

---

# *MicroVision*: AN OPEN DATASET AND BENCHMARK MODELS FOR DETECTING VULNERABLE ROAD USERS AND MICROMOBILITY VEHICLES

---

A PREPRINT

✉ **Alexander Rasch**

Chalmers University of Technology  
Gothenburg, Sweden  
alexander.rasch@chalmers.se

✉ **Rahul Rajendra Pai**

Chalmers University of Technology  
Gothenburg, Sweden  
rahul.pai@chalmers.se

March 20, 2026

## ABSTRACT

Micromobility is a growing mode of transportation, raising new challenges for traffic safety and planning due to increased interactions in areas where vulnerable road users (VRUs) share the infrastructure with micromobility, including parked micromobility vehicles (MMVs). Approaches to support traffic safety and planning increasingly rely on detecting road users in images—a computer-vision task relying heavily on the quality of the images to train on. However, existing open image datasets for training such models lack focus and diversity in VRUs and MMVs, for instance, by categorizing both pedestrians and MMV riders as “person”, or by not including new MMVs like e-scooters. Furthermore, datasets are often captured from a car perspective and lack data from areas where only VRUs travel (sidewalks, cycle paths). To help close this gap, we introduce the *MicroVision* dataset: an open image dataset and annotations for training and evaluating models for detecting the most common VRUs (pedestrians, cyclists, e-scooterists) and stationary MMVs (bicycles, e-scooters), from a VRU perspective. The dataset, recorded in Gothenburg (Sweden), consists of more than 8,000 anonymized, full-HD images with more than 30,000 carefully annotated VRUs and MMVs, captured over an entire year and part of almost 2,000 unique interaction scenes. Along with the dataset, we provide first benchmark object-detection models based on state-of-the-art architectures, which achieved a mean average precision of up to 0.723 on an unseen test set. The dataset and model can support traffic safety to distinguish between different VRUs and MMVs, or help monitoring systems identify the use of micromobility. The dataset and model weights can be accessed at <https://doi.org/10.71870/eepz-jd52>.

**Keywords** Micromobility · Open dataset · Traffic safety · Images · Bicycle · E-scooter · Object detection

## 1 Introduction

Micromobility refers to a mode of transport typically represented by compact and lightweight vehicles operating at lower speeds compared to more traditional modes of transport like passenger cars [1]. Micromobility has become a rising trend globally, particularly for personal transportation [2]. Its rise has been explained by its compactness and cost-effectiveness—often making it a popular last-mile option—but also by improvements to health and sustainability, as it is a more environmentally friendly mode of transportation [3, 2].

While the use of micromobility has increased over the last few years, so have interactions with other road users, as well as crashes [4]. As most users of micromobility vehicles (MMVs) fall within the category of vulnerable road users (VRUs), they are less protected than users of motorized vehicles such as cars or trucks. Consequently, VRUs are more likely to be injured in collisions, particularly when they do not wear helmets or collide with heavier motorized vehicles. As a response to critical interactions with motorized vehicles, VRUs may also become psychologically discouraged from continuing to travel, which may have negative long-term consequences for society. Therefore, improving the safety of micromobility is paramount to further promote and sustain its use.

Following the safe systems approach [5], initiatives to promote safe micromobility focus on road-user behavior, infrastructure, speed, and vehicle design. Initiatives focusing on behavior, for instance, study interactions between road users to understand what makes these interactions critical and to model behavior so that other initiatives can predict it and become more effective. For larger motorized vehicles, both active and passive safety systems have been developed to prevent crashes with VRUs or mitigate their consequences. However, for MMVs, such approaches have been scarce, partly due to tighter budget constraints and the limited space available for sensors and actuators.

A common challenge across all approaches is the need to accurately detect and distinguish micromobility vehicles and users in data. Such distinction is important because previous research has shown that e-scooterists, cyclists, and pedestrians have unique characteristics, such as differences in travel speed, route choice, and braking distance, for example during obstacle avoidance [6, 7, 8, 9, 10]. For instance, an e-scooterist may appear from afar like a pedestrian, but travel at speeds comparable to a cyclist [11]. Detecting VRUs has become particularly important for image data captured by cameras, which are central to modern safety systems due to their relatively low cost, small size, and ability to capture the environment in a manner similar to human vision. However, developing object-detection models with high accuracy requires large amounts of diverse training data.

## 2 Related Work

Previous work has provided multiple open traffic datasets that have proved useful to the research community developing models for detecting road users and vehicles [12, 13, 14]. However, these datasets have primarily been captured from a car perspective, thereby lacking data from cycle paths or sidewalks where only VRUs can travel. Furthermore, such datasets often miss finer distinctions between different types of micromobility, such as bicycles and e-scooters, or omit e-scooters entirely. E-scooters are known to perform differently compared to bicycles [15, 6]; therefore, distinguishing them from bicycles is important. In addition, detecting not only VRUs but also stationary MMVs has become increasingly important, as parked vehicles may pose obstacles to other road users.

Existing work on VRU detection has predominantly focused on pedestrians and cyclists. Most studies have developed annotated image datasets—some openly available—and trained deep-learning-based object-detection models. Model architectures have largely been based on *convolutional neural networks* (CNNs), which can be broadly categorized into two main classes: (1) region-based (two-stage) detectors, such as R-CNN, Fast R-CNN, and Faster R-CNN, which first generate region proposals that are subsequently classified; and (2) single-shot (one-stage) detectors, such as YOLO and SSD, which directly predict object locations and classes [16]. While two-stage detectors generally achieve higher accuracy at the cost of slower inference, single-stage detectors are often preferred for applications requiring faster inference, potentially at the expense of accuracy. With the advent of *transformers* in computer vision, transformer-based detectors such as DETR, RT-DETR, and RF-DETR have gained attention for their end-to-end formulation and strong performance in complex scenes, albeit with higher computational demands.

Previous research on cyclist detection has predominantly relied on CNN-based detection models. García-Venegas et al. [16], for example, developed an image dataset and detection model for cyclists and their orientations, concluding that region-based detectors yielded more accurate predictions, while single-shot detectors offered faster but less accurate results that could be sufficient when combined with object tracking. Xiaofei Li et al. [17] presented a public cyclist dataset from China containing more than 20,000 annotated instances and reported consistent detection performance across several then state-of-the-art architectures, including Fast R-CNN.

Due to the relatively recent appearance of e-scooters compared to bicycles on public roads, research on detecting e-scooters in images has been limited. Apurv et al. [18] were among the first to present

a public image dataset including e-scooters in the United States, consisting of approximately 10,000 images from 83 interaction scenes. Along with the dataset, they introduced a MobileNetV2-based benchmark model to distinguish e-scooter riders from pedestrians. Gilroy et al. [8] further improved rider classification by training a model on web-sourced images of partially occluded e-scooter riders. Chen et al. [19] published a dataset of approximately 2,000 images collected in Charlottesville, USA, with over 11,000 annotated objects, and presented benchmark detection models using one-stage YOLO architectures. Sabri et al. [20] introduced another public dataset of stationary MMVs, based on web-sourced videos featuring bicycles, e-scooters, and skateboards, and proposed a YOLOX-based detector augmented with temporal features, achieving improved detection accuracy.

