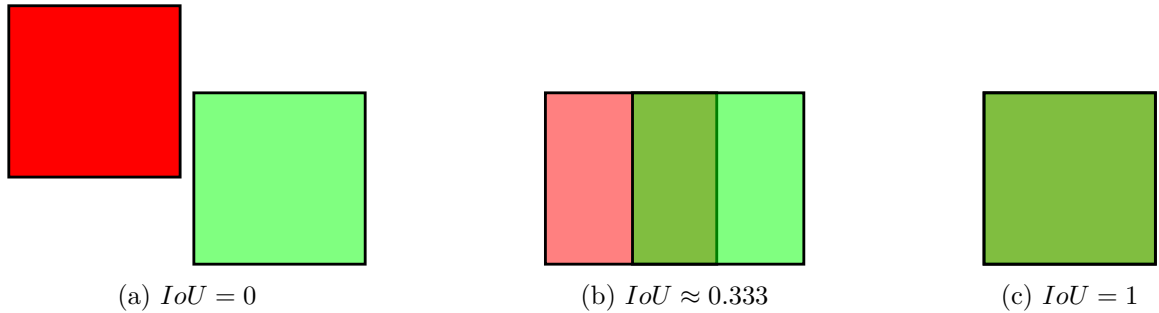


Figure 2.2: Visualization of the Intersection over Union calculation.

In this paper, we consider the COCO mAP specification and calculate the mAP as an average of AP calculated for all classes and over 10 IoU thresholds ranging from 0.50 to (and including) 0.95 with step 0.05. Similarly, values $mAP^{IoU:0.50}$ and $mAP^{IoU:0.75}$ denote mAP calculated with IoU thresholds equal to 0.50 and 0.75 respectively. Additionally, mAP^{small} is only calculated for objects of area smaller than 32^2 pixels, mAP^{medium} for objects larger than 32^2 px but smaller than 96^2 px and finally, mAP^{large} for objects larger than 96^2 px.

2.6 Vehicle Detection Based on Convolutional Neural Networks

[[tu by toho trebalo viac, nie? Ale nemám moc z čoho čerpať.. hlavne ohľadom embedded resp. optimalizovaných sú všetky články prosté zlé.]] [[Ale môžem písať o embedded a optimalizáciach všeobecne pre detektory, nemusí byť pre vozidlá]]

In this section, we discuss existing literature on vehicle detection using deep learning and convolutional neural networks. In our case, papers comparing the performance of different object detectors in the vehicle detection task were important for us to choose the right models to evaluate. Additional research was done on refining these detectors for detecting small vehicles, optimizing them for inference speed or vehicle detection accuracy. Of course, we advise the reader to be cautious when comparing their results because they vary based on the selected evaluation metrics, used datasets and devices used to evaluate the models.

Deep learning has revolutionized vehicle detection by significantly improving accuracy and robustness. In [34], the authors compare five mainstream deep learning-based object detection algorithms for vehicle detection in autonomous driving – Faster R-CNN, R-FCN, SSD, RetinaNet and YOLOv3. The results indicate that two-stage detectors generally have better detection accuracy than one-stage models, while SSD and YOLOv3 algorithms were found to have excellent real-time performance and generalization abilities.

In a more recent paper by Maity et al. [23], the authors present a comprehensive review of existing Faster R-CNN and YOLO-based vehicle detection and tracking methods, comparing them and highlighting their advantages over traditional vehicle detection methods. The authors also discuss several vehicle detection methods featuring modified YOLO models to improve detection performance.

In [38], the authors propose YOLOv5-Ghost – an improved neural network structure based on YOLOv5-small, to detect vehicles in the CARLA¹⁰ virtual environment. By replacing BottleneckCSP in the YOLOv5-small with Ghost Bottleneck, the detection speed is increased from 29 FPS (YOLOv5-small) to 47 FPS (YOLOv5-Ghost) while only reducing the mAP from 83.36 to 80.75 on their test dataset.

[[my sa zameriavame na porovnanie yolov8 na embedded zariadeniach a nič viac k tomu relevantné sme nenašli]]

2.7 Embedded Platforms for Machine Learning

Embedded platforms are a combination of hardware and software components that are designed to perform specific tasks. For object detection, or machine learning in general, these systems are optimized to be power efficient while providing significant processing capabilities for the development, deployment and execution of machine learning algorithms. In this section, we will discuss some of the most popular embedded devices and platforms used for machine learning, including object detection.

These embedded devices typically incorporate custom processors, microcontrollers or specialized accelerators specifically engineered for efficient execution of machine learning tasks, such as graphics processing units (GPUs), which can perform many computations in parallel, tensor processing units (TPUs) and many more.

2.7.1 Google Coral

Designed for TensorFlow Lite¹¹ models, the Google Coral platform offers an Edge TPU – a low-power, high-performance ASIC (application-specific integrated circuit) enabling on-

¹⁰CARLA (Car Learning to Act) is an open-source, realistic urban environment simulator designed for the development, testing and validation of autonomous driving systems in various traffic scenarios and conditions.

¹¹TensorFlow Lite is a lightweight machine learning framework designed for running TensorFlow models efficiently on devices with limited computational resources.

device machine learning inference¹². Coral provides development boards, USB accelerators along with a variety of modules and peripherals for edge AI applications.

2.7.2 Movidius Neural Compute Stick

Intel's Movidius Neural Compute Stick, commonly paired with a popular single-board computer Raspberry Pi, is a small, low-power USB-based hardware accelerator featuring a vision processing unit (VPU) designed to accelerate neural network computations.

2.7.3 NVIDIA Jetson

NVIDIA Jetson is a series of widely-used embedded computing platforms that feature powerful GPU accelerators and ARM-based CPUs while being energy-efficient.

[[Nejakú fotku jetsonov?]]

NVIDIA Jetson AGX Xavier

Jetson AGX Xavier is the flagship model in the NVIDIA Jetson family. It is a high-performance, energy-efficient platform designed for more demanding AI workloads. With an integrated NVIDIA Volta GPU with 512 CUDA cores and 64 Tensor cores, an 8-core NVIDIA Carmel ARM CPU and 16 GB of memory, it offers substantial computational capabilities for deploying state-of-the-art real-time object detectors. Its typical power consumption ranges from 20 W to 30 W.

NVIDIA Jetson Xavier NX

The NVIDIA Jetson Xavier NX features an NVIDIA Volta GPU with 384 CUDA cores and 48 Tensor cores, a 6-core NVIDIA Carmel ARM CPU and 8 GB of memory. Compared to Jetson AGX Xavier, it is more compact and power-efficient, with the typical consumption of 15 W, while of course offering lower performance and memory capacity.

NVIDIA Jetson Nano

The smallest and most popular module from the NVIDIA Jetson family is the Jetson Nano. Equipped with a 128-core NVIDIA Maxwell GPU, a quad-core ARM Cortex-A57 CPU and just 4 GB of memory, it serves as an entry-level, low-cost AI embedded platform, offering lower performance but lower power consumption (ranging from 5 to 10 W) compared to the previously mentioned Jetson boards.

Software for NVIDIA Jetson devices

Devices of the NVIDIA Jetson family support the Linux for Tegra (L4T) operating system, a customized Linux distribution designed specifically for the platform's unique capabilities. NVIDIA also offers the JetPack SDK, which includes the CUDA toolkit and the cuDNN library accelerating deep learning tasks on NVIDIA GPUs, the TensorRT library used for optimizing deep learning models to achieve higher inference speeds on NVIDIA GPUs, along with various multimedia and computer vision libraries.

¹²Inference in machine learning is the process of using a trained model to make predictions on new data, previously unseen by the model.

On all used NVIDIA Jetson platforms, JetPack SDK version 4.6.3 was used. It's important to note that using this version of the JetPack SDK means using older software versions, including Python version 3.6, TensorRT version 8.2.1 and CUDA version 10.2. The reason for not using a more recent JetPack SDK (even though it supports the Jetson AGX Xavier and Jetson Xavier NX), along with the complications it caused, will be discussed in [chapter 5](#).

2.8 Tools and Libraries

In this section, we will provide a brief overview of some crucial tools and libraries utilized throughout the project.

2.8.1 LabelBox annotation app

LabelBox [\[18\]](#) is a web-based application and data labeling platform widely used when training machine models. It provides an easy-to-use interface and a suite of tools to manage and annotate datasets efficiently. With the option of importing data, we have used the app to reannotate one of the datasets.

2.8.2 PyTorch

PyTorch [\[27\]](#) is a popular open-source, Python-based machine learning library designed to provide flexibility, ease of use, and high performance for deep learning applications. It offers a rich ecosystem of tools and libraries for various tasks, including computer vision. With support for GPU acceleration using NVIDIA's CUDA platform, PyTorch enables fast and efficient computation of large models, making it an ideal choice for high-performance deep learning tasks.

2.8.3 MMDetection Library

MMDetection [\[8\]](#) is a popular open-source deep learning toolbox for computer vision tasks, including object detection. Developed by the Multimedia Laboratory at the Chinese University of Hong Kong (OpenMMLab), MMDetection provides a flexible, extensible and modular framework that aims to simplify the process of training and deploying state-of-the-art models for various computer vision tasks, and provides a rich set of tools. Some of the key features of the library include:

- It is based on the PyTorch machine learning library.
- It builds upon and uses other open-source libraries from the OpenMMLab project, such as MMCV^{[13](#)} and MMEEngine^{[14](#)}.
- Library's modular design allows for easy customization and extension of the codebase.

¹³MMCV is a library for computer vision research including building blocks for convolutional neural networks, tools for image processing, transformations and much more.

¹⁴MMEEngine library serves as the training engine for all OpenMMLab codebases, supporting hundreds of algorithms frequently used in deep learning.

- It includes pre-trained models and their configurations which makes training and comparing different models fast and simple, while making it possible to use the transfer learning technique¹⁵ which significantly speeds up the process.

2.8.4 MMYOLO Library

Although the MMDetection library doesn't include the latest detectors of the YOLO family, such as the YOLOv8 detector which will be used in this project, the OpenMMLab project features a different library called MMYOLO [9] – an extension of the MMDetection library, which addresses this issue and focuses solely on detectors of the YOLO family. It contains implementations of YOLO-specific components, such as the CSPDarknet and PANet backbone networks, and YOLO-specific training techniques, like data augmentations or loss functions.

2.8.5 MMDeploy Library

The MMDeploy library [7], which is also a part of the OpenMMLab project, offers useful tools for deploying OpenMMLab models to a wide range of platforms and devices. It enables the conversion of PyTorch models trained with MMDetection or MMYOLO into backend models for execution on target devices. MMDeploy supports various backends, including ONNX, TensorRT, OpenVino, TorchScript and numerous others. In addition to streamlining the deployment process, the library also optimizes the converted models for their target platforms.

2.8.6 ONNX and ONNXRuntime

ONNX (Open Neural Network Exchange) [3] is an open-source project aimed at creating a consistent format for deep learning models. It was initially developed by Facebook and Microsoft, but several other companies and organizations joined later on. The primary goal of ONNX is to enable developers and researchers to easily switch between different machine learning frameworks without having to worry about model compatibility. ONNX defines a standard representation for neural network models, making it possible to train a model in one framework and use it for inference in another.

Additionally, ONNX provides a set of tools and libraries, such as ONNX Runtime [10], which is a high-performance inference engine for ONNX models.

2.8.7 TensorRT

TensorRT [[cite ale nenašiel som citáciu]] is a high-performance deep learning inference optimizer and runtime library developed by NVIDIA. It is designed to accelerate the deployment and inference of models on NVIDIA GPUs for various applications including computer vision.

In addition to optimizing the target models for improved inference performance and reduced memory footprint, TensorRT also supports multiple precision modes, including FP32, FP16 and INT8, allowing developers to choose the best balance between accuracy and performance.

¹⁵Transfer learning is a machine learning technique that leverages knowledge of a pre-trained model to more efficiently train a new model, reducing training time and data requirements.

Chapter 3

Datasets

Big and high-quality datasets are very important when training a CNN-based detector. In this section, used datasets are listed and analyzed. First, the criteria for recognizing an appropriate dataset for our task are explained. Each chosen dataset is then analyzed individually. Finally, a summary of all used datasets is provided along with details on how the dataset was split into training, validation and testing subsets.

[[Niekde popísať rozdiel medzi trénovacím, validačným a testovacím datasetom a na čo treba všetky?]]

3.1 Dataset Criteria

[[Nespísať to priamo do kapitoly? Zbytočný section je toto]]

Here, we briefly explain the most important criteria for selecting datasets to be used in a project like this.

Camera Angle

Since we're building a detector for a camera mounted on infrastructure, it is recommended to use datasets containing surveillance-type images. Datasets containing only images taken from, for example, a car dashboard camera, were therefore disregarded.

Classes

If we don't want to re-label the dataset manually, its classes must be mappable to *our* classes. In this work, 8 object classes are considered:

- Bicycle
- Motorcycle (any two-wheeled motorized vehicle)
- Passenger Car
- Transporter (or a van, pick-up truck, etc.)
- Bus (including a minibus)
- Truck
- Trailer

- Unknown

Many available datasets didn't annotate some of these classes or aggregated some of them into one and were therefore ignored.

Diversity of Images

For the trained detector to generalize well, a dataset needs to contain images with diverse camera angles, lighting conditions, weather, and other factors. For example, a 60 FPS continuous video does not bring much of an advantage.

Dataset Quality

We observed that many datasets contained incorrect annotations or classifications. It is important to check the dataset and either fix faulty annotations, ignore incorrectly annotated images or even disregard the whole dataset.

3.2 Individual Datasets

This section individually analyzes all datasets used in this work. We discuss their sizes, origins, advantages and drawbacks and also provide an example image from each of the datasets.

3.2.1 UA-DETRAC

The DETRAC dataset [36] provided by the University at Albany is the biggest and the most important dataset for this work, originally containing 1 274 055 annotations of 8250 vehicles in 138 252 images. The dataset is provided as frames from 100 video sequences of 25 fps with the resolution being 960×540 pixels. The length of these sequences varies, but is usually between 30 seconds and 90 seconds. The video is of a surveillance type and almost all sequences have a unique point of view, usually from a bridge. Many sequences were recorded in rain or at night and there is no lens flare from cars' headlights. One of the images from this dataset is shown in [Figure 3.1](#).

However, there are several issues with this dataset:

- Bicycles, motorcycles and a few vehicles that cannot be classified are not annotated at all.
- Bounding boxes are often loose and do not fit tightly to the objects.
- Vehicles near the edge of the frame, although fully visible, sometimes lack labels.
- In several sequences, vehicles are tracked and annotated even on frames on which they are fully occluded (by other vehicles or infrastructure).
- Vehicle annotations are inconsistent in relation to masks, as some are labeled even when masked, while others are left unlabeled even when already fully visible outside of a mask. This is likely due to the camera being hand-held and therefore moving while the mask is static throughout the sequence.



Figure 3.1: Example image from the UA-DETRAC dataset.

These problems were fixed by importing the dataset to the LabelBox labeling application, adjusting the masks (while also making them dynamic to compensate for camera movements), annotating or masking bicycles, motorcycles and “unknown” vehicles and finally, carefully repairing individual annotations if needed. Many sequences were hectic or were annotated so poorly that reannotating would be too time-consuming, so only 71 of 100 sequences were reannotated.

