# A Comprehensive Review on 3D Object Detection and 6D Pose Estimation With Deep Learning

**SABERA HOQUE**[1], **MD. YASIR ARAFAT**[2], **SHUXIANG XU**[1],
**ANANDA MAITI**[1], (Member, IEEE), AND **YUCHEN WEI**[1]
[1]School of Information and Communication Technology, University of Tasmania, Newnham, TAS 7248, Australia
[2]Bundle Australia Pty Ltd, Pialba 4655, Australia

Corresponding author: Sabera Hoque (sabera.hoque@utas.edu.au)

**ABSTRACT** Nowadays, computer vision with 3D (dimension) object detection and 6D (degree of freedom) pose assumptions are widely discussed and studied in the field. In the 3D object detection process, classifications are centered on the object's size, position, and direction. And in 6D pose assumptions, networks emphasize 3D translation and rotation vectors. Successful application of these strategies can have a huge impact on various machine learning-based applications, including the autonomous vehicles, the robotics industry, and the augmented reality sector. Although extensive work has been done on 3D object detection with a pose assumption from RGB images, the challenges have not been fully resolved. Our analysis provides a comprehensive review of the proposed contemporary techniques for complete 3D object detection and the recovery of 6D pose assumptions of an object. In this review research paper, we have discussed several proposed sophisticated methods in 3D object detection and 6D pose estimation, including some popular data sets, evaluation matrix, and proposed method challenges. Most importantly, this study makes an effort to offer some possible future directions in 3D object detection and 6D pose estimation. We accept the autonomous vehicle as the sample case for this detailed review. Finally, this review provides a complete overview of the latest in-depth learning-based research studies related to 3D object detection and 6D pose estimation systems and points out a comparison between some popular frameworks. To be more concise, we propose a detailed summary of the state-of-the-art techniques of modern deep learning-based object detection and pose estimation models.

**INDEX TERMS** Machine learning, deep neural network, computer vision, image processing, convolutional neural network, 3D object detection, 6D pose estimation.

## I. INTRODUCTION

Recently with the advancement of three-dimensional (3D) technology, the reconstruction of 3D models with pose assumptions has become a popular research topic. The main purpose of 3D model identification is to extract powerful features from RGB or RGBD images that can automatically improve the transportation system. Advanced models can make the map smarter and reduce vehicle costs. There are many challenges to this research concept, such as differentiation of perspectives, scaling, posture determination, illumination change, partial inclusion, adaptation detection, and background clutter.

Although many approaches and algorithms have been proposed and implemented for 2D image detection, the

challenges of retrieving 3D objects from 2D images are still being explored. Moreover, estimating poses from this model is also important for the robot industry. One of the core examples in the 3D object detection and pose estimation research sector is the autonomous vehicle, where image detection plays a vital role in recovering 3D objects from 2D images [109]. The modern world is automatically moving towards an intelligent transportation system that requires the successful implementation of autonomous vehicles. The most important issue for self-driving systems is how various modern technologies can be applied to enhance the efficiency of self-driving vehicles.

The great debate in smart car systems is which one works better for object detection, the LiDAR (Light Detection and Ranging) or camera. Also, it needs to be studied whether it is effective to use a combination of the LiDAR and camera systems. For example, both Waymo and Uber include LiDAR

The associate editor coordinating the review of this manuscript and approving it for publication was Claudio Cusano.

where Tesla only uses cameras in their smart car system. Yet no technology has been universally accepted as a final self-driving solution on the road [182].

LiDAR, a proven technique of measuring distance, applies light pulses to determine both the distance and range of the surrounding object to avoid a collision and reduce the vehicle's speed. The technology helps self-propelled vehicles create visual 3D maps using on-board software, sending millions of pulses per second based on readings from light pulses and providing the vehicle with information about its surroundings. LiDAR is used in conjunction with cameras that provide a 360-degree view of the surroundings in self-driving cars, so they are not a standalone solution in themselves.

The camera provides images in intelligent car software that can analyze with a high level of accuracy using AI (artificial intelligence). The autopilot system uses cameras to provide a 360-degree view of its surroundings. The system returns entirely visual data from the lens's optics to on-board software and does not rely on the range and detection like LiDAR for situation analysis. With the development of NNs (Neural Network) and CV (Computer Vision) algorithms, objects can be identified to provide surrounding information while driving. This helps the car avoid collisions, slow down or brake when there is traffic, change lanes safely, and read text from road or highway signs using OCR (Recognition of Optical Character).

Although LiDAR has been proven to see things even in dangerous or foggy weather, it is not always reliable, as it is affected by wavelength stability, temperature, and detective sensitivity. This difficulty makes LiDAR technology more expensive. Moreover, LiDAR requires more space to apply to cars, thus making self-driving cars look bulky and less attractive. On the other hand, cameras are better, easy to implement, and comparatively less expensive in visual recognition. The software requires more data processing to create images and identify objects for LiDAR data than visual data. Finally, the camera has been implemented with Tesla as a standalone system; however, other OEMs(original equipment manufacturer) believe that applying other sensors, including radar, to detect range and distance can improve the performance of self-driving.

The ultimate visual recognition system also required the accurate calculation of other vehicles pose on the road. Without predicting the actual pose of other vehicles, an autonomous car cannot make accurate decisions on whether to slow, brake or change direction. Recent state-of-the-art RGB-based 6 DoF (Degree of Freedom) pose estimation frameworks can be divided into two stages [51], [116], [241], including the object detection with 3D rotation by applying a trained framework and the estimation of 3D translation and 3D orientation (6D pose estimation) via relative distance estimation. Basically, the camera pose estimation is related to object localization, coordinates, and orientation. It is a crucial task not only for the autonomous car but also for the robot and navigation technology, the medical sector, and AR (augmented reality) [269]. In this review,

we will mainly focus on the papers that work on the autonomous car and predicting the position of on-road cars or obstacles.

The rest of the section is organized as follows: In I-A we present the contributions of this review article of deep learning for 3D Object Detection and 6D Pose Estimation. In I-B we have shown the difference between our review and other existing review articles. In I-C and I-D we have discussed the pervasiveness of both 3D object detection and 6D pose estimation. Finally, in I-E, we have briefly discussed the paper collection process.

### A. CONTRIBUTIONS OF THIS REVIEW TO DEEP LEARNING
The purpose of this thinking is to thoroughly review the advanced essays in the 3D learning object detection literature and the 6D pose assumptions from RGB and RGB-D images. It provides a brief overview of current research that is easily comprehensible, and anyone who is interested can grasp the basics of 3D object detection (3DOD) and a 6D pose aspiration (6DPE) system. Moreover, most importantly, this review provides explicit knowledge of 3DOD and 6DPE applications in the field of computer vision to encourage a whole new set of novel methods and ideas. This paper proposes a rich survey for academics interested in research, the autonomous industry and the 3DOD and 6DPE fields. The survey will provide rough guidelines and possible directions for 3D object detection and 6D pose estimation methods, where most of the paperwork relates to autonomous vehicles.

Altogether, the survey has several objectives, such as:
1) We have provided a comprehensive review for a 3D object detection and 6D pose estimation system based on deep learning.,
2) We have created an overview for advanced strategies,
3) We discussed the challenges, advantages, disadvantages of the various proposed strategies
4) We have identified and cited a significant number of innovative concepts and incoming directions in this research sector
5) We can detect vision and broaden the horizons of 3D object detection and research DL (Deep Learning) methods of 6D pose estimation research techniques,
6) In this review, we have tried to give a brief overview on some of the popular datasets available for computer vision.,
7) We have focused on a few popular assessment methods and created a shortlist.

### B. DIFFERENCE WITH OTHER FORMER REVIEWS
To date, much work has been done on 3D Object detection (3DOD) and 6D Pose estimation (6DPE), where most of them are deep learning-based. Nevertheless, the progress of a comprehensive review on the subject is still insufficient. This review sought to create a broad abstraction of modern research with DNN (Deep Neural Network) based 3D object detection 6D pose estimation systems and showed future directions. We can keep an eye on the paper by Mukhtar *et al.* [173], where they reviewed 194 documents

and worked on on-road-based vehicle detection and tracking systems for collision avoidance systems. This review is organized based on various vehicle detection processes, including car detection and tracking sensors.

Sahin *et al.* [209], presents a comprehensive and up-to-date review where authors discuss object detection, examine more than 200 documents, and pose recovery methods with some popular data sets. In addition, several evaluation methods, open issues, and future research directions have been discussed in the paper.

Sivaraman *et al.* [229], has also conducted a literary survey on the method of identifying, tracking and behaving on-road aspects of self-driving vehicles. This study focuses on the current literature related to vision and sensor-based vehicle detection techniques. It began with about 200 papers on environmental perception on the road from 2005. The review papers are mostly related to single vision, stereo vision, the combination of single and stereo vision and sensor-fusion methods for vehicle tracking, detailed image aircraft, 3D modelling, measurement and filtering. Finally, they have called for visionary vehicle identification, tracking, and behavioural analysis with future research directions.

Ioannidou *et al.* [105], discussed the various method of deep learning architecture on different types of 3D data and provided a classification of multiple approaches. Zhao *et al.* [295], provided a regular survey of DL-based object detection frameworks by reviewing a total of 194 research papers. This review begins with a brief history of deep learning with several DL type classifications. Generic Object Detection strategies are discussed here, along with some changes and improved detection performance concepts such as object detection, salient object detection, pedestrian detection, and face detection.

Zhou *et al.* [301], have conducted a review for aspect-based SFM (Structure Form Motion) method, VO (Visual Odometry), and SLAM (Simultaneous Localization and Mapping) based methods where the methods play an important role for support in autonomous driving systems. In their work, they focused on multiple sensor-based methods such as Internal Measurement Unit (IMU) sensors, LiDAR, GPS (global positioning system), monocular-based methods (depending on the height of the camera).

One of the latest online reviews of 3D object detection written by Liu [151], published in the science blog "Towards Data Science," has covered around 32 current state-of-the-art mono3DOD methods as of November 2019. This review did not focus on pose estimation and gave only a brief idea about it. This review is more organized (papers are grouped into several groups) than other previous surveys and gives a more accurate picture of the related article. Unfortunately, there are insufficient numbers of surveys on DL (deep learning) - which stem from the 6D pose estimation system, so researchers should focus on this.

Sahin *et al.* [210], wrote a review related to 6D pose hypotheses where they cover numerous research articles that analyze both object identification and pose hypotheses. Their review article mainly focuses on multiple dataset challenges such as occlusion, cluttered background, lighting conditions, symmetry, texture, illustration, and appearance. The reviewed datasets can be used to evaluate the effectiveness of methods that work in the RGB theme modality. According to the review, the 3D visual understanding is a challenge for complex interactions between objects in terms of perspective, fully or partially chaotic internal environments, and scale changes in different scenes.

Lateef *et al.*, [132] and Minaee *et al.* [167], have provided a comprehensive review of the literature of pioneering works for semantics and example level image division using over one hundred deep learning-based segmentation methods proposed in 2019 and 2020, respectively. Naseer *et al.* [177], created a review of advanced technology based on visual concepts, including visual classification, object identification, pose estimation, semantic segmentation, 3D reconstruction, salinity detection, physics-based reasoning and internal visual skills.

In addition, a recent comprehensive review was presented by Rahman *et al.* [199], where they reviewed the latest 3 DODTs (3D object detection technology). This review maintains some common steps, including descriptions of some popular public datasets, several performance appraisal metrics, and 3D BB techniques. They focused on cutting-edge technology in the 3DOD sector with their significance, contributions and future directional flaws. Zaixing *et al.* [89], discussed several approaches for 6D pose estimation in their review, including the advantages and disadvantages. A further up to date survey for 3D object understanding, classification, identification, defining size and shape, and tracking with 3D visualization and segmentation is present by Guo *et al.* [80].

Additionally, when listing recent approaches, we ignore traditional solutions to offer up-to-date reviews. Our survey paper looks back at later high profile research publications from a variety of perspectives on object detection and pose estimation. At the end of our survey, we proposed some new insights. In short, as of June 2021, this survey summarized and discussed more than 300 high profile states of art techniques (most of them the most recent). We have tried to make this review paper exceptional and comprehensive than other existing reviews by presenting the graphical outlines of the currently relevant papers. Also, we mentioned the future directions given by multiple authors and aim to make a decision based on them. This survey will help researchers (from start to end) who want to work with 3D object detection or 6D pose assessment.

## C. UNIVERSALITY AND UBIQUITY OF DEEP LEARNING IN 3D OBJECT DETECTION (3DOD) SYSTEMS

One of the critical and mandatory tasks for developing computer vision (CV) in the autonomous field is 3D object detection. Driving without a driver, for example, requires an authentic representation of 3D space around autonomous vehicles of various important categories (prediction, plan-

ning, detection and speed control). Although LiDAR point cloud has proven successful for accurate 3D object detection, it is weather sensitive and expensive. Although the concept of monocular 3D object detection (mono3DOD) from RGB or RGB-D image is not a fancy concept, it still differs immensely from the LiDAR-based approach.

In order to detect a 3D object and guess the pose, we need to fully understand an image, rather than just knowing the classification or image localization. 3DOD is a significant work that can be broken down into several subtasks to make important steps for accurate knowledge of images and videos, such as some other notable applications are classification [110], [123], human behavior analysis [26], pedestrian detection [53], skeleton detection [121], face recognition [299] and autonomous driving [33].

There are some significant hurdles in achieving the identification and object localization tasks such as occlusions, chaotic environment, lighting conditions, size differences and viewpoints. Due to the notable impact of accurate object detection in robotic and autonomous fields, more efforts are being made to identify a (3D / 2D) object more accurately with intense care and attention [76], [77], [202], [203]. 3D object identification can be divided into object localization (specific content located in a test image) and object classification (category by object). Conventional 3D object detection models can be divided into three main categories: informative zone selection, feature extraction, and classification.

Any given image can have multiple objects in different positions of the image with different aspect ratios or sizes; It is best to handle the whole image with a different image sliding window. Strategies attempt to identify all possible positions and orientations of objects. Due to a large number of test windows, the process is comparatively expensive and generates additional windows. Moreover, if some stable sliding window template is applied, it will create unsatisfactory areas.

Some important steps to detect object:

- **Feature extraction:** This step helps to identify diverse objects and reveal features with meaningful and strong representations about complex cells as neurons in the human brain [157] such as HOG [45], Haar-like [145] and SIFT [157]. However, due to the varied lighting conditions, it is challenging to accurately describe all kinds of things.
- **Classification**: A classifier has to differentiate the target object from different types to create recognition of more semantic, categorized, and informative ocular objects. The common classifier used for classifications is SVM [41], AdaBoost [68], DPM (Deformable Part-bas ed Model) [64] (more flexible for low level features).

A state of the art results has been achieved in the Pascal-VOC [62] object identification competition by applying the concept of describing local features. However, there were some issues with this model, such as inaccurate bounding boxes, inefficient and unwanted low-level descriptors, and

improperly trained models. These earlier object detection problems were overcome with the emergence of the Deep Neural Network (DNN) [123]. Eventually, identifying and detecting 3D objects from 2D images is a difficult task. The task becomes even more challenging as the level of depth of the 2D image during formation. Nevertheless, it is possible to identify 3D objects from 2D images with some efficient proposed methods.

### D. UNIVERSALITY AND UBIQUITY OF DEEP LEARNING IN 6D POSE ESTIMATION SYSTEMS

To detect 3D objects from monocular 2D RGB images, we need to create a 3D oriented BB (bounding box), while 3D reasoning from a single 2D input is a complex and difficult task. In the autonomous sector, other than object detection, pose estimation is a complex job that needs to be done. It is easier to predict the 6D pose in RGBD images than in RGB images because the 6D pose is a complex combination of 3D rotation of an object (raw, pitch, yaw) and 3D coordinates (X, Y, Z) at the camera focal point [152]. One significant step in identifying the 3D object and estimating the 6D pose of any object from the image can be divided into egocentric and allocentric positions [119]. In the context of autonomous driving, the orientation related to the camera is called egocentric, and the orientation related to an object is called allocentric. Also, full 6D pose estimation is required for successful implementation of AR (augmented reality) [163], robotics grasp [39], autopilot [33], and so on.

Recent improvements to visual depth sensors and the availability of low-cost depth data have significantly improved object pose estimation. In addition, successful implementation of 6D pose estimation method to solve some problems such as variability of viewpoint, similar objects, symmetrical property, occlusion and cluttered environment; All have been overcome due to the availability of RGB-D sensors and the recent improvement of the Convolutional Neural Networks (CNN).

Typically, the recovery of a 6D pose estimation depends on two factors, the familiar instances and the raw/unknown instance of an object. Moreover, some challenges such as shape mood, target domain, shift distribution between several sources, and classification of objects prevent calculating the pose accurately. These challenges have been widely studied in recent years because of their significance in augmented reality (AR) [163], robotics [251], and autonomous vehicles [74]. In the robotics and automated car industries, accurate object detection, the successful application of self-management of objects (robotic groups), and the assumption of 6D poses by robots play an important role in advancing the challenge of autonomous manipulation.