Few studies have addressed the direct detection of e-scooterists, defined as the combined entity of rider and vehicle, and most have focused instead on image-level classification to distinguish them from cyclists [18, 8]. Additionally, many datasets have been collected over relatively short time spans, potentially missing seasonal variations in environmental conditions and road-user appearance. Most available datasets including e-scooterists have also been recorded in North America [18, 19, 20]. To support better generalization of detection models, additional datasets from other regions, including Europe, are needed [12].

In this study, we address the lack of data for micromobility safety by introducing *MicroVision*, an open dataset accompanied by initial benchmark object-detection models. The dataset comprises over 8,000 high-resolution anonymized images with more than 30,000 annotated instances of VRUs and MMVs. Unlike existing car-centric datasets, the images are captured from the ego perspective of micromobility users, covering VRU infrastructure often absent from other datasets and spanning a full annual cycle in Gothenburg, Sweden. Furthermore, we introduce a state-aware classification scheme that explicitly distinguishes between active riders (e.g., e-scooterists) and stationary vehicles (e.g., parked e-scooters), a distinction critical for downstream safety tasks such as trajectory prediction and risk assessment. Finally, we provide benchmark models based on state-of-the-art object-detection architectures to support future research on micromobility detection, particularly for traffic-safety applications.

## 3 Method

### 3.1 Data Collection

The images were derived from video footage systematically collected within the city of Gothenburg, Sweden. Data acquisition spanned a full annual cycle (from July 2024 to June 2025) and multiple times of day, thereby capturing a wide range of seasonal and lighting conditions to ensure substantial environmental and temporal diversity. Footage was captured from a micromobility user’s perspective by recording the surroundings while (1) riding an e-scooter, (2) riding a bicycle, and (3) walking. A GoPro Max camera was used for all recordings. For vehicle-based recordings, the camera was mounted on the handlebar at heights varying between 1.0 and 1.5 meters above the ground (see Fig. 1); while walking, the camera was handheld. Videos were recorded at a resolution of  $1920 \times 1080$  pixels at 60 frames per second. The high frame rate ensured sharp captures of moving road users at close proximity. To minimize lens distortion and provide images suitable for standard object-detection architectures without intensive post-calibration, the camera was operated in a *linear* lens mode.

The systematic recording of public spaces raises ethical concerns regarding the incidental collection of personal data without explicit individual consent. Requesting consent from each road user is impractical and may affect the realism of captured scenes. To address this, the data collector displayed a clearly visible sign during all recording sessions informing nearby road users about the ongoing data collection and providing contact information for inquiries or data-removal requests. Prior to public release, all data underwent anonymization to obscure identifiable personal information and comply with local data-protection regulations. Specifically, faces and vehicle license plates were blurred using the Precision Blur service provided by brighter AI, following prior work [12]. The collection protocol and data provision were approved by the Swedish Ethical Review Authority (reference number 2024-00329-01).

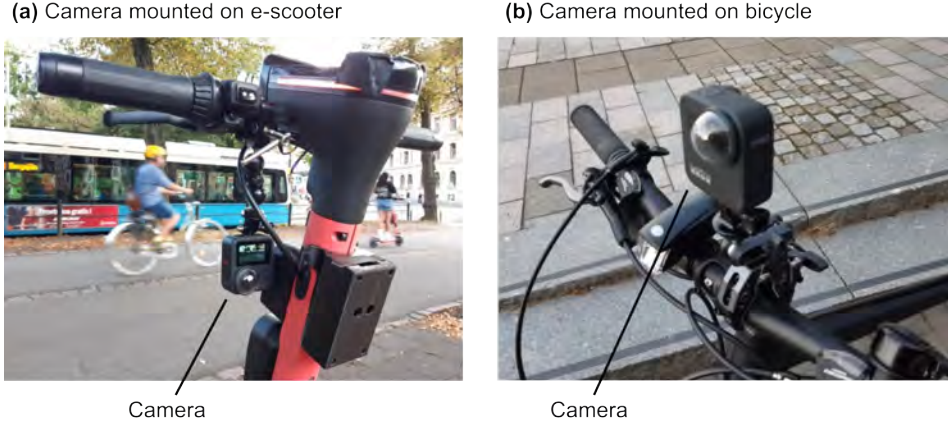


Figure 1: Placement of the GoPro Max camera on the e-scooter (a) and bicycle (b) used for data collection.

### 3.2 Frame Extraction

From the recorded videos, relevant frames were extracted for annotation. The extraction aimed to maximize environmental diversity while maintaining a balanced distribution of object classes. Only a limited number of frames depicting the same object from different angles were selected. This process was facilitated by pre-annotating the full videos with the current best model (YOLO11-based, fine-tuned from public COCO-pretrained weights). To maintain consistent object identities across frames, we employed the BoT-SORT tracker [21]. By tracking objects across frames and selecting frames sufficiently spaced in time, we ensured coverage of different viewpoints. Each candidate frame was then manually verified to confirm visual diversity and to exclude near-duplicate scenes.

### 3.3 Annotations

Selected frames were annotated with 2D rectangular bounding boxes for common VRUs and MMVs. Five classes were defined: (1) *pedestrian*, (2) *cyclist*, (3) *e-scooterist*, (4) stationary *bicycle*, and (5) stationary *e-scooter*. Consistent with the National Highway Traffic Safety Administration definition, a pedestrian was defined as a person who is walking or sitting, provided they were not on a vehicle [22]. For micromobility, we adopted a state-aware labeling strategy: a cyclist or e-scooterist was defined as the combination of the vehicle and at least one person traveling with it, including cases where the rider was stationary (e.g., waiting at a traffic light) or mounting the vehicle. Conversely, a person pushing a vehicle was labeled as two separate objects: a pedestrian and a vehicle. Stationary vehicles were annotated regardless of their orientation (upright or fallen).

Tricycles (e.g., bicycles with two front wheels) were included in the bicycle category in accordance with Swedish vehicle classification standards. Mopeds and large cargo cycles used for deliveries were excluded due to their rare appearance to avoid dataset imbalance. Smaller attachments (e.g., child seats or backpacks) were included within bounding boxes, while larger pushed or pulled objects (e.g., suitcases or child buggies) were excluded. Broken vehicles (e.g., missing wheels) were included in MMV categories. Following the COCO dataset conventions [23], occluded objects were annotated with boxes covering the visible extent; in cases of severe occlusion (e.g., dense bicycle racks), only the closest and most distinguishable object was annotated.

To maximize efficiency and consistency, a model-assisted semi-automated annotation workflow was employed (Fig. 2). First, an object-detection model (YOLO7; Wang et al. [24]) was fine-tuned on a small set of open-source web images; [25]. This model was used to pre-annotate the first batch of images. Three human annotators then reviewed and corrected labels using Label Studio v1.13.1 [26]. The corrected data were used to fine-tune a new model for pre-annotating subsequent batches. This iterative predict-correct-retrain cycle was repeated until the entire dataset was processed, progressively reducing manual effort.

