

Near-field Perception for Low-Speed Vehicle Automation using Surround-view Fisheye Cameras

Ciarán Eising^{1†}, Jonathan Horgan^{2†}, and Senthil Yogamani^{2†}

¹Department of Electronic and Computer Engineering, University of Limerick, Ireland

²Valeo Vision Systems, Tuam, County Galway, Ireland [†]co-first authors

Abstract—Cameras are the primary sensor in automated driving systems. They provide high information density and are optimal for detecting road infrastructure cues laid out for human vision. Surround-view camera systems typically comprise of four fisheye cameras with 190°+ field of view covering the entire 360° around the vehicle focused on near-field sensing. They are the principal sensors for low-speed, high accuracy, and close-range sensing applications, such as automated parking, traffic jam assistance, and low-speed emergency braking. In this work, we provide a detailed survey of such vision systems, setting up the survey in the context of an architecture that can be decomposed into four modular components namely Recognition, Reconstruction, Relocalization, and Reorganization. We jointly call this the *4R Architecture*. We discuss how each component accomplishes a specific aspect and provide a positional argument that they can be synergized to form a complete perception system for low-speed automation. We support this argument by presenting results from previous works and by presenting architecture proposals for such a system. Qualitative results are presented in the video at <https://youtu.be/ae8bCOF77uY>.

I. INTRODUCTION

Recently, Autonomous Driving (AD) gained huge attention with significant progress in deep learning and computer vision algorithms [1]. Within the next 5-10 years, AD is expected to be deployed commercially [2], with widespread deployment in the coming decades. Currently, most automotive original equipment manufacturers (OEMs) are working on development projects focusing on autonomous driving technology [3], with computer vision having high importance [4]. However, as more is asked from computer vision systems deployed for vehicle autonomy, the architectures of such systems become ever more complex. Thus, it is of advantage to take a step back, and consider the architectures at the highest level. While what we propose should be considered a general discussion on the structure of automotive computer vision systems, we will use specific examples of computer vision applied to Fisheye camera networks, such as surround-view/visual cocoon (see Figure 1). In this paper, we aim to provide two elements to the reader. Firstly, we provide a comprehensive survey of automotive vision, with specific focus on near-field perception for low-speed maneuvering. Secondly, we make a positional argument that considering perception system architectures as specializations of the 4Rs of automotive computer vision leads to significant synergies and efficiencies.

Several related surveys exist in the literature. Most notably, the work of Malik et al. [5], which we shall discuss in more detail soon. A prior survey from the authors [6] discussed in

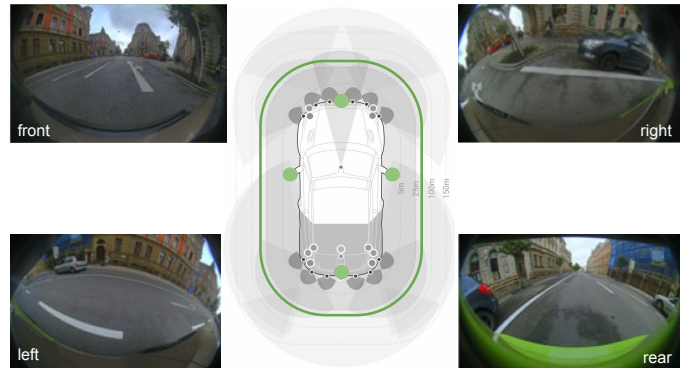


Fig. 1: Images from the surround-view camera network. Green perimeter shows 360° near-field sensing around the vehicle.