### E. THE PAPER COLLECTION PROCESS

Google Scholar is one of the primary sources of our paper collection. Also, the well-known Database "Web of Science" is another notable source through which we have introduced and collected a number of related papers.

In addition, we should mention "Wikipedia," which is an authentic source of information and documentation. YouTube plays a vital role in understanding any new concept in this case. UTAS (University of Tasmania) Open Access Repository [183], [184] is a great choice for collecting recent papers.

The keywords we have used for searching references include deep learning (DL), deep neural networks (DNN), Convolutional neural network (CNN), object localization, image processing, autonomous vehicle (AV), 2D/3D object detection, 2D/3D bounding box (BB), 2D/3D object proposal, 2D/3D object identification, and 6D pose estimation. In addition, ACM Digital Library, IEEE Explore, Scopus, ScienceDirect is a collection of the best research databases that make our survey resourceful. Last but not least, Researchgate is a legitimate source of information and paper. Most importantly, we have gone through some highly ranked conferences such as CVPR, ICCV, NIPS, AAI, ICLR, ECCV, ICRA, ICML, IV, IROS, ACM, ITSC, ICIP, TPAMI, IRS, WACV, ECCV, ACCV, and Sensors.

### F. OUTLINE OF THE REVIEW

This paper proposes a comprehensive study reviewing the current methods of object pose detection and recovery. Our contributions are as follows:

- Discussed computer vision and deep learning networks, autonomous car and its challenges in section II and III briefly.
- The datasets used for the 3D object detection and 6D pose estimation method were observed to identify its challenges, which are represented in Table 1.
- We Discussed the range of state-of-the-art (SOTA) technology from 3D BB detectors to full 6D pose guessers in IV section.
- In Table 3 where some of the SOTA 3D object detection methods are compared and in Table 5 some SOTA 6D poses estimation methods are compared.
- The Table 4 represents the 3D Object Detection Paper Collection and this table is represented graphically in Figure 2. The significant amount of paper collection on 6D pose estimation in recent years is shown in Table 6 and this table is represented figuratively in Figure 3.
- Open issues are discussed to identify potential future research directions in VI.
- Finally, section VII sums up the present situation of the field and concludes the review work.

### II. COMPUTER VISION AND DEEP LEARNING
### A. COMPUTER VISION

Computer Vision (CV) in artificial intelligence trains computers to interpret and understand the visual world, working with technologies where computers can achieve a high-level understanding of any digital image or video.

The CV system is a method of taking, processing, exploring and mastering digital images and using models built with the help of geometry, statistics, physics and some teaching theories to generate numerical or symbolic data from those images [108], [120], [170]. The CV process is gradually seeing new revolutionary concepts related to object detection where the main challenges are image processing and machine vision [193].

Moreover, a self-driving vehicle is a notable example where ANN and CV have been widely used. However, it is a big challenge for autonomous vehicles to accurately estimate the position of a 3D object from a 2D image. Although much progress has been made in identifying 2D objects from an image or video, identifying a 3D object and determining the 3D properties of an object from a single image is still a challenging problem.

*Typical Tasks of Computer Vision:* Content-based image retrieval [230], Pose estimation [269], Optical character recognition (OCR) [165], 2D code reading [206], Automatic face recognition, Recognition Features [67], Egomotion [16], [303], Optical flow [10]. Scene reconstruction [233], Image restoration [7], Image acquisition [46], Feature extraction [46], Detection/segmentation [155], High-level processing [46], Decision making [46].

### B. ARTIFICIAL NEURAL NETWORK (ANN)

The function of ANN is almost the same as that of the human brain, as knowledge is acquired through the network through a learning process from near it and stored using some synaptic weight neurons. To achieve the final design goal and change the synaptic weight of the network, NN has implemented a learning process known as a learning algorithm. Nowadays, ANN has been applied to multiple jobs, including computer vision, image recognition, speech recognition, social network filtering, machine translation, diagnostics, and video games [72], [178].

### C. DEEP NEURAL NETWORK (DNN)

Deep Neural Network (DNN), a section of a machine learning (ML) where the machine has to predict any output, can be supervised, semi-supervised or unsupervised [219]. Since traditional ML techniques cannot process natural data in their raw form, DL (Deep Learning), an advanced DNN technique, applies multiple layers to reveal high-level features from the raw data. For example, in image processing, where the lower layers of the DL model can recognize the edges only, the upper layers can detect a certain number of letters or objects or features of the object [239].

Eventually, DL processed unsorted/sorted, labelled or unlabelled data and construct a pattern to make a better prediction [120], [123], [176]. Though DL was popular since 1980-90s, offered the concept of the back-propagation classifier [207]; nonetheless, it soon lost its popularity due to over fitting, scarcity of big data, and poor computation capacity as compared to other ML tools.

The popularity of Deep learning algorithm has increased since 2006 [94] with the advancement in speech recognition [93] application. Convolutional Neural Networks (CNN), the popular DL framework, which is applied on

multiple sectors such as Natural Language Processing (NLP), Computer Vision, Speech Recognition, Audio Recognition, Machine Translation, Social Network Filtering, Bioinformatics, Medical image analysis, and much more [171].

### 1) CONVOLUTIONAL NEURAL NETWORK (CNN)

The Convolutional Neural Network (CNN) has a deep feed-forward architecture and remarkable ability to generalize to better networks with fully connected layers. CNN has largely applied to image analysis, especially pattern recognition, which can also be employed to solve other data analysis problems, such as classification problems. CNN is a deep learning network developed for image and video processing that has made significant progress since 2010 and is now widely used worldwide.

The two most notable qualities as classified composition and the ability to extract powerful features from an image prove that CNN is one of the most powerful object detection classifieds. Several important CNN architectures have been proposed times for image processing such as ImageNet [48], AlexNet [123], ConvNet [124], [134] LeNet [243], VGGNet [228], ResNet [87], ZFNet [284], GoogLeNet [243], GPU (Graphics Processing Unit) processor large-scale distributed clusters [47], and OverFeat [218].

On top of that, CNN is a powerful algorithm that is widely used for image classification and object detection [123], [284]. Because of the notable advantages, CNN has been widely applied in many research fields including image super-resolution reconstruction [179], [285], image classification, image retrieval [110], face recognition [299], pedestrian detection [249], [272], [294] and video analysis [228], [270], [284], car detection [33], and pose estimation [296].

Most importantly, CNN can be adequately trained that does not suffer from over fitting and is easy to apply to large networks [123]. However, CNN cannot provide accurate results when the length of the output level is variable and the presence of objects of interest is not fixed. Therefore, more sophisticated algorithms such as R-CNN, Fast R-CNN and YOLO have been developed to solve advanced image processing problems.

### 2) REGION BASED CONVOLUTIONAL NEURAL NETWORK (RCNN)

Girshick [76] proposed a method where a large number of regions were selected, and the Selective Search (SS) [253] method was applied to select only 2000 regions from an image, which he named the region proposals.

Since each region of the image is applied to CNN individually, the training time is about 84 hours, and the forecast time is about 47 seconds. As a result, the process becomes time-consuming because it has to classify 2000 region propositions for each image. Here, the CNN functions as a feature extractor and the revealed features are processed through an SVM [41] classifier to distribute the object inside the region proposal. Additionally, to anticipate the region proposals and increase the bounding box's precision quality, the algorithm creates

four offset values. The main problem with this classifier is time.

### 3) FAST RCNN

The algorithm previously proposed to create a quick object recognition classification updated some of the errors of R-CNN and renamed as Fast R-CNN [76]. This method is almost the same as the R-CNN classification. Since it uses CNN once, there is a significant gain over time. The training time is about 8.75 hours, and the estimated time is about 2.3 seconds.

### 4) FASTER RCNN

Both of the above algorithms (R-CNN and Fast R-CNN) use SS to determine region proposals. SS [253] is a slow and time-consuming process that over-segmenting the image affects network performance. Therefore, Shaoqing Ren *et al.* [203] proposed an object identification algorithm that removes the SS algorithm and allows the network to learn region recommendations. After the predicted regions are resized using the ROI(Region of Interest) pooling layer, which is then used to classify the image in the proposed region and predict the IoU (Intersection-over-Union) ratio of the bounding boxes.

### 5) SINGLE SHOT MultiBox DETECTOR (SSD)

Liu *et al.* [153] proposed SSD (Single Shot Multibox Detector), a single shot detector for multiple segments, applies an additional small conventional filter to maps that are faster and significantly more accurate than previous single shot detectors like YOLO.

### 6) MASK RCNN

He *et al.* [86] presents the concept of flexible structures called Mask R-CNN for object instance segmentation. This method effectively recognizes objects from an image while creates a high-quality segmentation mask for each instance at the same time. Mask R-CNN is a practical extension of Faster R-CNN, where an additional branch is added to predict an object mask parallel to an existing branch. Moreover, this method is a slightly improved version of R-CNN that runs at 5fps and can adapt quickly to predict human posture. Also, Mask R-CNN has won the COC 2016 Challenge by overcoming three key issues: Instant Segmentation, Bounding-Box Object Identification, and Individual Keypoint Identification.

### 7) YOLO

Redmon *et al.* [202] Proposed a novel object detection technique called YOLO (You Only Look Once), where the classifier does not process the whole image; Instead, it focuses partly on the image with a high probability of having the object in that part. This single convolutional network is faster than existing object detection algorithms. However, above all advantages, the YOLO algorithm struggles to detect small objects within the image. For example, the spatial limitations

of algorithms can make it difficult to identify flocks of birds. Some other notable DL structures are: RefineDet [287], Retina-Net [147], Deformable convolutional networks [44], Cascade R-CNN [21], 3D-RCNN [128], Libra R-CNN [186].

### 8) MESH R-CNN

Facebook introduced a novel RCNN method in artificial intelligence called Mesh R-CNN that can convert 2D objects to 3D shapes and mesh [267]. Facebook has highlighted its latest advances that can identify complex issues. This study has applied in-deep learning to understand the 3D shapes of complex objects and novel architectures such as Bounding Box, 3D Voxel Pattern, Point Cloud and Message for prediction and localization. Mesh R-CNN can effectively detect and classify objects in 3D form from chaotic 2D images and occluded objects and ultimately estimate their full 3D shape.

| No. | Algorithms | Features |
|-----|-----------|----------|
| 1 | CNN [284] | Impressive classification performance. |
| 2 | R-CNN [76] | Robust architectural structure [135] |
| 3 | Fast R-CNN [77] | Has additional sub-network. |
| 4 | Faster R-CNN [203] | Improved outcome by applying the RPN. |
| 5 | R-FCN [43] | Manages to address the dilemma. [87]. |
| 6 | Mesh R-CNN [267] | Novel architectures for 3D shape [29] . |
| 7 | YOLO [202] | Improved model for the fixed-grid regression. |
| 8 | 3D-RCNN [128] | 3D scene understanding to map image. |
| 9 | Libra-RCNN [128] | Re-balances by mixing : IoU, feature and L1 loss. |
| 10 | Cascade-RCNN [21] | It is extensively contextual for detection. |

## III. AUTONOMOUS VEHICLE (AV)

An autonomous vehicle (AV) is a combination of some actuators, sensors, complex analytical algorithms, machine learning methods, and high-speed processors that are needed to implement complex software. Self-driving vehicles create and maintain a map called Simultaneous Localization and Mapping (SLAM) of their surroundings by multiple sensors placed in different parts. Such as LiDAR or radar measures distances of other vehicles or obstacles, detect road edges. In addition, one or more cameras mounted on autonomous vehicles can detect traffic lights, road signs, lane signs, vehicles, obstacles and pedestrians [52].

During parking, ultrasonic sensors placed on wheels to detect obstructions and other vehicles. Advanced complex software [131] then processes all these sensory inputs,

generates outputs and sends commands to the vehicle actuator responsible for steering, braking and control acceleration. AV can be identified as a complete package of hard-coded rules, a complex algorithm and efficient predictive models, helping sophisticated software to run on the road smoothly [73], [192], [245].

To date, autonomous vehicles are equipped with two types of sensor such as active sensor: LiDAR [138], [139], [279], radar [277] and passive sensors: Single/Stereo cameras [32], [172], and their fused systems [33], [107], short-range sensor (Ultrasonic sensors) [122]. Veli *et al.* [104] made lots of progress in sensor technology and GNSS (Global Navigation Satellite Systems).

A research team from the Massachusetts Institute of Technology (MIT) [84] announced in May 2018 that they had successfully built a driverless car that could successfully navigate unmapped roads with a novel system known as MapLite [40]. This application enables the driverless car to drive on a completely new road without using pre-loaded 3D maps. The basic idea is to combine the vehicle's position with sensors that monitor the surrounding conditions, and OpenStreetMap (OSM) is used to detect the GPS of a vehicle [40].

Also, an AV has been divided into 5 levels such as level 1 - requires driver support, level 2 - partial automation phase, level 3 - limited driver support, level 4 - higher automation and level 5 - fully automated [225], [84] [190]. At present, level 3 autopilot is available on the road, as Level 4 and Level 5 autonomy require large-scale neural network training and visual recognition, including accurate pose estimations. Multiple companies produce intelligent vehicles and test them to drive autonomously in certain situations, such as Tesla Autopilot, Waymo, Uber, Volvo, Google, BMW, Mercedes Benz, Nissan and General Motors. However, they are still in the testing phase and unable to operate without assistance.

### A. TECHNICAL AND SOCIAL CHALLENGES OF AUTONOMOUS VEHICLES

Although the concept of autonomous or self-propelled vehicles has come a long way in recent years and numerous studies have been done in this sector, this technology is still not flawless. Lawmakers and consumers still feel confused and anxious about implementing self-driving cars and feel insecure and uncertain about the autonomous vehicle's ability to move freely. So self-propelled cars are still in the experimental stage, and more research is needed to perform them properly. One of the significant challenges of autonomous vehicles is accurately estimating the exact position and orientation of nearby vehicles. The five core reasons are classified as why the AV still are not on the roads are listed in below:

**Sensors:** An autonomous vehicle faces various challenges for smooth automation systems such as proper vehicle navigation system, GPS, environmental perception, LiDAR and radar, visual perception, speed and direct perception and

a vehicle control system [192], [293]. Furthermore, to be qualified as perfect autonomous vehicles, these sensors need to work in all weather conditions anywhere on the earth. Undoubtedly, the critical issue for driverless vehicles is that there should be a control system capable of automatically analyzing sensor data and making accurate estimates of vehicle postures, obstacles, pedestrians and road signs [306].

**Machine learning Algorithm:** At the moment, there is no widely approved and authorised ML algorithm to ensure that they are 100 % error-free, safe and secure for use in any driverless vehicle. The most popular algorithm applied to current driverless vehicles is SLAM, which integrates data from various sensing components and uses offline maps [266]. WAYMO has improved the performance of the algorithm SLAM and named DATMO (Detection and Tracking of Moving Objects), which can handle any curbs, including vehicles and pedestrians. Zhang *et al.* [291], proposed a concept that collaborated with the existing Visual odometry (VO) system such as SLAM and ORB-SLAM2 ( an updated version of the SLAM) [174].

**The open road:** When the AV drives on new roads, it should identify things that did not come before in the training process and may be subject to software updates. As a result, it would be not easy to ensure that the system is as secure as its former version.

**Regulation** To date, adequate standards and regulations for autonomous systems do not exist. There have been numerous high-profile accidents involving Tesla's current automobiles, as well as other automotive and autonomous vehicles [9].

**Social acceptability:** Applying an automatic car on the road is not only a problem for those who want to buy and use a driver-less vehicle, but also for others who share the road with them [9].

## IV. LITERATURE REVIEWS

An essential part of computer vision is the identification of objects from images or videos. Object detection helps in pose estimation, vehicle detection, pedestrian and other curb detection. Previously, the image was processed and classified using traditional machine learning (ML) algorithms such as colour histogram [220], SVM (Support Vector Machine) [41], logistic regression [202]. However, there are some differences between the recent object detection algorithms (CNN, R-CNN, YOLO) and traditional ML classification algorithms (SVM, logistic regression).

The definition of object identification problem determines where objects are located in a given image is called object localization, and what class each object belongs to is called object classification. Thus, the traditional thematic object detection model's pipeline is divided into three stages such as:

1) Informative region selection: Different objects can appear in any position of the image and have different aspect ratios or sizes, so scanning the entire image with a multi-scale sliding window is a natural choice.

2) Feature extraction: To identify different objects, we need to figure out visual features that can represent a semantic and robust.

3) Classification: A classifier needs to differentiate a target object from all other categories and further classify the presentation.

### A. 3D OBJECT RECOGNITION

Object recognition is one of the primary pillars of a computer's vision and is sometimes confused with the problem of object classification/shape retrieval. 3D object recognition methods can be divided into two main categories such as voting methods, Hough transform [6], and geometric hashing [130] and the correspondence based method, spin images [112], local feature histograms [90], 3D shape and harmonic shape context [69].