A final model-guided quality check was conducted to identify and correct potential human annotation errors. The best-performing model from the iterative loop was fine-tuned on the full dataset (90%/10% train/validation split) and used to predict labels for all images. Discrepancies between predictions and ground truth were manually reviewed, focusing on false negatives and false positives, and corrected as needed. This validation was repeated with varying splits until no labeling issues were observed.

To mitigate model bias toward object-dense scenes and reduce false positives in empty environments, a small number of background images (approximately 5% of the dataset) containing none of the defined classes were added. These images were identified by the model as having no detections and verified manually.

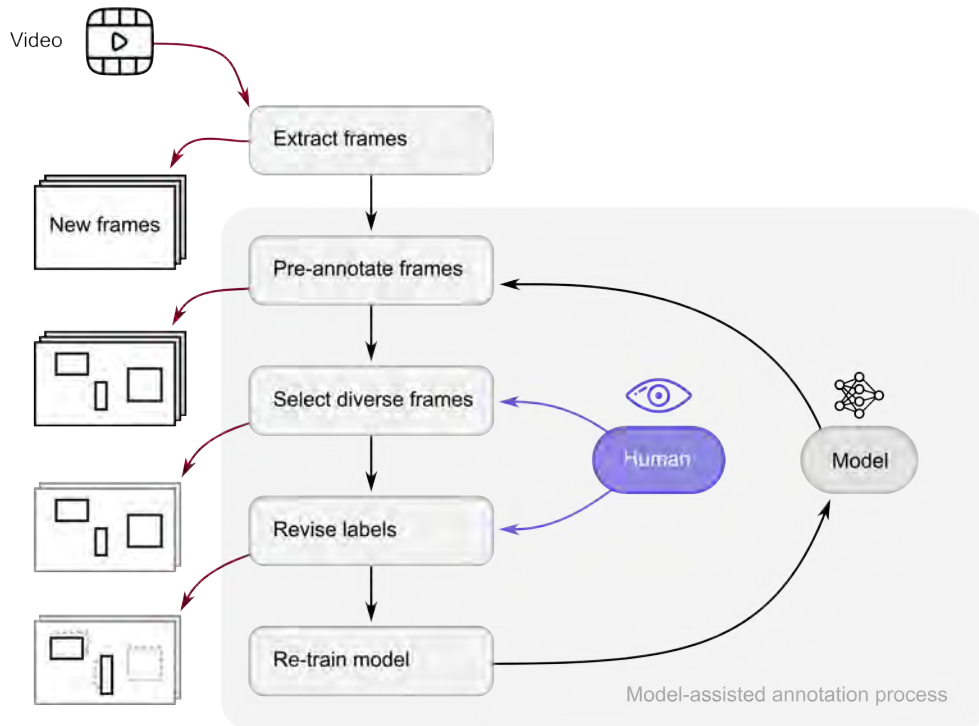


Figure 2: Overview of the data processing pipeline, from raw videos to revised labels used for model training.

### 3.4 Inter-Annotator Agreement

To assess reliability and consistency, an inter-annotator agreement analysis was conducted on a randomly selected subset of 258 images independently labeled by the three annotators. Pairwise agreement metrics were computed for all annotator pairs and averaged. For each pair, annotations were matched using the Hungarian algorithm to maximize total Intersection over Union (IoU) for bounding boxes of the same class [27]. Matches were considered valid only if IoU exceeded 0.5, consistent with common object-detection benchmarks such as Pascal VOC [28]. Based on valid matches, *box agreement* was computed as the average IoU to quantify spatial precision. *Class agreement* was assessed using Cohen’s Kappa [29], accounting for chance agreement. Metrics were stratified by object class and object size to ensure robustness across road-user types and distances.

### 3.5 Benchmark Object-Detection Models

To provide initial performance benchmarks, three state-of-the-art object-detection architectures representing distinct paradigms were trained and evaluated: (1) one-stage YOLO version 11 (YOLO11;

[30]), (2) two-stage Faster R-CNN [31], and (3) transformer-based RF-DETR [32]. Publicly available pretrained weights were fine-tuned for all models. To focus on architectural comparison rather than hyperparameter optimization, a uniform training setup was used. All models were trained for 100 epochs with an effective batch size of 32 and an input resolution of 1280 pixels (for RF-DETR, the nearest compatible resolution of 1232 pixels was used). YOLO11 (largest available model variant “X”) was trained using the Ultralytics package (v8.3.103). Faster R-CNN (largest variant “X101-FPN”) was trained using the Detectron2 package (v0.6; [33]). RF-DETR (largest variant “large”) was trained using the rfdetr package (v1.3.0) provided by Roboflow. Training was performed using up to four parallel NVIDIA A100 GPUs (80 GB VRAM each) on the Alvis compute cluster provided by the National Academic Infrastructure for Supercomputing in Sweden (NAISS).

To avoid data leakage, the dataset was split by *scenes* rather than by individual frames. Scenes were defined as temporal sequences containing the same objects, identified using YOLO11 tracking outputs; transitions between scenes were detected when the last tracked object disappeared for more than 10 s before the next appeared. Scenes were stratified into training (80%), validation (10%), and test (10%) sets.

Model performance was evaluated on the held-out test set using mean average precision (mAP), both per class and averaged across classes [34]. Following common practice, mAP was computed as the average over IoU thresholds from 0.5 to 0.95 in increments of 0.05 (mAP@0.5:0.95). In addition, performance was analyzed by object size following the COCO protocol [23]. Given the higher image resolution (1920 pixels wide) compared to COCO, size thresholds were scaled proportionally: *small* (area < 96<sup>2</sup> px<sup>2</sup>), *medium* (96<sup>2</sup> px<sup>2</sup> ≤ area < 288<sup>2</sup> px<sup>2</sup>), and *large* (area ≥ 288<sup>2</sup> px<sup>2</sup>).

