# Towards Reflected Object Detection: A Benchmark

Zhongtian Wang[1,†], You Wu[1,†], Hui Zhou[1], and Shuiwang Li[1,2]

[1] Guilin University of Technology, China
[2] Guangxi Key Laboratory of Embedded Technology and Intelligent System, China
`lishuiwang0721@163.com`

**Abstract.** Object detection has greatly improved over the past decade thanks to advances in deep learning and large-scale datasets. However, detecting objects reflected in surfaces remains an underexplored area. Reflective surfaces are ubiquitous in daily life, appearing in homes, offices, public spaces, and natural environments. Accurate detection and interpretation of reflected objects are essential for various applications. This paper addresses this gap by introducing a extensive benchmark specifically designed for Reflected Object Detection. Our Reflected Object Detection Dataset (RODD) features a diverse collection of images showcasing reflected objects in various contexts, providing standard annotations for both real and reflected objects. This distinguishes it from traditional object detection benchmarks. RODD encompasses 10 categories and includes 21,059 images of real and reflected objects across different backgrounds, complete with standard bounding box annotations and the classification of objects as real or reflected. Additionally, we present baseline results by adapting five state-of-the-art object detection models to address this challenging task. Experimental results underscore the limitations of existing methods when applied to reflected object detection, highlighting the need for specialized approaches. By releasing RODD, we aim to support and advance future research on detecting reflected objects. Dataset and code are available at: https://github.com/Tqybu-hans/RODD.
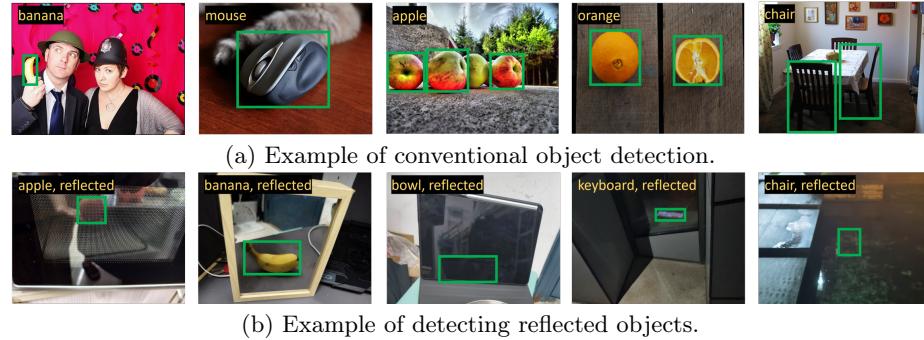
**Keywords:** Reflected object detection · Benchmark · Object detection

## 1 Introduction

The field of object detection has seen remarkable advancements over the past decade, driven by the development of deep learning techniques and the availability of large-scale datasets [16–18]. These advancements have significantly improved the accuracy and robustness of object detection systems in various applications [19]. However, one area that remains underexplored is the detection of objects reflected in surfaces, such as mirrors, glass windows, and other reflective materials. See Fig. 1 for an illustration of the difference between conventional object detection and reflected object detection.

---

† These authors contributed equally to this work.

(a) Example of conventional object detection.



(b) Example of detecting reflected objects.

**Fig. 1:** While previous object detection focused on the identification and localization of objects, this work focuses on information beyond that and concerns about the nature of objects in addition, as shown in (a) and (b), respectively. Note the nature of the objects (i.e., real or reflected) are marked in (b) additionally.

Reflective surfaces are ubiquitous in our daily lives, appearing in a wide array of environments and applications [20–22]. Mirrors, glass windows, water surfaces, and polished metals are just a few examples of materials that produce reflections. These reflective surfaces are prevalent in various settings, including homes, offices, public spaces, and natural environments, making the ability to detect and interpret reflected objects a crucial aspect of many technological applications [20, 23–25]. For instance, in surveillance, security systems can more effectively identify real intrusions or threats by differentiating reflections from genuine objects [32–34]. For autonomous driving, accurate identification of real objects versus reflections enables vehicles to navigate more safely and avoid accidents caused by misinterpretation [60–62, 65]. For service robots, robots can perform tasks with greater accuracy, such as picking and placing items, by correctly identifying real objects instead of their reflections [22, 67]. This improved object detection also facilitates better navigation in environments with reflective surfaces, such as warehouses [66, 68]. In smart homes, systems can provide more tailored responses by recognizing when a person is truly present rather than reacting to their reflection [47–49]. In medical applications, imaging and diagnostic tools yield more accurate results when they accurately interpret reflections, leading to better patient outcomes and more precise medical interventions.

Given the widespread presence of reflective surfaces in daily life, developing technologies that can effectively detect and interpret reflected objects is essential. This capability can enhance the performance and reliability of various applications, including smart home systems, surveillance, autonomous driving, and medical devices. However, to the best of our knowledge, there is currently no public benchmark for reflected object detection. This paper aims to address this gap by introducing a benchmark specifically designed for this purpose. We propose a comprehensive benchmark that includes a diverse set of images featuring reflected objects in various contexts. Our benchmark is designed to test the limits of current object detection methods and provide a standardized evaluation framework for developing and comparing new algorithms tailored to re-

flected object detection. The benchmark provides standard annotations used in object detection for identifying both actual (real) objects and their reflections. Additionally, it offers extra details that indicate whether an object is real or a reflection. This feature distinguishes it from traditional object detection benchmarks, which typically do not provide information about whether an object is a reflection. In addition to introducing the benchmark, this paper also presents baseline results by adapting several state-of-the-art object detection models. These results highlight the limitations of existing methods when applied to reflected object detection and underscore the need for specialized approaches. We analyze the performance of these models across different reflection scenarios and provide insights into the specific challenges posed by reflections.

## 1.1   Contribution

In this work, we make the first attempt to explore reflected object detection by introducing the RODD benchmark, which is specifically designed for detecting reflected objects. This benchmark provides a well-annotated dataset and robust evaluation metrics to facilitate research in this challenging area. The RODD benchmark fills a crucial gap in current object detection methods by focusing on reflected objects. It aims to provide researchers with a valuable resource to develop and test algorithms that handle the complexities of reflected objects. RODD comprises a diverse set of 10 classes of generic objects, totaling 21059 images annotated with axis-aligned bounding boxes, category labels, and object nature (real or reflected). Sample images from the RODD dataset are illustrated in Fig. 2. In addition, we developed five baseline detectors based on five state-of-the-art algorithms, namely RO-YOLOV8, RO-YOLOV10, RO-RTMDet, RO-YOLOX, and RO-PPYOLOE. These baselines serve to evaluate detectors' performance and provide benchmarks for future research on RODD. In summary, our contributions include:

- We make the first attempt to explore detecting reflected objects, a previously underexplored area in object detection. By focusing on this unique challenge, we hope to inspire further research and innovation in the detecting reflected objects.
- We introduce RODD, the first benchmark dedicated to detecting reflected objects, which consists of 10 classes of generic objects, with 21,059 images annotated with bounding boxes, object categories, and the nature of the objects. This dataset will enable detailed analysis and evaluation of algorithms developed for detecting reflected objects.
- To support further research on RODD, we develop five baseline detectors based on state-of-the-art models: RO-YOLOV8, RO-YOLOV10, RO-RTMDet, RO-YOLOX, and RO-PPYOLOE. These baseline models will provide initial performance metrics and serve as reference points for future studies.