The reannotated dataset contains 733 833 annotations in 99 771 images and these classes (the dataset contains no trailers):

- Bicycle
- Motorcycle
- Car
- Bus
- Van - mapped to *transporter*
- Others - mapped to *truck*
- Unknown

3.2.2 Miovision Traffic Camera Dataset

The MIO-TCD (Miovision traffic camera dataset) [1] is another huge and very important dataset. Images are taken at different times of day by thousands of traffic cameras in Canada and the United States. Roughly 79% of all images are of resolution 720×480 pixels, but the image quality seems to be lower. The rest is of resolution 342×228 pixels. Example images of both resolutions are shown in Figure 3.2.



(a) Resolution 720×480

(b) Resolution 342×228

Figure 3.2: Example images from the Miovision dataset.

The dataset consists of two parts: *Classification dataset* and *Localization dataset*. Only the *Train* subset of the *Classification* part is used here because the *Test* subset is not annotated.

The part used contains 351 549 objects of these classes:

- Pedestrian - ignored
- Bicycle
- Motorcycle
- Car
- Pickup truck - mapped to *transporter*
- Work van - mapped to *transporter*
- Bus
- Articulated truck - mapped to *truck*
- Single unit truck - mapped to *truck*
- Non-motorized vehicle - mapped to *trailer*
- Motorized vehicle - mapped to *unknown*

The processed dataset (without pedestrian annotations) contains 344 416 objects in 110 000 images.

3.2.3 AAU RainSnow Traffic Surveillance Dataset

Another important dataset is the *AAU RainSnow* dataset [2]. The authors mounted two synchronized (one RGB and one thermal) cameras on street lamps at seven different Danish intersections to take 5-minute long videos at different lighting and weather conditions - night and day, rain and snow. They then extracted 2200 frames from the videos and annotated them on a pixel-level. Several different types of masks were also created and included in the dataset.

In our work, we only use the annotated frames from the RGB camera, containing 13 297 annotations in 2200 frames (before processing) of resolution 640×480 pixels. See an example in Figure 3.3.



Figure 3.3: Example image from the AAU RainSnow dataset.

The dataset uses these 6 classes:

- Pedestrian
- Bicycle
- Motorbike
- Car
- Bus
- Truck

This introduces a problem - vehicles of our internal class *transporter* (van) don't have their own class in the dataset, but are classified as trucks. However, the dataset is small and it is impossible to perfectly divide transporters and trucks into two classes as there are many different models between which a line cannot be drawn. Ignoring this issue should therefore not cause any problems.

Several other minor problems were found when processing this dataset:

- Frames in groups Egensevej-1, Egensevej-3 and Egensevej-5 are hardly usable because of the low-quality camera and challenging weather and lighting conditions, so they were dismissed.
- Some frames had bounding boxes over the whole frame - this is most certainly an annotation error. These labels were ignored as well.

- The mask for Hadsundvej intersection didn't fully cover the area that should be ignored. This was fixed by simply editing the mask.

After processing, the dataset contains 10 545 objects in 1899 images.

3.2.4 Multi-View Traffic Intersection Dataset

For the MTID dataset, the authors [15] recorded one intersection from two points of view at 30 fps - one camera was mounted on existing infrastructure and one was attached to a hovering drone. The dataset contains 65 299 annotated objects in 5776 frames (equal share of frames for both cameras). An example from this dataset can be seen in Figure 3.4.



Figure 3.4: Example image from the drone subset of the Multi-View Traffic Intersection dataset.

All annotated objects fall into one of four classes:

- Bicycle
- Car
- Bus
- Lorry - mapped to *truck*

This, at first, might not seem like enough, but a closer inspection of the annotated frames reveals that there are no pedestrians or motorcycles. However, similarly to the AAU dataset in subsection 3.2.3, there are transporters in the frames that are classified as trucks. Again, this issue is simply ignored.

When processing this dataset, two other problems were encountered:

- Vehicles that are not on the road are not annotated, so they have to be masked out. This is not as easy for the drone video because the camera is moving, but it is still simple enough.

- Many frames of the drone footage lack some or all labels and have to be ignored. Ranges of images numbers which are ignored: [1, 31], [659, 659], [1001, 1318] and [3301, 3327].

The processed dataset contains 64 979 objects in 5399 frames.

3.2.5 Night and Day Instance Segmented Park Dataset

Another useful dataset is the *NDISPark* [22], which contains images of parked vehicles taken by a camera mounted on infrastructure. Although the annotated part of the dataset only contains 142 frames after processing, there are 3302 objects annotated in total. This still makes it a tiny dataset, but it provides images of vehicles from many different points of view and also contains many occluded vehicles. See Figure 3.5 for an example image from this dataset. Additionally, all frames are 2400 px in width and within 908 px and 1808 px in height.



Figure 3.5: Example image from the NDISPark dataset.

This dataset does not contain any classifications, but luckily, it only contains cars, *transporters* and a few unattached car trailers. All *transporters* and *trailers* were manually classified.

3.2.6 VisDrone Dataset

The VisDrone dataset [39] is very different from all the previous datasets since it doesn't just contain traffic surveillance images. Images are taken by a camera mounted on a drone, from many different points of view. After processing, there are 47 720 annotated objects in 1610 frames. An example is displayed in Figure 3.6.



Figure 3.6: Example image from the VisDrone dataset.

This dataset might be a helpful addition to our datasets as it contains useful negative images (many annotated people and new points of view) and it often captures vehicles from a bird's-eye view.

Additionally, objects in this dataset are classified into 12 classes, which can be easily mapped to our 8 *internal* classes.

3.3 Dataset processing

Before training, all datasets used in this project must be converted to a unified format, object classes need to be mapped to be the same in each dataset. Additionally, some datasets include images with regions that should be masked out, and certain subsets or images of some datasets need to be ignored.

Therefore, a Python script was developed for each dataset. It first loads the annotations from the original format – COCO¹ for AAU RainSnow, MTID and NDISPark, XML for DETRAC and different CSV formats for MIO-TCD and VisDrone. The script then maps the classes to ones shown in section 3.1 and if needed, removes the ignored subsets or images before saving the labels in the COCO format. The processing script for NDISPark dataset also corrects the object classes (since the original dataset annotates all objects with the

¹The widely used COCO annotation format defines how a dataset (its categories, list of images, annotations and other metadata) should be stored in the JSON (JavaScript Object Notation) format as a file.

same class) before saving, and therefore does not modify the original dataset file (class corrections are defined in the script itself).

MMDetection’s middle format was considered as it’s a more efficient alternative to the COCO format, but the COCO format is more popular and is supported by most relevant application, while also being human-readable (the MMDetection’s middle format is saved as a pickle² file), and most importantly, MMYOLO and the most recent version of MMDetection at the time of preparing the datasets (v0.4.0 and v3.0.0rc2) do not seem to fully support datasets in the middle format to be used in all of their dataset wrappers.

Several other Python processing scripts were developed, to apply masks, combine all ground truth files into one and to split the combined ground truth file into training, validation and testing subsets. Additionally, scripts to review individual datasets manually were created - one to visualize a dataset by simply adding bounding boxes (with class labels) to the images and one to convert the visualized images to video (or videos).

A few more scripts were created, of which two are worth mentioning – one for uploading the DETRAC dataset to the LabelBox application for reannotation and one for downloading the reannotated labels.

3.4 Summary of Datasets

In Table 3.1, we compare used datasets on a higher level and show the number of images and instances contained with additional comments. It is clear that the DETRAC dataset amounts for most of our data and might have been enough by itself, but to make this project as successful as possible, every available useful dataset should be used. The images from the selected datasets feature a great variety of lighting and weather conditions, points of view, object scales, occlusions and other relevant factors. Additionally, in Table 3.2 we show the number of annotated object instances per class in all used datasets combined.

3.5 Training, Validation and Testing Dataset Split

We chose to only include the Miovision and DETRAC-UA datasets for validation and testing, because they represent data from a typical traffic surveillance camera the best.

Because only the train subset of the Miovision dataset was used (only this subset contained annotations), the images used for validation and testing are chosen randomly. This, however, should not be a problem since the dataset contains many different camera angles and each image is very unique.

From the DETRAC-UA dataset, two sequences were selected for the validation subset (MVI_40201 and MVI_40244) and two for the test subset (MVI_40204 and MVI_40243). The chosen sequences are very different from the ones in the training subset, which is important for the evaluations to be accurate. However, the MVI_40201 sequence is recorded from the same angle as MVI_40204 and the same applies to sequences MVI_40244 and MVI_40243, so the validation and test subsets are alike, but of course, contain different data. Example images are shown in Figure 3.7.

The validation subset contains a total of 36 969 objects on 7770 images, of which 5500 images are from the Miovision dataset and 2270 from DETRAC-UA. Similarly, the test

²The Python’s pickle format refers to a binary serialization and deserialization method allowing for the conversion of Python objects into a byte stream.

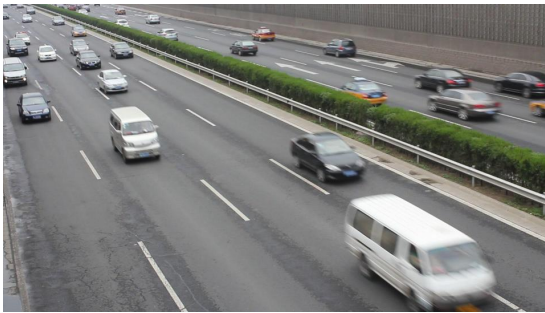
Dataset tag	# images	# instances	Comments
DETRAC	99 771	733 833	Large Continuous video High-quality camera Different lighting conditions
MIO-TCD	110 000	344 416	Large Low-quality images
AAU	1899	10 545	Small Different weather conditions
MTID	5399	64 979	Small Continuous video
NDISPark	142	3302	Small Occlusions
VisDrone	1610	47 720	Small Negative images New points of view
Total	218 821	1 204 795	-

Table 3.1: High-level comparison of used datasets – the number of images and object instances and additional comments for each of the datasets.

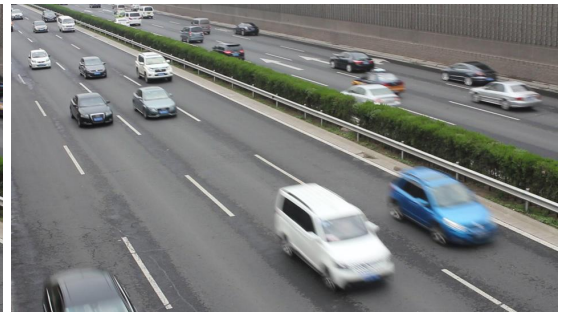
Class	# instances
Bicycle	14 036
Motorcycle	17 187
Passenger car	916 317
Transporter	123 585
Bus	64 529
Truck	38 069
Trailer	2360
Unknown	28 712
Total	1 204 795

Table 3.2: Numbers of class instances in all datasets combined.

subset contains a total of 50 357 objects on 7990 images – 5500 images from the Miovision dataset and 2490 from DETRAC-UA.



(a) MVI_40201 (validation subset)



(b) MVI_40204 (test subset)



(c) MVI_40241 (validation subset)



(d) MVI_40243 (test subset)

Figure 3.7: Example images from sequences from the DETRAC-UA dataset used for the validation and testing subsets.

Chapter 4

Object Detection Models

[[Nedať to k experimentom? Asi to nepotrebuje vlastnú kapitolu, no]]

In this chapter, the reader will find an overview of used object detection models, their training configurations and information about the training process.

4.1 Model Architectures

This section introduces all object detection architectures evaluated in this paper. The main focus is on the YOLOv8 object detector, which is currently the state-of-the-art real-time object detector and offers different model sizes for different applications.

Along with standard model sizes, YOLOv8-medium, YOLOv8-small and YOLOv8-nano, several others were trained and evaluated in this work: YOLOv8-pico and YOLOv8-femto model versions, which are simply smaller versions of the same YOLOv8 model. Finally, a YOLOv8-large model with the CSP Darknet backbone replaced by a popular MobileNetV2 backbone is introduced (hereafter referred to as YOLOv8 MobileNetV2).

In the YOLOv8 MobileNetV2 model, the output indices¹ selected were [2, 4, 6]. Several models with different `out_indices` settings were trained for just 75 epochs to compare their performance, however, none of the alternative settings led to improved performance (in terms of both inference speed and accuracy).

Normally, a square resolution is used for the model's input, but standard traffic cameras output a video of a rectangular shape, usually with the aspect ratio being 16:9. To optimize the inference speed, a rectangular input resolution is used for all trained detectors. The aim was to use aspect ratios as close to 16:9 as possible, but the widths and heights of the input resolution have to be multiples of 32, so some input resolutions are further from it than others.

For smaller models, several input resolutions were tested to provide insights into how a vehicle detection model can be optimized by decreasing the resolution of the input image. Although the inference speed should be (approximately) directly proportional to the number of pixels in the input image, we decided to test it and also find the limit – how small can the model's input resolution be for it to start performing poorly on our datasets.

For a summary of all used model architectures, their sizes, input resolutions and the amounts of their floating point operations (FLOPS) and parameters, see [Table 4.1](#). Quanti-

¹Output indices enable the selection of specific layers in the architecture, from which the output feature maps should be used. In this case, only these feature maps are then used as inputs to the object detector's neck.

ties of floating point operations and parameters were calculated using MMDetection’s analysis script `get_flops.py`. Although the script calculates the output using input resolution 720×480 , simply multiplying the output by the difference between the input resolutions provides an accurate result. This can be explained by the following equation:

$$N_{new} = N_{720 \times 480} \times \frac{new_width \times new_height}{720 \times 480} \quad (4.1)$$

[[napísať tu toho viac?]]

Architecture	Deepen factor	Widen factor	Input resolution	Number of floating point operations	Number of parameters
YOLOv8-medium	0.67	0.75	640×384	7.85 G	18.39 M
YOLOv8-small	0.33	0.5	640×384	2.63 G	3.73 M
YOLOv8 MobileNetV2	1	1	512×288	0.560 G	2.133 M
YOLOv8-nano	0.33	0.25	640×384	0.758 G	1.688 M
			512×288	0.455 G	1.013 M
			448×256	0.354 G	0.788 M
YOLOv8-pico	0.166	0.125	512×288	0.1510 G	0.3067 M
			448×256	0.1175 G	0.2386 M
			384×224	0.0881 G	0.1790 M
YOLOv8-femto	0.166	0.0625	512×288	0.0780 G	0.1288 M
			448×256	0.0607 G	0.1002 M
			384×224	0.0455 G	0.0751 M
			352×192	0.0358 G	0.0591 M

Table 4.1: Summary of different YOLOv8 model architectures used and comparison of amounts of their floating point operations and parameter with various model input resolutions.

4.2 Model Configurations

This section provides an overview of configurations used to train the vehicle detectors. A configuration of an object detection model when using the MMDetection or the MMYOLO library contains a set of parameters that define the model’s behavior and the training, validation and testing pipelines. These parameters can have a significant impact on the model’s speed and accuracy and tuning them is essential to achieve good results.