## 4 Results

### 4.1 The MicroVision Dataset

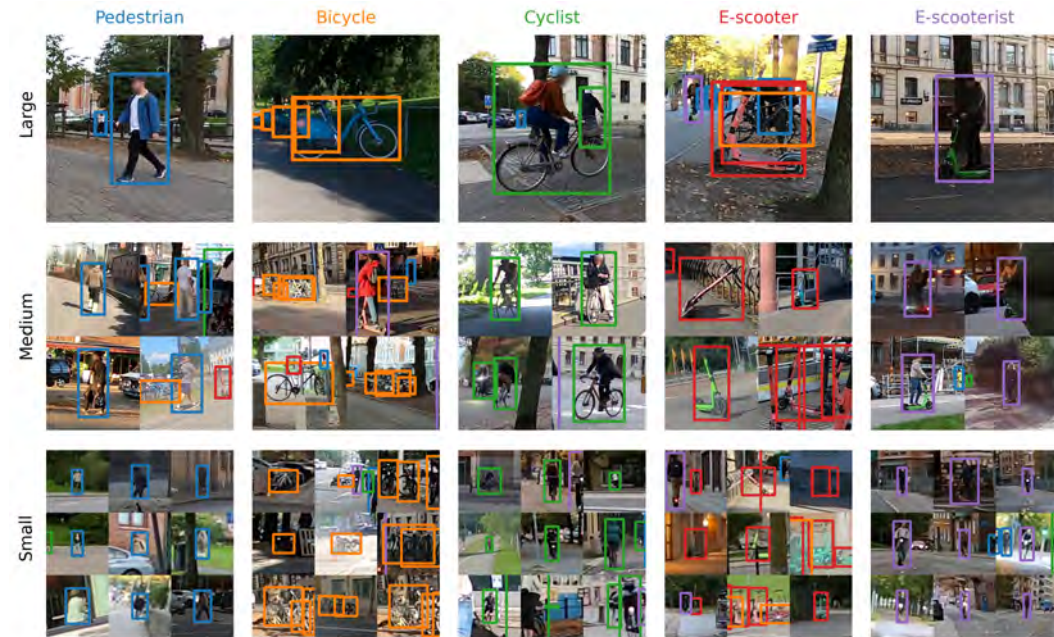


Figure 3: Example images and annotations from the MicroVision dataset for different object classes (columns) and object sizes (rows).

The MicroVision dataset contains 8,706 images, of which 594 are background images without any annotated objects. The remaining 8,113 images are annotated with a total of 34,866 objects using 2D bounding boxes. These include 17,032 pedestrians, 4,772 cyclists, 5,091 e-scooterists, 4,289 bicycles, and 3,682 e-scooters (see Fig. 3 for some examples). The images are organized into 1,984 unique scenes. Figure 4 shows the distributions of bounding-box widths and heights for all classes.

Overall, the dataset consists of 67% small, 25% medium, and 8% large objects. The comparatively large number of pedestrians, particularly small pedestrians, is a consequence of recording in dense urban environments, which increases background pedestrian frequency. While MMVs (bicycles and e-scooters) are predominantly represented by smaller and lower bounding boxes, VRU bounding boxes tend to be taller, as they combine both vehicle and rider.

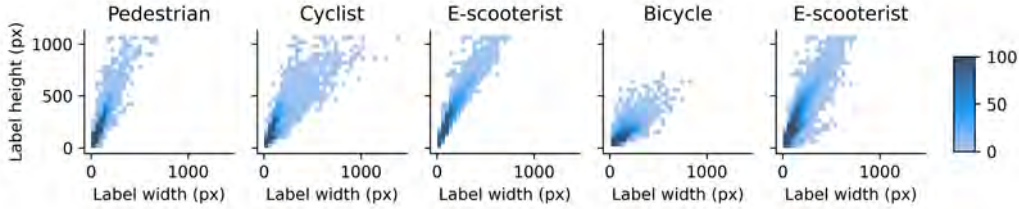


Figure 4: Bounding-box width and height distributions for the different object classes.

The inter-annotator agreement analysis yielded an overall Cohen’s Kappa of 0.824 and an average Intersection over Union (IoU) of 0.887. Agreement correlated positively with object size, increasing from small objects (Kappa = 0.783; IoU = 0.866) to large objects (Kappa = 0.963; IoU = 0.964). Across classes, active road users (cyclists and e-scooterists) exhibited higher annotation consistency than stationary vehicles, as summarized in Table 1.

Table 1: Inter-annotator agreement for different object classes and object sizes (S = small, M = medium, L = large). Class agreement is reported using Cohen’s Kappa, and box agreement using mean Intersection over Union (IoU).

Class	Class agreement (Cohen’s Kappa)				Box agreement (mean IoU)			
	S	M	L	All	S	M	L	All
Pedestrian	0.794	0.943	0.939	0.811	0.857	0.914	0.967	0.866
Bicycle	0.676	0.837	1.000	0.761	0.872	0.891	0.955	0.886
Cyclist	0.873	0.989	0.967	0.922	0.897	0.942	0.974	0.922
E-scooter	0.670	0.837	0.897	0.728	0.847	0.904	0.968	0.873
E-scooterist	0.877	0.987	0.978	0.938	0.911	0.953	0.957	0.937
All classes	0.783	0.919	0.963	0.824	0.866	0.922	0.964	0.887

## 4.2 Object-Detection Benchmark

Among the evaluated models, RF-DETR achieved the best overall performance on the held-out test set, with an overall mAP of 0.726, followed by YOLO11 (0.687) and Faster R-CNN (0.580), as reported in Table 2. While RF-DETR excelled in overall performance for each object size and class, and in particular for medium and large objects, YOLO11 outperformed the other models on small pedestrians and cyclists, small and medium-sized e-scooters, as well as small e-scooterists, almost closing the performance gap to RF-DETR for small objects overall. Faster R-CNN performed consistently worse across all object classes and object sizes. Image anonymization resulted in negligible performance differences across all models (YOLO11: 0.690 original vs. 0.687 anonymized; Faster R-CNN: 0.579 original vs. 0.580 anonymized; RF-DETR: 0.731 original vs. 0.726 anonymized).

Figure 5 illustrates ground-truth annotations and predictions from the two best-performing models (YOLO11 and RF-DETR) on representative test images using a confidence threshold of 0.5. Figure 6 highlights common failure cases. RF-DETR demonstrated superior performance for partially visible objects, such as occluded objects or objects near image boundaries. In some cases, both models incorrectly included adjacent objects (e.g., suitcases pulled by pedestrians) within bounding boxes; RF-DETR occasionally misclassified such instances. YOLO11, in contrast, correctly detected pedestrians pushing bicycles but sometimes produced redundant cyclist detections, indicating potential limitations in non-maximum suppression.

Table 2: Comparison of YOLO11, Faster R-CNN, and RF-DETR on the test set using COCO mAP@[0.5:0.95]. S, M, and L denote small, medium, and large objects, respectively. Bold values indicate the best performance per class and object size.

Class	Model	S	M	L	All
Pedestrian	YOLO11	<b>0.442</b>	0.645	0.769	0.597
	Faster R-CNN	0.279	0.559	0.726	0.499
	RF-DETR	0.438	<b>0.673</b>	<b>0.869</b>	<b>0.629</b>
Bicycle	YOLO11	0.110	0.373	0.724	0.518
	Faster R-CNN	0.131	0.280	0.633	0.431
	RF-DETR	<b>0.233</b>	<b>0.411</b>	<b>0.818</b>	<b>0.597</b>
Cyclist	YOLO11	<b>0.353</b>	0.700	0.883	0.766
	Faster R-CNN	0.093	0.570	0.818	0.669
	RF-DETR	0.332	<b>0.725</b>	<b>0.932</b>	<b>0.813</b>
E-scooter	YOLO11	<b>0.232</b>	<b>0.609</b>	0.837	0.692
	Faster R-CNN	0.052	0.408	0.694	0.510
	RF-DETR	0.226	0.583	<b>0.879</b>	<b>0.702</b>
E-scooterist	YOLO11	<b>0.574</b>	0.787	0.927	0.865
	Faster R-CNN	0.280	0.671	0.880	0.789
	RF-DETR	0.499	<b>0.816</b>	<b>0.950</b>	<b>0.889</b>
All classes	YOLO11	0.342	0.623	0.828	0.687
	Faster R-CNN	0.167	0.497	0.750	0.580
	RF-DETR	<b>0.346</b>	<b>0.641</b>	<b>0.890</b>	<b>0.726</b>

## 5 Discussion

### 5.1 Dataset and Benchmark Models

The MicroVision dataset provides a specialized resource for understanding complex interactions in urban micromobility environments. While its scale of 8,706 images is smaller than many established car-centric datasets [12, 13, 14], its scientific value lies in its qualitative uniqueness and technical specificity. Most existing open traffic datasets are recorded from the perspective of passenger cars and capture scenes primarily from the center of motorized traffic lanes. In contrast, the MicroVision dataset is captured directly from the perspective of micromobility users and pedestrians. This viewpoint shift is critical for developing safety systems intended for MMVs themselves, as it captures the specific visual context of sidewalks, cycle paths, and other shared urban spaces that are largely absent or underrepresented in car-centric data.