## 2   Related Work

### 2.1   Object Detection Algorithms

Object detection has been a critical area of research in computer vision, significantly advancing over the past few decades. Traditional object detection methods relied heavily on handcrafted features and shallow learning techniques. The advent of deep learning has revolutionized this field, leading to the development of more robust and accurate algorithms. Modern object detection methods are categorized into two types: two-stage detectors and one-stage detectors. Two-stage detectors, such as R-CNN [79], Fast R-CNN [78], Faster R-CNN [75], and Mask R-CNN [74], initially generate region proposals and then refine them through classification and bounding box regression, achieving high precision and efficiency. Variants like Cascade R-CNN [80] further enhance detection performance through multi-stage detection and regression. One-stage detectors, including SSD [81], YOLO, RetinaNet [82], and EfficientDet [83], predict object locations and categories in a single step, providing faster performance suitable for real-time applications. The YOLO series has evolved to YOLOv8 [85] and YOLOv10 [4], further optimizing speed and accuracy.

Despite these advancements, detecting objects reflected in surfaces such as mirrors and glass remains a challenging and underexplored problem. Most existing object detection algorithms are not specifically designed to differentiate between real objects and their reflections, leading to potential false positives and degraded performance in environments with reflective surfaces. Our work aims to address this gap by introducing a benchmark and developing specialized approaches for reflected object detection.

### 2.2   Object Detection Benchmarks

Object detection benchmarks play a crucial role in the development and evaluation of detection algorithms by providing standardized datasets and evaluation metrics that facilitate consistent and fair comparisons among different approaches. Over the years, several prominent benchmarks have emerged, each contributing uniquely to the field, such as PASCAL VOC [86], MS COCO [6], and ImageNet [84]. These benchmarks provide large-scale images and standardized evaluation metrics. For instance, PASCAL VOC comprises 20 categories with 11,530 images and 27,450 annotated bounding boxes. ImageNet covers 200 categories with approximately 500,000 annotated bounding boxes. MS COCO includes 91 categories, over 300,000 images, and 2.5 million annotated instances. These datasets have been instrumental in pushing the boundaries of object detection research, promoting the development of more accurate and robust models. In addition to these established benchmarks, several domain-specific benchmarks have emerged to address particular challenges in object detection. For instance, KITTI [87] focuses on autonomous driving scenarios, providing annotated data for detecting objects such as cars, pedestrians, and cyclists in street scenes.

UAVDT (UAV Detection and Tracking) [88] provides benchmarks for aerial object detection, emphasizing challenges unique to unmanned aerial vehicle (UAV) imagery, such as varying altitudes and viewpoints.

Despite significant advancements in object detection, no public benchmark specifically targets reflected object detection. This gap hinders the development and evaluation of algorithms for handling reflections. Reflective surfaces are common in real-world scenarios such as surveillance, autonomous driving, and smart homes. Accurate detection of reflected objects is crucial for enhancing performance and safety in these applications. This paper addresses this gap by introducing a benchmark specifically tailored for reflected object detection.

### 2.3 Dealing With Mirrors and Reflections in Vision

Mirrors or other reflective surfaces are common in natural images, and can cause false positive results in the tasks of detection, segmentation, counting, robotic navigation, scene reconstruction, and etc [20, 23–25, 69, 70]. Reflection detection focuses on identifying regions in an image that contain reflections. When we take a picture through glass windows, the photographs are often degraded by undesired reflections. One of the primary approaches to dealing with reflections involves removing or suppressing the reflections in images. For instance, Abiko et al. employed generative adversarial networks (GANs) to enhance the quality of reflection removal, yielding more natural and clear images [71]. Arvanitopoulos et al. propose a single image reflection suppression method based on a Laplacian data fidelity and an l-zero gradient sparsity regularization term [72]. Particularly, mirror surface detection aims to identify and segment mirror surfaces within a scene. For instance, Yang et al. proposed to address the mirror segmentation problem with a computational approach [21]. Since then, numerous methods have been developed to address mirror detection and segmentation [20, 23–25].
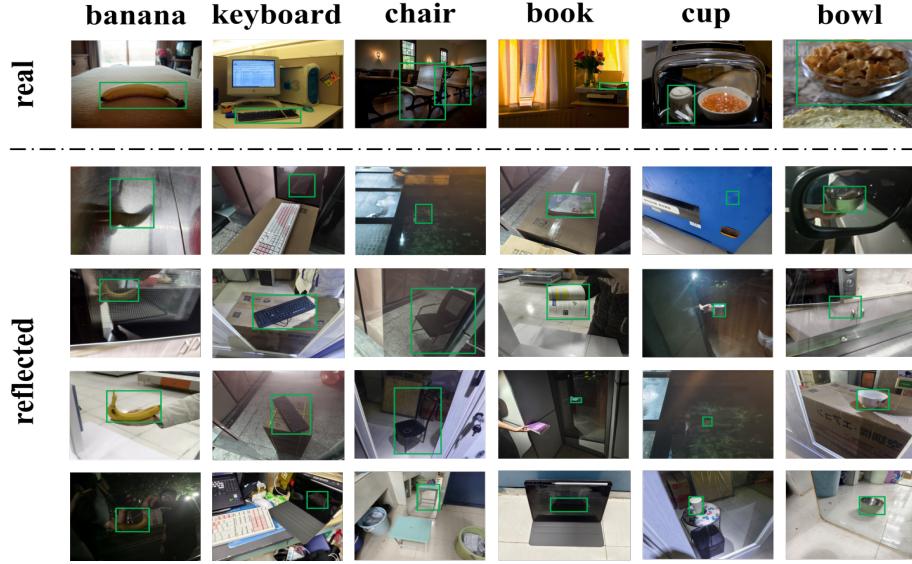
Despite extensive research efforts dedicated to dealing with mirrors and reflections in vision, most of these works focus primarily on identifying, localizing, segmenting, and suppressing reflective regions in images. In this work, we make the first attempt to differentiate reflected objects from real ones, a critical capability for various applications, including surveillance, autonomous driving, service robots, and smart homes.

## 3 Benchmark for Reflected Object Detection

We construct a dedicated dataset for Reflected Object Detection Dataset (RODD), which is a dataset that contains labels of both class and object nature, with prediction bounding-box labeled for each image.

### 3.1 Image Collection

For image collection, we selected 10 common objects in daily life, guided by the selection principles of PASCAL VOC [7] and COCO [6]. The chosen objects for
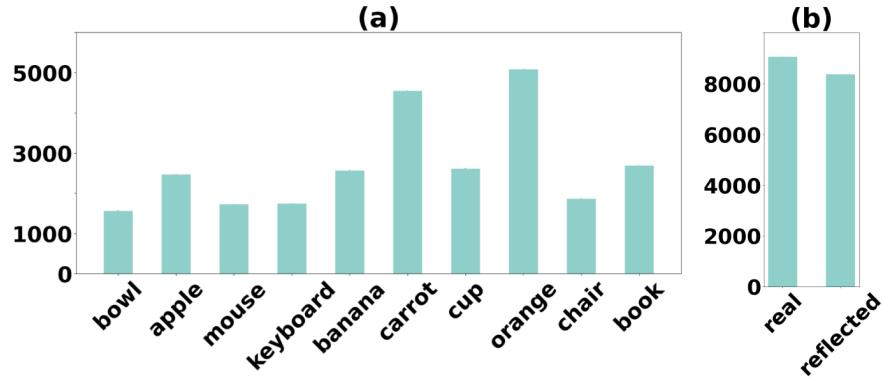
**Fig. 2:** Samples from six categories (i.e., 'banana', 'keyboard', 'chair','book', 'cup', and 'bowl', from left to right) and their corresponding natures (i.e., 'real' and 'reflected' from top to bottom) in the RODD dataset. Note that the objects have been marked with green bounding boxes.