Most of the YOLOv8 parameters are left unchanged from the default configuration of YOLOv8-m², like:

Optimizer: Stochastic Gradient Descent with momentum 0.937 and weight decay 0.0005

Parameter scheduler: Linear YOLOv5ParamScheduler with learning rate factor of 0.01

However, many parameters related to datasets and augmentations were adjusted and will be explained in the next subsections. Apart from those, the only relevant parameter that was changed is the batch size, which was set to the highest possible for every trained model. For the smallest one, YOLOv8-femto with 352×192 input resolution, the largest batch size of 760 was used. Because training batch sizes above 128 usually result in lower model precision [35], models with large batch sizes were trained for 500 epochs instead of

²The default YOLOv8-m configuration used, `yolov8_m_syncbn_fast_8xb16-500e_coco.py`, can be found at <https://github.com/open-mmlab/mmyolo/tree/v0.4.0/configs/yolov8>

the default 300 epochs. Along with other model-specific training parameters, including the number of warmup epochs³, these settings can be found in [Table 4.2](#).

Architecture	Input resolution	Learning rate	Batch size	Epochs	Warmup epochs
YOLOv8-medium	640×384	0.00125	46	300	5
YOLOv8-small	640×384	0.00125	76	300	5
YOLOv8 MobileNetV2	512×288	0.01	96	300	5
YOLOv8-nano	640×384	0.00125	112	300	5
	512×288	0.00125	192	300	5
	448×256	0.00125	256	300	5
YOLOv8-pico	512×288	0.01	224	500	10
	448×256	0.01	384	500	10
	384×224	0.01	512	500	10
YOLOv8-femto	512×288	0.01	380	500	10
	448×256	0.01	420	500	10
	384×224	0.01	640	500	10
	352×192	0.01	760	500	10

Table 4.2: Training configurations for individual models, including learning rate, batch size, number of training epochs and number of warmup epochs.

4.2.1 Dataset Wrappers

In the MMYOLO (and the MMDetection) model configurations, datasets to use for training, validation and testing are specified using a dataset wrapper, from which a `dataloader` (an object internally representing a dataset) is created. Because the datasets used in this project were in the COCO format, the `YOLOv5CocoDataset` wrapper was used for each of the 6 used datasets.

To compensate for some datasets being smaller than others while being important and of high quality, a dataset wrapper `RepeatDataset` is used, which makes the underlying dataset n-times more frequent when training. All datasets are finally concatenated into one by the `ConcatDataset` wrapper. The repetition factors of individual datasets are shown in [Table 4.3](#).

Dataset name	# images	Repetition factor	# images after over-sampling
DETRAC	99 771	1	99 771
MIO-TCD	110 000	1	110 000
AAU RainSnow	1899	3	5697
MTID	5399	5	26 995
NDISPark	142	25	3550
VisDrone	1610	4	6440

Table 4.3: Repetition factors for each dataset used when training.

³During warmup epochs, the learning rate is gradually increased from a lower value to the target learning rate.

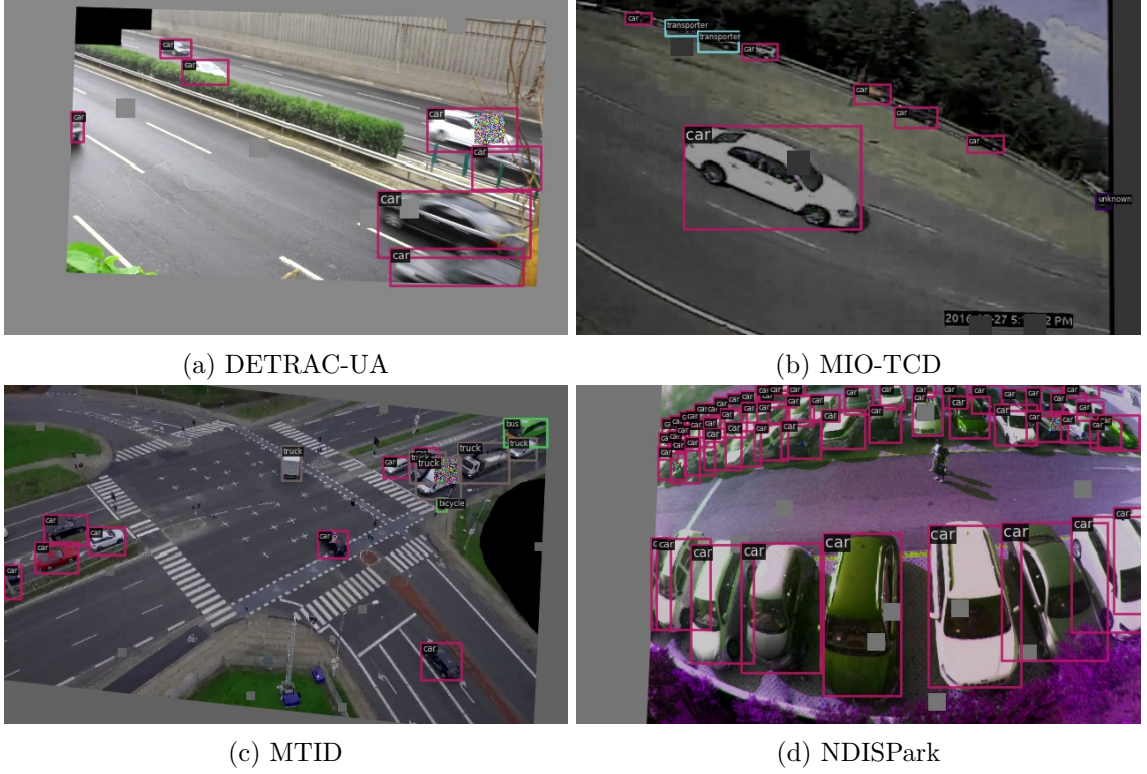


Figure 4.1: Examples of images augmented by the training augmentation pipeline, including the ground truth labels.

4.2.2 Training Augmentation Pipeline

In this subsection, the training augmentation pipeline is explained. Data augmentations are important when training an object detector, especially detectors of the YOLO family. An augmentation in this context refers to the process of applying a transformation to an image to artificially increase the size and diversity of the training dataset, which helps prevent overfitting and improves the generalization ability of the trained model. To a convolutional neural network, even a tiny rotation, translation, image flip, noise or color distortion makes an input image appear to be something completely different, so applying these transformations randomly to the input images is crucial when training a robust model. Examples of images augmented by the training augmentation pipeline can be seen in [Figure 4.1](#). Following are the transformations in the main training augmentation pipeline:

Resize

First, the image is resized to fit the model’s input resolution while, of course, keeping the aspect ratio unchanged.

Pad

If the aspect ratio of the input image is not the same as the model’s input, extra pixels around the input image must be added to adapt to the model input’s aspect ratio. The **color value** of the padded pixels is set to RGB(114, 114, 114).

Random Affine

The `YOLOv5RandomAffine` applies affine transformations to the image, while randomly selecting the values from configured ranges. Parameters:

Maximum translation ratio: 0.05

Maximum rotation degree: 5

Maximum shear degree: 3

Scaling ratio is set individually for each dataset as shown in [Table 4.4](#).

Dataset name	Minimum scaling ratio	Maximum scaling ratio
DETRAC	0.8	1.0
MIO-TCD	1.0	1.1
AAU RainSnow	0.9	1.1
MTID	0.9	2.0
NDISPark	0.9	1.5
VisDrone	1.5	2.5

Table 4.4: Image scaling ratios in the `RandomAffine` transformation for each dataset.

The **color value for padding** around the tranformed image (if needed) is set to RGB(114, 114, 114) to be the same as in the `Pad` transformation.

Cut-Out

The `CutOut` transformation randomly selects regions of the image and fills them with a single color. Again, parameters of this transformation are set individually for each dataset and are shown in [Table 4.5](#).

Dataset name	Number of regions	Size of a single region
DETRAC	6	22×22 px
MIO-TCD	4	26×26 px
AAU RainSnow	8	10×10 px
NDISPark	12	20×20 px
MTID	12	10×10 px
VisDrone	20	8×8 px

Table 4.5: `CutOut` transformation parameters for each dataset.

The **fill color value** is again set to RGB(114, 114, 114), same as for padding in the previous transformations.

Custom Cut-Out

A custom cut-out transformation was developed, similar to `CutOut` used in the previous step. Here, the regions are selected within each bounding box with a certain probability, rather than being chosen randomly within the entire image. Also, the region size is specified as a range of areas in relation to the bounding box area - if the upper value is 10% and a bounding box area is 100 pixels, the maximum cut-out region area can be 10 pixels.

With the boolean option `random_pixels` toggled, the color of each pixel of a cut-out region is generated randomly instead of filling it with a pre-defined color. However, it was found to have no effect.

The **probability** of a region being dropped from each bounding box is set to 0.05 (5%) and the **region area** is set to be randomly selected from interval [5%, 35%].

Albumentations

Albumentations [5] is a popular open-source library for data augmentation. The MMYOLO library provides the option to use its transformations in the data augmentation pipeline. Settings are left unchanged from the original YOLOv8-m configuration:

Blur probability: 0.01

Median blur probability: 0.01

Grayscale probability: 0.01

CLAHE probability: 0.01

HSV Random Augmentations

The YOLOv5HSVRandomAug simply adjusts the hue, saturation and value of the image randomly.

Random Flip

With **probability** of 0.5, the image is horizontally flipped using the **RandomFlip** augmentation.

Photometric Distortion

The **PhotoMetricDistortion** augmentation distorts an image sequentially, while each transformation is applied with a probability of 0.5. It modifies the brightness, contrast, converts color from BGR⁴ to HSV⁵, modifies the saturation, hue, converts from HSV to BGR, modifies the contrast and finally, randomly swaps the color channels.

Filter Annotations

As the last step in the pipeline, **FilterAnnotations** is called to remove bounding boxes with **width or height** lower than 8 pixels.

Fine-Tuning Augmentation Pipeline

Originally, MMYOLO's YOLOv8 models are configured to switch to a simplified augmentation pipeline for the last 10 training epochs. This model fine-tuning strategy is kept and in the augmentation pipeline and for the last 10 training epochs, cut-out augmentations are omitted and affine transformations are changed so that no rotation, translation or shear is applied to a sample.

However, this setting seemed to cause a consistent decrease in the validation mean Average Precision (mAP) metric during the fine-tuning phase (last 10 epochs).

[[že ktorú epochu som vybral pre každý model?? Aj graf validačného mAP? To asi nie]]

⁴BGR (Blue, Green, Red) only differs from the well known RGB format by the order of its channels. It is widely used in image processing for legacy and compatibility reasons.

⁵HSV (Hue, Saturation, Value) is a cylindrical-coordinate color model, where hue represents the color type, saturation refers to the color's intensity and value corresponds to brightness.

Chapter 5

Experiments

5.1 Devices Used in Experiments

While the main focus is on NVIDIA Jetson embedded devices, which were designed specifically for tasks like object detection, tests were run on several other devices for comparison of both ends of the performance gauge. In this section, details about each device that the models were evaluated on are provided, including details about the software and device configurations.

5.1.1 NVIDIA Jetson Platforms

Technical details of NVIDIA Jetson embedded platforms were discussed in [subsection 2.7.3](#) and software versions will be shown later in this section. However, one more important detail to note before reading about the experiments is the used power plan. All used Jetson devices feature several power plans to choose from to adjust the performance and power consumption for a specific task. For all experiments, these power plans were chosen:

NVIDIA Jetson AGX Xavier: 30 W power plan with 4 out of 8 cores running

NVIDIA Jetson Xavier NX: 20 W power plan with 4 out of 8 cores running

NVIDIA Jetson Nano: 10 W MAXN power plan with all 4 cores running

5.1.2 NVIDIA GeForce MX150

To be able to compare inference speeds on these embedded devices to ones on a regular GPU, tests were also done on an NVIDIA GeForce MX150 GPU with 2 GB of memory on a DELL Latitude 5401 laptop. Because of the GPU memory constraint, not all models can be evaluated with all inference batch sizes on this device, as will be pointed out in individual relevant experiments. Additionally, we were unable to perform any tests using the TensorRT library on this device because some of the packages required could not be installed on the system.

5.1.3 Intel Core i7-9850H

Models were also evaluated on a higher-end laptop CPU, Intel Core i7-9850H with the base frequency 2.6 GHz to demonstrate how the inference speeds of YOLOv8 object detection models differ between a GPU and a CPU.

Typically, the frequency at which a CPU operates is adjusted to accommodate the load and it often spikes up when a computation-hungry process starts. After a while, when

Name	All NVIDIA Jetson devices	NVIDIA MX150	Intel Core i7-9850H	Raspberry PI 4B
Operating system	Linux for Tegra (L4T OS 32.7.3)	Linux Debian 12		Linux Raspbian 11
JetPack SDK	4.6.3	-		
Python	3.6.9	3.10.0		3.7.0
CUDA	10.2	11.8	-	
TensorRT	8.2.1	-		
ONNX	1.13.1			1.12.0
ONNX Runtime	1.11.0 (ver. GPU)	1.12.0 (ver. GPU)	1.12.0	1.11.0
PyTorch	1.10.0	2.0.0		1.8.0
MMCV	2.0.0			
MMDeploy	1.0.0			1.0.0rc3
MMDetection	3.0.0			
MMEngine	0.7.2			
MMYOLO	0.4.0			

Table 5.1: Versions of relevant software installed on devices used to deploy and test the trained models.

the CPU temperature rises above a certain threshold, the CPU frequency has to drop to avoid overheating. This is called dynamic frequency scaling, also known as CPU throttling. To make the performance measurements as accurate as possible, the number of warmup samples when testing is increased from 10 to 100. This means that the inference speed (FPS) of these first 100 samples will be ignored when calculating the average FPS.

However, please note that the test results still might not be accurate enough and might depend on the underlying operating system, running applications and the environment.

5.1.4 Raspberry PI 4B

Finally, the popular Raspberry PI 4B single-board computer with ARM Cortex-A72 CPU with base frequency 1.8 GHz was used to test the trained models. Although it was not developed to run object detection models, it is perfect to test the smallest models and see how far go the possibilities of the real-time YOLOv8 object detector.

5.1.5 Software versions

Information about software installed on all previously mentioned devices can be found here. See [Table 5.1](#) for a compact table displaying versions of relevant software. Following are the reasons behind some odd choices or compatibility issues.

JetPack SDK

Despite both Jetson Xavier NX and Jetson AGX Xavier being supported by NVIDIA JetPack SDK version 5.1.1, deploying YOLOv8 models using this version often led to an untraceable fatal error. The error originated from one of NVIDIA’s proprietary libraries and provided limited information regarding its cause. Fortunately, downgrading the JetPack version to 4.6.3 resolved this issue.

Although the root cause of the error and the specific package (or library) that required downgrading were not fully determined, it is suspected that the problem was related to

TensorRT version 8.5.2, as the error originated from the TensorRT development library `libnvinfer`.

For the sake of providing further context to the reader, the error message received was `operation.cpp:203: DCHECK(!i->is_use_only()) failed`. No other relevant warnings or error messages preceded this one, making it challenging to pinpoint the exact cause.