Furthermore, the dataset introduces a state-aware classification strategy that explicitly distinguishes between active riders (e.g., cyclists or e-scooterists) and stationary vehicles (e.g., bicycles or e-scooters). This distinction is essential for traffic safety and trajectory prediction, as a vehicle with a rider represents a dynamic entity with immediate kinetic potential and intent, whereas a parked vehicle constitutes a static obstacle. By providing high-resolution data collected over a full annual cycle, MicroVision ensures that these classifications are robust to environmental and seasonal variations typical of a Northern European urban setting. The large number of interaction scenes (nearly 2,000) further contributes to diversity in road-user appearance, viewing angles, and surrounding environments.

The initial benchmarking of state-of-the-art object-detection models demonstrated promising performance across architectures on unseen scenes. The transformer-based RF-DETR model outperformed CNN-based architectures (YOLO11 and Faster R-CNN) in overall performance, achieving an mAP of 0.723 and excelling particularly at detecting partially visible and occluded objects. This performance gap suggests that attention mechanisms may offer superior scene comprehension in dense urban environments. However, this improved accuracy comes at increased computational cost: the evaluated RF-DETR variant contains approximately 135 million parameters and requires longer inference

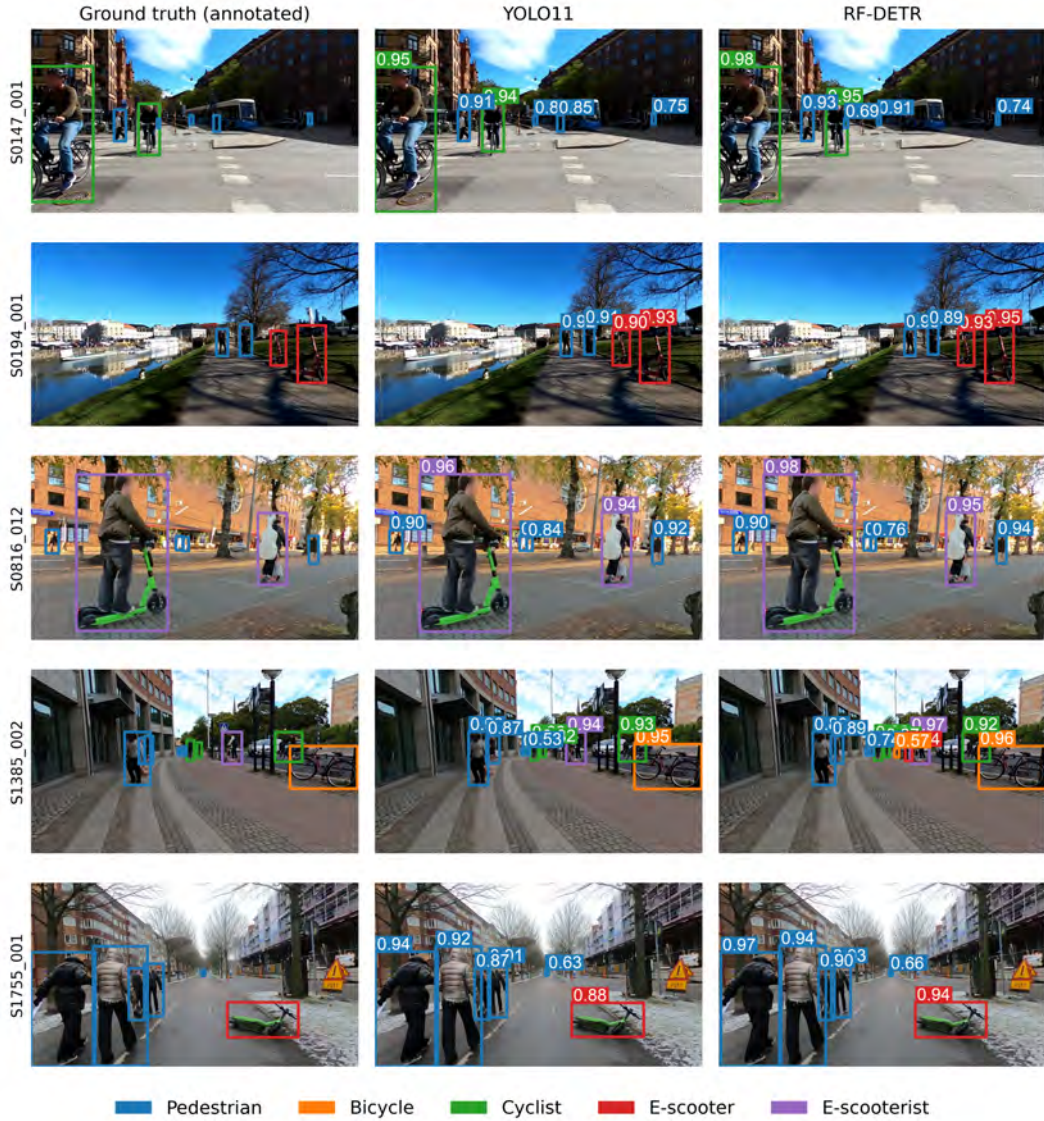


Figure 5: Successful example predictions on images from unseen test-set scenes for YOLO11 and RF-DETR, shown alongside ground-truth annotations. Numbers indicate predicted confidence scores (threshold = 0.5).

times, making it less suitable for latency-critical real-time applications compared to more lightweight models such as YOLO11 (57 million parameters). All evaluated models struggled with reliably detecting stationary MMVs, largely due to dense clustering in parking zones where vehicles occlude one another. Such cases pose challenges not only for automated detection but also for consistent human annotation, as reflected in the inter-annotator agreement analysis.

Model performance further degraded for small and distant objects, such as far-away pedestrians, which is a well-known limitation in object detection. From a traffic-safety perspective, immediate collision risks are typically lower for such distant objects; however, applications requiring long-range perception may benefit from complementary techniques. Object tracking can mitigate missed detections over time, and Slicing Aided Hyper Inference (SAHI) can improve detection of small objects by performing inference on image subregions [35]. Additional challenging cases included fallen or deformed (e.g., folded) MMVs, which were comparatively rare in the dataset and therefore more difficult for models to learn reliably.