David *et al.* [156] developed an object recognition system using local image features in cluttered real-world scenarios. Cordelia *et al.* Schmid *et al.* [214], has shown that recognition of successful objects can often be achieved by applying a sample local image descriptor to a large number of repetitive locations. Papazov *et al.* [187] proposed the recognition of a 3D object, especially for noisy and scattered data in cluttered and occluded environments. This proposed concept applies a combination of strong geometric descriptors, a hashing technique and a sampling technique - RANSAC [65].

### B. 3D OBJECT DETECTION FROM RGB AND RGB-D IMAGE

3D object detection is a significant key part of the visual perception system of robotic and autonomous technologies. It has many applications with different category some of them described in FIGURE 1.



**FIGURE 1.** The application domain of object detection.

In generic object detection, object instances are identified by applying predefined sections/categories. It has some challenges such as the immense range of inner-class variations and the large-scale object categories [150]. Salient object detection detects the most significant and notable object in an image, and then it segments the whole area of that object [12].

### 1) FEATURE EXTRACTION, SEGMENTATION AND MATCHING

Rapid and accurate image segmentation with feature extraction is the primary task of the computer vision field. Lowe *et al.* [156] proposed SIFT (Scale Invariant Feature Transform), an object recognition system for image scaling, translation, matching and rotation, and a partial constant for illumination changes, including 3D projection. The images here have been converted to a wide collection of local feature vectors and can generate approximately 1000 SIFT keys in 1k ms during each image count by applying classification. Although occlusion may be present in the image, SIFT can provide a high level of accuracy.

Kang *et al.*, [114] Created a structure called DaSNet-V2 that matches identification, category, localization, and object instances. A method capable of achieving real-time performance by adopting PWP 3D (count per pixel) and applying the region-based simultaneous strategy of 2D partitioning using the NVIDIA CUDA framework is largely developing parallel algorithms [194]. Fu *et al.* [70] introduced DORN (Deep Ordinary Regression Network), a multi-scale network framework that achieves a spacing-increasing discretion (SID) strategy to rebuild depth and depth networks to reduce the complexity of existing feature maps.

### 2) SHAPE VARIATION

Xiang and Dollar offered 3DVP (3D Voxel pattern), which uses ACF (Aggregate Channel Features) detectors to find out the basic features of each object such as shape, appearance, aspect and curbs [54], [271]. In addition, the 3D pose of a vehicle can be accurately localized from the context of this method and can detect other vehicles and guess the pose [74].

Novotny *et al.*, [181] have created the C3DPO (Canonical 3D Pose) Network for non-rigid structure motion where no training images and messes are available. It has partially reconstructed a 3D object from a monochromatic RGB image to change perspectives and distort the object. It has also emphasized the mandatory presence of certain canonicalization functions of reconstituted size and shape. The input depth proposes objects pose according to the classification of convex hulls that align the clusters of convex sections drawn from the images. This is an example of a highly efficient size identification pipeline that uses the CHAL (convex hull alignment) algorithm for hypothesis generation and is used to identify objects in complex scenes with multiple objects [42].

Qian *et al.*, [196] presented a method for evaluating individual 3D sizes, where there was a balance and robustness between the accuracy and efficiency of the conventional stage recovery method, significant measurement limits and high-frequency fringe patterns. Chabot [28] made a framework called Deep MANTA for 3D object detection based on a single-dimensional image in an end-to-end fashion network, determining the object class, 2D region proposal generation, 2D location, orientation, dimension and 3D position. This model has implemented a 3D vehicle dataset featuring 3D mesh with real size to match vehicle parts (wheels, headlights,

mirrors) and defines a 3D shape for each 3D model. Zhou *et al.*, [304] has built the CenterNet framework, which is simpler, faster, and more accurate than traditional BB detectors and poses estimators.

### 3) 3D PROJECTIONS OF THE 3D BOUNDING BOX VERTICES

Chen *et al.*, [34] proposed 3DOP (3D Object Proposal) for accurate object class identification in the context of autonomous driving. 3DOP produced several sets of 3D candidate boxes to identify almost every object in 3D space. This method has featured object size, ground plane, different depths, spaces, the density of points inside the box, visibility and soil distance.

Mono 3D (Monocular 3D Object Detection) [32] uses ground planes and some segmentation features to generate 3D proposals from monocular images in the context of autonomous driving. In addition, both 3DOP and Mono3D methods applied some common hand-crafted features. This technique applies several intuitive potentials to each candidate box expected in the image plane encoding synthetic segmentation, relevant information, size and location pre-requisites, and ideal object sizes. Also, the S-SVM [111], structured SVM [252], parallel cutting plane [228] and IoU has been implemented with a comprehensive search model.

The proposed DSS (Deep Sliding Shapes) [236] is a 3D convergent formulation that takes 3D volumetric views as input from an RGB-D image and then outputs a 3D object bounding boxes. In addition, this method proposes the first 3D Region Proposal Network (RPN) to learn objects from geometric shapes and the first Joint Object Recognition Network (ORN) to extract geometric features in colour properties in 2D.

Ding *et al.*, [49] proposed a fancy wire-frame model called the CPO (Cross Projection Optimization Method) that can detect both 3D pose and shape estimation of a vehicle for an autonomous vehicle. The CPO method applies a simple wire-frame model combined with the Hierarchical Wire-frame Constant (HWC) method instead of bounding box annotation to shape detection for 3D pose and accurate 3D localization [33].

The solution provided is primarily based on local properties, especially for matching objects in a 2D image of a rigid 3D object [79]. This method creates an accurate 3D model of the object with the locations of its features and then places it in an image to identify new features. Finally, the position, orientation, and shape of the virtual object are defined concerning the object's coordinates.

Rad [198] has created a framework where a total of 8 corners of the bounding box are applied to the multiple-input image called BB8. This method is trained to predict their poses in the form of 2D projections of the corners of their 3D bounding boxes and calculates 3D poses from this 2D-3D correspondence with a PNP algorithm [136].

Another strategy called Mono3OD [227] where a single RGB image uniquely transformed to reduce object detection and increase the credit count for 3D BBs. Li [142]

suggested RTM3D (Real-time Monocular 3D Detection), the first real-time 3D identification method for autonomous driving, predicting the nine point-of-view of the 3D BB in place of the image and using 3D and 2D perspective geometry to restore orientation, location and dimensions in 3D space.

Liu [150] has claimed a deep fitting degree scoring network for mono 3DOD, which focuses on the active fitting degree among proposals and objects. It is discrete from other existing monocular frameworks by attaining localization by computing the visible degree of calculation among the 3D project proposals and the object. A concept named FQNet, [150], can assume the 3D IoU (Intersection over Union) among the 3D proposals and the object.

Zhang *et al.*, [291] proposed a framework for 3D object detection by determining object class, 2D region proposition production, 2D position, position, dimension and 3D positioning based on a single image in an end-to-end fashion network. Furthermore, Bao *et al.* [8], recently introduced Mono-Fenet, a compelling feature enhancement method for the 3D object detection, which includes the ROI Mean Pooling layer, the PointFE network, and feature enhancement networks using 3D-NMS and exclusive RGB imagery.

The full 3D poses and dimensions of an object from a 2D BB by applying some restrictions to calculate the orientation and volume of the object using DCNN, where the novel DCNN method known as MultiBin regression is used to estimate the orientation of the object [172]. SS3D, a single-phase monocular 3D object detector where the 3D representation is returned by a representative and uses for the geometric shapes of the 3D box with autonomous driving [113].

Hu *et al.*, [102], introduced a complete 3D vehicle bounding box tracking information method from exclusive videos and a method for dealing with 3D vehicle detection guesswork. A new pipeline based on LSTM [254] is designed to collect large-sized 3D trajectories from real-world driving environments and track 3D vehicles within 30 meters. The method called M3D-RPN implemented exclusive 3D identification and 3D zone proposal networks and lifting the geometric relationship of 2D and 3D perspectives, including 3D boxes [15].

### 4) 3D OBJECT DETECTION IN POINT CLOUD

Scientists proposed a method for identifying free-form 3-dimensional objects in point clouds with global representations [56]. The basic idea of the model is to create a universal approach statement based on the point pair factor. Free-form objects in 3D datasets can be achieved by a number of sensors, such as a laser scan, a TOF (Time of Flight) camera, which has been widely disseminated from a computer perspective [25], [160].

Chen *et al.*, [33] introduced accurate 3D object detection for individual behaviour, known as Multi-View 3D Network (MV3D), which works with multimodal datasets. MV3D framework creates efficient 3D candidate boxes from a 3D point cloud BEV (Bird's Eye View) [154] image, and the main goal of this method is to identify 3D objects using both LiDAR and image data. Current LiDAR-based methods set 3D windows in 3D Voxel Grid [58], [260] or apply convolutional networks [139] to front viewpoint maps.

On the other hand, a hybrid method has introduced that combined both LiDAR and camera data for 2D detection to get accurate results [59], [78]. Qi *et al.*, [195] offered a fancy concept called "Frastum PointNets" based on RGB-D data in a point cloud and expects a semantic class for each point in that point cloud. A method named PV-RCNN provides accurate 3D object detection from point clouds that deeply integrates 3D visualization with point-to-point set-based abstraction with a 3D visual convoluted neural network and multiple receiving fields [221]. Finally, a novel method called SAANet (Special Adaptive Alignment) uses an "SAA" module that addresses fusion-based deep structures that combine clouds and images for 3D object detection with complements cloud properties and image properties [31].

### 5) SPEED / ACCURACY TRADE-OFF

Huang *et al.*, [103] introduced a process that helps determine the speed and accuracy of the calculation and also recommend which method is better suited for a specific application. Shrivastava *et al.*, [224] proposed a TDM (top-down modulation) approach to include image quality for a ConvNet architecture such as VGGNet [228], ResNet [87], and Inception-Resnet [242]. Song [236] proposed Deep Sliding Shapes (DSS) that convert an RGB-D image into a point cloud and then slides a 3D detection window into 3D space. Luo [158] made a concept that identifies 3D objects and accurately predicts the position, size, orientation and division of objects in 3D space at very fast speeds. Li [140] has come up with an idea called GS3D, a 3DOD method based on an RGB (single) image in autonomous driving.

### 6) OBJECT DETECTION BY KEY POINT ESTIMATION.

The most famous classifier that detected an object using key-point inference (identifies the object as a point to the key) is Cornernet [133], ExtremeNet [305], and CenterNet [304]. In CornerNet, the corners of 2D BB are used as semantic key points. ExtremeNet, on the other hand, highlights all points, including the top, left, bottom, right, and centre of the bounding box. Compared to these classifiers, the Centernet is much faster, which only chooses the object's centre.

## V. LITERATURE REVIEW OF 6D/6DoF (DEGREE OF FREEDOM) POSE ESTIMATION

In the computer vision sector, guessing a 6D pose of an object is a significant problem that needs to detect both 3D orientation and 3D position of an object in the case of the camera centred coordinates [116]. In short, the three factors for the 6D pose estimation are the critical role of rotating left and right on the X-axis (roll) side as well as on the Y-axis

**TABLE 1.** Data-sets used for multiple 3D object detection and pose estimation method.

| No. | Benchmarks | Frame | Format | Categories | Features |
|---|---|---|---|---|---|
| 1 | KITTI [74] | 14999 | RGB | 3 | >100 gb data |
| 2 | COCO [148] | 328000 | RGB | 11 | 2,500,000 labeled instances. |
| 3 | SUN RGB-D [55] | 10335 | RGB- D | 800 | Combination of NYU, B3DO and SUN3D. |
| 4 | NYU [226] | 1449 | RGB | 19 | 795 trained data. |
| 5 | PASCAL-VOC[62] | 14999 | RGB | 12 | Extension of PASCAL3D+ |
| 6 | ImageNet [48] | >14 million | RGB | 21 | 7481 trained and 80256 labeled objects. |
| 7 | RGB-D Object Dataset [129] | 250,000 | RGB-D | 51 | 300 physically distinct objects. |
| 8 | Parsing IKEA Objects [146] | >1k | RGB | NA | 3D models of IKEA furniture. |
| 9 | CVonline [66] | All | NA | NA | A rich database for CV, ML, and IP. |
| 10 | Yu Xiang et al. [275] | 2640 | RGB | 4 | Ratio between test and training images is 50% |
| 11 | NYC3DCars [164] | 2299 | RGB | 20 | Over 2000 annotated photos from New York. |
| 12 | EPFL Cars [185] | 2000 | RGB | 20 | Acquired from a car show. |
| 13 | ApolloCar3D (Kaggle) [238] | 5,277 | RGB (Monocular) | 3 | 5 GB Data.. |
| 14 | LINEMOD [91] | 30899 | RGB | 13 | 7481 trained and 80256 labeled objects. |
| 15 | MULT-I [247] | 1,100 | RGB | 13 | Both cluttered and occluded. |
| 16 | OCC [14] | 10k | RGB-D | 20 | cluttered, textured and texture-less, rigid and non-rigid objects. |
| 17 | BIN-P [55] | NA | RGB-D | 20 | The first fully annotated bin picking dataset |
| 18 | T-LESS [96] | > 49K | 3 | 50 | 39K training and 10K test images. |
| 19 | RU-APC [204] | nearly 10,000 images | RGB-D | 25 | Low illumination |
| 20 | BOP [98] | Nearly 340k images | RGB-D | 89 | Combination of RU-APC,TUD-L, and TYO-L Datasets. |
| 21 | YouTube-BoundingBoxes (YT-BB) [201] | 380,000 video segments | videos | 23 | YTBB is the largest human-annotated detection data set |
| 22 | ShapeNet [29] | 3,000,000 | 3D models | 3,135 | Autonomous robots and vehicles. |
| 23 | YCB-Video dataset [276] | 133K images | video dataset | 21 | Live RGB-D camera. |
| 24 | Yale-CMU-Berkeley dataset [22, 23] | 2K images | RGB + RGB-D | NA | BigBIRD Object . |
| 25 | JHUScene-50 [37] | 22520 images | RGB-D images | 50 | >20K labeled poses. |
| 26 | NOCS-REAL275 [261] | (275K training, 25K testing) images | RGB-D images | 6 | 18 different scenes and 42 unique instances. |
| 27 | Falling Things (FAT) [250] | 60k images | RGB | 21 | 2D/3D bounding box coordinates for all objects |
| 29 | ObjectNet3D [273] | 90,127 i images | 2D images | 100 categories | 201,888 objects and 44,147 3D shapes |
| 30 | nuScenes [20] | 1,400,000 images | 390,000 LiDAR sweeps | 23 | Contains 100 times images than KITTI dataset. |
| 31 | RobotP [281] | 4200000 | RGBD | NA | synthesized photo-realistic color-and-depth images |

(pitch) and tilting backwards on the Z-axis (Yaw). Thus, these features encourage concentration on the recovery of vehicle posture and size estimates to enhance the intelligence of the intelligent transport system and the robotic sector. Therefore, the conventional states of industrial techniques of 6D pose estimation are discussed here in the context of the autonomous car.

### A. 6D POSE ESTIMATION DIRECTLY FROM RGB IMAGES

Wu *et al.*, [269] proposed an algorithm named 6D-VNet, and won the first place in the "Apolloscape Challenge 3D Car Instance" competition. It is an abstract structure for autonomous vehicles assuming 6 DOF object poses that can detect all aspects of traffic in a single RGB image while rotating vectors and 3D translation. The basic technique of

**TABLE 2.** Evaluation metrics: Different types of evaluation metrics to identify and measure the performance of proposed classifiers.

| No. | Evaluation metrics | Estimation | Function | |
|-----|--------------------|------------|----------|---|
| 1 | Intersection over Union (IoU) [61, 236] | Assesses the 2D space performance. | $$W_{I_o}U_{2D} = \dfrac{B \cap \overline{B}}{B \cup \overline{B}}$$ | (1) |
| 2 | Average Precision (AP) [61] | Assesses the shape of the Precision-Recall (PR) curve. | $$AP = \dfrac{1}{11} \sum_{r \in 0,0.1,..,1} pinterp(r)$$ | (2) |
| 3 | Average Orientation Similarity (AOS) [61, 74] | Applies the AP metric for object detection. | $$AOS = \dfrac{1}{11} \sum_{r \in 0,0.1,..,1} max_{\widetilde{r}:\widetilde{r} \geq r} s(\widetilde{r})$$ | (3) |
| 4 | Average Viewpoint Precision (AVP) [274] | Generats a VPR (Viewpoint Precision-Recall) curve. | NA | |
| 5 | Translational and Angular Error [56] | Measure the ground truth angel and translations and rotation of an object. | $$W_{TE} =\| X - \overline{X} \|_2$$ | (4) |
| 6 | 2D Projection. [274] | Estimated poses by projecting the model's vertices onto the image plane. | $$W_{2D}proj = \dfrac{1}{|V|} \sum_{v \in V} \|C\overline{R}_v - CR_v\|$$ | (5) |
| | | | $$W_{RE} = \arccos\left(Tr(R\overline{R}^{-1} - 1)/2\right)$$ | (6) |
| 7 | Average Distance (AD). [91] | Remove obscurities from the symmetry property and occluded background. | $$W_{AD} = avg_{s \in M}\|(\overline{R}s + \overline{T}) - (Rs + T)\|$$ | (7) |
| | | | $$W_{AD} = avg_{s1 \in M} min_{s2 \in M}\|(\overline{R}s + \overline{T}) - (Rs + T)\|$$ | (8) |
| 8 | Visible Surface Discrepancy (VSD) [97] | A method that remove ambiguities. | $$W_{VSD} = \begin{cases} 0 & \text{if } p \in \overline{V} \cap V \wedge |\overline{D}(p) - D(p)| \\ 1 & \text{otherwise} \end{cases}$$ | (9) |
| 9 | Sym Pose Distance. [97] | Fixed the ambiguity due to the symmetrical shape of an object. | $$W_{Sym} = min_{G_1,G_2 \in G} \sqrt{\dfrac{1}{S} \int_s \|T_2 \circ G_2(X) - T_1 \circ G_1(X)\|^2 \times ds}$$ | (10) |
| 10 | Average 3D precision (A3DP) [237] | Inspired by the AVP metric ( for both 3D shape and 3D pose ) | NA | |

this method is to control the 6D position of the vehicle using the outputs from the RPN (Region Proposal Network) [77] and 2D object detection network (Mask R-CNN) [86] that can learn both rotation and translation by outlining a loss function model.