Old Python Version on Jetson Devices

All NVIDIA Jetson devices used to evaluate the trained models utilize the same version of the JetPack SDK, 4.6.3, which includes Python version 3.6.9. However, several essential Python packages—specifically, `protobuf` version 3.20.2, `MMCV` version 2.0.0 and `MMEEngine` version 0.7.2—require a Python version of 3.7.0 or higher. As these packages (and the specified versions) were crucial for the model deployment and testing to be possible, we had to manually modify their requirements to allow for installation with Python version 3.6.9.

Of course, this approach is not an ideal solution to the problem and its success was not guaranteed. Fortunately, no indications of compatibility issues were discovered during model deployment or testing.

One might suggest that upgrading to a newer Python version would be the most appropriate solution. However, due to compatibility and dependency constraints on Jetson devices, this is not feasible.

JetPack version 5.1.1 includes a more recent Python version but does not support Jetson Nano, and when installed on Jetson Xavier NX or Jetson AGX Xavier, the deployment of the trained models fails, as explained earlier in this section. Although installing a different Python version than the one provided with the JetPack SDK is possible, installing other necessary packages for the newer Python is not. This is because many such packages were developed specifically for Jetson devices and only support a certain Python version—the one pre-installed with the JetPack SDK.

5.2 Model Deployment and Optimizations

[[niekde napísať ktoré epochy som nakoniec vybral pre evaluation a že to je čisto na základe validačných mAP?]]

For model deployment, we have created an automated script utilizing the `MMDeploy` library's deployment script. Our script uses all available deploy configurations¹ to deploy all trained PyTorch models to a specified backend (or backends).

There are only two deploy configurations for the ONNX Runtime backend—one with a static model shape² and one with a dynamic shape in the batch size dimension³. The dynamic models were deployed to accept a maximum batch size of 32, but were optimized for a batch size of 1. For the TensorRT backend, two additional deploy configurations were created, both for a model with a dynamic shape in the batch dimension—one for weight quantization to the `FP16` representation, and one for weight quantization to the `INT8` representation including weight calibration using the validation dataset.

¹When deploying models using `MMDeploy`, deploy configuration files are used to specify the parameters of the deployment process, including target backend or whether to apply post-training quantization

²A model shape can be static or dynamic. A static model can only receive inputs of a certain shape specified when deploying (batch size and input resolution—width and height), whereas a dynamic model accepts various input shapes.

³A dynamic model shape in the batch size dimension means the model accepts different input shapes but still requires a fixed input image resolution.

To optimize inference on the Raspberry PI 4B, we aimed to use the NCNN inference framework designed for mobile and embedded devices with limited computing resources. However, deploying a YOLOv8 model to the NCNN format is not yet possible, as the YOLOv8 model contains operations that are not yet supported by the framework and the MMYOLO library does not yet support converting the unsupported operations to supported ones.

Devices Used for Deployment

The models in the ONNX format were all deployed on a single device (with ONNX package version 1.13.1) and distributed to all devices. The TensorRT backend is only used on NVIDIA Jetson devices and all models were deployed to the TensorRT engine⁴ individually on each NVIDIA Jetson device, because they are platform-specific and transferring them across different devices is not recommended.

Weight Quantization

The MMDeploy library supports post-training quantization to FP16 and INT8 representations during the process of model deployment to TensorRT. To preserve the model’s accuracy after quantization to INT8 precision, weight calibration was done using the validation dataset.

Although the ONNX Runtime framework also supports quantization (to both FP16 and INT8), the process is not as straightforward and doesn’t seem to be supported by the MMDeploy library. Although evaluating models quantized to INT8 representation would be beneficial, the quantization to the FP16 precision would probably have a little effect on performance on CPUs as they do not usually support operations with this precision and calculate them using the same operators as numbers in the FP32 representation.

Additionally, TensorRT on the NVIDIA Jetson Nano with Maxwell GPU doesn’t support fast inference of models quantized to INT8, so quantization to INT8 won’t be performed on this device.

5.3 Experiments and Evaluation

In this section, we will discuss additional details of all performed experiments – mainly the combinations of models, devices, inference backends, deployment configurations, optimization techniques and batch sizes of individual tests. Insights gained from these experiments and additional experiment-specific details will be discussed in [chapter 6](#).

All experiments will come from data collected by testing. Testing was done using a Python script included as a tool in the MMDeploy library, `test.py`, which uses the library to create a wrapper for a deployed model. It then uses the test dataset provided to evaluate the model’s performance, calculating the inference speed in FPS and mAP metrics. The inference speed is calculated including the input pre-processing and non-maximum suppression⁵ (NMS).

To test all possible combinations of models, backends, deployment configurations and batch sizes on individual devices automatically, a new Python script `test_all.py` was

⁴A TensorRT engine is a deployed model in the TensorRT format.

⁵Non-maximum suppression in object detection is a method that selects a single prediction bounding box out of several overlapping bounding boxes, which are likely generated for the same object.

developed, which uses the MMDeploy's `test.py` script and executes it for each possible test combination (unless specified differently by the user).

The MMDeploy's `test.py` saves the test logs, from which the mAP metrics and inference speed can be read. This is also done automatically by our `collect_test_results.py` Python script, which reads all available test logs to output all metrics structured in JSON format to a file. Although the `test.py` script also outputs the final test metrics in JSON format to a folder, the name of the folder cannot be specified and is named after the test start time (eg. `20230429_043829/`), so reading the data from the log file is simpler and less error-prone.

When testing models on CPUs in the ONNX format using the MMDeploy library, the library doesn't control how many processes the ONNX Runtime library uses to run inference. To provide accurate results, the source code was modified to create the ONNX Runtime inference session with the option of only running in a single thread. This means that the inference speeds on CPUs are not the highest possible. However, simply multiplying the FPS metric by the number of CPU cores to get the maximum inference speed possible on the device is generally incorrect, because running a multi-threaded inference is not as efficient as running it on a single thread. On the other hand, if running n separate inference processes on a CPU with n cores, it is theoretically possible to achieve inference speeds n times higher than with a single-threaded inference, although the inference duration would remain the same.

[[Možno tabuľka všetkých testov alebo tak? Alebo, ako na začiatku sľubujem, napísať že testujem všetky modely atď.]]

Chapter 6

Results

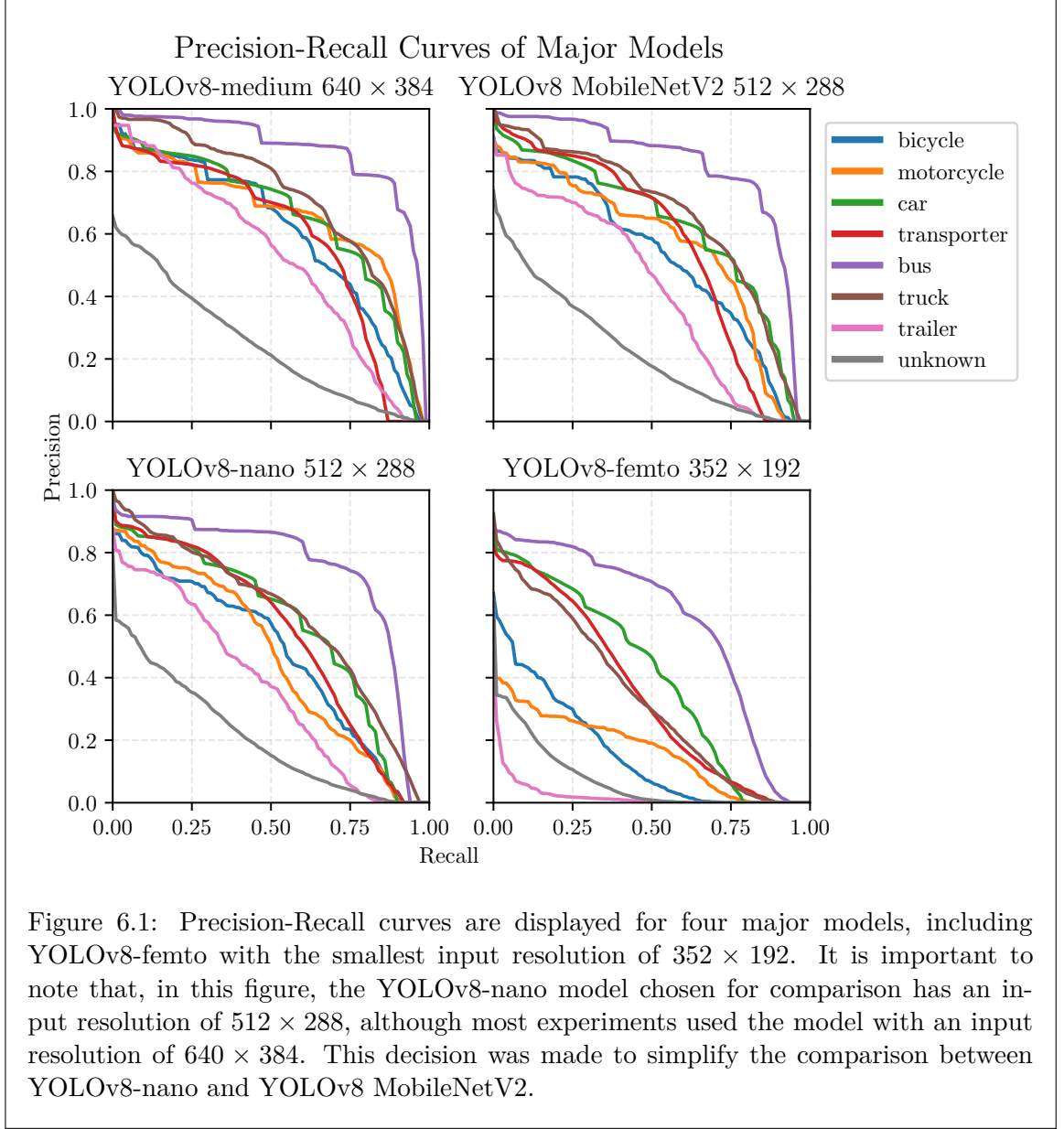
povedať krátko že v akom poradí to budem evaluovať a napr. že graf FPS/mAP bude až na konci alebo tak

6.1 Precision-Recall Curves of Major Models

In this experiment, we measured the precision and recall values of four of the trained models on the test dataset to plot the precision-recall (PR) curves in [Figure 6.1](#). We could not show the curves of all the models because the figure would simply not fit on a single page, so instead, only the four most representative models were selected. A more complete figure with PR curves for the rest of the models can be seen in [Appendix A](#). The precision and recall values were calculated by MMDeploy’s test script, to which we have inserted a few lines of code to save the values to a file.

In the plots, the performance for the bus class appears abnormally high. We hypothesize that the reason behind it is the presence of a few easily-detectable, large buses in numerous images from the DETRAC-UA dataset within the testing set. Additionally, the low performance of models in detecting vehicles belonging to the *unknown* category should not be a major concern, as the diversity of the vehicles in this category makes it difficult to accurately detect and classify them.

As could be expected, the YOLOv8-medium model demonstrated superior performance, achieving the highest precision and recall values. Furthermore, the performance of the YOLOv8 MobileNetV2 model is comparable to that of the YOLOv8-nano, although slightly superior. In contrast, the YOLOv8-femto with an input resolution of 352×192 could be an example of a model with considerably lower performance.



6.2 Effect of Model's Input Resolution on Inference Speed

Smaller models – YOLOv8-nano, YOLOv8-pico and YOLOv8-femto were trained with different input resolutions to compare how the image resolution affects the model's performance on our test dataset.

We expect that the reduction in the number of pixels in a model's input resolution should be directly proportional to the increase in its inference speed (FPS). For example, we anticipate that a detector with the input resolution of 640×320 would be twice as fast compared to a detector with the input resolution 640×640 . However, this hypothesis needed to be tested to bIn this experiment, we measured the precision and recall values of four of the trained models on the test dataset to plot the precision-recall (PR) curves

in Figure 6.1. We could not show the curves of all the models because the figure would simply not fit on a single page, so instead, only the four most representative models were selected. A more complete figure with PR curves for the rest of the models can be seen in Appendix A.e proven correct In Figure 6.2, we compare the inference speed of models with different input resolutions. The biggest input resolution is selected as the base to which other input resolutions will be compared to, and expected FPS is calculated using this simple equation:

$$\text{Expected FPS} = \text{Base FPS} \times \frac{\text{Base input width} \times \text{Base input height}}{\text{Compared input width} \times \text{Compared input height}} \quad (6.1)$$

While our experiments showed that the relationship between input resolution and FPS cannot be exactly captured by the above equation, the results generally supported our hypothesis. Following this experiment, we will suppose that the input resolution reduction is indeed directly proportional to the inference speed increase, and in some of the following experiments, we will only compare the biggest input resolutions of each model to make the results more informative and concise.

6.3 Inference Speeds: ONNX Runtime vs. TensorRT

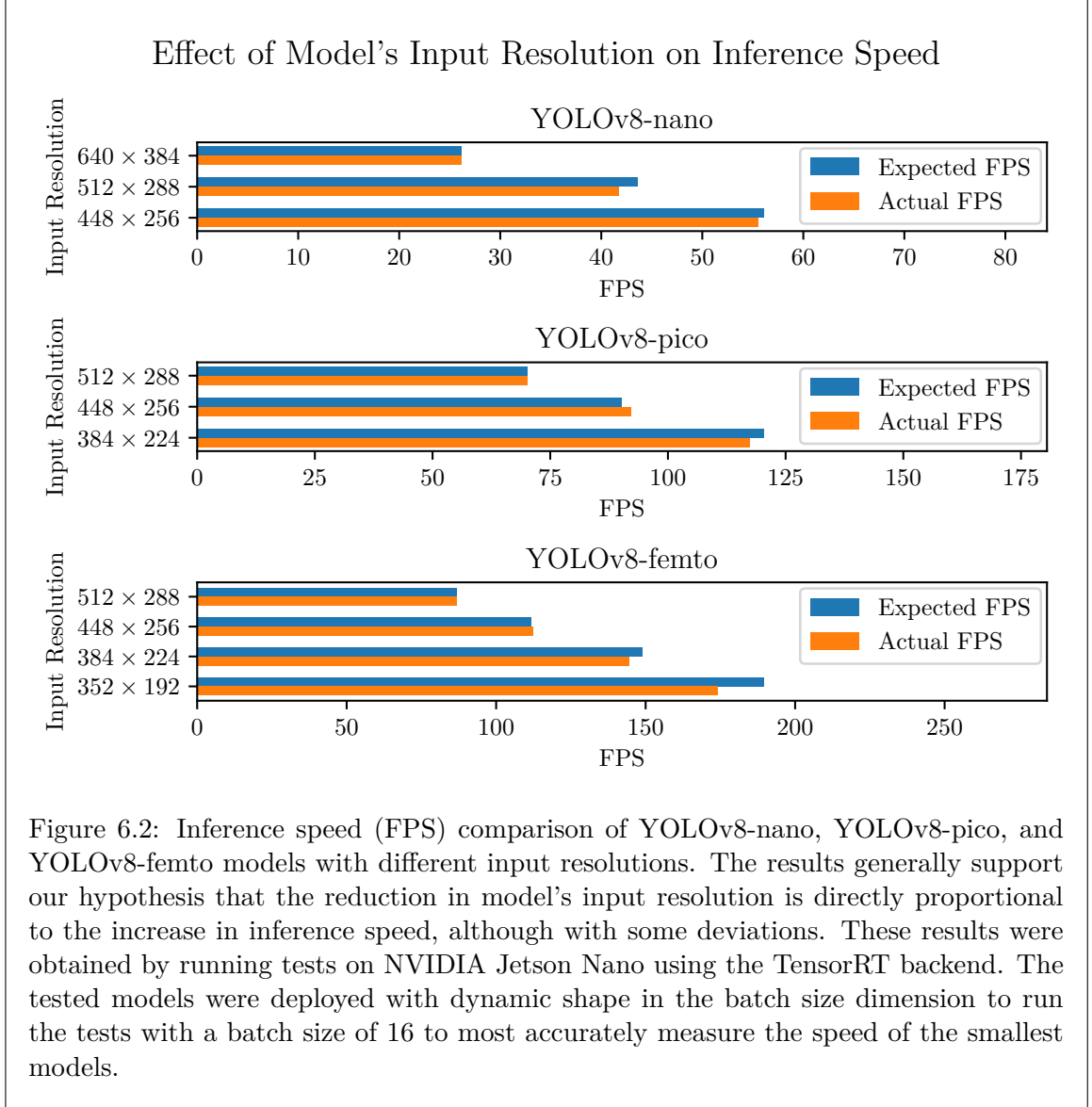
In this experiment, we compare how the inference speed is increased by using the TensorRT backend for inference on NVIDIA Jetson devices instead of the ONNX Runtime backend. The inference speeds (FPS) measured by our test script are compared on both the lowest and the highest-performing NVIDIA Jetson devices with different batch sizes to be able to draw general conclusions from the results. For each model architecture, only the biggest input resolution is selected for testing, following the insights gained by the experiment in section 6.2. Additionally, all tested models were dynamic in shape.