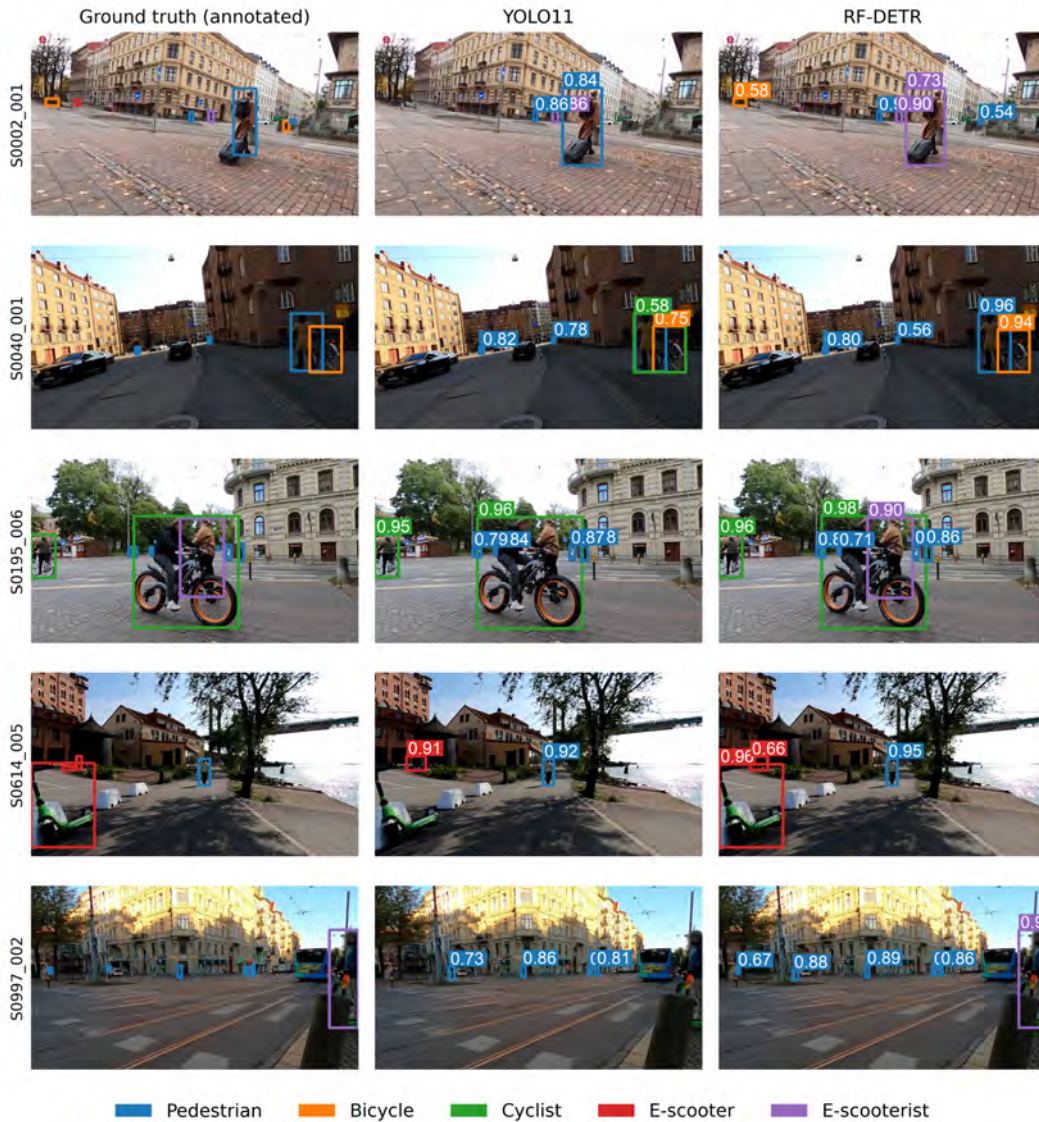


Figure 6: Representative examples of prediction errors from YOLO11 and RF-DETR, shown alongside ground-truth annotations.

## 5.2 Implications and Applications

Within the transportation domain, the MicroVision dataset enables a wide range of applications. Traffic-safety research and development can leverage the dataset and benchmark models to build advanced driver or rider assistance systems for both motorized vehicles and MMVs, aiming to prevent collisions with VRUs and MMVs. Detection outputs can support downstream tasks such as trajectory prediction that account for distinct behavioral patterns, for example, the greater unpredictability often associated with e-scooters compared to bicycles [15].

Beyond real-time systems, the benchmark models can support retrospective traffic-safety research by enabling automated identification of road users in large-scale naturalistic video datasets [36, 37, 10, 38]. Such datasets are typically annotated manually, a process that is both time-consuming and error-prone. The MicroVision dataset and models may facilitate behavioral research by enabling efficient discovery of specific road-user interactions, appearances, and conflict scenarios.

Although primarily designed for safety-related research, the dataset may also benefit traffic management and urban planning. Automated detection and counting of VRUs can support infrastructure assessment, demand modeling, and policy evaluation [39, 40]. However, the relatively low mounting height of the camera, corresponding to a road-user perspective, may limit the direct applicability of trained models to imagery captured from elevated viewpoints, such as building-mounted cameras.

### 5.3 Limitations and Future Work

While the dataset captures a broad range of objects and environmental conditions, it is geographically limited to Gothenburg, Sweden. As a result, the generalization of trained models to regions with different infrastructure designs, traffic regulations, or vehicle types remains an open question. Moreover, most images were collected during daytime and under favorable weather conditions, reflecting periods of higher VRU activity. Future work should investigate whether additional data are required to improve generalization to nighttime or adverse weather conditions. Nonetheless, the dataset provides a strong foundation for future model-assisted annotation workflows, enabling efficient extension to new domains.

Regarding scope, the annotation effort focused on the most prevalent micromobility forms currently observed in Sweden. Less common VRUs and MMVs, such as low-speed mopeds or monowheels, were not included, nor were fine-grained distinctions within categories (e.g., bicycles with trailers). Future dataset expansions could address these gaps by extending the temporal and spatial coverage. Additionally, extending annotations beyond 2D bounding boxes to include instance segmentation could enable more precise pixel-level scene understanding and support tasks such as drivable-space estimation.

## 6 Conclusion

We present an open image dataset with annotations and a first object-detection benchmark to advance the detection of vulnerable road users and micromobility vehicles. Comprising over 8,000 high-resolution images from nearly 2,000 unique interaction scenes and more than 30,000 annotations, the dataset addresses a critical gap by capturing the visual contexts of sidewalks and cycle paths and by employing a state-aware annotation strategy that distinguishes active riders from stationary vehicles.

Benchmarking state-of-the-art object-detection architectures shows that transformer-based models such as RF-DETR achieve superior accuracy in complex and occluded scenes, while efficient single-stage detectors like YOLO11 offer a favorable trade-off for real-time applications. Together, the dataset and benchmark models provide a foundation for next-generation traffic systems that can improve the safety and comfort of micromobility users, both through real-time assistance systems and by enabling partial automation of video-based behavioral research. By releasing these resources openly, we aim to support the development of models that generalize across the diverse and evolving landscape of urban transportation with a focus on micromobility.