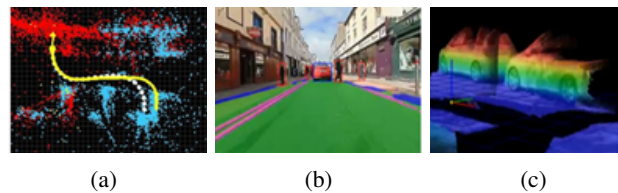


Fig. 2: Examples of (a) *Relocalization*, (b) *Recognition* and (c) *Reconstruction* in action on a surround-view camera system.

detail the role of computer vision in automated parking applications. However, vision components were discussed mostly in isolation, with interdependencies between them not being discussed. In this work, while we briefly mention some embedded considerations, it is only in the context of architecture design. A more detailed overview of embedded considerations for driver assistance systems is provided in [7]. Loce et al. [8] provide a very high-level overview of vision systems for automotive, including traffic management, driver monitoring and security and law enforcement. Their treatment of on-vehicle vision systems is limited to lane keeping, pedestrian detection and driver monitoring. In [9], one of the few surveys that discusses automotive surround-view systems is provided, though the focus is entirely on blind-zone viewing. The output of any perception system used in autonomous driving will be passed to the decision making, trajectory planning, and vehicle control system(s). While in this paper we focus solely on the visual perception elements, the reader is referred to [10] for a useful overview of autonomous driving systems.

The work of this paper is inspired, in part, by the work of Malik et al. in [5]. The authors of that work propose

that the core problems of computer vision are reconstruction, recognition, and reorganization, what they dub as the 3Rs of Computer Vision. Here, we propose to extend and specialize the 3Rs of Computer Vision to the 4Rs of Automotive Computer Vision: *Reconstruction*, *Recognition*, *Reorganization* and *Relocalization*. Figure 2 shows examples of the first three Rs.

As with [5], *Reconstruction* means inferring scene geometry from a video sequence, including the position of the vehicle within the scene. The importance of this should be obvious, as it is central to problems in scene mapping, obstacle avoidance, maneuvering and vehicle control. Malik et al. extend this beyond just geometric inference to include properties such as reflectance and illumination. However, these additional properties are not (currently, at least) significant in the context of automotive computer vision, and so we define *Reconstruction* in the more traditional sense of meaning 3D geometry recovery.

Recognition is the term used for attaching semantic labels to aspects of a video image or scene. As in [5], hierarchies are included in recognition. For example, a cyclist has a spatial hierarchy, as it can be divided into the subsets of bicycle and rider, and a vehicle category can have taxonomic subcategories of car, lorry, bicycle, etc. This can continue as far as is useful for an autonomous driving system. Lights can be categorized by the type (vehicle light, streetlights, stop lights, etc.), color (red, yellow, green), and their importance to the autonomous vehicle (need to respond, can ignore), which infers higher level reasoning of the system.

Relocalization is place recognition and metric localization of a vehicle relative to its surroundings. Relocalization can happen against a pre-recorded trajectory in the host vehicle, for example, for trained parking [6], or against a map that is transferred from the infrastructure, for example, HD Maps [11]. It is highly related to loop closure in SLAM [12], though rather than consider just the problem of loop closure, we consider the broader problem of the localization of the vehicle against one or many pre-defined maps.

Reorganization is the approach of combining information from the previous three components of computer vision into a unified representation. In the work of Malik et al., reorganization is derived from the term “perceptual organization”, and roughly equates it with segmentation. Image segmentation approaches are now dominated by CNN-based semantic segmentation and instance segmentation and therefore fall into the domain of recognition. In this paper, we use the term to equate with “late fusion”, which is the manipulation, filtering and reorganization of inputs into a unified output. This is an important step in the context of vehicle automation, as a unified representation of the sensor outputs is required for vehicle control. This also admits the fusion of the outputs of multiple cameras at a late stage and can be a pre-filter for automotive sensor fusion [13].

The useful information flow paths are shown in Figure 3. As we shall discuss, each of the three of Relocalization, Recognition and Reconstruction have useful information for the other, and all three feed the Reorganization component. It could be argued that even Reorganization has useful information for the other three components. While recurrent systems

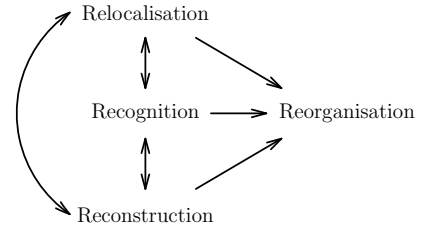


Fig. 3: The 4Rs of automotive vision, and their information flow paths.

have seen some traction in specific automotive vision problems [14], [15], this is not commonplace at an architecture level in automotive visual systems.