RODD are bowl, apple, mouse, keyboard, banana, carrot, cup, orange, chair, and book, all of which are categories included in the COCO dataset. However, gathering varied images of these objects or their reflected ones in different scenes can be challenging. To address this, we initially sourced images using web crawlers and online repositories that focus on real-world scenarios with reflective surfaces. Additionally, we conducted field photography sessions in various environments such as homes, offices, and public spaces to capture images that include mirrors and other reflective surfaces. To ensure that the dataset was representative of real-world conditions, we made sure to capture images under various lighting conditions and from different angles. The final collection consists of 21,059 images, covering 10 objects (bowl, apple, mouse, keyboard, banana, carrot, cup, orange, chair, and book) and 2 attribute indicating the nature of the objects (i.e., real or reflected). Fig. 2 presents some sample images from RODD, demonstrating that each object category is captured in multiple scenes.

### 3.2   Annotation

This section provides a detailed introduction to the image annotation process, covering three aspects: category, bounding box, and the nature of the object, as follows:
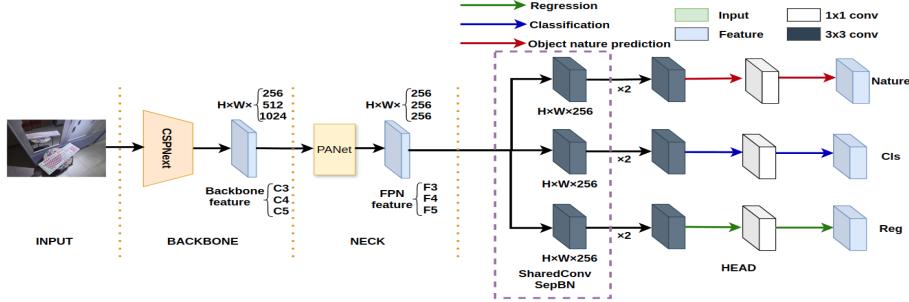
**Fig. 3:** (a) Number of images per category in RODD. (b) Number of images containing real or reflected objects in RODD.

- **Category:** one of: bowl, apple, mouse, keyboard, banana, carrot, cup, orange, chair, and book.
- **Bounding box:** an axis-aligned bounding box that encloses the visible part of the object in the image.
- **Nature of the object:** a real or reflected object.

We follow three steps, i.e., manual annotation, visual inspection, and box refinement, to complete the annotation of images, guided by the annotation guidelines proposed in [7] and [6]. Specifically, all the images are first annotated by an expert, i.e., a student engaged in object detection, during the initial stage. Manual annotation can lead to occasional errors or inconsistencies, prompting the verification team to carefully review the annotated files in the second step. Annotation errors identified by the validation team in the third stage will be sent back to the initial annotation stage for refinement. By employing this three-stage strategy, the dataset ensures its contained objects have high-quality annotation. Fig. 2 displays five examples of box annotations from RODD.

### 3.3   Dataset Statistics

The statistics of the RODD dataset are summarized in Fig. 3. Fig. 3 (a) presents a histogram showing the number of images in the dataset for each category. As observed, the 'orange' category is the most frequent, with 5,086 images. Fig. 3 (b) displays a histogram that illustrates the number of images containing real or reflected objects. This detailed breakdown highlights the distribution and prevalence of each object category within the dataset, providing insight into the dataset's composition and the representation of reflections. To facilitate training and evaluation, the RODD dataset is split into two primary subsets: the training set and the test set, with a ratio of 7:3.

**Fig. 4:** The network structure of the RO-RTMDet detector, inherited from RTMDet, is different from the addition of an additional branch head for the object nature.

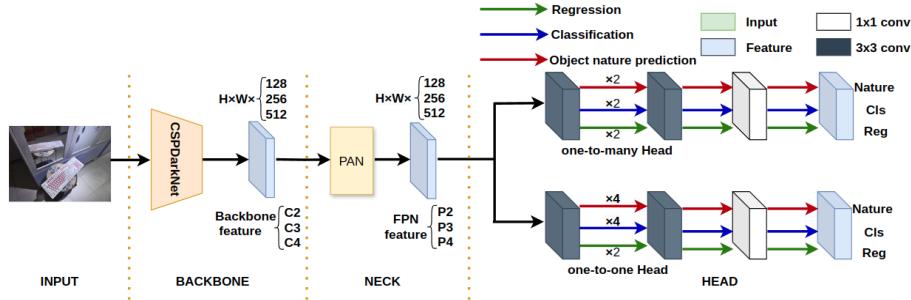## 4    Baseline Detectors for Detecting Reflected Objects

We develop five baseline detectors based on five state-of-the-art object detection algorithms, i.e., RTMDet [5], YOLOv10 [4], YOLOV8 [85], YOLOX [1], and PPYOLOE [2], to facilitate the development of detecting reflected objects. For each model, we add an additional head or branch to predict the nature of the objects without altering the overall framework. The resulting baseline detectors are named RO-RTMDet, RO-YOLOv10, RO-YOLOV8, RO-YOLOX, and RO-PPYOLOE, respectively. Given space constraints and the fact that YOLOv10, YOLOV8, YOLOX, and PPYOLOE are all YOLO variants, we detail only RO-RTMDet and RO-YOLOv10 in the following sections. The extension to YOLOV8, YOLOX, and PPYOLOE is straightforward and will not be elaborated upon here.

### 4.1    RO-RTMDet

The network architecture of the proposed RO-RTMDet is shown in Fig. 4. CSP-Net [3] serves as the backbone, generating output features C3, C4, and C5 with 128, 256, and 512 channels, respectively. These features are fused into CSP-PAFPN [5], the neck of RO-RTMDet, which employs the same block as the backbone. The classification head and the regression head are two parallel components used for classification and regression, respectively, forming the head of the original RTMDet. Building upon the original RTMDet model, we introduce a new classification head to predict the nature of objects (i.e., real or reflected). During RO-RTMDet training, the overall loss of the model is defined as follows:

$$L_{total} = L_{cls} + L_{reg} + \lambda L_{nat}, \tag{1}$$

where $L_{cls}$, $L_{reg}$, and $L_{nat}$ represent the losses for classification, regression, and object nature prediction, respectively. $\lambda$ is a constant that weights the loss for

**Fig. 5:** The network structure of the RO-YOLOv10 detector is inherited from YOLOv10, except for the addition of an additional reflected nature branch head.

the reflected objects prediction head. Below are their specific definitions:

$$L_{cls} = \frac{1}{N_{pos}} \sum_{n=1}^{N_{pos}} \sum_{cls \in classes} -|y_n^{cls} - p_n^{cls}|^\beta ((1 - y_n^{cls})log(1 - p_n^{cls}) + y_n^{cls}log(p_n^{cls})),$$

$$L_{reg} = \frac{1}{N_{pos}} \sum_{n=1}^{N_{pos}} 1 - (\text{IOU}(b_n^t, b_n^p) - \frac{|C - b_n^t \bigcup b_n^p|}{|C|}),$$

$$L_{nat} = \frac{1}{N_{pos}} \sum_{n=1}^{N_{pos}} \sum_{nat \in natures} -|y_n^{nat} - p_n^{nat}|^\beta ((1 - y_n^{nat})log(1 - p_n^{nat}) + y_n^{nat}log(p_n^{nat})),$$

$$(2)$$

where $y_n^{cls}$ and $y_n^{nat}$ are the labeled value of classification and the object nature, $p_n^{cls}$ and $p_n^{nat}$ are the corresponding predictions, $N_{pos}$ is the number of positive anchor, $\beta$ is the hyperparameter for the dynamic scale factor, which is set to 2, $b_n^t$ and $b_n^p$ represent the ground truth bounding boxes and the prediction, respectively; IOU and $C$ are the IOU loss function and the smallest enclosing convex box of these two bounding boxes. We utilize the same training pipeline as RTMDet for training RO-RTMDet.