Brachmann *et al.*, [13] created a template-based model for calculating 6D pose for a specific object from a single RGB image. The algorithm optimizes the power following the RANSAC concept for a large and uninterrupted 6D pose space. The technical feasibility of classification is using a new composite dense 3D object coordinate form, including object class labelling. Kehl *et al.*, [116] developed SSD-6D, a CNN method to detect the 3D object and accurately guess the 6D pose from an RGB image. It is a unique detector method for relevant training on synthetic model information, which

applies to the collection of small objects and objects with many conceptual and practical advantages.

Inspired by BB8 [198] method Zhang *et al.*, [289] re-imposed the coordinates of the image and applied the Perspective-n-Point (PNP) [136] algorithm without any post-refinement. Similar to recent work, the method uses 2D BB to calculate the coordinate regression of images based on their centres, focusing on the gap between image classification and pose estimates [248]. Deep-6DPose is an end-to-end deep learning solution, which finds objects and compresses them and retrieves instances of 6D objects from single RGB images [50]. It consists of two main components, such as RPN and a mask R-CNN, including Lie algebra.

Billings *et al.*, [11] has developed a new proposal to predict 6D object poses from monocular RGB images

**TABLE 3.** Advantages and disadvantages of some state-of-the-art 3D object detection techniques.

| No. | Algorithms | Advantages | Disadvantages | DataSets |
|-----|-----------|-----------|---------------|----------|
| 1 | **3DVP [271]** | Has strong accuracy in occluded and complex scenes. | Its overlapping detection graph is often very complex. | KITTI [74] |
| 3 | MV3D [33] | High precision accuracy. | Computationally expensive. | KITTI |
| 4 | 3DOP [34] | Highly accurate object proposal skills. | Average performance on cluttered environment. | KITTI |
| 5 | Mono3D [32] | High-performance in monocular imagery with better recognition rate. | Detection GRAPH is often very complex | KITTI |
| 7 | DSS [236] | Reveal powerful 3D and color features from the data. | NA | KITTI |
| 10 | SAANet [31] | Higher Precision and inference time in LiDAR + Img-based methods | Fails to localize and explore local orientation information. | KITTI |
| 11 | Deep MANTA [28] | Less time consuming and overcomes loss of information. | NA | KITTI |
| 12 | MonoDepth [291] | Perform very good especially for finding smaller objects | NA | KITTI |
| 13 | Mousavian et al [172] | Better performance for large dataset. | NA | KITTI and Pascal 3D+ |
| 2 | **C3DPO** [181] | 3D reconstruction and object deformations capacity. | Requires expensive hardware. | PASCAL3D+ |
| 6 | CPO [49] | High precision accuracy. | Overlapping in 3D shape and pose. | 4 Real World Benchmark. |
| 8 | SSD [153] | High-accuracy, high speed and very robust. | Confused with similar categories and worse performance on smaller objects. | SUN RGB-D andNYUv2 |
| 9 | 3D-SSD [31] | Significant accuracy and computation efficiency. | Receptive field is limited. | SUN RGB-D andNYUv2 |

**TABLE 4.** 3D object detection paper collection.

| No. | Approach | References |
|-----|----------|-----------|
| 1 | Classification | SVM [41], S-SVM [111], [111], [252] , [100], [4], |
| 1 | Representation Transformation (Pseudo-LiDAR, BEV) | MLF [28] , Pseudo-LiDAR [264], Pseudo-LiDAR++ [280], Pseudo LiDAR-e2e [265], pseudo LiDAR color [159], ForeSeE [263], BEV-IPM [118], [17], SAANet [31], PIXOR [279], [82] . |
| 2 | Keypoints and Shape | Deep MANTA [28] , 3D-RCNN [128], Mono3D [32], ROI-10D [161], MonoGRNet [197], ApolloCar3D [237], RTM3D [142], [215], [274], [35], [300], [42], [259]. |
| 3 | Distance via 2D/2.5D/3D constraint | [159] , 3D-RCNN [128], Mono3D [32], ROI-10D [161], MonoGRNet [197], ApolloCar3D [237], 6D-VNet [269], GPP [200], RTM3D [142], [49] , [79], Deep3dBox [172], Shift R-CNN [175], GS3D [140], [278], [169], [103]. |
| 4 | Feature extraction and Matching | [28] , [79], [13], Deep Sliding Shape [236], [137], MVTec ITODD [57], [14], [156], [205], [114], [302]. |
| 5 | 3D object proposal methods | 3DVP [271], 3DOP [34], Mono3D [32] , MV3D [33] , Deep Sliding Shapes [236], 3D voxel grids [58, 260, 271] , CPO [49], Joint mono3D [102], SS [253], SS3D [113], CasGeo[63], M3D-RPN [15], MonoDIS [227], [102], CenterNet[304], MLF-Mono [8], MonoDepth [291], [144]. |
| 7 | 3D Object Detection in Point Cloud. | [56], [162], [14], [160], [221], [138], [195], PV-RCNN++ [222], SE-SSD [298], SVGA-Net [88], CIA-SSD [297], PC-RGNN [292], [262]. |
| 8 | 3D Object Detection from RGB-D IMAGE. | DSS [236], SS [235]. |
| 9 | Speed /Accuracy Trade-off | [103], [224], [244], SSD [153], 3D-SSD [158]. |

by applying the CNN pipeline with the ROI proposal. It predicts the presence of intermediate outlines for 3D objects, 3D orientation and 3D translation vectors. For the

6-D category level pose estimation, two-level BB-based alternative methods have been developed that directly output the 6D pose without the use of any PNP but consist of

ResNet (Residual Neural Network), RPN, and FCN (Fully Convolutional Network) [149].

### B. POSE FROM DEPTH / POINT CLOUD METHOD

Mitash *et al.*, [168] advocated a concept for efficient object 6D pose estimation in cluttered scenes, where the Cartesian product of the candidate's post for interactive objects is used to identify the best view and create an efficient search, and the candidate post clusters for each object. The MCTS (Monte Carlo Tree Search) technique is applied to conduct tradeoffs in fine-tuning and explore new instances.

Xiang *et al.*, [276] has created a generic structure called POSNN that calculates the 3D translation of an object in the image and predicts its distance from the camera. Furthermore, this method reduces the ShapeMatch-Loss function and enables POSNN to handle symmetrical objects where the VGG16 backbone is used to extract features.

The PointPoseNet classifier for 6DoF objects gives the idea of inference of rigid objects using deep learning in point clouds. A point-to-point correspondence assignment is performed with a joint classification and segmentation within a point cloud system [83]. Capellen [27] suggested that ConvPoseCNN has evolved from the concept of PoseCNN but can avoid cutting individual objects. Instead, it offers accurate predictions for pixel-based translation of object poses and orientation modules and has been replaced with a complete CNN prediction network. Also, [191] recently removed the ROI pooled orientation layer and introduced PVNet (Pixelwise Voting Network) to deny pixel-based vectors and use them for key-point positions.

### C. 6D POSE ESTIMATION DIRECTLY FROM RGB-D IMAGES

A scene coordinate regression (SCoRe) forest is used, trained in a specific scene, employs only RGB-D image pixel comparison features and has fast calculation accuracy. The proposed method is an RNSAC-based pose optimization algorithm where SCoRe Forest is evaluated by the RNSAC algorithm and makes accurate posture estimates [223]. Since the additional depth channel of the RGB-D image helps extracts the entire 6D pose (3D rotation and 3D translation) of rigid object instances present in the scene. The core objective of the approach is the intermediate representation of the form of a dense 3D object coordinate labelled and paired with a dense class.

On the other hand, Taylor [246] did not predict 6 DoF directly from an RGB image but instead followed the object's coordinates in that image. Each pixel in this image points to a coordinate of the canonical body in a canonical position called VM (Vitruvian Manifold). The popular RF (Random Forest) [3] classifier is used to vote here, and geometric validity is used.

Brachmann *et al.*, [14] provided an idea that is both an extension and combination of [223] and [246]. This hybrid concept estimates the 6D pose of a specific object from a single RGB-D image. Wang *et al.*, [258] initiated a compelling method in cluttered scenes, which can successfully predict the object's posture. It has mixed colour and depth data from the RGB-D image and then integrates repetitive refinement methods into neural network architectures.

Li *et al.*, [141] applied CNN to process the depth image as an additional image channel. However, the built-in 3D structure in the depth channel was neglected. In contrast, the geometric features of the dense fusion method convert pixels into sectional depths into 3D point clouds by applying built-in cameras. The proposed DPOD (Dense Pose Object Detector) applies PNP and RANSAC to compute an input image and a map of dense multi-class 2D / 3D correspondence between available 3D models [282].

### D. INSTANCE-LEVEL 6 DoF POSE ESTIMATION

Collet *et al.*, [38] created 3D object metric models using local descriptors of different images. Each model was optimized to easily fit a sequential training image, resulting in the best possible alignment between the 3D model and the original object. It combines the well-known RANSAC [65] and Mean Shift algorithm [36] to register multiple instances of each object that can successfully guess the 6-DOF pose for any complex and chaotic scene. In addition, it can handle randomly complex non-planning objects, powerful to handle outliers and occlusions, and able to control illumination, scale and rotation change.

The vision-based system, which is actually an extension of Gordon's method [79], enables the accurate localization initialization step called POSESEQ and enables full pose inference in object recognition in a complete cluttered environment. Thanh *et al.* presented LieNet [51], as a unique template-based pose estimation method that uses the Lie algebraic rotation matrix to estimate the rotation matrix of an object. It estimates the translation vector by predicting the distance of the object from the centre of the camera. This method takes the input of an image and then outputs the object's identification with a 6D pose, including a bounding box, label, and segmentation mask.

Vidal *et al.*, [256] developed a method that followed the basic structure of the point pair feature (PPF) method introduced by Drost [56], which is a combination of two levels, such as global modelling and local matching. The main structure identifies the rotation points, model points and angles of each scene. The expansion of Vidal's work is the concept of the posture of free-form objects, critical work in favour of a highly confused autonomous system. A novel pre-processing step has been added here, transforming the classification into a better efficient feature matching method.

### E. CATEGORY-LEVEL 6 DoF POSE ESTIMATION

Sahin *et al.*, [208] covers various challenges for 6D pose estimation such as inconsistency of viewpoint, objects (both texture and texture-less), curbs, cluttered scene and identical objects. Wang *et al.*, [261] has created a method that assumes both 6D poses of hidden object instances without an object CAD model in an RGB-D image. Furthermore, a novel concept called NOCS (Normalised Object Coordinate Space)

**FIGURE 2.** Increasing amount of efforts in literature on monocular 3D object detection in recent years.



**FIGURE 3.** Increasing amount of efforts in literature on 6D pose estimation in recent years.

has been introduced here, representing a partnership principle for all possible instances of an object.

Schuster *et al.*, [216] evaluates dense 3D data located in multiple light situations and applies online graph SLAM to generate a dense 3D composite map and estimates 6D poses. This technique also creates a fancy graph topology for incorporating the results of local reference filters and overall high-bandwidth sensor data into sub-maps.

### F. FEATURE MATCHING METHODS
To solve the 6D object pose hypothesis and ensure the best possible accuracy, Krull [126] successfully applied Reinforcement Learning to the pose agent classification for the first time. Each decision here follows the potential distribution of a stochastic policy gradient approach that takes a direct gradient in terms of the expected loss of interest.

### G. TEMPLATE-MATCHING TECHNIQUES
Hinterstoisser *et al.*, [91] built a framework called LineMod for automatic detection and tracking of 3D objects based on the latest template-based approach that uses both depth and colour images to capture the object's presence and 3D shape on a set of templates with different aspects of an object. Also, the 3D model can be used for the accurate estimation of the position of the object. Tejani *et al.*, [247] developed a novel patch-based framework where a Latent-Class Hough Forests method for 3D object detection was introduced, and estimations were made in a heavily cluttered and occluded environment. This method absorbs the classification labels during training, and as a by-product, it creates the right image-ground mask.

### H. CNN/ DEEP LEARNING - BASED APPROACHES
Krull [125] presented a model for 6D pose estimation, which applied a CNN to map images and revealed that training on a single object was sufficient and that CNN successfully generalized all the different objects and backgrounds of an image. Rangesh *et al.*, [200] applied an exclusive idea for a 3D identification box suitable for the object on the ground

to combine 2D visual context, 3D dimension and ground plane. Eppner *et al.*, [60] presented and evaluated the winning system for the 2015 Amazon Picking Challenge, where they created four key aspects of system building: integration, manipulation, manipulation planning, and estimation.

Google has announced the release of MediaPipe, a 3D object detection pipeline that identifies objects in 2D images on everyday objects and estimates their poses and sizes. MediaPipe is a cross-platform structure that builds ML pipelines and creates 3D bounding boxes with augmented reality (AR) [5] and identifies additional information such as camera pose, 3D point cloud, lighting and planar surfaces [85], [268]. Basically, MediaPipe performs object detection, face detection, hand tracking, hair segmentation with ML frameworks called Tensorflow and Tensorflow Lite [1].

A novel model [101] designed to predict the pose and size of an object from a monocular RGB image has applied a multi-task-learning approach named MobileNetv2 [212] and predicts object size. The Gaussian regression task applies a pose estimation algorithm (EPnP) [136] to the final 3D coordinates for the bounding box. A novel model [101] designed to predict the pose and size of an object from a monocular RGB image has applied a multi-task-learning approach named MobileNetv2 [212] and predicts object size. The Gaussian regression task applies a pose estimation algorithm (EPnP) [136] to the final 3D coordinates for the bounding box.

Tremblay *et al.*, [251] introduced the first one-shot deep neural network for robotic manipulation trained only on synthetic data capable of achieving 6-DoF object pose estimates of 3D objects. The system is called DOPE (Deep Object Pose Estimation), which applies the Perspective-N-Point (PNP) algorithm, which combines 3D bounding boxes with 2D images. Li [141] has proposed a pose correction algorithm where the solution is to correct the pose because the object is being observed from the centre line of the camera. It is a multi-philosophy fusion framework with a single philosophical ambiguity and quick guess selection based on a voting scheme.