Malik et al. [5] describe that in early years of computer vision, the concept of advantageous division of computer vision into low, medium, and high-level tasks was popular, but later became redundant. Later in this paper, we will argue that there is still merit to that mode of thinking, particularly when considering real-time embedded computer vision systems. For example, accessing the pixel information from an image is memory bandwidth intensive, and so doing all intensive pixel processing in a single step is logical. Therefore, we also propose a sub-division of computer vision into *pipeline stages*, being *Pixel Processing*, *Intermediate Processing* and *Object stages*, wrapped by *pre-* and *post-processing* stages. We will generally refer to this as the computer vision *pipeline* or *pipeline stages*.

Additionally, while what we discuss can, in the general sense, be applied to the entire field of automotive computer vision, we must acknowledge that convolutional neural networks (CNNs) have become the standard building block for many visual perception tasks in vehicle autonomy. Thus, many of our examples will focus on the role of neural networks. Bounding boxes for object detection is one of the first successful applications of CNNs for detecting not only pedestrians and vehicles, but also their positions. Recently semantic segmentation is becoming more mature [16], [17], starting with detection of roadway composition like road surface, lanes, road markings, curbs, etc. CNNs are also becoming competitive for geometric vision tasks like depth estimation [18] and Visual SLAM [19].

The rest of the paper is structured as follows. In Section II, we provide some background information on pertinent low-speed, near-field sensing use cases, we give an overview of fisheye cameras, and we provide a brief overview of the *Wood-Scape* dataset [20]. This is the dataset we use to provide some results later in the paper. Section III provides an overview of a surround-view sensing system architecture. Section IV discusses the components of 4R in detail individually by providing a survey of work in each of the 4R areas. In Section V, we discuss the interactions between the 4R components, provide architecture proposals and give a set of results from prior works that support the architectural arguments.

II. OVERVIEW OF NEAR-FIELD SENSING

In this section, we will provide some background material that will give the reader a deeper understanding of what will

come in the next sections. As most research in the autonomous driving perception field focuses on far-field perception using standard or narrow field of view cameras, it is pertinent to give some details here on near-field perception. We give an overview of near-field perception use cases, of how a fisheye camera differs from a standard camera, and of the WoodScape dataset that we use throughout the paper.

A. Near-field sensing use cases

Here we will discuss a few of the most pertinent use cases in vehicle autonomy for surround-view computer vision systems.

1) *Automated Parking Systems*: Automated parking systems are one of the primary use cases for short range sensing [6], with some typical parking use cases described in Figure 4. As early as 1992, prototypes of semi-automated parking systems using radar systems were proposed, though not produced commercially [21]. Early commercial partially automated parking systems employed either ultrasonic sensors or radar [22], [23]. However, more recently, surround-view cameras are becoming one of the primary sensors for automated parking [6], [24]. A major limitation of ultrasonic and radar sensors for automated parking is that parking slots can only be identified based on the presence of other obstacles (Figure 5). Extending this, surround-view camera systems allow for parking in the presence of visual parking slot markings such as painted line markings, while also being seen as a key enabling technology for Valet Parking systems to become a reality [25], [26].

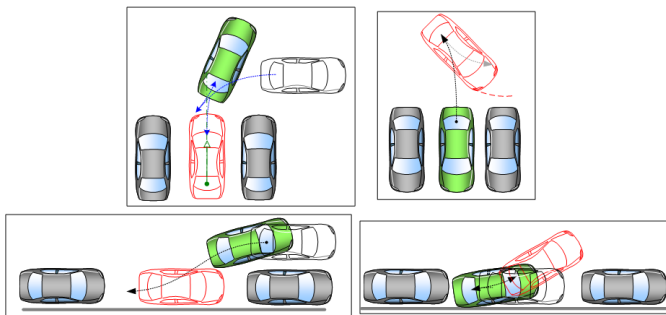


Fig. 4: Typical Parking use cases. Top row: Perpendicular backward park in and perpendicular backward park out. Bottom row: Parallel backward parking in and park Out.