### 4.2 RO-YOLOv10

The network architecture of the proposed RO-YOLOv10 detector is shown in Fig. 5. RO-YOLOv10 uses a modified CSPDarknet as backbone. It replaces the C2f module used in YOLOv8 [85] with a compact inverted block (CIB) module and introduces an efficient partial self-attention (PSA) module [4]. These features are input into the neck to enhance feature representation, which is made up of the PAN (Path Aggregation Network). The original YOLOv10 model has two types of heads: (1) a one-to-many (o2m) head for regression and classification tasks, and (2) a one-to-one (o2o) head for precise localization. In RO-YOLOv10, we add object nature prediction branch into both of these two heads. During RO-YOLOv10 training, the overall loss of the model is defined as follows:

$$L_{total} = L_{o2m-head} + L_{o2o-head},$$
$$L_{o2m-head} = L_{o2m-cls} + \lambda L_{o2m-nat} + L_{o2m-reg} + L_{o2m-dfl}, \quad (3)$$
$$L_{o2o-head} = L_{o2o-cls} + \lambda L_{o2o-nat} + L_{o2o-reg} + L_{o2o-dfl}$$

In the o2m head, $L_{o2m-cls}$ and $L_{o2m-nat}$ represent the losses for classification and object nature prediction, respectively, while $L_{reg}$ and $L_{dfl}$ indicate the CompleteIntersectin over Union (CIoU) Loss [76] and the Distribution Focal loss (DFL) [77]. Similarly, each loss function in the o2o head carries the same meaning as in the the o2m head. $\lambda$ is a constant that weights the loss for the object nature prediction branch. Below, the $L_{cls}$ and $L_{nat}$ in the o2m head are used as examples to provide their specific definitions. The specific definition of $L_{reg}$ and $L_{dfl}$ are omitted, as it is too intricate to elaborate on here and may divert from the main focus of our discussion. For a comprehensive understanding of $L_{reg}$ and $L_{dfl}$, we recommend referring to the detailed explanations provided in the original documentation by Zheng et al. [76] and Li et al. [77].

$$L_{cls} = \frac{1}{N_{pos}} \sum_{n=1}^{N_{pos}} \sum_{cls \in classes} y_n^{cls} log(p_n^{cls}) + (1 - y_n^{cls})log(1 - p_n^{cls}),$$

$$(4)$$

$$L_{nat} = \frac{1}{N_{pos}} \sum_{n=1}^{N_{pos}} \sum_{nat \in natures} y_n^{nat} log(p_n^{nat}) + (1 - y_n^{nat})log(1 - p_n^{nat}),$$

where $y_n^{cls}$ and $y_n^{nat}$ are the labeled value of classification and the object nature, $p_n^{cls}$ and $p_n^{nat}$ are the corresponding predictions, $N_{pos}$ is the number of positive anchor. We utilize the same training pipeline as YOLOv10 for training RO-YOLOv10.

## 5    Evaluation

### 5.1    Evaluation Metrics

In the experiment, the proposed baseline detectors are evaluated for the performance by using two common metrics, i.e., average precision (AP) and mean average precision (mAP). IOU (Intersection over Union) measures the overlap between the predicted bounding box (bbox) and the ground truth bbox. In object detection tasks, a complete prediction comprises two main components: first, the model must identify specific objects within a given image, and second, it needs to accurately determine their respective locations. Specifically, precision is the proportion of objects predicted by the model that match the real objects, whereas recall measures the proportion of real objects detected by the model. These two measures are combined in mAP, which highlights the significance of properly balancing each during the evaluation process.

Guided by the COCO evaluation [6], three IoU thresholds are used: fixed thresholds at 0.5 and 0.75 and a range threshold from 0.5 to 0.95 with a step size of 0.05. The corresponding average precisions (APs) are evaluated under
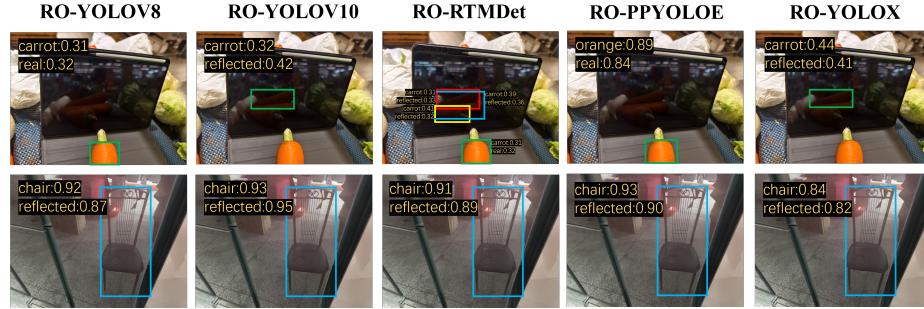
**Table 1:** The evaluation results of the five proposed baseline detectors, i.e., RO-YOLOv8, RO-YOLOv10, RO-RTMDet, RO-YOLOX, and RO-PPYOLOE, on the RODD dataset. It is important to note that $AP_c$, $AP_n$, and $AP_{cn}$ represent the precision metrics for predicting the object's category, the object's nature, and the combination of both.

| | $\{AP_c, AP_n, AP_{cn}\}$@0.5 | $\{AP_c, AP_n, AP_{cn}\}$@0.75 | $\{mAP_c, mAP_n, mAP_{cn}\}$ |
|---|---|---|---|
| RO-YOLOv8 | (0.693, 0.736, 0.583) | (0.663, 0.700, 0.561) | (0.638, 0.671, 0.540) |
| RO-YOLOv10 | (**0.766**, **0.789**, **0.585**) | (**0.726**, **0.751**, **0.562**) | (**0.697**, **0.723**, **0.542**) |
| RO-RTMDet | (0.650, 0.706, 0.526) | (0.620, 0.681, 0.506) | (0.591, 0.654, 0.485) |
| RO-YOLOX | (0.668, 0.700, 0.497) | (0.594, 0.633, 0.447) | (0.527, 0.558, 0.387) |
| RO-PPYOLOE | (0.579, 0.654, 0.532) | (0.550, 0.621, 0.512) | (0.526, 0.590, 0.486) |

these IoU thresholds, denoted as AP@0.5, AP@0.75, and AP@[.50:.05:.95], respectively. In the experiment, COCO mAP is employed to evaluate the performance of detectors in detecting reflected objects. Following [8–12], we use $AP_c$, $AP_n$, and $AP_{cn}$ to represent the precision metrics for predicting the object's category, the object's nature, and their combination, respectively. Additionally, an extra prefix 'm' is added to represent mean AP, i.e., mAP.