The results in Figure 6.3 clearly show the superiority of the TensorRT inference backend on NVIDIA Jetson devices for all tested models and different batch sizes. Naturally, the choice of the inference backend has no effect on the mean Average Precision of the tested models. Therefore, none of the following experiments will feature the ONNX Runtime backend on NVIDIA Jetson devices.

6.4 Effect of Model’s Shape on Inference Speed

While a model with a static shape has a fixed input resolution and a fixed batch size, it can be set to accept input images of different resolutions and different batch sizes. In this paper, we used the MMDeploy library to deploy the trained models into both static and dynamic shapes. However, the deployment process was configured so that dynamic models would only be flexible in the batch size dimension while the model’s input resolution stays static. Anyways, we will call these models dynamic.

In this experiment, we compare inference speeds achieved by static models and dynamic models to learn whether deploying a model to a dynamic shape (in the batch size dimension) decreases the inference speed. Only the most representative models were selected for comparison to keep the plots concise. Naturally, all tests were conducted with a batch size of a single image. Tests were conducted on the NVIDIA Jetson Nano with the TensorRT inference backend and on the Raspberry PI 4B with ONNX Runtime to make the measurements as accurate as possible by using the least-performing devices for both inference



backends. The comparison of inference speeds on all model architectures is presented in [Figure 6.4](#).

Because the inference speed is not consistently faster for static models, we conclude that it does not significantly decrease if a dynamic shape is used for the model's input in the batch size dimension. Therefore, we will only evaluate dynamic models in future experiments.

6.5 Effect of Batch Size on Inference Speed

An object detector can achieve higher inference speeds (FPS) if a larger batch size is used – if the detector is given more than one frame in a single input. Although rarely suitable for real-time object detection, a larger batch size enables more efficient utilization of GPU resources at the cost of higher memory usage. In this experiment, we will discuss how

the inference speed increases with larger batch sizes used when performing inference with different models and on different devices.

In [Figure 6.5](#) (please note the logarithmic horizontal axis), results of tests on Raspberry PI 4B, NVIDIA Jetson Nano and NVIDIA Jetson AGX Xavier are shown. Although tests were run with six different batch sizes, to make the figure more compact, only batch sizes 1, 2, 8 and 32 were selected for comparison.

The results show that using a larger batch size for inference does indeed often increase the inference speed, but the increase is generally only significant when performing inference on smaller models, where the inference speed is higher. This is well illustrated by the NVIDIA Jetson AGX Xavier, where the increase in inference speed with larger batch sizes is relatively subtle for the largest detector, YOLOv8-medium 640×384 . In contrast, for the smallest model, YOLOv8-femto 352×192 , raising the batch size from 1 to 32 results in an almost ten-fold increase in inference speed. Therefore, using larger batch sizes is generally only suitable in cases where high inference speeds (higher than real-time) are desired, for example in a multi-camera real-time vehicle detection system.

6.6 Mean Average Precision and Inference Speed Benchmark

In this experiment, we compare the trained and models in terms of their mAP metric (mean Average Precision) and their inference speed in FPS.

For the tests, we wanted to choose two least-performing devices, with the Raspberry PI 4B being an obvious choice. Additionally, the NVIDIA Jetson Nano would be a great choice, but since the device is not optimized for INT8 precision and the models were therefore not quantized on this device, NVIDIA Jetson Xavier NX was selected instead. For the tests to be as representative as possible even for the smallest models, tests on Jetson Xavier NX were run with a large batch size of 32, while on the Raspberry PI, where the batch size doesn't really matter, a batch size of 1 was used instead. Naturally, ONNX Runtime backend was used on the Raspberry PI and TensorRT on the Jetson Xavier NX.

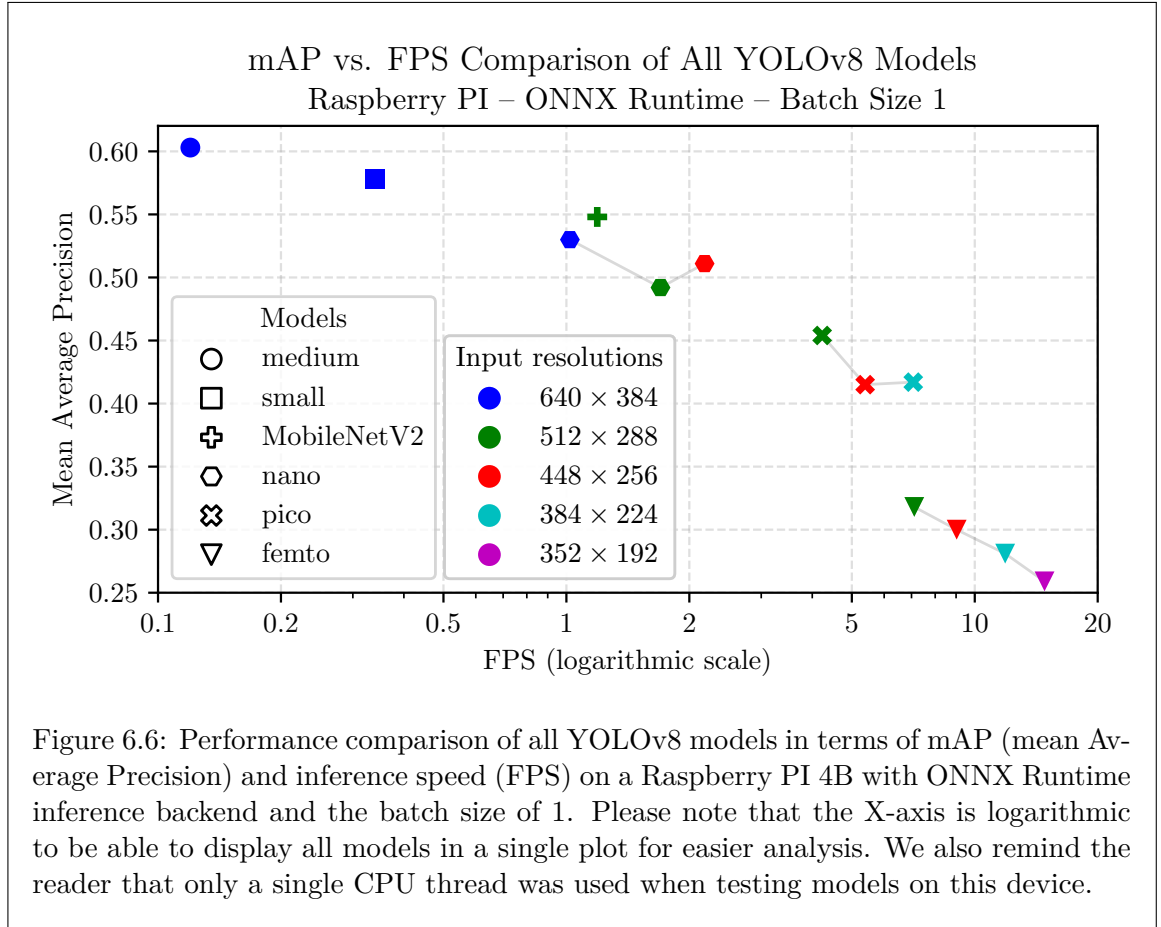
We therefore present two plots, comparing the mAP and FPS of all models. On the Raspberry PI, ONNX Runtime backend was used to test the models on the CPU. For the NVIDIA Jetson Xavier NX with TensorRT backend, the plot also features results of models quantized to both FP16 and INT8 precisions and individual models with different precisions were connected in the plot. The mAP and FPS comparison for the Raspberry PI can be found in [Figure 6.6](#), while the corresponding results for the NVIDIA Jetson Xavier NX are shown in [Figure 6.7](#). Please note that in both plots, a logarithmic scale was chosen for the X-axis (FPS) to include all models (large and tiny) in the figure. For a complete benchmark of all trained and quantized models on all devices, see [Appendix B](#).

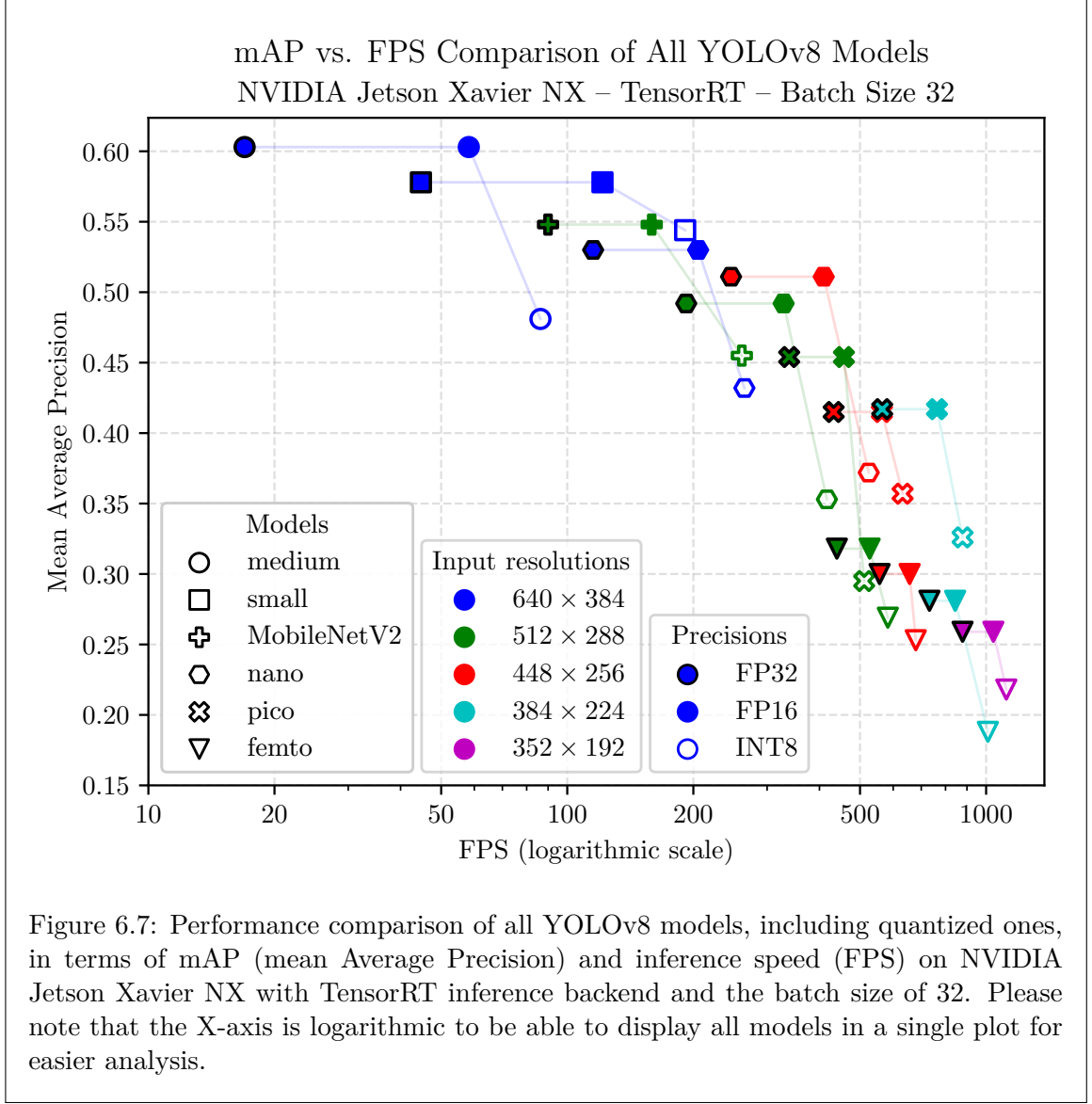
An important thing to notice is the effect of quantization on the mAP and FPS. Although quantization to INT8 precision seems to rarely result in better performance than training a smaller model, it is clear that quantizing to FP32 results in a significantly higher inference speed with a little or no decrease in the mAP.

The performance of the YOLOv8 MobileNetV2 model can also be discussed here, as tests on the Raspberry PI show the superiority of the YOLOv8 model with the MobileNetV2 backbone compared to the YOLOv8-nano models. However, tests conducted on the NVIDIA Jetson Xavier NX with the TensorRT backend show that the opposite is true and the MobileNetV2 brings little to no advantage when compared to YOLOv8-small or YOLOv8-nano models on this device. To draw further conclusions about this matter, the reader is encouraged to look at the tables in [Appendix B](#), which present the complete mAP

and FPS values. Additionally, all 6 COCO mAP metrics of different models are compared in the next experiment.

In both plots resulting from this experiment, a strange phenomenon can be seen in which the mean Average Precision of a model with a higher input resolution – YOLOv8-nano 512×288 is lower than the one of a model with a lower input resolution of 448×256 . The same seems to apply for YOLOv8-pico models, where the mAP with input resolution 448×256 is lower (although not as significantly) than with 384×224 , instead of being higher. However, this problem doesn't seem to affect the YOLOv8-femto models, of which the mAP and FPS values are as expected. We tried examining the problem, mainly by double-checking the results and training configurations of these models, but could not find the root cause and did not have enough resources to conduct further experiments to provide more insights into this phenomenon.