## Acknowledgments

The authors thank Shiyi Qiu, Mahin Garg, and Anton Broman (Chalmers University of Technology) for their assistance with data processing and annotation, and Marco Dozza (Chalmers University of Technology) for valuable discussions and funding acquisition. The authors also thank the Chalmers Data Office for support with the practical and administrative aspects of data set preparation and publication.

The computations were enabled by resources provided by the National Academic Infrastructure for Supercomputing in Sweden (NAISS), partially funded by the Swedish Research Council through grant agreement no. 2022-06725.

This work was carried out within the *MicroVision* project, funded by Vinnova (Sweden’s innovation agency), the Swedish Energy Agency, and Formas (the Swedish Research Council for Sustainable Development) through the DriveSweden program (reference number 2023-01047).

## Data Availability

The MicroVision dataset and benchmark model weights are publicly available at <https://doi.org/10.71870/eepez-jd52>, hosted on the Data Organisation and Information System (DORIS) by the Swedish National Data Service (SND). Code to reproduce the analyses and data-processing pipeline is available at <https://github.com/microlab-chalmers/microvision>.

## Declarations

### Conflict of Interest

The authors declare no conflict of interest.

## References

- [1] SAE International. *Taxonomy and Classification of Powered Micromobility Vehicles*, SAE Standard J3194\_201911, November 2019. URL [https://doi.org/10.4271/J3194\\_201911](https://doi.org/10.4271/J3194_201911).
- [2] A. G. Olabi, Tabbi Wilberforce, Khaled Obaideen, Enas Taha Sayed, Nabila Shehata, Abdul Hai Alami, and Mohammad Ali Abdelkareem. Micromobility: Progress, benefits, challenges, policy and regulations, energy sources and storage, and its role in achieving sustainable development goals. *International Journal of Thermofluids*, 17(January):100292, 2023. ISSN 26662027. doi:10.1016/j.ijft.2023.100292. URL <https://doi.org/10.1016/j.ijft.2023.100292>.
- [3] Michael McQueen, Gabriella Abou-Zeid, John MacArthur, and Kelly Clifton. Transportation Transformation: Is Micromobility Making a Macro Impact on Sustainability? *Journal of Planning Literature*, 36(1):46–61, 2021. ISSN 15526593. doi:10.1177/0885412220972696.
- [4] George Yannis, Virginia Petraki, and Philippe Crist. Safer Micromobility: Technical Background Report The International Transport Forum. Technical Report March, 2024. URL <https://www.itf-oecd.org/sites/default/files/safer-micromobility-technical-report.pdf>.
- [5] Md Nasim Khan and Subasish Das. Advancing traffic safety through the safe system approach: A systematic review. *Accident Analysis & Prevention*, 199(February):107518, may 2024. ISSN 00014575. doi:10.1016/j.aap.2024.107518. URL <https://doi.org/10.1016/j.aap.2024.107518><https://linkinghub.elsevier.com/retrieve/pii/S0001457524000630>.
- [6] Marco Dozza, Alessio Violin, and Alexander Rasch. A data-driven framework for the safe integration of micro-mobility into the transport system: Comparing bicycles and e-scooters in field trials. *Journal of Safety Research*, 81:67–77, jun 2022. ISSN 00224375. doi:10.1016/j.jsr.2022.01.007. URL <https://doi.org/10.1016/j.jsr.2022.01.007><https://linkinghub.elsevier.com/retrieve/pii/S002243752200007X>.
- [7] Marco Dozza, Tianyou Li, Lucas Billstein, Christoffer Svernlöv, and Alexander Rasch. How do different micro-mobility vehicles affect longitudinal control? Results from a field experiment. *Journal of Safety Research*, 84:24–32, feb 2023. ISSN 00224375. doi:10.1016/j.jsr.2022.10.005. URL <https://doi.org/10.1016/j.jsr.2022.10.005><https://linkinghub.elsevier.com/retrieve/pii/S0022437522001591>.
- [8] Shane Gilroy, Darragh Mullins, Edward Jones, Ashkan Parsi, and Martin Glavin. E-Scooter Rider detection and classification in dense urban environments. *Results in Engineering*, 16 (October):100677, 2022. ISSN 25901230. doi:10.1016/j.rineng.2022.100677. URL <https://doi.org/10.1016/j.rineng.2022.100677>.
- [9] Khashayar Kazemzadeh, Milad Haghani, and Frances Sprei. Electric scooter safety: An integrative review of evidence from transport and medical research domains. *Sustainable Cities and Society*, 89(November 2022):104313, feb 2023. ISSN 22106707. doi:10.1016/j.scs.2022.104313. URL <https://doi.org/10.1016/j.scs.2022.104313><https://linkinghub.elsevier.com/retrieve/pii/S22106707220006175>.
- [10] Rahul Rajendra Pai and Marco Dozza. Understanding factors influencing e-scooterist crash risk: A naturalistic study of rental e-scooters in an urban area. *Accident Analysis and Prevention*, 209(June 2024):107839, 2025. ISSN 00014575. doi:10.1016/j.aap.2024.107839. URL <https://doi.org/10.1016/j.aap.2024.107839>.