Fig. 5: Parking with cameras can be guided by road markings and not just objects, unlike Ultrasonic and Radar.

2) *Traffic Jam Assistance Systems*: As a substantial proportion of accidents are rear-end collisions at low speed [27], traffic jam situations are considered to be one of the areas of driving that automation can give benefit in the short term

[28], though current systems perhaps lack robustness [29]. In automated traffic jam assistance systems, the vehicle assumes control of the longitudinal and lateral position while in the traffic jam scenario (Figure 6). This functionality is typically used in low speed environments, with maximum speeds of ~ 60 kph [29], though even lower maximum speeds of 40kph are suggested [30]. While typically highway scenarios are considered for traffic jam assistance [28], there has been investigation into the urban traffic jam assistance systems [31]. Given the low-speed nature of this application, surround-view cameras are an ideal sensor, particularly in urban settings where, for example, pedestrians can attempt to cross from areas that are outside the field of view of traditional forward-facing cameras or radar systems. Figure 7 shows examples of using surround-view cameras for traffic jam assist. In addition to detecting other road users and markings, features such as a depth estimation [18] and SLAM [19] are also important for inferring distances to objects and controlling the vehicle position.

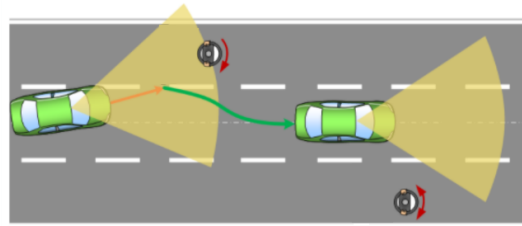


Fig. 6: In traffic jam assist systems, the vehicle can assume both lateral and longitudinal control, including stopping and starting the motion of the vehicle.



Fig. 7: Surround-view cameras can be used for traffic jam assistance systems by detection leading vehicles and lane markings, for example. In addition, vulnerable road users (e.g., crossing pedestrians) that are outside the angular range of other sensors can be detected.

3) *Low Speed Braking*: Protection of vulnerable road users in low speed reversing situations has become a focus of legislation in some jurisdictions [32], with initial efforts to simply display the rearward portion of the vehicle to the driver [9]. It is shown in one study that automatic rearward braking significantly reduced collision claim rates [33], with vehicles equipped with rear camera, parking assistance and automatic braking showing a 78% reduction of reported collisions. Surround-view camera systems are extremely useful for low-speed braking, as the combinations of depth estimation and object detection are building blocks for this functionality.

B. Fisheye cameras

Fisheye cameras offer a distinct advantage for automotive applications. Given their extremely wide field of view, they can observe the full surrounding of a vehicle with a minimal number of sensors. Typically four cameras is all that is required for full 360° coverage of a car (Figure 1). However, this advantage comes with a cost given the significantly more complex projection geometry. Several papers in the past have provided reviews of how to model fisheye geometry, e.g., [34]. We do not aim to repeat this here and will rather focus on the problems that the use of fisheye camera technology brings to automotive visual perception.

In standard field of view cameras, the principles of rectilinear projection and perspective are closely approximated, with the usual perspective properties, i.e., straight lines in the real world are projected as straight lines on the image plane. Parallel sets of straight lines are projected as a set of lines that are convergent on a single vanishing point on the image plane. Deviations from this through optical distortions are easily corrected. Many automotive datasets provide image data with optical distortions removed [35], with an easy means for correction [36], or with almost imperceptible optical distortion [37]. As such, most research in automotive vision makes an implicit assumption of rectilinear projection. Fish-eye perspective differs significantly from rectilinear perspective. A straight line in the camera scene is projected as a curved line on the fish-eye image plane, and parallel sets of lines are projected as a set of curves that converge at two vanishing points [38]. This distortion is immediately obvious in the examples presented in this paper (e.g., Figures 1, 2 and 5). However, distortion is not the only effect. Figure 8 shows an image from a typical mirror mounted camera in a surround-view system. In a fisheye camera, the orientation in the image of objects depends on their location in the image. In this example, the vehicle on the left is rotated almost 90° compared to the vehicle on the right. This has an impact on the translation invariance assumed in convolutional approaches to object detection. In standard cameras, translation invariance is an acceptable assumption. However, this is not the case in fisheye imagery, as evidenced in Figure 8. One must carefully consider how to handle this in any computer vision algorithm design.