6.7 Evaluating Models Across Multiple Mean Average Precision Metrics

In other experiments, we compare the performance of models just by using a single metric – the mean Average Precision (mAP) for objects of all sizes and over ten IoU thresholds ranging from 0.50 to 0.95. However, models can also be compared using five additional, more specific mAP metrics as explained in [section 2.5](#) – mAP with the IoU threshold equal to 0.50, with IoU threshold 0.75, and mAP for small, medium-sized or large objects.

In [Table 6.1](#), we display these metrics for all trained models. However, we omit the quantized models from the table because the mAP values tend to stay the same after reducing precision to FP16 and quantizing to INT8 is rarely beneficial, as illustrated by the experiment in [section 6.6](#). For the complete table featuring all precisions of all models, see [Appendix C](#).

Model (YOLOv8)	Input resolution	mAP					
		mAP	IoU:50	IoU:75	small	medium	large
medium	640×384	0.603	0.805	0.695	0.316	0.606	0.770
small	640×384	0.578	0.786	0.660	0.273	0.585	0.744
MobileNetV2	512×288	0.548	0.759	0.618	0.235	0.550	0.715
nano	640×384	0.530	0.749	0.608	0.229	0.529	0.685
	512×288	0.492	0.699	0.569	0.210	0.483	0.673
	448×256	0.511	0.723	0.592	0.206	0.508	0.681
pico	512×288	0.454	0.664	0.502	0.151	0.448	0.614
	448×256	0.415	0.621	0.468	0.136	0.408	0.552
	384×224	0.417	0.615	0.465	0.128	0.405	0.589
femto	512×288	0.318	0.488	0.351	0.072	0.298	0.416
	448×256	0.300	0.477	0.318	0.070	0.259	0.429
	384×224	0.281	0.452	0.302	0.070	0.263	0.386
	352×192	0.259	0.430	0.275	0.054	0.227	0.370

Table 6.1: Comparison of multiple mean Average Precision metrics for each trained model and input resolution. The mAP metrics displayed in this table are explained in [section 2.5](#).

In the results, we can observe that smaller models detect large and even medium-sized objects reasonably well. However, when it comes to detecting small objects (of area lower than 32×32 px), the performance of these models drops significantly, even at higher input resolutions. However, the smaller models might still be highly suitable for cases in which detecting small objects is not a priority, as they generally perform well for larger objects while being dramatically faster.

6.8 Inference Speeds on Different Devices

In this experiment, we benchmark all available devices by measuring the inference speeds of all model architectures. For each architecture, only one input resolution was selected, however, differently than in the other experiments: YOLOv8-nano 512×288 was chosen to better compare with YOLOv8 MobileNetV2 and YOLOv8-femto 352×192 to include the smallest model. A second set of measurements was created to only benchmark models quantized to FP16 precision on the NVIDIA Jetson devices. Both plots can be seen in [Figure 6.8](#). Please note that the X-axis is in a logarithmic scale in both plots because of the performance differences between individual devices.

We decided to measure the inference speed with higher batch sizes to get the highest possible inference speeds possible, but since larger models do not fit into memory on lower-performing devices when a large batch size was used¹, we could not benchmark these models with the largest batch size of 32. However, thanks to insights provided by [section 6.5](#), we consider it safe to use a smaller batch size of 8 for YOLOv8-medium, YOLOv8-small and YOLOv8 MobileNetV2 on all devices. For the rest of the models, a batch size of 32 was used. This applies to both plots resulting from this experiment.

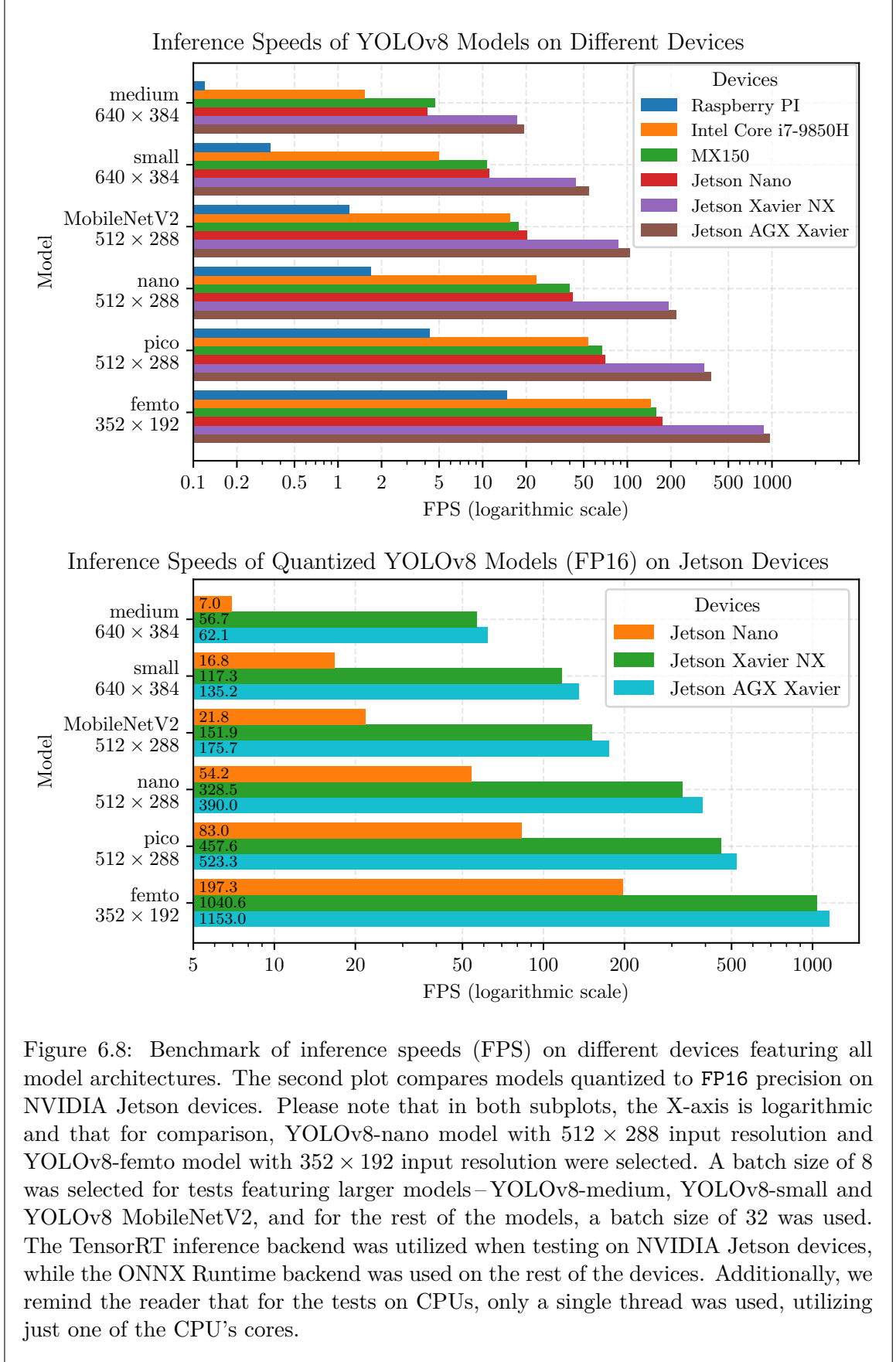
What’s interesting to see from the first plot is that the NVIDIA Jetson Nano is comparable to the NVIDIA MX150, while consuming just a fraction of the power. Additionally,

¹List of tests that failed because the model could not fit into the memory of a device when performing inference with a specific batch size: YOLOv8-medium on NVIDIA Jetson Nano with batch size 32, YOLOv8-medium on NVIDIA MX150 with batch sizes 16 and 32, YOLOv8-small on NVIDIA MX150 with batch size 32.

there is a significant difference between the performance of the Jetson Nano compared to Jetson Xavier NX, while the inference speed measured on Jetson AGX Xavier is only slightly higher than on Jetson Xavier NX.

Another notable difference lies in how weight quantization to FP16 increases the inference speeds on NVIDIA Jetson devices. The highest difference in FPS between precisions FP32 and FP16 on Jetson nano was for the YOLOv8-small model, where the inference speed after quantization increased by 50.13 % from 11.19 FPS to 16.79 FPS. On the other hand, the inference speed for the same model on Jetson Xavier NX increased by 165.84 % from 44.12 FPS to 117.29 FPS. Furthermore, the difference is even higher for YOLOv8-medium, where the quantization to FP16 resulted in a 228.33 % increase in inference speed. However, it's important to note that for smaller models, these gains are lower and more comparable between the two devices.

For a complete table of inference speeds of all models (including all quantized models) on all devices with a batch size of 1, please see [Table B.1](#) in [Appendix B](#).



6.9 TODO Confusion matrix?

Asi netreba keď mám PR krivky. Možno v prílohách ak bude čas

6.10 TODO DETRAC vs MIO-TCD?

NEDÔLEŽITÉ!

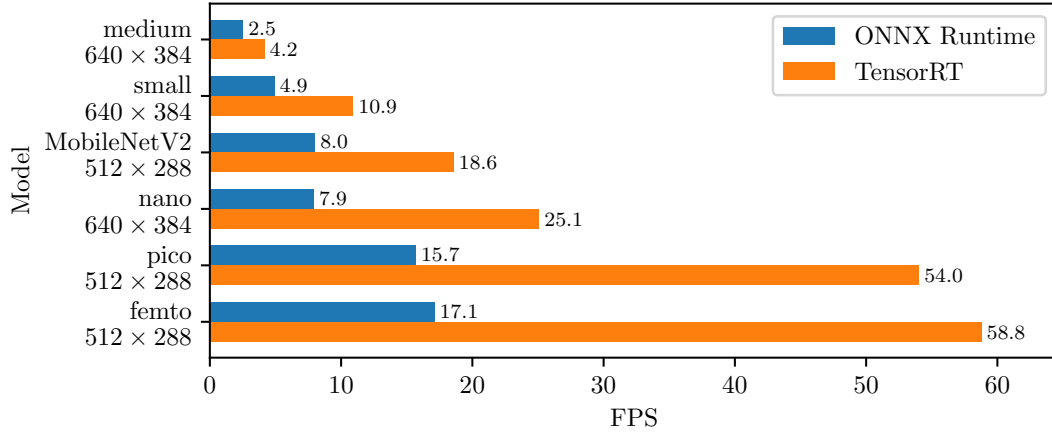
porovnať mAP niektorých, možno všetkých detektorov, medzi dvoma test subsetmi. Nedôležité!

Že detrac je ľahký do istej miery ale pri rohoch je anotovaný inak ako je detektor schopný sa naučiť; mio-tcd má veľa objektov náročných

ako graf, vľavo modely a resolution, x os mAP a pre každý model tri bary - all, detrac a mio-tcd