## 5.2 Evaluation Results

**Overall performance.** We extensively evaluate RO-YOLOv8, RO-YOLOv10, RO-RTMDet, RO-YOLOX, and RO-PPYOLOE, the five baseline detectors proposed in this paper, on the RODD dataset. Table 1 presents the evaluation results using the three precision metrics defined in Section 5.1, i.e., $AP_c$, $AP_n$, and $AP_{cn}$. As can be seen, RO-YOLOv10 is the best detector, consistently achieving highest Average Precisions (APs) compared to other detectors. Besides, the evaluation results of object categories in the five baseline detectors show lower Average Precision (AP) at fixed IoUs, specifically at 0.5 and 0.75, as well as lower mean AP compared to the evaluation results of the object's nature. For all these detectors, the differences between the APs for categories and the APs for nature exceed 2%. Notably, the largest gap is up to 7.5%, observed in the RO-PPYOLOE detector. This indicates that identifying and localizing the object itself is more challenging than recognizing the object's nature in the images within RODD. This discrepancy may be explained by the fact that the number of natures (i.e., 2) is significantly smaller than the number of categories (i.e., 10) in RODD. The smaller number of natures means that there is more sufficient training data for each nature class, leading to a more effective training process. In contrast, the larger number of categories results in a more dispersed dataset for each object type, making it harder for the model to learn and generalize effectively. This imbalance in the training data distribution can lead to more robust performance in recognizing the object's nature compared to identifying and localizing the object itself. This suggests that as the number of categories increases, it becomes necessary to implement effective methods that take this factor into account. It is noteworthy that when conventional object detection and the prediction
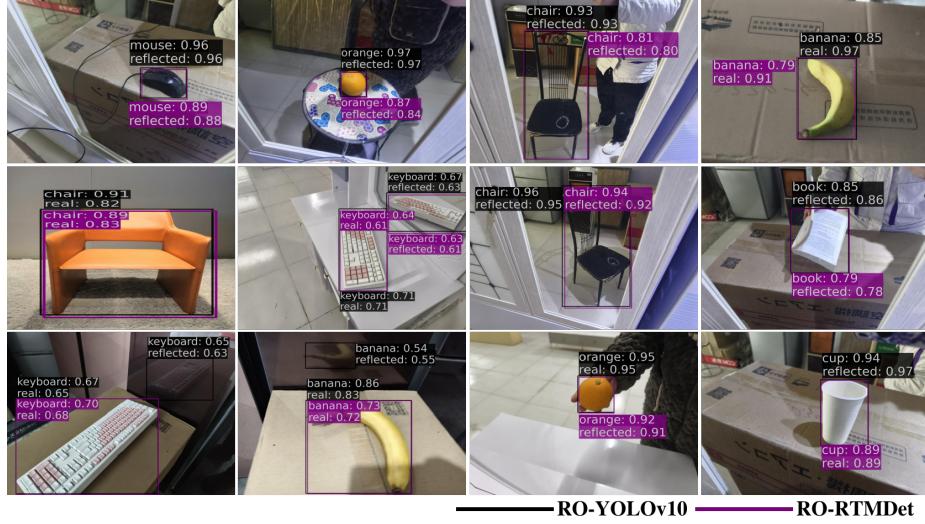
**Fig. 6:** A qualitative comparison of the five detectors on 2 samples from the carrot and the chair category, respectively. Note that all these detectors successfully detect the chair but fail to detect the carrots correctly. The predicted bounding box, object category, object nature, and the corresponding scores have been marked in the images.

of objects' nature are integrated into a composite task, namely detecting reflected objects, the detectors exhibit lower performance compared to handling each task independently. This observation suggests that the task of detecting reflected objects we propose in this work is more challenging than conventional object detection tasks. This lower performance highlights the importance of developing specialized algorithms and training strategies that can better manage the intricacies of this composite task.

**Table 2:** Comparison of the $mAP_{cn}$ of the five baseline detectors on the RODD dataset. It is important to note that $mAP_{cn}$ is the mAP for prediction of the composite of the object's category and its nature.

|  | bowl | apple | mouse | keyboard | banana | carrot | cup | orange | chair | book |
|---|---|---|---|---|---|---|---|---|---|---|
| $mAP_{cn}$(RO-YOLOv8) | 0.691 | 0.561 | 0.744 | 0.780 | **0.683** | **0.484** | 0.787 | 0.522 | 0.837 | 0.623 |
| $mAP_{cn}$(RO-YOLOv10) | **0.718** | **0.592** | 0.761 | 0.775 | 0.679 | 0.483 | **0.847** | **0.650** | **0.841** | 0.629 |
| $mAP_{cn}$(RO-RTMDet) | 0.673 | 0.478 | **0.781** | **0.788** | 0.645 | 0.260 | 0.736 | 0.387 | 0.829 | 0.338 |
| $mAP_{cn}$(RO-YOLOX) | 0.629 | 0.526 | 0.641 | 0.612 | 0.492 | 0.259 | 0.660 | **0.650** | 0.632 | 0.554 |
| $mAP_{cn}$(RO-PPYOLOE) | 0.690 | 0.472 | 0.763 | 0.787 | 0.622 | 0.254 | 0.752 | 0.387 | 0.825 | 0.349 |

**Performance on per Category.** To get a deeper analysis and understanding of the performance in detecting the nature of objects using our proposed baseline detectors, we further conduct performance evaluations on each category. Table 2 presents the $mAP_{cn}$ of the five detectors evaluated on RODD. As observed, the five detectors demonstrate their best performance on the chair category and their worst on the carrot category. The $mAP_{cn}$ values for the chair category are all above 80%, except for the RO-YOLOX detector. In contrast, for the carrot category, the $mAP_{cn}$ values are all below 50%, with RO-RTMDet, RO-YOLOX, and RO-PPYOLOE even falling below 30%. This disparity could be explained by the fact that images typically contain a single chair target, often presented at

**Fig. 7:** A qualitative evaluation was conducted on 12 samples from RODD. The first two rows display examples accurately predicting the nature of objects using RO-YOLOv10 and RO-RTMDet detectors, while last row shows error detection results generated by these two detectors. Note that the predicted bounding box, object category, object nature, and the corresponding scores have been marked in the images.

standard sizes. Conversely, images frequently contain numerous carrot targets, leading to clusters and occlusions. See Fig. 6 for an intuitive comparison of examples of these two categories, the first row is the qualitative results of the carrot obtained by the five detectors. The challenges in detecting these carrots on the screen are compounded by their size and the properties of the reflective surfaces themselves. The screen has a lower reflectivity coefficient compared to mirrors, which obscures the appearance of the carrots when reflected. This reduced reflectivity makes it difficult for detectors such as RO-YOLOV8 and RO-PPYOLOE to accurately identify the characteristics of carrots in their reflected versions. Moreover, carrots, being similar in color and appearance to oranges, can further confuse detectors like RO-PPYOLOE, leading to misjudgments in object categorization. In contrast, these detectors succeed in detecting the chair due to the high reflectivity of the mirror and the absence of cluttered backgrounds. In addition, the experimental results in Table 2 also demonstrate that the same detectors will achieve varying performance across different categories. This variation may be attributed to inherent differences in object characteristics, such as size, shape, texture, and context within the images, as well as the unbalanced distribution of categories. These results underscore the importance of considering object-specific challenges in detecting reflected objects.