**TABLE 5.** Advantages and disadvantages of some state-of-the-art 6D pose estimation methods.

| No. | Algorithms | Advantages | Disadvantages | DataSets |
|---|---|---|---|---|
| 1 | LieNet [51] | Quite simple and inexpensive pose refiner | Poor performance for rotational symmetry objects.(e.g., coffee mug) | LINEMOD and Tejani's dataset |
| 2 | DenseFusion[258] | Shows robustness when the scene is extremely cluttered. | Unable to grasp a specific type of object. | YCB-Video and LineMOD |
| 3 | Brachmann [13] | Swift, scalable, powerful and exceptionally perfect method for generic objects. | Direct using an auto-context random forest hampered test performance. | RGB - D Images [92] |
| 4 | POSESEQ [38] | Handles complex non-planar objects with great speed. | Poor latency scale into household environments. | RGB - Images |
| 5 | MOPED [162] | Improved scalability and optimized speed with high robustness, accuracy and low latency. | Computationally expensive and large Image show low performance. | 4 Real World Benchmark |
| 6 | Brachmann [14] | Covers the textured and textureless, rigid and non-rigid objects, symmetrical objects. | The pipeline is linear in the number of objects. | RGB - D Images |
| 7 | LINEMOD [91] | Easy to deploy, reliable, and fast. | Poor performances and suffers for false positives data. | LINEMOD [91] |
| 8 | SSD-6D [116] | Stable training and robust prediction. | Computationally expensive. | LINEMOD and OCCLUSION |
| 9 | BB8 [198] | Perform better on rotational symmetrical and similar objects. | Performs badly for Texture-LESS dataset. | LINEMOD, OCC, T-LESS. |
| 10 | Iryna [79] | An efficient, incremental and jitters reduction method. | Unscalable for operation in large environments | Live video sequence |
| 11 | Posecnn [276] | ShapeMatch-Loss for symmetric objects pose estimation | Estimate pose for only color images. | YCB-Video, LINEMOD and OCC |
| 12 | PoseAgent [126] | Dramatically reduces computation time and variance. | Estimate pose for only color images. | Krull et al [127] |
| 13 | Deep-6DPose [50] | Better trade-off between speed and accuracy. | Not very strong to nearly rotational symmetric objects. | LINEMOD [91] |
| 14 | DeepIM [143] | Accurately manage poses and shapes for texture-less and unseen objects. | Estimate pose for color images only. | LINEMOD and OCCLUSION |
| 15 | BB8 [198] | Effective method. | Relatively Slow for large image. | LINEMOD and OCCLUSION |

**TABLE 6.** 6D full object pose estimation paper collection.

| No. | Approach | References |
|---|---|---|
| 1 | Pose from RGB images. | [208], DenseFusion[258], MPOED[39], SSD-6D[116], BB8 [198], [248], [101], [19], [26], 6D-VNet [269], GPP [200], [289], [11], Deep-6DPose [50], RePOSE [106], DPOD [283], E2E6DoF [81] |
| 2 | Pose from depth / point cloud method. | [290], [168], PointPoseNet [83], [24], [255], [286]. |
| 2 | Pose from RGB-D data. | [217], [99], [97], [96], [98], [282], [246], [55], [204], [37]. [211], [180], [231]. |
| 3 | Instance-Level 6 DoF Pose Estimation | POSESEQ [38], MOPED [162], LieNet [51] , [24], [208], [256], [232] [2], GSNet [115]. |
| 4 | Category-Level 6 DoF Pose Estimation | [79], [185], C3DPO [181], [208], [213], [261], [216], [149], Pix2Pose [188]. |
| 5 | Feature matching methods | [157], [39], [257], [14], [126], [166], [240], EfficientPose [18], |
| 7 | Template-matching techniques | [91], [117], [205], [247], [189], |
| 8 | CNN/ Deep Learning -based approaches. | [125],[276], GPP [200], [27], Deep3DBox [172], [60], [95], [251], [141], [75], DeepIM [143], DeepHMap++ [71], HRNet [30]. |
| 9 | Template-clustering approaches. | PVNet [191], CT-LineMod[288], HybridPose [234] |

On the other hand, a model called DeepIM [143] is able to predict a relational pose transformation by applying 3D location and 3D orientation and a repetitive training process. The network FlowNetSimple architecture uses the backbone network project as DeepHMap++, which centres a two-stage pipeline and integrates two learning concepts to estimate 6D poses of invisible objects in challenging scenes [71].

## I. TEMPLATE-CLUSTERING APPROACHES

Zhang *et al.*, [288] has proposed City-LineMod (advanced to Cognitive Template-Clustering Line mode) method. The technique applies a 7D (4D geometry + 3D texture) cognitive feature vector to restore the standard 3D spacing points in the patch-linemode clustering method. Moreover, the distance of different 3D spatial points will also be affected by the 4D additional information regarding the direction and width of the features.

## J. HYBRID POSE METHOD

Martinez *et al.*, [162] presented a hybrid GPU / CPU architecture that uses parallelism at all levels named MOPED (Multiple Object Position Estimation and Detection), a bright and measurable perception concept for both object recognition and fracture estimation. Furthermore, a mode based on another object recognition algorithm known as POS-ESEQ [38] showed a massive increase in scalability and accuracy and optimizes the algorithm's speed. Technically, MOPED has employed a new feature-matching algorithm that optimizes databases to handle complexity and a robust pose merge algorithm capable of efficiently rejecting outsiders with matching K-NN (K-Nearest Neighbour) method where $k > 2$ [245]. The default classification algorithm and SIFTGPU is the MOPED feature extraction algorithm.

## VI. FUTURE RESEARCH DIRECTIONS

From the above discussion, it is clear that plenty of work has been done on 3D object detection, which forms a solid foundation for this field. Nevertheless, further research is needed as 6D pose estimation systems have not yet performed adequately. Therefore, this part of the article will give some possible ideas of the future directions for both sectors, which will help understand the status and involvement of 3DOD and 6 DPE.

## A. FUTURE RESEARCH DIRECTIONS FOR 3D OBJECT DETECTION

### 1) DETECTING A RIGID OBJECT

Several existing work [32], [33], [151], [271] showed the efficacy of deep learning in detecting a rigid object. Even though the classifiers are mainly focused on the ''car'' category, the concept of these methods can be contextual to all other solid and inflexible types of objects. The accurate detection of the 3D rigid object is a complex job and is very significant in the domain of computer vision. Currently, using some proposed deep learning techniques, we can detect inflexible objects, but still, lots of works need to do to make the process flawless. In developing an autonomous car, accurately identifying rigid objects can be a significant research idea.

### 2) HANDING ROTATIONALLY SYMMETRIC OBJECTS

Identifying half and full symmetry objects like a coffee mug or glass is a confusing and complex matter that classifiers usually fail to give accurate results. Do *et al.*, [51] has suggested an idea to overcome this complicated problem, but not entirely successful. Although the work on symmetric object detection is starting to get deeper, not much work has been done to date. More attention and research needs to be done on this case of symmetrical object identification.

### 3) TRACKING A OBJECT FROM VIDEO

One possible incoming direction is to simultaneously explore its application as part of a system that uses repetitive neural networks to detect and track objects in video [116]. In addition, the intense colour variation between the CAD model and the visual avatar is a significant work. Another potential research aspect in this context is online model learning and relocalization [223]. A hypothesis can be developed to represent both single and multiview with an extended update of a new frame [141]. In addition, to avoid the problem of proper loss, the term-balancing required for upcoming potential research direction [116].

## B. FUTURE RESEARCH DIRECTIONS FOR 6D POSE ESTIMATION

### 1) IMPROVING THE VO (VISUAL ODOMETRY)

Appropriate VO (Visual Odometry) is mandatory in the context of autonomous driving; Thus, the future design of both automated car and street scene construction needs to be improved [291]. It is not possible to apply driverless cars without the proper implementation of VO.

### 2) IMPROVING THE 6D POSE ESTIMATE ACCURACY

One can improve the version of the DeepIM method [143] for autonomous applications to produce accurate 6D pose estimates from high-resolution camera images (colour only) at high frame rates with a large field view. The authors also mentioned using stereotype camera images as input to improve the quality of this method. Another work can be done by combining the advanced two-step method to transform it into a new pose tracking framework where the pose parameters from the previous frame can be reused to replace the pose detection step in DeepHMap [71]. Adding a branch to the back for object segmentation in DeepHMap may provide some additional regularization.

### 3) IMPROVING THE MAP OR VPS

It is an open challenge to efficiently and consistently merge sub-maps into multi-robot systems to create a long-term mapping system, aiming at improving the algorithm that matches the map [216]. The globalization strategy combines visual positioning services (VPS), street view, and machine learning for more accurate location and adaptation detection. Mutual technology is essential to enhance the correct positioning and orientation of blue dots on digital maps in our cars, smartphones, and up-to-date interactions.

### 4) IMPROVING THE POSE OF SYMMETRICAL OBJECTS

Since managing the poses of symmetrical or symmetrical objects is a complex task, relevant classifiers should be improved to accomplish the task efficiently [276]. In order to properly manage symmetrical properties, methods need to learn the symmetry of objects and update their capabilities [27]. One of the notable tasks may be to manage the symmetry property of objects and pose estimation automatically.

### 5) IMPROVING THE FUNCTION OF MOBILE MANIPULATORS

The efficiency of the POSSEQ [38] classifier can be enriched by enabling mobile manipulators to work more perfectly to communicate with the crowd's internal environment. Some hypothetical 6DoF pose [269] reprocessing techniques will be filtered using repetitive closet point-based algorithms or repetitive retrieval networks. Also, classifiers need to successfully model and recognize scenes of different sizes and complexities in large environments [79] (campus, laboratory, shopping centre or a museum).

### 6) IMPROVING THE 3D POINT CLOUD NETWORKS

A number of 3D point cloud networks can be replaced directly by the PointNet network [83] for potential improvement in accurate 3D object detection and 6DF pose estimation. A computational budget can be created to know the appropriate time for the softer version of the PoseAgent [126] classification. For parallelism, multiple computational cores can be applied by advanced PoseAgent. In addition, training can be provided to replace the processing steps of an existing CNN method and improve the results by observing and predicting updated postures from the given images [125].

### 7) IMPROVING THE DATASET

To deal with the common challenges of objects, such as reflective and texture-less objects, and the adverse conditions, such as occlusion and changing lighting conditions, we can integrate some multi-dimensional object models into the dataset packages. To facilitate the reconstruction of indoor and outdoor dynamic scenes, 4D or 5D models can be added to the dataset, which can play an important role in any visual applications such as navigational systems for moving objects (for example: autonomous car) [281].

### 8) REMOVING THE VISIBLE GAP BETWEEN MACHINE PERFORMANCE AND THAT OF HUMAN's

In AppolloCar3D, researchers [238] mentioned four visible surfaces and manually defines a correspondence between critical points and surfaces. They suggested that a total of 66 key points were assigned to every single car model (for both SUVs and Sedans). According to [238], since people cannot memorize the semantic meaning of 66 key points correctly, there is a noticeable gap ($\sim$ 10 %) in between algorithms/machines with humans. Henceforth, correctly

resolving visible gaps between machines and humans can be a future inspiration for research.

### 9) EXPLORE GEOMETRIC PROPERTIES

Estimating the 6DoF pose of an object from a single RGB image is a significant and challenging task, especially under heavy occlusion and for the Texture-Less object. In such a case, the exploring of geometric features needs to be improved to estimate the 6 DOF object more efficiently [81].

## VII. CONCLUSION

This review paper studies the state-of-the-art deep learning techniques for 3D object detection and 6D pose estimation. Most current object detection methods identify images with a 2D bounding box technique that can recognize both the position and range of the objects in the image. However, recognizing a vehicle as a 2D BB is not always sufficient for perfect autonomous driving. Therefore, predicting the position of the 3D object from the images is just as important as determining the 2D position of the vehicle. For 3D object detection, current works report sophisticated results using RGB / RGB-D imagery, point cloud, and fusion-based techniques.

Here, with the help of this review, we have addressed the advantages and disadvantages of each of the basic techniques, both 3D object detection and 6D pose estimation techniques. We have also mentioned some traditional theoretical evaluation metrics and summarised the popular Big Image datasets applied by well-known object identification and pose estimation methods. Since the deep learning method of 3D object detection and 6D pose estimation are not as mature as 2D object detection, research is needed for real-time operation. From now on, a significant improvement needs to be made to manage a fast and reliable 3 DOD and 6 DPE system across a broad set of real-time practical applications. Although RGB-D is much simpler than RGB, it faces problems for some depth issues, such as not being able to recognize small objects properly.

Several classifications have been proposed in the 6D Pose estimation functions, such as the point addition method, the template matching method, the Hough forest method, and the deep learning method. However, the effectiveness of the proposed classifiers is still far from the level of actual application, which should be able to successfully predict 6D poses of multi-objects, including severe occurrence and chaos scene situations. Therefore, this article presents an in-depth review of the most significant work to date on in-depth learning-based 3D object detection and 6D pose estimation systems. Until then, we believe that this review article can be cited and used as a sample source of reference and forms an important endorsement to the research community.

### REFERENCES

[1] *An End-to-End Open Source Machine Learning Platform.*
[2] G. N. Albanis, N. Zioulis, A. Chatzitofis, A. Dimou, D. Zarpalas, and P. Daras, "On end-to-end 6DOF object pose estimation and robustness to object scale," in *Proc. ML Reproducibility Challenge*, 2021, pp. 1–9.

[3] M. Y. Arafat, S. Hoque, and D. M. Farid, "Cluster-based under-sampling with random forest for multi-class imbalanced classification," in *Proc. 11th Int. Conf. Softw., Knowl., Inf. Manage. Appl. (SKIMA)*, Dec. 2017, pp. 1–6.

[4] M. Y. Arafat, S. Hoque, S. Xu, and D. M. Farid, "An under-sampling method with support vectors in multi-class imbalanced data classification," in *Proc. 13th Int. Conf. Softw., Knowl., Inf. Manage. Appl. (SKIMA)*, Aug. 2019.

[5] T. R. Azuma, *A Survey of Augmented Reality*, vol. 6. Hughes Research Laboratories, 1997.

[6] D. H. Ballard, "Generalizing the Hough transform to detect arbitrary shapes," *Pattern Recognit.*, vol. 13, no. 2, pp. 111–122, 1981.

[7] M. R. Banham and A. K. Katsaggelos, "Digital image restoration," *IEEE Signal Process. Mag.*, vol. 14, no. 2, pp. 24–41, Mar. 1997.

[8] W. Bao, B. Xu, and Z. Chen, "MonoFENet: Monocular 3D object detection with feature enhancement networks," *IEEE Trans. Image Process.*, vol. 29, pp. 2753–2765, 2020.

[9] I. Barabás, A. Todoruţ, N. Cordoş, and A. Molea, "Current challenges in autonomous driving," *IOP Conf. Ser., Mater. Sci. Eng.*, vol. 252, Oct. 2017, Art. no. 012096.

[10] J. L. Barron, D. J. Fleet, and S. S. Beauchemin, "Performance of optical flow techniques," *Int. J. Comput. Vis.*, vol. 12, no. 1, pp. 43–77, 1994.

[11] G. Billings and M. Johnson-Roberson, "SilhoNet: An RGB method for 6D object pose estimation," *IEEE Robot. Autom. Lett.*, vol. 4, no. 4, pp. 3727–3734, Oct. 2019.

[12] A. Borji, M.-M. Cheng, H. Jiang, and J. Li, "Salient object detection: A benchmark," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5706–5722, Dec. 2015.

[13] E. Brachmann, F. Michel, A. Krull, M. Y. Yang, S. Gumhold, and C. Rother, "Uncertainty-driven 6D pose estimation of objects and scenes from a single RGB image," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 3364–3372.

[14] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and A. C. Rother, "Learning 6D object pose estimation using 3D object coordinates," in *Computer Vision—ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham, Switzerland: Springer, 2014, pp. 536–551.

[15] G. Brazil and X. Liu, "M3D-RPN: Monocular 3D region proposal network for object detection," Tech. Rep., 2019.

[16] M. Bruijning, M. D. Visser, C. A. Hallmann, and E. Jongejans, "Trackdem: Automated particle tracking to obtain population counts and size distributions from videos in R," *Methods Ecol. Evol.*, vol. 9, no. 4, pp. 965–973, Apr. 2018.

[17] F. Bu, T. Le, X. Du, R. Vasudevan, and M. Johnson-Roberson, "Pedestrian planar LiDAR pose (PPLP) network for oriented pedestrian detection based on planar LiDAR and monocular images," *IEEE Robot. Autom. Lett.*, vol. 5, no. 2, pp. 1626–1633, Apr. 2020.

[18] Y. Bukschat and M. Vetter, "EfficientPose: An efficient, accurate and scalable end-to-end 6D multi object pose estimation approach," in *Proc. CVPR*, 2020.

[19] M. Burenius, J. Sullivan, and S. Carlsson, "3D pictorial structures for multiple view articulated pose estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2013, pp. 3618–3625.

[20] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "NuScenes: A multimodal dataset for autonomous driving," 2019, *arXiv:1903.11027*. [Online]. Available: http://arxiv.org/abs/1903.11027

[21] Z. Cai and N. Vasconcelos, "Cascade R-CNN: Delving into high quality object detection," in *Proc. Comput. Vis. Pattern Recognit.*, 2017.

[22] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M. Dollár, "Yale-CMU-Berkeley dataset for robotic manipulation research," *Int. J. Robot. Res.*, vol. 36, no. 3, pp. 261–268, Mar. 2017.

[23] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollár, "The YCB object and model set: Towards common benchmarks for manipulation research," in *Proc. Int. Conf. Adv. Robot. (ICAR)*, Jul. 2015.

[24] D. Campbell, L. Petersson, L. Kneip, and H. Li, "Globally-optimal inlier set maximisation for camera pose and correspondence estimation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 2, pp. 328–342, Feb. 2020.

[25] R. J. Campbell and P. J. Flynn, "A survey of free-form object representation and recognition techniques," *Comput. Vis. Image Understand.*, vol. 81, no. 2, pp. 166–210, 2001.

[26] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2D pose estimation using part affinity fields," Tech. Rep., 2016.

[27] C. Capellen, M. Schwarz, and S. Behnke, "ConvPoseCNN: Dense convolutional 6D object pose estimation," in *Proc. 15th Int. Joint Conf. Comput. Vis., Imag. Comput. Graph. Theory Appl.*, 2020.