Inference Speeds: ONNX Runtime vs. TensorRT on YOLOv8 Models

NVIDIA Jetson Nano – Batch Size 1



NVIDIA Jetson AGX Xavier – Batch Size 32

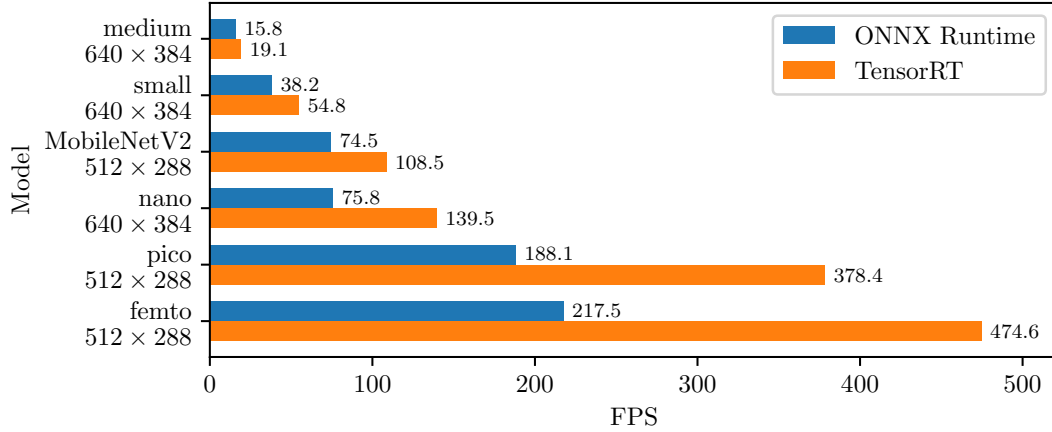
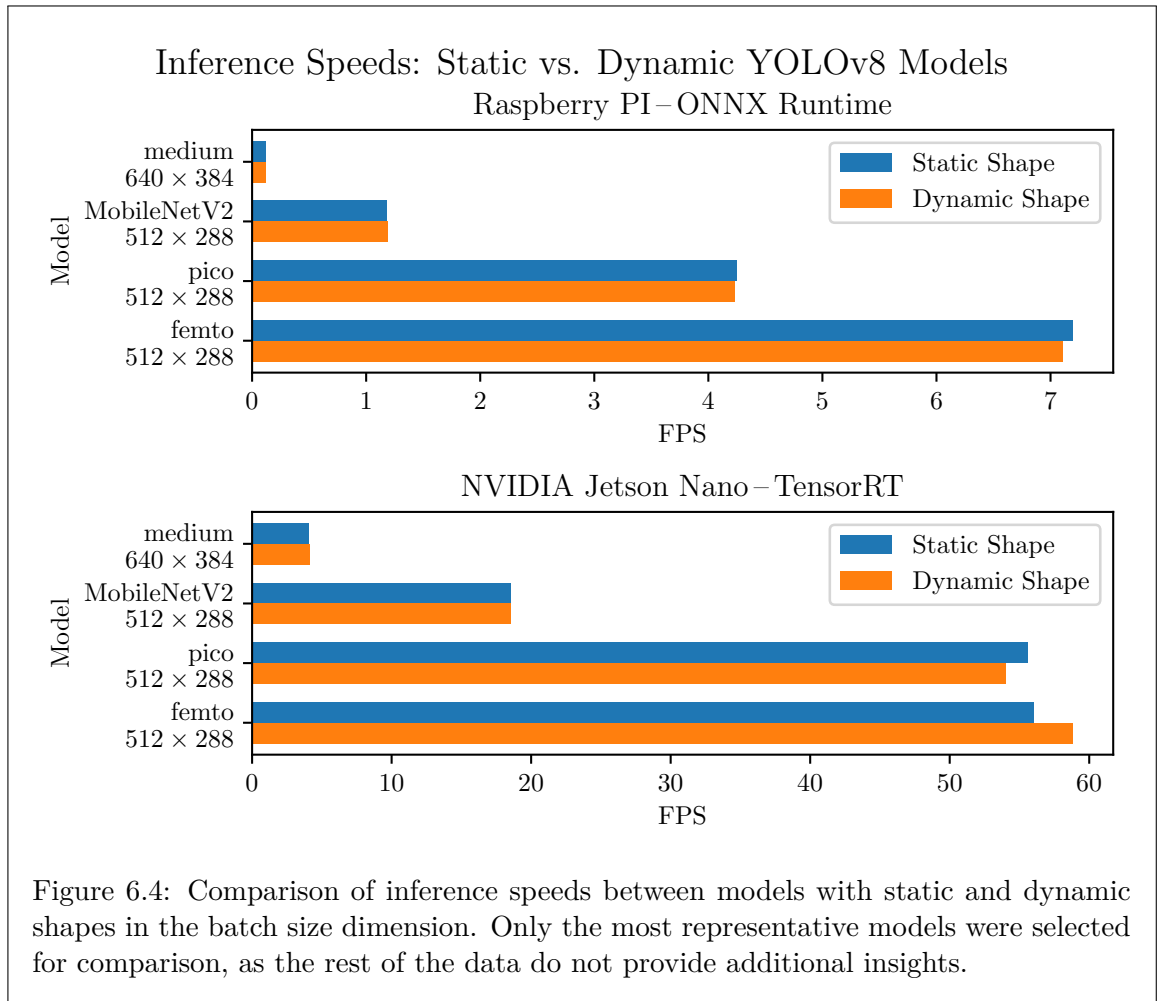
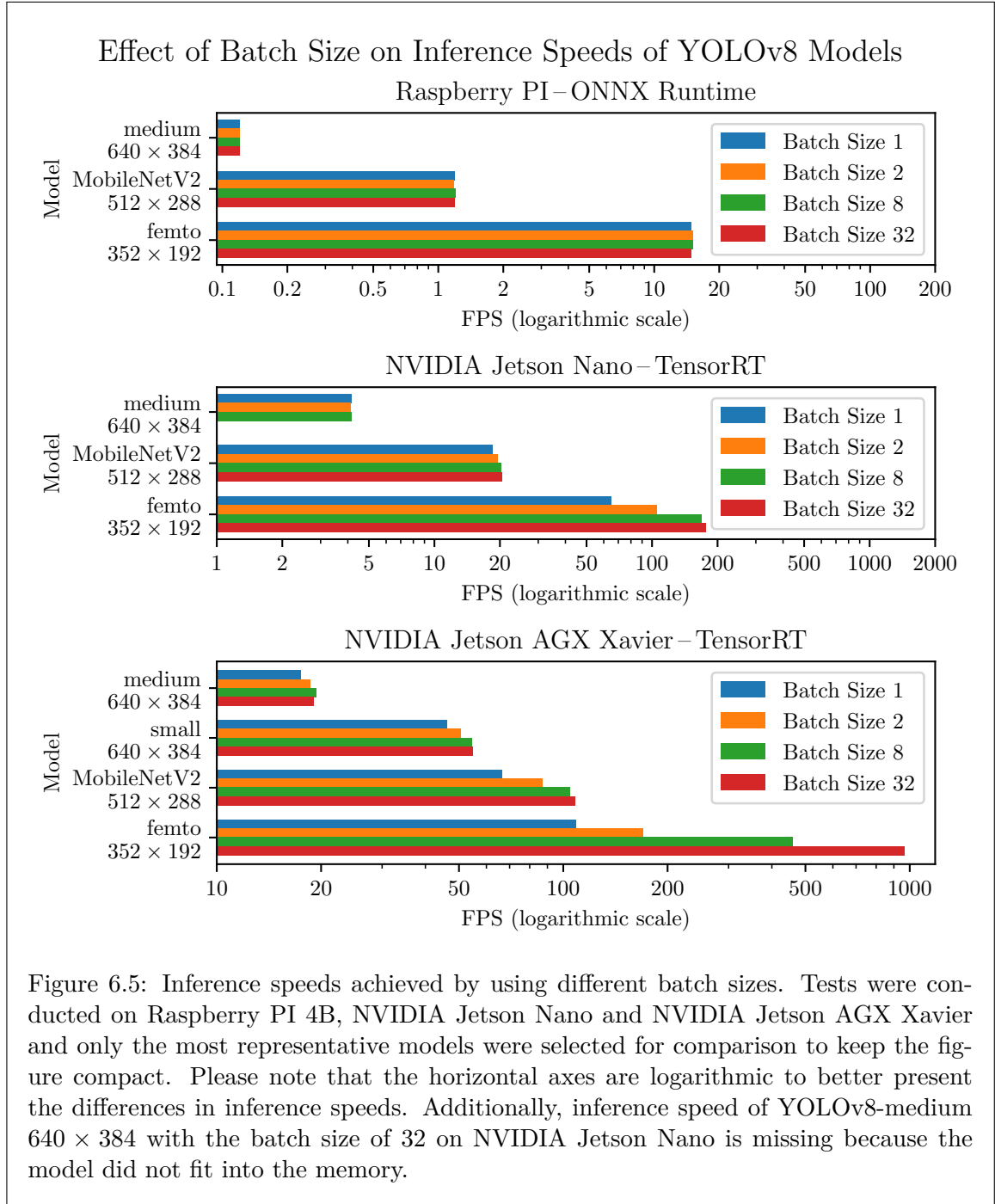


Figure 6.3: Comparison of inference speeds of all trained YOLOv8 model architectures between the TensorRT and the ONNX Runtime inference backends (of course, both utilizing the devices' GPU). Tests were conducted on NVIDIA Jetson Nano with a batch size of 1 and on NVIDIA Jetson AGX Xavier with a batch size of 32. Please note that for each model architecture, only the model with the largest input resolution was tested and all tested models were deployed as dynamic in the batch size dimension.





Chapter 7

Future Work

TODO TODO TODO

knowledge distillation, pruning,... kvantizácia onnx

Porovnať onnxruntime a tensorRT na bežnej GPU

NCNN alebo ARM NN pre RPI

SSD porovnať tiež

TFLite

Chapter 8

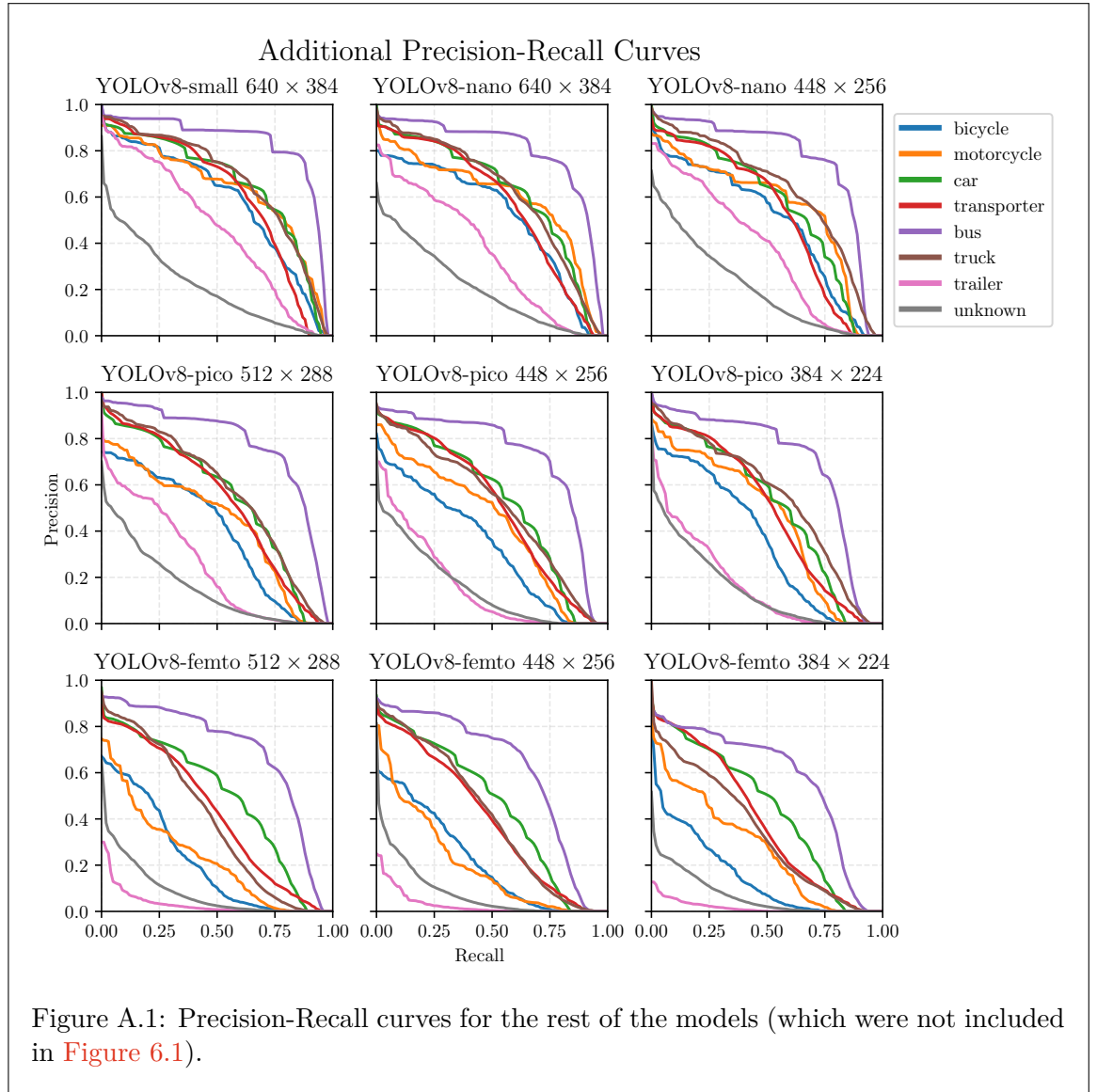
Conclusion

TODO TODO TODO

Appendices

Appendix A

Additional Precision-Recall Curves



Appendix B

Complete Benchmark of All Models on All Devices

[[pripomenúť že CPU len na jednom threade]]

Model (YOLOv8)	Input Resolution	Preci- sion	mAP	Inference Speed with a batch size of 1 (FPS)					
				Rasp- berry PI	Intel Core i7 9850H	NVIDIA MX150	NVIDIA Jetson		
							Nano	Xavier NX	AGX Xavier
medium	640×384	FP32	0.603	0.12	1.62	4.04	4.15	16.36	17.48
		FP16	0.603	-	-	-	6.82	47.49	49.29
		INT8	0.471	-	-	-	-	48.62	53.37
small	640×384	FP32	0.578	0.34	5.12	9.54	10.91	39.48	46.18
		FP16	0.578	-	-	-	16.27	52.01	57.44
		INT8	0.544	-	-	-	-	70.68	64.25
MobileNet V2	512×288	FP32	0.548	1.19	15.82	15.34	18.55	51.59	66.72
		FP16	0.548	-	-	-	19.93	53.30	58.66
		INT8	0.459	-	-	-	-	69.95	55.24
nano	640×384	FP32	0.530	1.02	14.78	20.40	25.08	54.30	57.50
		FP16	0.530	-	-	-	31.90	63.13	55.33
		INT8	0.430	-	-	-	-	59.78	64.21
	512×288	FP32	0.492	1.70	23.67	30.58	37.56	57.25	63.53
		FP16	0.492	-	-	-	48.29	70.01	71.38
		INT8	0.356	-	-	-	-	81.91	83.87
	448×256	FP32	0.511	2.18	29.55	36.21	49.25	67.73	58.17
		FP16	0.511	-	-	-	54.04	77.80	83.35
		INT8	0.372	-	-	-	-	61.81	91.60
pico	512×288	FP32	0.454	4.23	50.48	43.96	54.02	69.78	71.11
		FP16	0.454	-	-	-	57.14	59.80	87.29
		INT8	0.297	-	-	-	-	64.16	94.32
	448×256	FP32	0.415	5.39	61.38	52.89	57.37	78.47	81.29
		FP16	0.415	-	-	-	59.12	66.16	98.80
		INT8	0.356	-	-	-	-	71.28	101.98
	384×224	FP32	0.417	7.07	77.38	59.04	60.39	62.16	88.68
		FP16	0.418	-	-	-	61.85	77.94	108.55
		INT8	0.325	-	-	-	-	80.42	113.24
femto	512×288	FP32	0.318	7.11	62.40	55.16	58.78	55.81	84.13
		FP16	0.318	-	-	-	60.28	91.62	93.83
		INT8	0.300	-	-	-	-	65.07	97.91
	448×256	FP32	0.300	9.03	81.19	63.29	61.44	62.33	95.11
		FP16	0.300	-	-	-	58.12	70.06	104.58
		INT8	0.264	-	-	-	-	72.38	107.26
	384×224	FP32	0.281	11.87	103.49	76.33	66.01	72.29	103.47
		FP16	0.281	-	-	-	57.66	82.98	113.06
		INT8	0.235	-	-	-	-	83.54	115.56
	352×192	FP32	0.259	14.82	121.83	84.24	64.84	78.18	109.02
		FP16	0.259	-	-	-	65.91	89.41	120.63
		INT8	0.214	-	-	-	-	88.15	123.08

Table B.1: Benchmark of all trained and quantized YOLOv8 models on all available devices. For these tests, a batch size of 1 was used and all models were of a dynamic shape in the batch size dimension. Additionally, the TensorRT inference backend was used on NVIDIA Jetson devices, while the ONNX Runtime backend was used for tests on the rest of the devices. Models were only quantized to FP16 precision on NVIDIA Jetson devices and INT8 precision was only used on Jetson Xavier NX and Jetson AGX Xavier.

Model (YOLOv8)	Input Resolution	Preci- sion	mAP	Inference Speed with a batch size of 32* (FPS)					
				Rasp- berry PI	Intel Core i7 9850H	NVIDIA MX150	NVIDIA Jetson		
							Nano	Xavier NX	AGX Xavier
medium	640×384	FP32	0.603	0.12	1.57	4.68*	4.15*	16.97	19.11
		FP16	0.603	-	-	-	6.96	58.20	63.91
		INT8	0.471	-	-	-	-	86.30	102.17
small	640×384	FP32	0.578	0.34	4.98	10.70*	11.23	44.73	54.78
		FP16	0.578	-	-	-	16.79	121.36	140.62
		INT8	0.544	-	-	-	-	191.18	206.82
MobileNet V2	512×288	FP32	0.548	1.19	15.46	17.17	20.36	89.98	108.48
		FP16	0.548	-	-	-	22.00	159.10	185.10
		INT8	0.459	-	-	-	-	261.10	290.35
nano	640×384	FP32	0.530	1.03	14.38	21.22	26.30	115.28	139.49
		FP16	0.530	-	-	-	33.69	205.36	251.42
		INT8	0.430	-	-	-	-	264.79	305.43
	512×288	FP32	0.492	1.70	23.71	40.13	41.86	192.31	219.17
		FP16	0.492	-	-	-	54.16	328.51	390.02
		INT8	0.356	-	-	-	-	416.75	466.65
	448×256	FP32	0.511	2.21	30.16	51.74	55.73	246.03	283.35
		FP16	0.511	-	-	-	71.49	409.57	483.80
		INT8	0.372	-	-	-	-	524.52	575.99
pico	512×288	FP32	0.454	4.30	53.92	67.29	70.51	338.69	378.36
		FP16	0.454	-	-	-	83.02	457.63	523.29
		INT8	0.297	-	-	-	-	511.57	596.65
	448×256	FP32	0.415	5.48	68.66	83.70	92.82	432.56	484.74
		FP16	0.415	-	-	-	107.79	563.96	641.45
		INT8	0.356	-	-	-	-	631.91	740.18
	384×224	FP32	0.417	7.18	89.89	105.58	118.47	564.88	648.23
		FP16	0.418	-	-	-	138.83	763.25	888.08
		INT8	0.325	-	-	-	-	878.58	1016.77
femto	512×288	FP32	0.318	7.10	70.84	85.43	87.42	440.21	474.64
		FP16	0.318	-	-	-	98.72	527.16	579.90
		INT8	0.300	-	-	-	-	582.56	589.16
	448×256	FP32	0.300	9.04	93.97	106.62	113.09	556.71	600.81
		FP16	0.300	-	-	-	127.27	656.74	682.84
		INT8	0.264	-	-	-	-	679.27	726.82
	384×224	FP32	0.281	11.92	125.96	133.62	145.79	732.24	813.68
		FP16	0.281	-	-	-	158.38	843.47	1004.03
		INT8	0.235	-	-	-	-	1009.65	1034.69
	352×192	FP32	0.259	14.80	145.94	158.72	175.98	878.53	968.72
		FP16	0.259	-	-	-	197.29	1040.57	1152.97
		INT8	0.214	-	-	-	-	1117.36	1237.66

Table B.2: Benchmark of all trained and quantized YOLOv8 models on all available devices. For these tests, a batch size of 32 was used. However, for special cases marked by *, where the model did not fit into the device’s memory with the large batch size, a batch size of 8 was used instead. Models tested were of a dynamic shape in the batch size dimension. Additionally, the TensorRT inference backend was used on NVIDIA Jetson devices, while the ONNX Runtime backend was used for tests on the rest of the devices. Models were only quantized to FP16 precision on NVIDIA Jetson devices and INT8 precision was only used on Jetson Xavier NX and Jetson AGX Xavier.