- [11] Matteo della Mura, Serena Failla, Nicolò Gori, Alfonso Micucci, and Filippo Paganelli. E-Scooter Presence in Urban Areas: Are Consistent Rules, Paying Attention and Smooth Infrastructure Enough for Safety? *Sustainability (Switzerland)*, 14(21), 2022. ISSN 20711050. doi:10.3390/su142114303.
- [12] Mina Alibeigi, William Ljungbergh, Adam Tonderski, Georg Hess, Adam Lilja, Carl Lindstrom, Daria Motorniuk, Junsheng Fu, Jenny Widahl, and Christoffer Petersson. Zenseact open dataset: A large-scale and diverse multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [13] A Geiger, P Lenz, C Stiller, and R Urtasun. Vision meets robotics: The KITTI dataset. *The International Journal of Robotics Research*. *The International Journal of Robotics Research*, (October):1–6, 2013.
- [14] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurélien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2443–2451, 2020. ISSN 10636919. doi:10.1109/CVPR42600.2020.00252.
- [15] Natalia Distefano, Salvatore Leonardi, Mariusz Kieć, and Carmelo D’Agostino. Comparison of E-Scooter and Bike Users’ Behavior in Mixed Traffic. *Transportation Research Record*, 2024. ISSN 21694052. doi:10.1177/03611981241263339.
- [16] M. García-Venegas, D. A. Mercado-Ravell, L. A. Pinedo-Sánchez, and C. A. Carballo-Monsivais. On the safety of vulnerable road users by cyclist detection and tracking. *Machine Vision and Applications*, 32(5):1–17, 2021. ISSN 14321769. doi:10.1007/s00138-021-01231-4. URL <https://doi.org/10.1007/s00138-021-01231-4>.
- [17] Xiaofei Li, Fabian Flohr, Yue Yang, Hui Xiong, Markus Braun, Shuyue Pan, Keqiang Li, and Dariu M. Gavrilă. A new benchmark for vision-based cyclist detection. In *2016 IEEE Intelligent Vehicles Symposium (IV)*, volume 2016-Augus, pages 1028–1033. IEEE, jun 2016. ISBN 978-1-5090-1821-5. doi:10.1109/IVS.2016.7535515. URL <http://ieeexplore.ieee.org/document/7535515/>.
- [18] Kumar Apurv, Renran Tian, and Rini Sherony. Detection of e-scooter riders in naturalistic scenes, 2021. URL <https://arxiv.org/abs/2111.14060>.
- [19] Dong Chen, Arman Hosseini, Arik Smith, Amir Farzin Nikkhah, Arsalan Heydarian, Omid Shoghli, and Bradford Campbell. Performance Evaluation of Real-Time Object Detection for Electric Scooters, 2024. URL <https://arxiv.org/abs/2405.03039>.
- [20] Khalil Sabri, Céilia Djilali, Guillaume-Alexandre Bilodeau, Nicolas Saunier, and Wassim Bouachir. Detection of Micromobility Vehicles in Urban Traffic Videos. *Proceedings of the Conference on Robots and Vision, (Mmv)*, may 2024. doi:10.21428/d82e957c.abc3243f. URL <https://crv.pubpub.org/pub/du7cg0ee>.
- [21] Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. Bot-sort: Robust associations multi-pedestrian tracking, 2022. URL <https://arxiv.org/abs/2206.14651>.
- [22] National Center for Statistics and Analysis. Pedestrians: 2019 data (Traffic Safety Facts. Report No. DOT HS 813 079). Technical Report May, National Highway Traffic Safety Administration, 2021. URL <https://crashstats.nhtsa.dot.gov/Api/Public/Publication/813079>.
- [23] Tsung Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 8693 LNCS(PART 5):740–755, 2014. ISSN 16113349. doi:10.1007/978-3-319-10602-1\_48.
- [24] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors, 2022. URL <https://arxiv.org/abs/2207.02696>.
- [25] Luhan Fang and Yahui Wu. Threat assessment from naturalistic video-data: How to detect, classify, and estimate the position of multiple road users from cameras, 2024. URL <https://odr.chalmers.se/items/1f59fb35-7b6a-49b6-958e-d0fe540917ee>.

- [26] Maxim Tkachenko, Mikhail Malyuk, Andrey Holmanyuk, and Nikolai Liubimov. Label Studio: Data labeling software, 2025. URL <https://github.com/HumanSignal/label-studio>.
- [27] H W Kuhn. The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2(1-2):83–97, 1955. doi:<https://doi.org/10.1002/nav.3800020109>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/nav.3800020109>.
- [28] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–338, jun 2010. ISSN 0920-5691. doi:10.1007/s11263-009-0275-4. URL <http://link.springer.com/10.1007/s11263-009-0275-4>.
- [29] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20:37–46, 1960. ISSN 1552-3888(Electronic),0013-1644(Print). doi:10.1177/001316446002000104.
- [30] Glenn Jocher and Jing Qiu. Ultralytics YOLO11, 2024. URL <https://github.com/ultralytics/ultralytics>.
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In C Cortes, N Lawrence, D Lee, M Sugiyama, and R Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL [https://proceedings.neurips.cc/paper\\_files/paper/2015/file/14bfa6bb14875e45bba028a21ed38046-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/14bfa6bb14875e45bba028a21ed38046-Paper.pdf).
- [32] Isaac Robinson, Peter Robicieux, Matvei Popov, Deva Ramanan, and Neehar Peri. RF-DETR: Neural Architecture Search for Real-Time Detection Transformers, 2025. URL <https://arxiv.org/abs/2511.09554>.
- [33] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2, 2019. URL <https://github.com/facebookresearch/detectron2>.
- [34] Steven M Beitzel, Eric C Jensen, and Ophir Frieder. *MAP*, pages 1691–1692. Springer US, Boston, MA, 2009. ISBN 978-0-387-39940-9. doi:10.1007/978-0-387-39940-9\_492. URL [https://doi.org/10.1007/978-0-387-39940-9\\_492](https://doi.org/10.1007/978-0-387-39940-9_492).
- [35] Fatih Cagatay Akyon, Sinan Onur Altinuc, and Alptekin Temizel. Slicing Aided Hyper Inference and Fine-Tuning for Small Object Detection. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 966–970. IEEE, oct 2022. ISBN 978-1-6654-9620-9. doi:10.1109/ICIP46576.2022.9897990. URL <http://arxiv.org/abs/2202.06934><http://dx.doi.org/10.1109/ICIP46576.2022.9897990><https://ieeexplore.ieee.org/document/9897990/>.
- [36] Ben Beck, Derek Chong, Jake Olivier, Monica Perkins, Anthony Tsay, Adam Rushford, Lingxiao Li, Peter Cameron, Richard Fry, and Marilyn Johnson. How much space do drivers provide when passing cyclists? Understanding the impact of motor vehicle and infrastructure characteristics on passing distance. *Accident Analysis & Prevention*, 128: 253–260, jul 2019. ISSN 00014575. doi:10.1016/j.aap.2019.03.007. URL <https://www.sciencedirect.com/science/article/abs/pii/S0001457518309990><https://linkinghub.elsevier.com/retrieve/pii/S0001457518309990>.
- [37] Marco Dozza and Julia Werneke. Introducing naturalistic cycling data: What factors influence bicyclists’ safety in the real world? *Transportation Research Part F: Traffic Psychology and Behaviour*, 24:83–91, may 2014. ISSN 1369-8478. doi:10.1016/J.TRF.2014.04.001. URL <https://www.sciencedirect.com/science/article/pii/S1369847814000394>.
- [38] K. Schleinitz, T. Petzoldt, L. Franke-Bartholdt, J. Krems, and T. Gehlert. The German Naturalistic Cycling Study – Comparing cycling speed of riders of different e-bikes and conventional bicycles. *Safety Science*, 92:290–297, feb 2017. ISSN 0925-7535. doi:10.1016/J.SSCI.2015.07.027. URL <https://www.sciencedirect.com/science/article/pii/S0925753515001976>.
- [39] Eduardo Peixoto, João Moutinho, and Rui José. A Low-Cost Video-Based Solution for City-Wide Bicycle Counting in Starter Cities. In Henrique Santos, Gabriela Viale Pereira, Matthias Budde, Sérgio F Lopes, and Predrag Nikolic, editors, *Science and Technologies for Smart Cities*, pages 139–150, Cham, 2020. Springer International Publishing. ISBN 978-3-030-51005-3.
- [40] Francisco Vacalebri, Sara Moll, Griselda López, and Alfredo García. AI-Enhanced Road Safety: Real-Time Information on Cycle Traffic on Rural Roads. In Angel A Juan, Javier Faulin, and

David Lopez-Lopez, editors, *Decision Sciences*, pages 281–288, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-78241-1.