[28] F. Chabot, M. Chaouch, J. Rabarisoa, C. Teulière, and A. T. Chateau, "Deep MANTA: A coarse-to-fine many-task network for joint 2D and 3D vehicle analysis from monocular image," Tech. Rep., 2017.

[29] X. A. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "ShapeNet: An information-rich 3D model repository," in *Proc. Comput. Vis. Pattern Recognit.*, 2015.

[30] B. Chen, A. Parra, J. Cao, N. Li, and T.-J. Chin, "End-to-end learnable geometric vision by backpropagating PnP optimization," Tech. Rep., 2020.

[31] J. Chen and T. Bai, "SAANet: Spatial adaptive alignment network for object detection in automatic driving," *Image Vis. Comput.*, vol. 94, Feb. 2020, Art. no. 103873.

[32] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3D object detection for autonomous driving," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 2147–2156.

[33] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3D object detection network for autonomous driving," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017.

[34] X. Chen, K. Kundu, Y. Zhu, A. G. Berneshawi, H. Ma, S. Fidler, and R. Urtasun, "3D object proposals for accurate object class detection," in *Advances in Neural Information Processing Systems*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Red Hook, NY, USA: Curran Associates, 2015, pp. 424–432.

[35] M.-M. Cheng, Y. Liu, W.-Y. Lin, Z. Zhang, L. P. Rosin, and S. P. H. Torr, "BING: Binarized normed gradients for objectness estimation at 300fps," *Comput. Vis. Media*, vol. 5, no. 1, pp. 3–20, 2019.

[36] Y. Cheng, "Mean shift, mode seeking, and clustering," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 17, no. 8, pp. 790–799, Aug. 1995.

[37] C. Li, J. Bohren, E. Carlson, and G. D. Hager, "Hierarchical semantic parsing for object pose estimation in densely cluttered scenes," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2016, pp. 5068–5075.

[38] A. Collet, D. Berenson, S. S. Srinivasa, and D. Ferguson, "Object recognition and full pose registration from a single image for robotic manipulation," in *Proc. IEEE Int. Conf. Robot. Autom.*, May 2009, pp. 48–55.

[39] A. Collet, M. Martinez, and S. S. Srinivasa, "The MOPED framework: Object recognition and pose estimation for manipulation," *Int. J. Robot. Res.*, vol. 30, no. 10, pp. 1284–1306, Apr. 2011.

[40] A. Conner-Simons and R. Gordon, "Self-driving cars for country roads," Tech. Rep., May 2018.

[41] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, 1995, doi: 10.1023/A:1022627411411.1995.

[42] R. Cupec, I. Vidović, D. Filko, and P. Đurović, "Object recognition based on convex hull alignment," *Pattern Recognit.*, vol. 102, Jun. 2020, Art. no. 107199.

[43] J. Dai, Y. Li, K. He, and J. Sun, "R-FCN: Object detection via region-based fully convolutional networks," 2016, arXiv:1605.06409. https://arxiv.org/abs/1605.06409

[44] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," in *Proc. Comput. Vis. Pattern Recognit.*, 2017.

[45] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 1, Jun. 2005, pp. 886–893.

[46] E. R. Davies, *Machine Vision: Theory, Algorithms, Practicalities*. San Mateo, CA, USA: Morgan Kaufmann, 2004.

[47] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, V. Q. Le, and Y. A. Ng, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Red Hook, NY, USA: Curran Associates, 2012, pp. 1223–1231.

[48] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009.

[49] W. Ding, S. Li, G. Zhang, X. Lei, and H. Qian, "Vehicle pose and shape estimation through multiple monocular vision," Tech. Rep., 2018.

[50] T.-T. Do, M. Cai, T. Pham, and I. Reid, "Deep-6DPose: Recovering 6D object pose from a single RGB image," 2018, *arXiv:1802.10367*. [Online]. Available: http://arxiv.org/abs/1802.10367

[51] T.-T. Do, T. Pham, M. Cai, and D. I. Reid, "LieNet: Real-time monocular object instance 6D pose estimation," in *Proc. BMVC*, 2018.

[52] J. Dokic, B. Müller, and G. Meyer, "European roadmap smart systems for automated driving," in *Proc. Eur. Technol. Platform Smart Syst. Integr. (EPoSS)*, Apr. 2015.

[53] P. Dollár, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: An evaluation of the state of the art," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 4, pp. 743–761, Apr. 2012.

[54] P. Dollár, R. Appel, S. Belongie, and P. Perona, "Fast feature pyramids for object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 8, pp. 1532–1545, Aug. 2014.

[55] A. Doumanoglou, R. Kouskouridas, S. Malassiotis, and T.-K. Kim, "Recovering 6D object pose and predicting next-best-view in the crowd," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016.

[56] B. Drost, M. Ulrich, N. Navab, and S. Ilic, "Model globally, match locally: Efficient and robust 3D object recognition," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 998–1005.

[57] B. Drost, M. Ulrich, P. Bergmann, P. Hartinger, and C. Steger, "Introducing MVTec ITODD—A dataset for 3D object recognition in industry," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops (ICCVW)*, Oct. 2017.

[58] M. Engelcke, D. Rao, D. Z. Wang, C. H. Tong, and I. Posner, "Vote3Deep: Fast object detection in 3D point clouds using efficient convolutional neural networks," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2017.

[59] M. Enzweiler and D. M. Gavrila, "A multilevel mixture-of-experts framework for pedestrian classification," *IEEE Trans. Image Process.*, vol. 20, no. 10, pp. 2967–2979, Oct. 2011.

[60] C. Eppner, S. Höfer, R. Jonschkowski, R. Martín-Martín, A. Sieverling, V. Wall, and O. Brock, "Lessons from the Amazon picking challenge: Four aspects of building robotic systems," in *Proc. 12th Robot., Sci. Syst.*, 2017, pp. 4831–4835.

[61] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal visual object classes challenge: A retrospective," *Int. J. Comput. Vis.*, vol. 111, no. 1, pp. 98–136, Jan. 2015.

[62] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. (2007). *The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results*. [Online]. Available: http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html

[63] J. Fang, L. Zhou, and G. Liu, "3D bounding box estimation for autonomous vehicles by cascaded geometric constraints and depurated 2D detections using 3D results," Tech. Rep., 2019.

[64] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, Sep. 2010.

[65] M. A. Fischler and R. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Commun. ACM*, vol. 24, no. 6, pp. 381–395, Jun. 1981.

[66] R. Fisher, "CVonline: The evolving, distributed, non-proprietary, on-line compendium of computer visio," School Inform., Univ. Edinburgh, Edinburgh, U.K., Tech. Rep., 2019.

[67] A. D. Forsyth and J. Ponce, *Computer Vision: A Modern Approach*, no. 792. London, U.K.: Pearson, 2011.

[68] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *J. Comput. Syst. Sci.*, vol. 55, no. 1, pp. 119–139, Aug. 1997.

[69] A. Frome, D. Huber, R. Kolluri, T. Bülow, and J. Malik, "Recognizing objects in range data using regional point descriptors," in *Computer Vision—ECCV 2004*, T. Pajdla and J. Matas, Eds. Berlin, Germany: Springer, 2004, pp. 224–237.

[70] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018.

[71] M. Fu and W. Zhou, "DeepHMap++: Combined projection grouping and correspondence learning for full DoF pose estimation," *Sensors*, vol. 19, no. 5, p. 1032, Feb. 2019.

[72] L. A. Gatys, A. S. Ecker, and M. Bethge, "A neural algorithm of artistic style," 2015, *arXiv:1508.06576*. [Online]. Available: http://arxiv.org/abs/1508.06576

[73] S. K. Gehrig and F. J. Stein, "Dead reckoning and cartography using stereo vision for an autonomous car," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Systems. Hum. Environ. Friendly Robots High Intell. Emotional Quotients*, vol. 3, Oct. 1999, pp. 1507–1512.

[74] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? The KITTI vision benchmark suite," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2012, pp. 3354–3361.

[75] N. Gessert, M. Schlüter, and A. Schlaefer, "A deep learning approach for pose estimation from volumetric OCT data," *Med. Image Anal.*, vol. 46, pp. 162–179, May 2018.

[76] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2014, pp. 580–587.

[77] R. Girshick, "Fast R-CNN," Tech. Rep., 2015.

[78] A. González, D. Vázquez, A. M. López, and J. Amores, "On-board object detection: Multicue, multimodal, and multiview random forest of local experts," *IEEE Trans. Cybern.*, vol. 47, no. 11, pp. 3980–3990, Nov. 2017.

[79] I. Gordon and D. Lowe, "What and where: 3D object recognition with accurate pose," in *Toward Category-Level Object Recognition*, vol. 4170, J. Ponce, M. Hebert, C. Schmid, and A. Zisserman, Eds. Berlin, Germany: Springer, Jan. 2006, pp. 67–82.

[80] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Bennamoun, "Deep learning for 3D point clouds: A survey," in *Proc. CVPR*, 2020.

[81] A. Gupta, J. Medhi, A. Chattopadhyay, and V. Gupta, "End-to-end differentiable 6DoF object pose estimation with local and global constraints," Tech. Rep., 2020.

[82] F. K. Gustafsson, M. Danelljan, and T. B. Schon, "Accurate 3D object detection using energy-based models," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2021.

[83] F. Hagelskjaer and A. Buch, "PointVoteNet: Accurate object detection and 6 DOF pose estimation in point clouds," in *Proc. Comput. Vis. Pattern Recognit. (CVPR)*, Dec. 2019.

[84] J. A. Hawkins, "Mit built a self-driving car that can navigate unmapped country roads," Tech. Rep., May 2015.

[85] M. Hays and T. Mullen, "Mediapipe on the web, blog," Tech. Rep., Jan. 2020.

[86] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017.

[87] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016.

[88] Q. He, Z. Wang, H. Zeng, Y. Zeng, S. Liu, and B. Zeng, "SVGA-Net: Sparse voxel-graph attention network for 3D object detection from point clouds," in *Proc. CVPR*, 2020.

[89] Z. He, W. Feng, X. Zhao, and Y. Lv, "6D pose estimation of objects: Recent technologies and challenges," *Appl. Sci.*, vol. 11, no. 1, p. 228, Dec. 2020.

[90] G. Hetzel, B. Leibe, P. Levi, and B. Schiele, "3D object recognition from range images using local feature histograms," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, vol. 2, Dec. 2001, p. 2.

[91] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab, "Model based training, detection and pose estimation of texture-less 3D objects in heavily cluttered scenes," in *Computer Vision—ACCV 2012*, K. M. Lee, Y. Matsushita, J. M. Rehg, and Z. Hu, Eds. Berlin, Germany: Springer, 2013, pp. 548–562.

[92] S. Hinterstoisser, V. Lepetit, N. Rajkumar, and K. Konolige, *Going Further With Point Pair Features* (Lecture Notes in Computer Science). 2016, pp. 834–848.

[93] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, and B. Kingsbury, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal Process. Mag.*, vol. 29, no. 6, pp. 82–97, Nov. 2012.

[94] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

[95] D.-C. Hoang, A. J. Lilienthal, and T. Stoyanov, "Panoptic 3D mapping and object pose estimation using adaptively weighted semantic information," *IEEE Robot. Autom. Lett.*, vol. 5, no. 2, pp. 1962–1969, Apr. 2020.

[96] T. Hodan, P. Haluza, S. Obdrzalek, J. Matas, M. Lourakis, and X. Zabulis, "T-LESS: An RGB-D dataset for 6D pose estimation of texture-less objects," in *Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, Mar. 2017.

[97] T. Hodan, E. J. S. Matas, and S. Obdržálek, "On evaluation of 6D object pose estimation," in *Proc. ECCV Workshops*, 2016.

[98] T. Hodan, F. Michel, E. Brachmann, W. Kehl, A. G. Buch, D. Kraft, B. Drost, J. Vidal, S. Ihrke, X. Zabulis, C. Sahin, F. Manhardt, F. Tombari, T.-K. Kim, J. Matas, and C. Rother, "Bop: Benchmark for 6D object pose estimation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 19–34.

[99] T. Hodaň, X. Zabulis, M. Lourakis, Š. Obdržálek, and J. Matas, "Detection and fine 3D pose estimation of texture-less objects in RGB-D images," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Sep. 2015, pp. 4421–4428.

[100] S. Hoque, D. M. Farid, and M. Y. Arafat, "Advanced data balancing method with SVM decision boundary and bagging," in *Proc. 6th IEEE CSDE*. Melbourne, VIC, Australia: CQUniv. Australia, 2019.

[101] T. Hou, A. Ahmadyan, L. Zhang, J. Wei, and M. Grundmann, "MobilePose: Real-time pose estimation for unseen objects with weak shape supervision," in *Proc. Comput. Vis. Pattern Recognit. (CVPR)*, 2020.

[102] H.-N. Hu, Q.-Z. Cai, D. Wang, J. Lin, M. Sun, P. Krähenbühl, T. Darrell, and F. Yu, "Joint monocular 3D vehicle detection and tracking," Tech. Rep., 2018.

[103] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy, "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017.

[104] V. Ilci and C. Toth, "High definition 3D map creation using GNSS/IMU/LiDAR sensor integration to support autonomous vehicle navigation," *Sensors*, vol. 20, no. 3, p. 899, Feb. 2020.

[105] A. Ioannidou, E. Chatzilari, S. Nikolopoulos, and I. Kompatsiaris, "Deep learning advances in computer vision with 3D data: A survey," *ACM Comput. Surveys*, vol. 50, no. 2, pp. 1–38, Jun. 2017.

[106] S. Iwase, X. Liu, R. Khirodkar, R. Yokota, and M. K. Kitani, "Repose: Real-time iterative rendering and refinement for 6D object pose estimation," Tech. Rep., 2021.

[107] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander, "Joint 3D proposal generation and object detection from view aggregation," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Oct. 2018, pp. 1–8.

[108] B. Jähne and H. Haußecker, *Computer Vision and Applications a Guide for Students and Practitioners*. New York, NY, USA: Academic, 2000.

[109] J. Janai, F. Güney, A. Behl, and A. Geiger, "Computer vision for autonomous vehicles: Problems, datasets and state of the art," Tech. Rep., 2017.

[110] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," Tech. Rep., 2014.

[111] T. Joachims, T. Finley, and C.-N. J. Yu, "Cutting-plane training of structural SVMs," *Mach. Learn.*, vol. 77, no. 1, pp. 27–59, Oct. 2009.

[112] A. E. Johnson and M. Hebert, "Using spin images for efficient object recognition in cluttered 3D scenes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 5, pp. 433–449, May 1999.

[113] E. Jörgensen, C. Zach, and F. Kahl, "Monocular 3D object detection and box fitting trained end-to-end using intersection-over-union loss," Tech. Rep., 2019.

[114] H. Kang and C. Chen, "Fruit detection, segmentation and 3D visuali-sation of environments in apple orchards," *Comput. Electron. Agricult.*, vol. 171, Apr. 2020, Art. no. 105302.

[115] L. Ke, S. Li, Y. Sun, Y.-W. Tai, and C.-K. Tang, "GSNet: Joint vehicle pose and shape reconstruction with geometrical and scene-aware supervision," in *Proc. CVPR*, 2020.

[116] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab, "SSD-6D: Making RGB-based 3D detection and 6D pose estimation great again," Tech. Rep., 2017.

[117] W. Kehl, F. Tombari, N. Navab, S. Ilic, and V. Lepetit, "Hashmod: A hashing method for scalable 3D object detection," in *Proc. Brit. Mach. Vis. Conf.*, 2015.

[118] Y. Kim and D. Kum, "Deep learning based vehicle position and orientation estimation via inverse perspective mapping image," in *Proc. IEEE Intell. Vehicles Symp. (IV)*, Jun. 2019, pp. 317–323.

[119] L. R. Klatzky, *Allocentric and Egocentric Spatial Representations: Definitions, Distinctions, and Interconnections* (Lecture Notes in Computer Science), vol. 1404. Springer, 1998.

[120] R. Klette, *Concise Computer Vision, An Introduction Into Theory and Algorithms*. Springer, 2014.

[121] H. Kobatake and Y. Yoshinaga, "Detection of spicules on mammogram based on skeleton analysis," *IEEE Trans. Med. Imag.*, vol. 15, no. 3, pp. 235–245, Jun. 1996.

[122] L. Koval, J. Vaňuš, and P. Bilík, "Distance measuring by ultrasonic sensor," *IFAC-PapersOnLine*, vol. 49, no. 25, pp. 153–158, 2016.

[123] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Neural Inf. Process. Syst.*, vol. 25, Jan. 2012.

[124] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Commun. ACM*, vol. 60, no. 6, pp. 84–90, May 2017.

[125] A. Krull, E. Brachmann, F. Michel, M. Y. Yang, S. Gumhold, and C. Rother, "Learning analysis-by-synthesis for 6D pose estimation in RGB-D images," Tech. Rep., 2015.