A natural approach to address these issues is to rectify the images in some manner. We can immediately discard rectifying to a single planar image, as firstly, too much of the field of view would necessarily be lost thus negating the advantage of fisheye imagery and secondly interpolation and perspective artefacts would quickly dominate the rectified output. A common approach is use multi-planar rectification, whereby different portions of the fisheye image are warped to different planar images. For example, we can define a cube, and warp the image on to the surfaces of the cube. Figure 9 shows warpings on to two such surfaces. Even here, interpolation and perspective effects are visible, and one must deal with the complexity of surface transitions. Another rectification approach is to consider warping to a cylindrical surface, as per Figure 10. In such a warping, the axis of the cylinder is configured such that it is vertical to



Fig. 8: Significant deviation from translation invariance is evident in fisheye cameras.

the ground. The observation is that most objects of interest in an automotive scene lie and move on an approximately horizontal plane, being the road surface. Therefore, we wish to retain the horizontal field of view, while allowing some vertical field of view to be sacrificed. This brings about an interesting combination of geometries. The vertical is per a linear perspective projection, and as such vertical lines in the scene are projected as vertical lines in the image. Objects that are distant or small in the image are visually like a perspective camera. It is even proposed that, through this warping, you can train a network using standard perspective cameras, and use them on fisheye imagery directly without training [39]. However, in the horizontal, distortion exists in the new image. Large, close objects exhibit strong distortion, sometimes even greater than in the original fisheye image. It is also interesting to consider what induces a translation in the resulting cylindrical image. As described in Figure 11, when we are dealing with a perspective camera, translation is induced when the object moves with a constant Z -distance from the camera. That is, on a plane that is parallel to the image plane. However, in a cylindrical image, the distance over the horizontal plane must remain constant to induce an image translation. That is, the object must undergo a rotation about the cylinder axis. In contrast, it is not clear in a raw fisheye image what object motion, if any, would induce an image translation.

The possible warpings are essentially infinite, as we can consider any viable surface, and different warpings may prove better for certain applications and processing techniques than others. This will naturally prove to be an issue as we try to unify automotive vision techniques. With that in mind, we should start thinking instead about how to natively process fisheye imagery without warping. This is non-trivial. For convolutional techniques, for example, this necessitates revisiting how convolutions work. Work has proceeded to address this in the similar field of spherical imagery [40], [41], [42], but such spherical geometry is simpler to consider than native fisheye, and spherical images are, in fact, often warpings of multiple fisheye images [43].



Fig. 9: Warping on to multiple planes. Perspective distortion is evident in the left-hand image, and interpolation artefacts (blurring) is visible in both images.



Fig. 10: Warping to a cylindrical surface. In the automotive scene, vertical objects are vertical in the cylindrical image, and if they are not too close to the camera, appear like a perspective camera (observe pedestrians and cars in the top row). However, closer, and larger objects can be even more distorted in the cylindrical mapping.

C. WoodScape dataset

Even though this paper should be considered as part review and part positional, throughout we will offer insights into results to support our arguments. Mostly, though not ubiquitously, we will use the WoodScape Dataset [20]. It is therefore pertinent to give a description of the dataset, which was captured in two distinct geographical locations: USA and Europe. While most data were obtained from saloon vehicles there is a significant subset from a sports utility vehicle ensuring a strong mix in sensor mechanical configurations. Driving scenarios are divided across highway, urban driving, and parking use cases. Intrinsic and extrinsic calibrations are provided for all sensors as well as timestamp files to allow synchronization of the data. Relevant vehicle’s mechanical data (e.g., wheel circumference, wheelbase) are included. High-quality data are ensured via quality checks at all stages of the data collection process. Annotation data undergo a rigorous quality assurance by highly skilled reviewers. The sensors recorded for this dataset are listed below:

- 4x 1MPx RGB fisheye cameras (190° horizontal FOV)
- 1x LiDAR rotating at 20Hz (Velodyne HDL-64E)
- 1x GNSS/IMU (NovAtel Propak6 & SPAN-IGM-A1)

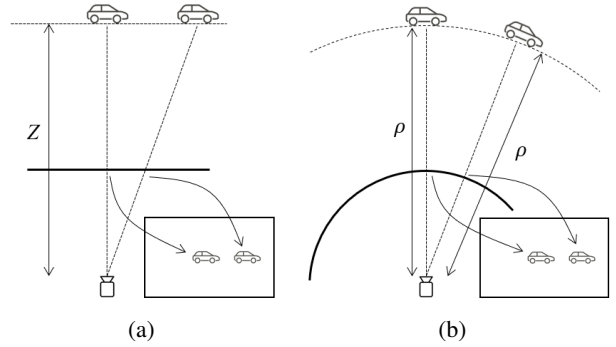


Fig. 11: In a perspective image (a), a translation with a constant Z -depth in the scene leads to a translated image. In contrast, with the cylindrical projection, the horizontal distance ρ must be constant (that is, the object must undergo rotation around the cylinder’s axis) to induce a translation in the image.

TABLE I: Summary of baseline results of 4Rs on our WoodScape dataset.

Task	Model	Metric	Value
Recognition			
Segmentation	ENet [44]	IoU	51.4
2D Bounding Box	Faster R-CNN [45]	mAP (IoU>0.5)	31
Soiling Detection	ResNet10 [46]	Category (%)	84.5
Reconstruction			
Depth Estimation	Eigen [47]	RMSE	7.7
Motion Segmentation	MODNet [48]	IoU	45
Visual Odometry	ResNet50 [46]	Translation (<5mm)	51
		Rotation (<0.1°)	71
Relocalization			
Visual SLAM	LSD SLAM [49]	Relocalization (%)	61

- 1x GNSS Positioning with SPS (Garmin 18x)
- Odometry signals from the vehicle bus.

Table I gives an overview of the baseline results for this dataset for three of the four Rs. Most of the metrics are standard. The relocalization metric is the percentage of instances in which the estimated pose is within a tolerance of 2° in orientation and 0.05m in position.

III. SYSTEM ARCHITECTURE CONSIDERATIONS

A significant consideration in the design of automotive computer vision, in particular the pipelining, is the constraints of embedded systems in which multiple cameras and multiple computer vision algorithms must run in parallel. It is therefore useful to give a brief overview to understand the constraints better, though readers are referred to [6] for a more detailed review of these considerations. Perhaps the most important component to consider is the System-on-Chip (SoC), and typically the first step in designing a commercial automotive camera system is the selection of the SoC for embedded systems, based on criteria including performance (Tera Operations Per Second (TOPS), utilization, bandwidth), cost, power consumption, heat dissipation, high to low end scalability and programmability. The SoC choice provides the computational bounds in the design of algorithms. As computer vision algorithms are compute intensive, Automotive SoCs have a lot of dedicated hardware accelerators for image

signal processing, lens distortion correction, dense optical flow, stereo disparity, etc. In computer vision, deep learning is playing a dominant role in various recognition tasks and gradually for geometric tasks, like depth [18] and motion estimation [50]. The progress in CNN has also led to the hardware manufacturers typically including a custom hardware intellectual property core to provide a high throughput of over 10 TOPS [51].

A. Pipeline Stages

To maximize the performance of processing hardware, it is advantageous to consider embedded vision in terms of processing stages, and to consider shared processing at each processing stage. In this manner, expensive early operations are shared amongst later processing stages that are closer to application layers. In Figure 12, we show an example of a 4R architecture split into pipeline stages.