**Qualitative Evaluation.** Given the potential for overwhelming viewers with too many methods in a single image, Fig. 7 presents qualitative detection results from just the RO-YOLOv10 and RO-RTMDet detectors. The first two rows dis-

play eight correctly predicted samples, while the third row shows examples where the detectors inaccurately predicted the object's nature. In these cases, reflected objects might blend into low-light backgrounds or lack distinct texture features (i.e., the first and second samples), or their mirrored background may resemble the real background (i.e., the third and fourth samples), leading to missed or inaccurate detections. This evaluation highlights that in complex scenes, the detectors are prone to struggle with accurately identifying the nature of objects.

**Table 3:** The ablation study of the RO-YOLOv10 model is conducted on RODD using various weighting coefficients.

| $\lambda$ | $\{AP_c, AP_n, AP_{cn}\}$@0.5 | $\{AP_c, AP_n, AP_{cn}\}$@0.75 | $\{mAP_c, mAP_n, mAP_{cn}\}$ |
|---|---|---|---|
| 0.2 | (0.744, 0.765, 0.552) | (0.706, 0.731, 0.526) | (0.682, 0.706, 0.514) |
| 0.4 | (0.745, 0.779, 0.568) | (0.710, 0.742, 0.549) | (0.682, 0.716, 0.530) |
| 0.6 | (0.742, 0.776, 0.553) | (0.708, 0.738, 0.533) | (0.679, 0.710, 0.515) |
| 0.8 | (0.751, 0.776, 0.567) | (0.716, 0.739, 0.544) | (0.684, 0.708, 0.526) |
| 1.0 | (**0.766**, 0.789, **0.585**) | (**0.726**, 0.751, **0.562**) | (**0.697**, 0.719, **0.542**) |
| 1.2 | (0.749, 0.772, 0.560) | (0.708, 0.729, 0.532) | (0.682, 0.705, 0.518) |
| 1.4 | (0.743, 0.771, 0.563) | (0.705, 0.738, 0.540) | (0.672, 0.705, 0.519) |
| 1.6 | (0.758, **0.790**, 0.576) | (0.721, **0.754**, 0.557) | (0.688, **0.723**, 0.534) |
| 1.8 | (0.747, 0.775, 0.524) | (0.712, 0.739, 0.545) | (0.680, 0.707, 0.535) |
| 2.0 | (0.743, 0.784, 0.574) | (0.707, 0.746, 0.549) | (0.674, 0.710, 0.525) |

### 5.3   Ablation Study

We train the RO-YOLOv10 model on RODD using different weighting coefficients, i.e., $\lambda$ in Eq. (3), which varies from 0.2 to 2.0 in steps of 0.2, in order to study the impact of the coefficient for weighting the loss of predicting the nature of objects. This experiment aims to determine how different weightings influence the model's ability to balance the two tasks: detecting the objects and predicting their nature. By adjusting $\lambda$, we can observe how the model prioritizes the nature prediction task relative to the conventional object detection task. Table 3 presents the experimental results for the mAP and the AP at fixed IoUs (0.5 and 0.75). As can be seen, RO-YOLOv10 can obtain part of the best APs when $\lambda$ is set to 1.0 or 1.6. An obvious partial difference between $AP_c$ and $AP_n$ is evident. In general, as $\lambda$ ranges from 0.2 to 1.6, $AP_c$ initially increases and then decreases, while $AP_n$ consistently rises, reaching their optimal values at 1.0 and 1.6, respectively. A compromise is achieved when $\lambda$ is set to 1.0, serving as the default setting, where $AP_{cn}$ reaches its highest value. Specifically, the highest $mAP_c$, which equals 0.697, occurs at $\lambda = 1.0$; the highest $mAP_n$, which equals 0.723, occurs at $\lambda = 1.6$; and the highest $mAP_{cn}$, which equals 0.542, occurs at $\lambda = 1.0$. The results suggest that there may be a counteracting impact between object localization and prediction of objects' nature when these two tasks are done concurrently as a composite task. More effective methods for mitigating this counteracting effect are needed.

# 6    Conclusions

In this paper, we investigated the underexplored challenge of reflective object detection and introduced the Reflective Object Detection Dataset (RODD), an extensive benchmark specifically designed for this task. RODD includes 10 categories and 21,059 images of real or reflected objects in various backgrounds, accompanied by standard annotations of bounding boxes and the nature of the objects (real or reflected), distinguishing it from traditional object detection benchmarks. In addition to introducing RODD, we adapted five state-of-the-art object detection models to this challenging task and presented baseline results. The experimental findings reveal the limitations of current methods when applied to reflected object detection, underscoring the necessity for specialized approaches. By releasing RODD, we aim to foster and advance future research in detecting reflected objects. This dataset provides a valuable resource for developing and evaluating new methods, ultimately contributing to improved performance in applications such as surveillance, autonomous driving, service robots, smart homes, and medical imaging.

# References

1. Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "Yolox: Exceeding yolo series in 2021," *arXiv preprint arXiv:2107.08430*, 2021. 8
2. S. Xu, X. Wang, W. Lv, Q. Chang, C. Cui, K. Deng, G. Wang, Q. Dang, S. Wei, Y. Du *et al.*, "Pp-yoloe: An evolved version of yolo," *arXiv preprint arXiv:2203.16250*, 2022. 8
3. C.-Y. Wang, H.-Y. M. Liao, Y.-H. Wu, P.-Y. Chen, J.-W. Hsieh, and I.-H. Yeh, "Cspnet: A new backbone that can enhance learning capability of cnn," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 390–391. 8
4. A. Wang, H. Chen, L. Liu, K. Chen, Z. Lin, J. Han, and G. Ding, "Yolov10: Real-time end-to-end object detection," *arXiv preprint arXiv:2405.14458*, 2024. 4, 8, 9
5. C. Lyu, W. Zhang, H. Huang, Y. Zhou, Y. Wang, Y. Liu, S. Zhang, and K. Chen, "Rtmdet: An empirical study of designing real-time object detectors," *arXiv preprint arXiv:2212.07784*, 2022. 8
6. T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*.    Springer, 2014, pp. 740–755. 4, 5, 7, 10
7. M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, pp. 303–338, 2010. 5, 7
8. L. Qin, H. Zhou, Z. Wang, J. Deng, Y. Liao, and S. Li, "Detection beyond what and where: a benchmark for detecting occlusion state," in *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*.    Springer, 2022, pp. 464–476. 11
9. Y. Guo, Y. Chen, J. Deng, S. Li, and H. Zhou, "Identity-preserved human posture detection in infrared thermal images: A benchmark," *Sensors*, vol. 23, no. 1, p. 92, 2022. 11

10. Y. Wu, H. Ye, Y. Yang, Z. Wang, and S. Li, "Liquid content detection in transparent containers: A benchmark," *Sensors*, vol. 23, no. 15, p. 6656, 2023. 11

11. Y. Li, Y. Wu, X. Chen, H. Chen, D. Kong, H. Tang, and S. Li, "Beyond human detection: A benchmark for detecting common human posture," *Sensors*, vol. 23, no. 19, p. 8061, 2023. 11

12. H. Zhou, Y. Wu, J. Li, L. Pan, H. Ye, and S. Li, "Beyond animal detection: a benchmark for detecting animal age group," in *Fifth International Conference on Artificial Intelligence and Computer Science (AICS 2023)*, vol. 12803.  SPIE, 2023, pp. 506–515. 11

13. C.-T. Chien, R.-Y. Ju, K.-Y. Chou, C.-S. Lin, and J.-S. Chiang, "Yolov8-am: Yolov8 with attention mechanisms for pediatric wrist fracture detection," *arXiv preprint arXiv:2402.09329*, 2024.