[126] A. Krull, E. Brachmann, S. Nowozin, F. Michel, J. Shotton, and C. Rother, "PoseAgent: Budget-constrained 6D object pose estimation via reinforcement learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017.

[127] A. Krull, F. Michel, E. Brachmann, S. Gumhold, S. Ihrke, and A. C. Rother, "6-DOF model based tracking via object coordinate regression," in *Computer Vision—ACCV 2014*, D. Cremers, I. Reid, H. Saito, M.-H. Yang, Eds. Cham, Switzerland: Springer, 2015, pp. 384–399.

[128] A. Kundu, Y. Li, and J. M. Rehg, "3D-RCNN: Instance-level 3D object reconstruction via render-and-compare," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 3559–3568.

[129] K. Lai, L. Bo, X. Ren, and D. Fox, "A large-scale hierarchical multi-view RGB-D object dataset," in *Proc. IEEE Int. Conf. Robot. Autom.*, May 2011, pp. 1817–1824.

[130] Y. Lamdan and H. J. Wolfson, "Geometric hashing: A general and efficient model-based recognition scheme," in *Proc. 2nd Int. Conf. Comput. Vis.*, 1988, pp. 238–249.

[131] T. Lassa, "The beginning of the end of driving," Tech. Rep., Nov. 2012.

[132] F. Lateef and Y. Ruichek, "Survey on semantic segmentation using deep learning techniques," *Neurocomputing*, vol. 338, pp. 321–348, Apr. 2019.

[133] H. Law and J. Deng, "CornerNet: Detecting objects as paired keypoints," Tech. Rep., 2018.

[134] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Comput.*, vol. 1, no. 4, pp. 541–551, Dec. 1989.

[135] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.

[136] V. Lepetit, F. Moreno-Noguer, and P. Fua, "EPnP: An accurate O(*n*) solution to the PnP problem," *Int. J. Comput. Vis.*, vol. 81, no. 2, pp. 155–166, Feb. 2009.

[137] V. Lepetit, L. Vacchetti, D. Thalmann, and P. Fua, "Fully automated and stable registration for augmented reality applications," in *Proc. 2nd IEEE ACM Int. Symp. Mixed Augmented Reality*, Nov. 2003, pp. 93–102.

[138] B. Li, "3D fully convolutional network for vehicle detection in point cloud," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Sep. 2017, pp. 1513–1518.

[139] B. Li, T. Zhang, and T. Xia, "Vehicle detection from 3D lidar using fully convolutional network," Tech. Rep., 2016.

[140] B. Li, W. Ouyang, L. Sheng, X. Zeng, and X. Wang, "GS3D: An efficient 3D object detection framework for autonomous driving," Tech. Rep., 2019.

[141] C. Li, J. Bai, and D. G. Hager, "A unified framework for multi-view multi-class object pose estimation," Tech. Rep., 2018.

[142] P. Li, H. Zhao, P. Liu, and F. Cao, "RTM3D: Real-time monocular 3D detection from object keypoints for autonomous driving," Tech. Rep., 2020.

[143] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, "DeepIM: Deep iterative matching for 6D pose estimation," *Int. J. Comput. Vis.*, vol. 128, no. 3, pp. 657–678, Nov. 2019.

[144] M. Liang, B. Yang, Y. Chen, R. Hu, and R. Urtasun, "Multi-task multi-sensor fusion for 3D object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019.

[145] R. Lienhart and J. Maydt, "An extended set of Haar-like features for rapid object detection," in *Proc. Int. Conf. Image Process.*, vol. 1, 2002, pp. 1–900.

[146] J. J. Lim, H. Pirsiavash, and A. Torralba, "Parsing IKEA objects: Fine pose estimation," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 2992–2999.

[147] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," Tech. Rep., 2017.

[148] 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*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham, Switzerland: Springer, 2014, pp. 740–755.

[149] F. Liu, P. Fang, Z. Yao, R. Fan, Z. Pan, W. Sheng, and H. Yang, "Recovering 6D object pose from RGB indoor image based on two-stage detection network with multi-task loss," *Neurocomputing*, vol. 337, pp. 15–23, Apr. 2019.

[150] L. Liu, J. Lu, C. Xu, Q. Tian, and J. Zhou, "Deep fitting degree scoring network for monocular 3D object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019.

[151] P. L. Liu, "Monocular 3D object detection in autonomous driving—A review," *Blog*, 2019.

[152] P. L. Liu, "Orientation estimation in monocular 3D object detection," *Blog*, Oct. 2019.

[153] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and C. A. Berg, *SSD: Single Shot MultiBox Detector* (Lecture Notes in Computer Science). 2016, pp. 21–37.

[154] Y.-C. Liu, K.-Y. Lin, and Y.-S. Chen, "Bird's-eye view vision system for vehicle surrounding monitoring," in *Robot Vision*, G. Sommer and R. Klette, Eds. Berlin, Germany: Springer, 2008, pp. 207–218.

[155] Z. Liu, L. Wang, G. Hua, Q. Zhang, Z. Niu, Y. Wu, and N. Zheng, "Joint video object discovery and segmentation by coupled dynamic Markov networks," *IEEE Trans. Image Process.*, vol. 27, no. 12, pp. 5840–5853, Dec. 2018.

[156] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proc. IEEE Int. Conf. Comput. Vis.*, vol. 2, Sep. 1999, pp. 1150–1157.

[157] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.

[158] Q. Luo, H. Ma, Y. Wang, L. Tang, and R. Xiong, "3D-SSD: Learning hierarchical features from RGB-D images for amodal 3D object detection," Tech. Rep., 2017.

[159] X. Ma, Z. Wang, H. Li, P. Zhang, X. Fan, and W. Ouyang, "Accurate monocular object detection via color-embedded 3D reconstruction for autonomous driving," Tech. Rep., 2019.

[160] G. Mamic and M. Bennamoun, "Representation and recognition of 3D free-form objects," *Digit. Signal Process.*, vol. 12, pp. 47–76, Jan. 2002.

[161] F. Manhardt, W. Kehl, and A. Gaidon, "ROI-10D: Monocular lifting of 2D detection to 6D pose and metric shape," Tech. Rep., 2018.

[162] M. Martinez, A. Collet, and S. S. Srinivasa, "MOPED: A scalable and low latency object recognition and pose estimation system," in *Proc. IEEE Int. Conf. Robot. Autom.*, May 2010, pp. 2043–2049.

[163] E. Marchand, H. Uchiyama, and F. Spindler, "Pose estimation for augmented reality: A hands-on survey," *IEEE Trans. Vis. Comput. Graph.*, vol. 22, no. 12, pp. 2633–2651, Dec. 2016.

[164] K. Matzen and N. Snavely, "NYC3DCars: A dataset of 3D vehicles in geographic context," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 761–768.

[165] J. Memon, M. Sami, and R. A. Khan, "Handwritten optical character recognition (OCR): A comprehensive systematic literature review (SLR)," Tech. Rep., 2020.

[166] F. Michel, A. Kirillov, E. Brachmann, A. Krull, S. Gumhold, B. Savchynskyy, and C. Rother, "Global hypothesis generation for 6D object pose estimation," Tech. Rep., 2016.

[167] S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and A. D. Terzopoulos, "Image segmentation using deep learning: A survey," in *Proc. CVPR*, 2020.

[168] C. Mitash, A. Boularias, and K. E. Bekris, "Improving 6D pose estimation of objects in clutter via physics-aware Monte Carlo tree search," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2018, pp. 1–8.

[169] F. Mokhtarian, N. Khalili, and P. Yuen, "Multi-scale free-form 3D object recognition using 3D models," *Image Vis. Comput.*, vol. 19, no. 5, pp. 271–281, 2001.

[170] T. Morris, *Computer Vision and Image Processing*. Red Globe Press, 2003.

[171] T. Morris, *Enlarge Computer Vision and Image Processing*, no. 320. Red Globe Press, 2004.

[172] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, "3D bounding box estimation using deep learning and geometry," Tech. Rep., 2016.

[173] A. Mukhtar, L. Xia, and T. B. Tang, "Vehicle detection techniques for collision avoidance systems: A review," *IEEE Trans. Intell. Transp. Syst.*, vol. 16, no. 5, pp. 2318–2338, May 2015.

[174] R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: An open-source slam system for monocular, stereo, and RGB-D cameras," *IEEE Trans. Robot.*, vol. 33, no. 5, pp. 1255–1262, Oct. 2017.

[175] A. Naiden, V. Paunescu, G. Kim, B. Jeon, and M. Leordeanu, "Shift R-CNN: Deep monocular 3D object detection with closed-form geometric constraints," Tech. Rep., 2019.

[176] V. Nair and E. G. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proc. 27th Int. Conf. Int. Conf. Mach. Learn. (ICML)*. Madison, WI, USA: Omnipress, 2010, pp. 807–814.

[177] M. Naseer, S. Khan, and F. Porikli, "Indoor scene understanding in 2.5/3D for autonomous agents: A survey," *IEEE Access*, vol. 7, pp. 1859–1887, 2019.

[178] C. Nicholson, "A beginner's guide to important topics in AI, machine learning, and deep learning," Tech. Rep., 2019.

[179] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," Tech. Rep., 2015.

[180] D. Nospes, K. Safronov, S. Gillet, K. Brillowski, and U. E. Zimmermann, "Recognition and 6D pose estimation of large-scale objects using 3D semi-global descriptors," in *Proc. 16th Int. Conf. Mach. Vis. Appl. (MVA)*, May 2019, pp. 1–6.

[181] D. Novotny, N. Ravi, B. Graham, N. Neverova, and A. Vedaldi, "C3DPO: Canonical 3D pose networks for non-rigid structure from motion," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019.

[182] Oceanservice.noaa.gov, "What is lidar," *Nat. Ocean. Atmos. Admin.*, Feb. 2021.

[183] *International Theses*, Univ. Tasmania, Hobart, TAS, Australia.

[184] *Open Access Repository*, Univ. Tasmania, Hobart, TAS, Australia.

[185] M. Ozuysal, V. Lepetit, and P. Fua, "Pose estimation for category specific multiview object localization," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 778–785.

[186] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra R-CNN: Towards balanced learning for object detection," in *Proc. Comput. Vis. Pattern Recognit.*, 2019.

[187] C. Papazov and D. Burschka, "An efficient RANSAC for 3D object recognition in noisy and occluded scenes," in *Proc. 10th Asian Conf. Comput. Vis. (ACCV)*, Jan. 2010, pp. 135–148.

[188] K. Park, T. Patten, and M. Vincze, "Pix2Pose: Pixel-wise coordinate regression of objects for 6D pose estimation," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019.

[189] N. Payet and S. Todorovic, "From contours to 3D object detection and pose estimation," in *Proc. Int. Conf. Comput. Vis.*, Nov. 2011, pp. 983–990.

[190] S. Pendleton, H. Andersen, X. Du, X. Shen, M. Meghjani, Y. Eng, D. Rus, and M. Ang, "Perception, planning, control, and coordination for autonomous vehicles," *Machines*, vol. 5, no. 1, p. 6, Feb. 2017.

[191] S. Peng, Y. Liu, Q. Huang, H. Bao, and X. Zhou, "PVNet: Pixel-wise voting network for 6DoF pose estimation," Tech. Rep., 2018.

[192] K. Piper, "It's 2020. Where are our self-driving cars?" Tech. Rep., 2020.

[193] D. Forsyth and J. Ponce, *Computer Vision: A Modern Approach*. 2003.

[194] V. Prisacariu and I. Reid, "PWP3D: Real-time segmentation and tracking of 3D objects," *Int. J. Comput. Vis.*, vol. 98, no. 3, pp. 335–354, 2012.

[195] R. C. Qi, W. Liu, C. Wu, H. Su, and J. L. Guibas, "Frustum PointNets for 3D object detection from RGB-D data," Tech. Rep., 2017.

[196] J. Qian, S. Feng, T. Tao, Y. Hu, Y. Li, Q. Chen, and C. Zuo, "Deep-learning-enabled geometric constraints and phase unwrapping for single-shot absolute 3D shape measurement," *APL Photon.*, vol. 5, no. 4, Apr. 2020, Art. no. 046105.

[197] Z. Qin, J. Wang, and Y. Lu, "MonoGRNet: A geometric reasoning network for monocular 3D object localization," Tech. Rep., 2018.

[198] M. Rad and V. Lepetit, "BB8: A scalable, accurate, robust to partial occlusion method for predicting the 3D poses of challenging objects without using depth," in *Proc. ICCV*, 2017, pp. 3828–3836.

[199] M. M. Rahman, Y. Tan, J. Xue, and K. Lu, "Recent advances in 3D object detection in the era of deep neural networks: A survey," *IEEE Trans. Image Process.*, vol. 29, pp. 2947–2962, 2020.

[200] A. Rangesh and M. M. Trivedi, "Ground plane polling for 6DoF pose estimation of objects on the road," Tech. Rep., 2018.

[201] E. Real, J. Shlens, S. Mazzocchi, X. Pan, and V. Vanhoucke, "YouTube-BoundingBoxes: A large high-precision human-annotated data set for object detection in video," Tech. Rep., 2017.

[202] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," Tech. Rep., 2015.

[203] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015.

[204] C. Rennie, R. Shome, E. K. Bekris, and F. A. D. Souza, "A dataset for improved RGBD-based object detection and pose estimation for warehouse pick-and-place," Tech. Rep., 2015.

[205] R. Rios-Cabrera and T. Tuytelaars, "Discriminatively trained templates for 3D object detection: A real time scalable approach," in *Proc. IEEE Int. Conf. Comput. Vis.*, Dec. 2013, pp. 2048–2055.

[206] J. Rouillard, "Contextual QR codes," in *Proc. 3rd Int. Multi-Conf. Comput. Global Inf. Technol. (ICCGI)*, Jul. 2008, pp. 50–55.

[207] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, *Learning Internal Representations by Error Propagation*. Cambridge, MA, USA: MIT Press, 1986, pp. 318–362.

[208] C. Sahin, G. Garcia-Hernando, J. Sock, and T.-K. Kim, "Instance- and category-level 6D object pose estimation," Tech. Rep., 2019.

[209] C. Sahin, G. Garcia-Hernando, J. Sock, and T.-K. Kim, "A review on object pose recovery: From 3D bounding box detectors to full 6D pose estimators," Tech. Rep., 2020.

[210] C. Sahin and T.-K. Kim, "Recovering 6D object pose: A review and multi-modal analysis," in *Computer Vision—ECCV 2018 Workshops*, L. Leal-Taixé and S. Roth, Eds. Cham, Switzerland: Springer, 2019, pp. 15–31.

[211] H. Sahloul, S. Shirafuji, and J. Ota, "3D affine: An embedding of local image features for viewpoint invariance using RGB-D sensor data," *Sensors*, vol. 19, no. 2, p. 291, Jan. 2019, doi: 10.3390/s19020291.

[212] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," Tech. Rep., 2018.

[213] S. Savarese and L. Fei-Fei, "3D generic object categorization, localization and pose estimation," in *Proc. IEEE 11th Int. Conf. Comput. Vis.*, Oct. 2007, pp. 1–8.

[214] C. Schmid and R. Mohr, "Local grayvalue invariants for image retrieval," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 5, pp. 530–535, May 1997.

[215] H. Schneiderman and T. Kanade, "A statistical method for 3D object detection applied to faces and cars," in *Proc. Comput. Vis. Pattern Recognit. (CVPR)*, vol. 1, Feb. 2000, pp. 746–751.

[216] M. J. Schuster, K. Schmid, C. Brand, and M. Beetz, "Distributed stereo vision-based 6D localization and mapping for multi-robot teams," *J. Field Robot.*, vol. 36, no. 2, pp. 305–332, Mar. 2019.

[217] M. Schwarz, H. Schulz, and S. Behnke, "RGB-D object recognition and pose estimation based on pre-trained convolutional neural network features," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2015.

[218] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and A. LeCun, "OverFeat: Integrated recognition, localization and detection using convolutional networks," Tech. Rep., 2013.

[219] L. G. Shapiro and G. C. Stockman, *Computer Vision*. Upper Saddle River, NJ, USA: Prentice-Hall, 2001.

[220] L. G. Shapiro and G. C. Stockman, *Computer Vision*. London, U.K.: Pearson, 2001.

[221] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and A. Li, "PV-RCNN: Point-voxel feature set abstraction for 3D object detection," in *Proc. Comput. Vis. Pattern Recognit.*, 2019.

[222] S. Shi, L. Jiang, J. Deng, Z. Wang, C. Guo, J. Shi, X. Wang, and H. Li, "PV-RCNN++: Point-voxel feature set abstraction with local vector representation for 3D object detection," Tech. Rep., 2021.

[223] J. Shotton, B. Glocker, C. Zach, S. Izadi, A. Criminisi, and A. Fitzgibbon, "Scene coordinate regression forests for camera relocalization in RGB-D images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2013.

[224] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta, "Beyond skip connections: Top-down modulation for object detection," in *Proc. CVPR*, 2016.

[225] J. Shuttleworth, "Automated driving levels of driving automation are defined in new SAE international standard j3016," SAE Int. J3016, 2018, p. 2.