Appendix C

Complete Mean Average Precision Metrics for YOLOv8 Models

Model (YOLOv8)	Input resolution	Precision	mAP					
			mAP	IoU:50	IoU:75	small	medium	large
medium	640×384	FP32	0.603	0.805	0.695	0.316	0.606	0.770
		FP16	0.603	0.805	0.693	0.316	0.607	0.770
		INT8	0.471	0.672	0.556	0.164	0.476	0.663
small	640×384	FP32	0.578	0.786	0.660	0.273	0.585	0.744
		FP16	0.578	0.786	0.660	0.274	0.585	0.742
		INT8	0.544	0.763	0.625	0.245	0.549	0.703
MobileNetV2	512×288	FP32	0.548	0.759	0.618	0.235	0.550	0.715
		FP16	0.548	0.759	0.619	0.236	0.549	0.718
		INT8	0.459	0.655	0.523	0.212	0.453	0.607
nano	640×384	FP32	0.530	0.749	0.608	0.229	0.529	0.685
		FP16	0.530	0.749	0.607	0.230	0.529	0.683
		INT8	0.430	0.650	0.484	0.182	0.453	0.530
	512×288	FP32	0.492	0.699	0.569	0.210	0.483	0.673
		FP16	0.492	0.698	0.569	0.209	0.482	0.669
		INT8	0.356	0.535	0.398	0.150	0.359	0.500
	448×256	FP32	0.511	0.723	0.592	0.206	0.508	0.681
		FP16	0.511	0.723	0.592	0.205	0.508	0.683
		INT8	0.372	0.555	0.417	0.140	0.371	0.516
pico	512×288	FP32	0.454	0.664	0.502	0.151	0.448	0.614
		FP16	0.454	0.664	0.500	0.151	0.448	0.612
		INT8	0.297	0.446	0.329	0.065	0.260	0.497
	448×256	FP32	0.415	0.621	0.468	0.136	0.408	0.552
		FP16	0.415	0.621	0.468	0.136	0.409	0.551
		INT8	0.356	0.552	0.396	0.122	0.355	0.474
	384×224	FP32	0.417	0.615	0.465	0.128	0.405	0.589
		FP16	0.418	0.615	0.465	0.128	0.405	0.588
		INT8	0.325	0.490	0.361	0.101	0.308	0.486
femto	512×288	FP32	0.318	0.488	0.351	0.072	0.298	0.416
		FP16	0.318	0.488	0.350	0.072	0.298	0.417
		INT8	0.300	0.466	0.326	0.071	0.270	0.420
	448×256	FP32	0.300	0.477	0.318	0.070	0.259	0.429
		FP16	0.300	0.477	0.319	0.070	0.259	0.430
		INT8	0.264	0.413	0.287	0.061	0.219	0.404
	384×224	FP32	0.281	0.452	0.302	0.070	0.263	0.386
		FP16	0.281	0.453	0.302	0.070	0.264	0.387
		INT8	0.235	0.409	0.242	0.059	0.211	0.331
	352×192	FP32	0.259	0.430	0.275	0.054	0.227	0.370
		FP16	0.259	0.430	0.275	0.054	0.226	0.370
		INT8	0.214	0.348	0.235	0.044	0.182	0.323

Table C.1: Comparison of multiple mean Average Precision metrics for each trained model and input resolution, including models quantized to FP16 and INT8.

Bibliography

- [1] MIO-TCD: A New Benchmark Dataset for Vehicle Classification and Localization. *IEEE Transactions on Image Processing*. Institute of Electrical and Electronics Engineers (IEEE). oct 2018, vol. 27, no. 10, p. 5129–5141. DOI: 10.1109/tip.2018.2848705.
- [2] BAHNSEN, C. H. and MOESLUND, T. B. Rain Removal in Traffic Surveillance: Does it Matter? *IEEE Transactions on Intelligent Transportation Systems*. Institute of Electrical and Electronics Engineers (IEEE). aug 2019, vol. 20, no. 8, p. 2802–2819. DOI: 10.1109/tits.2018.2872502.
- [3] BAI, J., LU, F., ZHANG, K. et al. *ONNX: Open Neural Network Exchange* [<https://github.com/onnx/onnx>]. 2019.
- [4] BUGDOL, M., MIODONSKA, Z., KRECICHWOST, M. and KASPEREK, P. Vehicle detection system using magnetic sensors. *Transport Problems*. march 2014, vol. 9.
- [5] BUSLAEV, A., PARINOV, A., KHVEDCHENYA, E., IGLOVIKOV, V. I. and KALININ, A. A. Albumentations: fast and flexible image augmentations. *ArXiv e-prints*. 2018.
- [6] CHOUDHARY, T., MISHRA, V., GOSWAMI, A. and SARANGAPANI, J. A comprehensive survey on model compression and acceleration. *Artificial Intelligence Review*. Springer Science and Business Media LLC. feb 2020, vol. 53, no. 7, p. 5113–5155. DOI: 10.1007/s10462-020-09816-7.
- [7] CONTRIBUTORS, M. *OpenMMLab’s Model deployment toolbox*. December 2021. Available at: <https://github.com/open-mmlab/mmdploy>.
- [8] CONTRIBUTORS, M. *OpenMMLab Detection Toolbox and Benchmark*. August 2018. Available at: <https://github.com/open-mmlab/mmdetection>.
- [9] CONTRIBUTORS, M. *MMYOLO: OpenMMLab YOLO series toolbox and benchmark* [<https://github.com/open-mmlab/mmyolo>]. 2022.
- [10] DEVELOPERS, O. R. *ONNX Runtime* [<https://onnxruntime.ai/>]. 2021.
- [11] FANG, J., MENG, H., ZHANG, H. and WANG, X. A Low-cost Vehicle Detection and Classification System based on Unmodulated Continuous-wave Radar. In: *2007 IEEE Intelligent Transportation Systems Conference*. 2007, p. 715–720. DOI: 10.1109/ITSC.2007.4357739.
- [12] GIRSHICK, R., DONAHUE, J., DARRELL, T. and MALIK, J. Rich feature hierarchies for accurate object detection and semantic segmentation. *arXiv*. november 2013. DOI: 10.48550/ARXIV.1311.2524.

- [13] HE, K., ZHANG, X., REN, S. and SUN, J. Deep Residual Learning for Image Recognition. arXiv. december 2015. DOI: 10.48550/ARXIV.1512.03385.
- [14] HOWARD, A. G., ZHU, M., CHEN, B., KALENICHENKO, D., WANG, W. et al. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv. april 2017. DOI: 10.48550/ARXIV.1704.04861.
- [15] JENSEN, M. B., MOGELMOSE, A. and MOESLUND, T. B. Presenting the Multi-view Traffic Intersection Dataset (MTID): A Detailed Traffic-Surveillance Dataset. In: *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, Sep 2020. DOI: 10.1109/itsc45102.2020.9294694.
- [16] JOCHER, G., CHAURASIA, A. and QIU, J. YOLO by Ultralytics. [online]. 2023. Accessed: May 5, 2023. Available at: <https://github.com/ultralytics/ultralytics>.
- [17] KRIZHEVSKY, A., SUTSKEVER, I. and HINTON, G. E. ImageNet Classification with Deep Convolutional Neural Networks. In: PEREIRA, F., BURGESS, C., BOTTOU, L. and WEINBERGER, K., ed. *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2012, vol. 25. Available at: https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf.
- [18] LABELBOX. *Labelbox* [online]. 2022. Accessed: May 5, 2023. Available at: <https://labelbox.com>.
- [19] LECUN, Y., BOTTOU, L., BENGIO, Y. and HAFFNER, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*. Institute of Electrical and Electronics Engineers (IEEE). 1998, vol. 86, no. 11, p. 2278–2324. DOI: 10.1109/5.726791.
- [20] LI, Z., LIU, F., YANG, W., PENG, S. and ZHOU, J. A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects. *IEEE Transactions on Neural Networks and Learning Systems*. Institute of Electrical and Electronics Engineers (IEEE). dec 2022, vol. 33, no. 12, p. 6999–7019. DOI: 10.1109/tnnls.2021.3084827.
- [21] LIU, W., ANGUELOV, D., ERHAN, D., SZEGEDY, C., REED, S. et al. SSD: Single Shot MultiBox Detector. Springer International Publishing. december 2015, p. 21–37. DOI: 10.1007/978-3-319-46448-0_2.
- [22] LUCA, C., CARLOS, S., PAULO, C. J., CLAUDIO, G. and GIUSEPPE, A. *Night and Day Instance Segmented Park (NDISPark) Dataset: a Collection of Images taken by Day and by Night for Vehicle Detection, Segmentation and Counting in Parking Areas*. Zenodo, 2022. DOI: 10.5281/ZENODO.6560822.
- [23] MAITY, M., BANERJEE, S. and CHAUDHURI, S. S. Faster R-CNN and YOLO based Vehicle detection: A Survey. In: *2021 5th International Conference on Computing Methodologies and Communication (ICCMC)*. IEEE, Apr 2021. DOI: 10.1109/iccmc51019.2021.9418274.
- [24] MURALI, S. and V K, G. Shadow Detection and Removal from a Single Image Using LAB Color Space. *Cybernetics and Information Technologies*. march 2013, vol. 13. DOI: 10.2478/cait-2013-0009.

- [25] OPENMMLAB. Dive into YOLOv8: How does this state-of-the-art model work? [online]. 2023. Accessed: May 5, 2023. Available at: <https://openmmlab.medium.com/dive-into-yolov8-how-does-this-state-of-the-art-model-work-10f18f74bab1>.
- [26] O'SHEA, K. and NASH, R. An Introduction to Convolutional Neural Networks. arXiv. november 2015. DOI: 10.48550/ARXIV.1511.08458.
- [27] PASZKE, A., GROSS, S., CHINTALA, S., CHANAN, G., YANG, E. et al. Automatic differentiation in PyTorch. 2017.
- [28] REDMON, J., DIVVALA, S., GIRSHICK, R. and FARHADI, A. You Only Look Once: Unified, Real-Time Object Detection. arXiv. june 2015. DOI: 10.48550/ARXIV.1506.02640.
- [29] SANDLER, M., HOWARD, A., ZHU, M., ZHMOGINOV, A. and CHEN, L.-C. MobileNetV2: Inverted Residuals and Linear Bottlenecks. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4510-4520. arXiv. january 2018. DOI: 10.48550/ARXIV.1801.04381.
- [30] SHAH, D. *Mean Average Precision (mAP) Explained: Everything You Need to Know* [online]. 2022. Accessed: May 5, 2023. Available at: <https://www.v7labs.com/blog/mean-average-precision>.
- [31] SIMONYAN, K. and ZISSERMAN, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. arXiv. september 2014. DOI: 10.48550/ARXIV.1409.1556.
- [32] STIAWAN, R., KUSUMADJATI, A., AMINAH, N. S., DJAMAL, M. and VIRIDI, S. An Ultrasonic Sensor System for Vehicle Detection Application. *Journal of Physics: Conference Series*. IOP Publishing. apr 2019, vol. 1204, no. 1, p. 012017. DOI: 10.1088/1742-6596/1204/1/012017. Available at: <https://dx.doi.org/10.1088/1742-6596/1204/1/012017>.
- [33] TANG, C., OUYANG, K., WANG, Z., ZHU, Y., WANG, Y. et al. Mixed-Precision Neural Network Quantization via Learned Layer-wise Importance. arXiv. march 2022. DOI: 10.48550/ARXIV.2203.08368.
- [34] WANG, H., YU, Y., CAI, Y., CHEN, X., CHEN, L. et al. A Comparative Study of State-of-the-Art Deep Learning Algorithms for Vehicle Detection. *IEEE Intelligent Transportation Systems Magazine*. Institute of Electrical and Electronics Engineers (IEEE). 2019, vol. 11, no. 2, p. 82-95. DOI: 10.1109/its.2019.2903518.
- [35] WANG, T., ZHU, Y., ZHAO, C., ZENG, W., WANG, Y. et al. Large Batch Optimization for Object Detection: Training COCO in 12 minutes. In: *Computer Vision – ECCV 2020*. Springer International Publishing, 2020, p. 481-496. DOI: 10.1007/978-3-030-58589-1_29.
- [36] WEN, L., DU, D., CAI, Z., LEI, Z., CHANG, M.-C. et al. UA-DETRAC: A New Benchmark and Protocol for Multi-Object Detection and Tracking. arXiv. november 2015. DOI: 10.48550/ARXIV.1511.04136.

- [37] WU, K., XU, T., ZHANG, H. and SONG, J. Overview of video-based vehicle detection technologies. *ICCSE 2011 - 6th International Conference on Computer Science and Education, Final Program and Proceedings*. august 2011. DOI: 10.1109/ICCSE.2011.6028764.
- [38] WU, T.-H., WANG, T.-W. and LIU, Y.-Q. Real-Time Vehicle and Distance Detection Based on Improved Yolo v5 Network. In: *2021 3rd World Symposium on Artificial Intelligence (WSAI)*. IEEE, Jun 2021. DOI: 10.1109/wsai51899.2021.9486316.
- [39] ZHU, P., WEN, L., DU, D., BIAN, X., FAN, H. et al. Detection and Tracking Meet Drones Challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. Institute of Electrical and Electronics Engineers (IEEE). nov 2022, vol. 44, no. 11, p. 7380–7399. DOI: 10.1109/tpami.2021.3119563.