14. Z. Zheng, P. Wang, D. Ren, W. Liu, R. Ye, Q. Hu, and W. Zuo, "Enhancing geometric factors in model learning and inference for object detection and instance segmentation," *IEEE transactions on cybernetics*, vol. 52, no. 8, pp. 8574–8586, 2021.

15. H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 658–666.

16. Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," *Proceedings of the IEEE*, vol. 111, pp. 257–276, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:152282225 1

17. Z.-Q. Zhao, P. Zheng, S. tao Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, pp. 3212–3232, 2018. [Online]. Available: https://api.semanticscholar.org/CorpusID:49862415 1

18. A. B. Amjoud and M. AMROUCH, "Object detection using deep learning, cnns and vision transformers: A review," *IEEE Access*, vol. 11, pp. 35 479–35 516, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:258077107 1

19. L. Liu, W. Ouyang, X. Wang, P. W. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, "Deep learning for generic object detection: A survey," *International Journal of Computer Vision*, vol. 128, pp. 261 – 318, 2018. [Online]. Available: https://api.semanticscholar.org/CorpusID:52177403 1

20. D. Owen and P.-L. Chang, "Detecting reflections by combining semantic and instance segmentation," *ArXiv*, vol. abs/1904.13273, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:140289466 2, 5

21. X. Yang, H. Mei, K. Xu, X. Wei, B. Yin, and R. W. H. Lau, "Where is my mirror?" *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8808–8817, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:201666156 2, 5

22. J. Lin, X. Y. Tan, and R. W. H. Lau, "Learning to detect mirrors from videos via dual correspondences," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9109–9118, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:261081000 2

23. G.-P. Ji, K. Fu, Z. Wu, D.-P. Fan, J. Shen, and L. Shao, "Full-duplex strategy for video object segmentation," *Computational Visual Media*, vol. 9, pp. 155–175, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:236950747 2, 5

24. F. Liu, Y. Liu, J. Lin, K. Xu, and R. W. H. Lau, "Multi-view dynamic reflection prior for video glass surface detection," in *AAAI Conference on Artificial Intelligence*, 2024. [Online]. Available: https://api.semanticscholar.org/CorpusID:268678293 2, 5

25. J. Lin, G. Wang, and R. W. H. Lau, "Progressive mirror detection," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3694–3702, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:219964489 2, 5

26. Y. Wu, Y. Wang, Y. Liao, F. Wu, H. Ye, and S. Li, "Tracking transforming objects: A benchmark," *ArXiv*, vol. abs/2404.18143, 2024.

27. W. H. Ittelson, L. Mowafy, and D. F. Magid, "The perception of mirror-reflected objects," *Perception*, vol. 20, pp. 567 – 584, 1991.

28. E. Gregory and M. McCloskey, "Mirror-image confusions: Implications for representation and processing of object orientation," *Cognition*, vol. 116, no. 1, pp. 110–129, 2010.

29. I. Bianchi and U. Savardi, "The relationship perceived between the real body and the mirror image," *Perception*, vol. 37, no. 5, pp. 666–687, 2008.

30. O. H. Turnbull and R. A. Mccarthy, "Failure to discriminate between mirror-image objects: A case of viewpoint-independent object recognition?" *Neurocase*, vol. 2, no. 1, pp. 63–72, 1996.

31. H. Yoshimura and T. Tabata, "Relationship between frames of reference and mirror-image reversals," *Perception*, vol. 36, no. 7, pp. 1049–1056, 2007.

32. B. R. Ray, M. Aalsma, N. D. Zaller, E. B. Comartin, and E. Sightes, "The perpetual blind spot in public health surveillance." *Journal of correctional health care : the official journal of the National Commission on Correctional Health Care*, 2022. 2

33. Y. Shen and W. Q. Yan, "Blind spot monitoring using deep learning," *2018 International Conference on Image and Vision Computing New Zealand (IVCNZ)*, pp. 1–5, 2018. 2

34. C. Singhal and S. Barick, "Ecms: Energy-efficient collaborative multi-uav surveillance system for inaccessible regions," *IEEE Access*, vol. 10, pp. 95 876–95 891, 2022. 2

35. G. Sasikala and V. R. Kumar, "Development of advanced driver assistance system using intelligent surveillance," *International Conference on Computer Networks and Communication Technologies*, 2018.

36. M. S. A. Latif, A. A. Ismail, and A. Zariman, "Smart mirror for home automation," 2020.

37. M. A. Hossain, P. K. Atrey, and A. El Saddik, "Smart mirror for ambient home environment," in *2007 3rd IET International Conference on Intelligent Environments*. IET, 2007, pp. 589–596.

38. W. M. Freysteinson, "Assessing the mirrors in long-term care homes: A preliminary survey," *Journal of Gerontological Nursing*, vol. 36, no. 1, pp. 34–40, 2010.

39. D. Kilic and N. Sailaja, "User-centred repair: From current practices to future design," in *International Conference on Human-Computer Interaction*. Springer, 2024, pp. 52–71.

40. P. Sareen, K. A. Ehinger, and J. M. Wolfe, "Through the looking-glass: Objects in the mirror are less real," *Psychonomic Bulletin & Review*, vol. 22, pp. 980–986, 2015.

41. S. J. Nightingale, K. A. Wade, H. Farid, and D. G. Watson, "Can people detect errors in shadows and reflections?" *Attention, Perception & Psychophysics*, vol. 81, pp. 2917 – 2943, 2019.

42. C. D. Mole and R. M. Wilkie, "Looking forward to safer hgvs: The impact of mirrors on driver reaction times," *Accident Analysis & Prevention*, vol. 107, pp. 173–185, 2017.

43. I. Bianchi, M. Bertamini, and U. Savardi, "Differences between predictions of how a reflection behaves based on the behaviour of an object, and how an object behaves based on the behaviour of its reflection," *Acta Psychologica*, vol. 161, pp. 54–63, 2015.

44. D. Baysal, "Assessment of vehicular vision obstruction due to driver-side b-pillar and remediation with blind spot eliminator," *arXiv preprint arXiv:2302.07088*, 2023.

45. B. L. R. Stojkoska and K. V. Trivodaliev, "A review of internet of things for smart home: Challenges and solutions," *Journal of cleaner production*, vol. 140, pp. 1454–1464, 2017.

46. A. Saad al sumaiti, M. H. Ahmed, and M. M. Salama, "Smart home activities: A literature review," *Electric Power Components and Systems*, vol. 42, no. 3-4, pp. 294–305, 2014.

47. D. Marikyan, S. Papagiannidis, and E. Alamanos, "A systematic review of the smart home literature: A user perspective," *Technological Forecasting and Social Change*, vol. 138, pp. 139–154, 2019. 2

48. B. K. Sovacool and D. D. F. Del Rio, "Smart home technologies in europe: A critical review of concepts, benefits, risks and policies," *Renewable and sustainable energy reviews*, vol. 120, p. 109663, 2020. 2

49. A. Chakraborty, M. Islam, F. Shahriyar, S. Islam, H. U. Zaman, and M. Hasan, "Smart home system: a comprehensive review," *Journal of Electrical and Computer Engineering*, vol. 2023, no. 1, p. 7616683, 2023. 2

50. M.-Z. Poh, D. McDuff, and R. Picard, "A medical mirror for non-contact health monitoring," in *ACM SIGGRAPH 2011 Emerging Technologies*, 2011, pp. 1–1.