[226] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgbd images," in *Computer Vision—ECCV 2012*, A. Fitzgibbon, S. Lazebnik, P. Perona, Y. Sato, and C. Schmid, Eds. Berlin, Germany: Springer, 2012, pp. 746–760.

[227] A. Simonelli, S. R. R. Bulò, L. Porzi, M. López-Antequera, and P. Kontschieder, "Disentangling monocular 3D object detection," Tech. Rep., 2019.

[228] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. ICLR*, 2015.

[229] S. Sivaraman and M. M. Trivedi, "Looking at vehicles on the road: A survey of vision-based vehicle detection, tracking, and behavior analysis," *IEEE Trans. Intell. Transp. Syst.*, vol. 14, no. 4, pp. 1773–1795, Dec. 2013.

[230] A. W. M. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based image retrieval at the end of the early years," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 12, pp. 1349–1380, Dec. 2000.

[231] J. Sock, S. H. Kasaei, L. S. Lopes, and T.-K. Kim, "Multi-view 6D object pose estimation and camera motion planning using RGBD images," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops (ICCVW)*, Oct. 2017, pp. 2228–2235.

[232] J. Sock, P. Castro, A. Armagan, G. Garcia-Hernando, and T.-K. Kim, "Tackling two challenges of 6D object pose estimation: Lack of real annotated RGB images and scalability to number of objects," Tech. Rep., Mar. 2020.

[233] A. Soltani, H. Haibin, J. Wu, T. Kulkarni, and J. Tenenbaum, "Synthesizing 3D shapes via modeling multi-view depth maps and silhouettes with deep generative networks," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, Jul. 2017, pp. 1511–1519.

[234] C. Song, J. Song, and Q. Huang, "HybridPose: 6D object pose estimation under hybrid representations," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020.

[235] S. Song and J. Xiao, "Sliding shapes for 3D object detection in depth images," in *Computer Vision—ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham, Switzerland: Springer, 2014, pp. 634–651.

[236] S. Song and J. Xiao, "Deep sliding shapes for amodal 3D object detection in RGB-D images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016.

[237] X. Song, P. Wang, D. Zhou, R. Zhu, C. Guan, Y. Dai, H. Su, H. Li, and R. Yang, "ApolloCar3D: A large 3D car instance understanding benchmark for autonomous driving," Tech. Rep., 2018.

[238] X. Song, P. Wang, D. Zhou, R. Zhu, C. Guan, Y. Dai, H. Su, H. Li, and R. Yang, "ApolloCar3D: A large 3D car instance understanding benchmark for autonomous driving," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 5452–5462.

[239] M. Sonka, V. Hlavac, and R. Boyle, *Image Processing, Analysis, and Machine Vision*, 2nd ed., no. 3216-7. Boston, MA, USA: Springer, 1993.

[240] G. Spampinato, J. Lidholm, C. Ahlberg, F. Ekstrand, M. Ekstrom, and L. Asplund, "An embedded stereo vision module for 6D pose estimation and mapping," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Sep. 2011, pp. 1626–1631.

[241] M. Sundermeyer, Z.-C. Marton, M. Durner, M. Brucker, and A. Triebel, "Implicit 3D orientation learning for 6D object detection from RGD images," in *Proc. CVPR*, 2019.

[242] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proc. CVPR*, 2016.

[243] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," Tech. Rep., 2014.

[244] C. Szegedy, S. Reed, D. Erhan, D. Anguelov, and S. Ioffe, "Scalable, high-quality object detection," 2014, *arXiv:1412.1441*. [Online]. Available: http://arxiv.org/abs/1412.1441

[245] A. Taeihagh and H. S. M. Lim, "Governing autonomous vehicles: Emerging responses for safety, liability, privacy, cybersecurity, and industry risks," *Transp. Rev.*, vol. 39, no. 1, pp. 103–128, Jul. 2019.

[246] J. Taylor, J. Shotton, T. Sharp, and A. Fitzgibbon, "The vitruvian manifold: Inferring dense correspondences for one-shot human pose estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2012.

[247] A. Tejani, D. Tang, R. Kouskouridas, and T.-K. Kim, "Latent-class Hough forests for 3D object detection and pose estimation," in *Computer Vision—ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, T. Tuytelaars, Eds. Cham, Switzerland: Springer, 2014, pp. 462–477.

[248] B. Tekin, S. N. Sinha, and P. Fua, "Real-time seamless single shot 6D object pose prediction," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018.

[249] D. Tomè, F. Monti, L. Baroffio, L. Bondi, M. Tagliasacchi, and S. Tubaro, "Deep convolutional neural networks for pedestrian detection," *Signal Process., Image Commun.*, vol. 47, pp. 482–489, Sep. 2016.

[250] J. Tremblay, T. To, and S. Birchfield, "Falling things: A synthetic dataset for 3D object detection and pose estimation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2018.

[251] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and A. Birchfield, "Deep object pose estimation for semantic robotic grasping of household objects," Tech. Rep., 2018.

[252] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun, "Support vector machine learning for interdependent and structured output spaces," in *Proc. 21st Int. Conf. Mach. Learn. (ICML)*, New York, NY, USA, 2004, p. 104.

[253] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders, "Selective search for object recognition," *Int. J. Comput. Vis.*, vol. 104, no. 2, pp. 154–171, Apr. 2013.

[254] G. Van Houdt, C. Mosquera, and G. Nápoles, "A review on the long short-term memory model," *Artif. Intell. Rev.*, vol. 53, no. 8, pp. 5929–5955, Dec. 2020.

[255] J. Vidal, C.-Y. Lin, X. Lladó, and R. Martí, "A method for 6D pose estimation of free-form rigid objects using point pair features on range data," *Sensors*, vol. 18, no. 8, p. 2678, Aug. 2018.

[256] J. Vidal, C.-Y. Lin, and R. Marti, "6D pose estimation using an improved method based on point pair features," in *Proc. 4th Int. Conf. Control, Autom. Robot. (ICCAR)*, Apr. 2018, pp. 405–409.

[257] D. Wagner, G. Reitmayr, A. Mulloni, T. Drummond, and D. Schmalstieg, "Pose tracking from natural features on mobile phones," in *Proc. 7th IEEE/ACM Int. Symp. Mixed Augmented Reality*, Sep. 2008, pp. 125–134.

[258] C. Wang, D. Xu, Y. Zhu, R. Martín-Martín, C. Lu, L. Fei-Fei, and S. Savarese, "Densefusion: 6D object pose estimation by iterative dense fusion," Tech. Rep., 2019.

[259] C.-Y. Wang, I.-H. Yeh, and H.-Y. M. Liao, "You only learn one representation: Unified network for multiple tasks," Tech. Rep., 2021.

[260] D. Z. Wang and I. Posner, "Voting for voting in online point cloud object detection," in *Robotics: Science and Systems*. 2015.

[261] H. Wang, S. Sridhar, J. Huang, J. Valentin, S. Song, and J. L. Guibas, "Normalized object coordinate space for category-level 6D object pose and size estimation," Tech. Rep., 2019.

[262] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: Fully convolutional one-stage object detection," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019.

[263] X. Wang, W. Yin, T. Kong, Y. Jiang, L. Li, and C. Shen, "Task-aware monocular depth estimation for 3D object detection," Tech. Rep., 2019.

[264] Y. Wang, W.-L. Chao, D. Garg, B. Hariharan, M. Campbell, and Q. K. Weinberger, "Pseudo-LiDAR from visual depth estimation: Bridging the gap in 3D object detection for autonomous driving," Tech. Rep., 2018.

[265] X. Weng and K. Kitani, "Monocular 3D object detection with pseudo-lidar point cloud," Tech. Rep., 2019.

[266] H. D. Whyte and T. Bailey, "Simultaneous localization and mapping," *IEEE Robot. Autom. Mag.*, vol. 13, no. 2, pp. 99–110, Jun. 2006.

[267] K. Wiggers, "Facebook highlights AI that converts 2D objects into 3D shapes," Online Blog, Oct. 2019.

[268] K. Wiggers, "Google brings cross-platform AI pipeline framework," Tech. Rep., Jan. 2020.

[269] D. Wu, Z. Zhuang, C. Xiang, W. Zou, and X. Li, "6D-VNet: End-to-end 6DoF vehicle pose estimation from monocular RGB images," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2019.

[270] Z. Wu, X. Wang, Y.-G. Jiang, H. Ye, and X. Xue, "Modeling spatial-temporal clues in a hybrid deep learning framework for video classification," Tech. Rep., 2015.

[271] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Data-driven 3D voxel patterns for object category recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 1903–1911.

[272] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Subcategory-aware convolutional neural networks for object proposals and detection," Tech. Rep., 2016.

[273] Y. Xiang, W. Kim, W. Chen, J. Ji, C. Choy, H. Su, R. Mottaghi, L. Guibas, and S. Savarese, "ObjectNet3D: A large scale database for 3D object recognition," in *Proc. Eur. Conf. Comput. Vis.*, vol. 9912, Oct. 2016, pp. 160–176.

[274] Y. Xiang, R. Mottaghi, and S. Savarese, "Beyond Pascal: A benchmark for 3D object detection in the wild," in *Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, Mar. 2014, pp. 75–82.

[275] Y. Xiang and S. Savarese, "Estimating the aspect layout of object categories," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2012, pp. 3410–3417.

[276] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "PoseCNN: A convolutional neural network for 6D object pose estimation in cluttered scenes," in *Proc. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018.

[277] X. Mao, D. Inoue, S. Kato, and M. Kagami, "Amplitude-modulated laser radar for range and speed measurement in car applications," *IEEE Trans. Intell. Transp. Syst.*, vol. 13, no. 1, pp. 408–413, Mar. 2012.

[278] B. Xu and Z. Chen, "Multi-level fusion based 3D object detection from monocular images," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 2345–2353.

[279] B. Yang, W. Luo, and R. Urtasun, "PIXOR: Real-time 3D object detection from point clouds," Tech. Rep., 2019.

[280] Y. You, Y. Wang, W.-L. Chao, D. Garg, G. Pleiss, B. Hariharan, M. Campbell, and K. Q. Weinberger, "Pseudo-LiDAR++: Accurate depth for 3D object detection in autonomous driving," Tech. Rep., 2019.

[281] H. Yuan, T. Hoogenkamp, and R. C. Veltkamp, "RobotP: A benchmark dataset for 6D object pose estimation," *Sensors*, vol. 21, no. 4, p. 1299, Feb. 2021.

[282] S. Zakharov, I. Shugurov, and S. Ilic, "DPOD: 6D pose object detector and refiner," Tech. Rep., 2019.

[283] S. Zakharov, I. Shugurov, and S. Ilic, "DPOD: 6d pose object detector and refiner," Tech. Rep., 2019.

[284] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," Tech. Rep., 2013.

[285] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus, "Deconvolutional networks," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 2528–2535.

[286] H. Zhang and Q. Cao, "Fast 6D object pose refinement in depth images," *Int. J. Speech Technol.*, vol. 49, no. 6, pp. 2287–2300, Jun. 2019.

[287] S. Zhang, L. Wen, X. Bian, Z. Lei, and Z. S. Li, "Single-shot refinement neural network for object detection," Tech. Rep., 2017.

[288] T. Zhang, Y. Yang, Y. Zeng, and Y. Zhao, "Cognitive template-clustering improved LINEMOD for efficient multi-object pose estimation," *Cognit. Comput.*, pp. 1–10, Mar. 2020.

[289] X. Zhang, Z. Jiang, and H. Zhang, "Real-time 6D pose estimation from a single RGB image," *Image Vis. Comput.*, vol. 89, pp. 1–11, Sep. 2019.

[290] X. Zhang, Z. Jiang, and H. Zhang, "Out-of-region keypoint localization for 6D pose estimation," *Image Vis. Comput.*, vol. 93, Jan. 2020, Art. no. 103854.

[291] Y. Zhang, H. Zhang, G. Wang, J. Yang, and J.-N. Hwang, "Bundle adjustment for monocular visual odometry based on detections of traffic signs," *IEEE Trans. Veh. Technol.*, vol. 69, no. 1, pp. 151–162, Jan. 2020.

[292] Y. Zhang, D. Huang, and Y. Wang, "PC-RGNN: Point cloud completion and graph neural network for 3D object detection," in *Proc. CVPR*, 2020.

[293] J. Zhao, B. Liang, and Q. Chen, "The key technology toward the self-driving car," *Int. J. Intell. Unmanned Syst.*, vol. 6, no. 1, pp. 2–20, Jan. 2018.

[294] Z.-Q. Zhao, H. Bian, D. Hu, W. Cheng, and H. Glotin, "Pedestrian detection based on fast R-CNN and batch normalization," in *Proc. Intell. Comput. Theories Appl. (ICIC)*, Jul. 2017, pp. 735–746.

[295] Z.-Q. Zhao, P. Zheng, S. T. Xu, and X. Wu, "Object detection with deep learning: A review," Tech. Rep., 2018.

[296] Z. Cao, Y. Sheikh, and N. K. Banerjee, "Real-time scalable 6DOF pose estimation for textureless objects," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2016, pp. 2441–2448.

[297] W. Zheng, W. Tang, S. Chen, L. Jiang, and C.-W. Fu, "CIA-SSD: Confident IoU-aware single-stage object detector from point cloud," in *Proc. CVPR*, 2020.

[298] W. Zheng, W. Tang, L. Jiang, and C.-W. Fu, "SE-SSD: Self-ensembling single-stage object detector from point cloud," in *Proc. CVPR*, 2021.

[299] Z. Yang and R. Nevatia, "A multi-scale cascade fully convolutional network face detector," in *Proc. 23rd Int. Conf. Pattern Recognit. (ICPR)*, Dec. 2016, pp. 633–638.

[300] Y. Zhong, "Intrinsic shape signatures: A shape descriptor for 3D object recognition," in *Proc. IEEE 12th Int. Conf. Comput. Vis. Workshops (ICCV Workshop)*, Sep. 2009, pp. 689–696.

[301] D. Zhou, Y. Dai, and H. Li, "Ground-plane-based absolute scale estimation for monocular visual odometry," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 2, pp. 791–802, Feb. 2019.

[302] D. Zhou, J. Fang, X. Song, L. Liu, J. Yin, Y. Dai, H. Li, and R. Yang, "Joint 3D instance segmentation and object detection for autonomous driving," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 1836–1846.

[303] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised learning of depth and ego-motion from video," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 6612–6619.

[304] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," Tech. Rep., 2019.

[305] X. Zhou, J. Zhuo, and P. Krähenbühl, "Bottom-up object detection by grouping extreme and center points," Tech. Rep., 2019.

[306] W. Zhu, J. Miao, J. Hu, and L. Qing, "Vehicle detection in driving simulation using extreme learning machine," *Neurocomputing*, vol. 128, pp. 160–165, Mar. 2014.

**SHUXIANG XU** received the bachelor's degree in applied mathematics from the University of Electronic Science and Technology of China, China, in 1986, the master's degree in applied mathematics from Sichuan Normal University, China, in 1989, and the Ph.D. degree in computing from Western Sydney University, Australia, in 2000. He is currently working as a Lecturer and a Ph.D. Student Supervisor with the Discipline of Information and Communication Technology, School of Technology, Environments and Design, University of Tasmania, Australia. Much of his work is focused on developing new machine learning algorithms and using them to solve problems in various application fields. His research interests include artificial intelligence, machine learning, and data mining. He received an Overseas Postgraduate Research Award from the Australian Government, in 1996, to research his Ph.D. degree in computing.

**SABERA HOQUE** received the bachelor's degree in computer science and engineering from Northern University Bangladesh, in 2007, and the master's degree in computer science and engineering from United International University, in 2017. She is currently pursuing the Ph.D. degree with the School of Information and Communication Technology, University of Tasmania, Australia. Her research interests include artificial intelligence, machine learning, data mining, image processing, and software development.

**ANANDA MAITI** (Member, IEEE) received the Ph.D. degree from the University of Southern Queensland, in 2016. He is currently an Early Career Researcher. His current research interests include computer networking, and algorithms along with the Internet-of-Things and its various applications in agriculture and remote laboratories. He is also interested in augmented and virtual reality and their application in e-learning.

**MD. YASIR ARAFAT** received the bachelor's degree in computer science and engineering from Northern University Bangladesh, in 2007, and the master's degree in computer science and engineering from United International University, in 2017. Throughout his career, he has been working as a Software Engineer. He is currently working as a Lead Developer at Bundle Australia Pty Ltd. His research interests include artificial intelligence, machine learning, data mining, image processing, and software development.

**YUCHEN WEI** received the B.E. degree in information engineering from China University of Mining and Technology, Xuzhou, China, in 2012, and the M.Sc. degree from Tongji University, Shanghai, China, in 2015. He is currently pursuing the Ph.D. degree with the School of Technology, Environments and Design, University of Tasmania. His research interests include artificial intelligence, machine learning, and image processing.

• • •