51. V. M. Soppimath, M. G. Hudedmani, M. Chitale, M. Altaf, A. Doddamani, and D. Joshi, "The smart medical mirror-a review," *International Journal of Advanced Science and Engineering*, vol. 6, no. 1, pp. 1244–1250, 2019.

52. C. Bichlmeier, S. M. Heining, M. Feuerstein, and N. Navab, "The virtual mirror: a new interaction paradigm for augmented reality environments," *IEEE Transactions on Medical Imaging*, vol. 28, no. 9, pp. 1498–1510, 2009.

53. J. L. JASMINE, M. A. JINU, T. DHANALAKSHMI, B. RAKSHAMBIGAI, and S. MADHUMITHA, "Medical mirror for health care," in *National Conference on Recent Advancements in Communication*, vol. 7, no. 08, 2020.

54. C. Bichlmeier, T. Sielhorst, and N. Navab, "The tangible virtual mirror: New visualization paradigm for navigated surgery," in *International Workshop on Augmented Reality Environments for Medical Imaging and Computer-Aided Surgery, Copenhagen, Denmark*, 2006.

55. A. Roy and M. Mukhopadhyay, "Situs inversus totalis: operating on the mirror image," *Hellenic Journal of Surgery*, vol. 92, pp. 150–152, 2020.

56. J. Jansen, L. Dubois, R. Schreurs, P. J. Gooris, T. J. Maal, L. F. Beenen, and A. G. Becking, "Should virtual mirroring be used in the preoperative planning of an orbital reconstruction?" *Journal of Oral and Maxillofacial Surgery*, vol. 76, no. 2, pp. 380–387, 2018.

57. J. Abi-Rafeh, D. Zammit, M. M. Jaberi, B. Al-Halabi, and S. Thibaudeau, "Nonbiological microsurgery simulators in plastic surgery training: a systematic review," *Plastic and Reconstructive Surgery*, vol. 144, no. 3, pp. 496e–507e, 2019.

58. S. Dhalwar, S. Ruby, S. Salgar, and B. Padiri, "Image processing based traffic convex mirror detection," in *2019 Fifth International Conference on Image Information Processing (ICIIP)*.   IEEE, 2019, pp. 41–45.

59. S. Sadrhaghighi, M. Dolati, M. Ghaderi, and A. Khonsari, "Monitoring openflow virtual networks via coordinated switch-based traffic mirroring," *IEEE Transactions on Network and Service Management*, vol. 19, no. 3, pp. 2219–2237, 2022.

60. J. Díaz, E. Ros, S. Mota, G. Botella, A. Cañas, and S. Sabatini, "Optical flow for cars overtaking monitor: the rear mirror blind spot problem," *Ecovision (European research project)*, 2003. 2

61. C. Zhang, F. Steinhauser, G. Hinz, and A. Knoll, "Traffic mirror-aware pomdp behavior planning for autonomous urban driving," in *2022 IEEE Intelligent Vehicles Symposium (IV)*.   IEEE, 2022, pp. 323–330. 2

62. D. Li, H. Hagura, T. Miyabashira, Y. Kawai, and S. Ono, "Traffic mirror detection and annotation methods from street images of open data for preventing accidents at intersections by alert," in *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*.   IEEE, 2023, pp. 3256–3262. 2

63. D. Mittal, V. Verma, and R. Rastogi, "A comparative study and new model for smart mirror," *International Journal of Scientific Research in Research Paper, Computer Science and Engineering*, vol. 5, no. 6, pp. 58–61, 2017.

64. O.-S. Loizides, M. Kastrinakis, G. Badawy, M. N. Smadi, D. Murray, and P. Koutsakis, "Efficient policing for screen mirroring traffic," *Human-centric Computing and Information Sciences*, vol. 8, pp. 1–14, 2018.

65. I. Noriaki, O. Shintaro, S. Yoshihiro, O. Kazuya, and R. Grewe, "Collision risk prediction utilizing road safety mirrors at blind intersections," in *27th International Technical Conference on the Enhanced Safety of Vehicles (ESV) National Highway Traffic Safety Administration*, no. 23-0164, 2023. 2

66. V. N. Lu, J. Wirtz, W. H. Kunz, S. Paluch, T. Gruber, A. Martins, and P. G. Patterson, "Service robots, customers and service employees: what can we learn from the academic literature and where are the gaps?" *Journal of Service Theory and Practice*, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:212964611 2

67. D. Park and Y. H. Park, "Identifying reflected images from object detector in indoor environment utilizing depth information," *IEEE Robotics and Automation Letters*, vol. 6, pp. 635–642, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:231714728 2

68. D. Damodaran, S. Mozaffari, S. Alirezaee, and M. J. Ahamed, "Experimental analysis of the behavior of mirror-like objects in lidar-based robot navigation," *Applied Sciences*, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:257188940 2

69. R. Bajcsy, S. W. Lee, and A. Leonardis, "Detection of diffuse and specular interface reflections and inter-reflections by color image segmentation," *International Journal of Computer Vision*, vol. 17, pp. 241–272, 1996. [Online]. Available: https://api.semanticscholar.org/CorpusID:6799933 5

70. A. DelPozo and S. Savarese, "Detecting specular surfaces on natural images," *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8, 2007. [Online]. Available: https://api.semanticscholar.org/CorpusID:16144262 5

71. R. Abiko and M. Ikehara, "Single image reflection removal based on gan with gradient constraint," *IEEE Access*, vol. 7, pp. 148 790–148 799, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:204820165 5

72. N. Arvanitopoulos, R. Achanta, and S. Süsstrunk, "Single image reflection suppression," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1752–1760, 2017. [Online]. Available: https://api.semanticscholar.org/CorpusID:13095034 5

73. A. Warren, K. Xu, J. Lin, G. K. Tam, and R. W. Lau, "Effective video mirror detection with inconsistent motion cues," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 17 244–17 252.

74. K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969. 4

75. S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015. 4

76. Z. Zheng, P. Wang, D. Ren, W. Liu, R. Ye, Q. Hu, and W. Zuo, "Enhancing geometric factors in model learning and inference for object detection and instance segmentation," *IEEE Transactions on Cybernetics*, vol. 52, pp. 8574–8586, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:218538057 10

77. X. Li, W. Wang, L. Wu, S. Chen, X. Hu, J. Li, J. Tang, and J. Yang, "Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection," *ArXiv*, vol. abs/2006.04388, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:219531292 10

78. R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448. 4

79. R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587. 4

80. Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6154–6162. 4

81. W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*.   Springer, 2016, pp. 21–37. 4

82. T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988. 4

83. M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10 781–10 790. 4

84. A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012. 4

85. Ultralytics, "Yolov8: Real-time object detection and image segmentation," 2023, accessed: 2024-06-27. [Online]. Available: https://github.com/ultralytics/ultralytics 4, 8, 9

86. M. Everingham, L. V. Gool, C. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," 2012, accessed: 2024-06-27. [Online]. Available: http://host.robots.ox.ac.uk/pascal/VOC/ 4

87. A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 4

88. D. Du, Y. Qi, H. Yu, Y. Yang, K. Duan, G. Li, W. Zhang, Q. Huang, and Q. Tian, "The unmanned aerial vehicle benchmark: Object detection and tracking," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 370–386. 5