1) *Pre-processing*: The pre-processing stage of the pipeline can be thought of as the processing that prepares the data for computer vision. This consists of the Image Signal Processing (ISP) steps, such as White Balance, Denoise, Color Correction and Color Space Conversion. For a full discussion on ISP and the tuning of ISP for computer vision tasks in an automotive setting, the reader is referred to [52]. ISP is usually done by hardware engines, e.g., as part of the primary SoC. It is rarely done in software, as there is a vast amount of pixel level processing to be completed. Methods are being proposed to automatically tune the hyperparameters of ISP pipelines to optimize the performance of computer vision algorithm [52], [53]. It should be noted that methods are being proposed to simplify the ISP pipeline for visual perception [54].

2) *Pixel Processing Stage*: Pixel processing can be considered as those parts of a computer vision architecture that *touch the image* directly. In classical computer vision, these would be algorithms such as edge detection, feature detection, descriptors, morphological operations, image registration, stereo disparity and so on. Readers are referred to any one of a number of textbooks for further information on these steps, such as [55]. In neural networks, this would equate with the early layers of a CNN encoder. The processing at this stage is dominated by relatively simple algorithms that must run on potentially millions of pixels many times each second. That is, the computational cost is associated with the fact that these algorithms may run many millions of times each second, rather than the complexity of the algorithms themselves. Processing hardware at this stage is typically dominated by Hardware Accelerators and GPUs, though some elements may be suitable for DSP.

3) *Intermediate Processing Stage*: As the name suggests, the intermediate processing stage bridges the gap from the pixels to the object detection stage. Here, the amount of data to process is still high, but significantly lower than the pixel processing stage. This may include steps such as estimating vehicle motion through visual odometry, stereo triangulation of a disparity map, and the general feature-wise reconstruction of a scene. We would also include CNN decoders at this stage of the pipeline. Processing hardware at this stage would typically be digital signal processors.

4) *Object Processing Stage*: The object processing stage is where higher level reasoning is incorporated. It is here that we may cluster point clouds to create objects, where objects are classified, and where, through said reasoning, we can apply algorithms to suppress relocalization on movable objects. The processing at this stage is dominated by more complex algorithms but operating on fewer data points. In terms of hardware, it is often suitable to run these on general purpose processing units, such as ARM, though digital signal processors would commonly be utilized as well.

5) *Post-processing*: Finally, we have the post processing stage, which could also be termed the *global stage* of processing. It is here that we persist data temporally and spatially. As we can have long temporal persistence and large spatial maps, the overall goal of the preceding stages is to minimize the amount of data reaching this stage while maintaining all the pertinent information that will finally be used for vehicle control. In this stage, we would include steps such as bundle adjustment, map building, high level object tracking and prediction and fusion of the various computer vision inputs. As we are dealing with the highest level of reasoning in the system, and ideally with the lowest amount of data points, general purpose processing units are typically desirable here.

IV. 4R COMPONENTS

It is useful to begin the detailed description of the 4Rs with an example, and so, in Figure 12, we show at a high level what a computer vision architecture for autonomous driving might look like. Four video streams from an automotive surround-view system (Figure 1) are passed to an SoC. Each algorithmic block is mapped to one of a set of available processing units that are typically available on high-end automotive vision SoCs, being hardware accelerator or graphical processing unit, digital signal processor, and general-purpose processing core (e.g., ARM). Each of the 4R pipelines shows one of the standard algorithms, being visual SLAM for relocalization, CNN for object detection and motion stereo for reconstruction. Reorganization shows a typical map, tracking and fusion pipeline.

What is of interest is the possible links between the 4Rs, even in such a standard system. For example, Visual Odometry is a natural part of any SLAM pipeline but can be reused in a motion stereo [56] context for dense triangulation, and equally, the stereo triangulation can be used for the scene reconstruction component of visual SLAM, which is really a seed for bundle adjustment. Robust SLAM can only be achieved in scenes that are dominated by dynamic objects if motion segmentation is readily available [57]. CNNs provide options on moving object detection [58], potentially as part of a multi-task network [59], and thus can be used to suppress issues in the relocalization pipeline associated with dynamic objects. Motion segmentation can also be used to suppress issues with incorrect triangulation in reconstruction, and as such can be included in the clustering stage of the reconstruction pipeline to suppress false detections. Finally, naturally, the output of the object detector can be directly input into the sensor map.

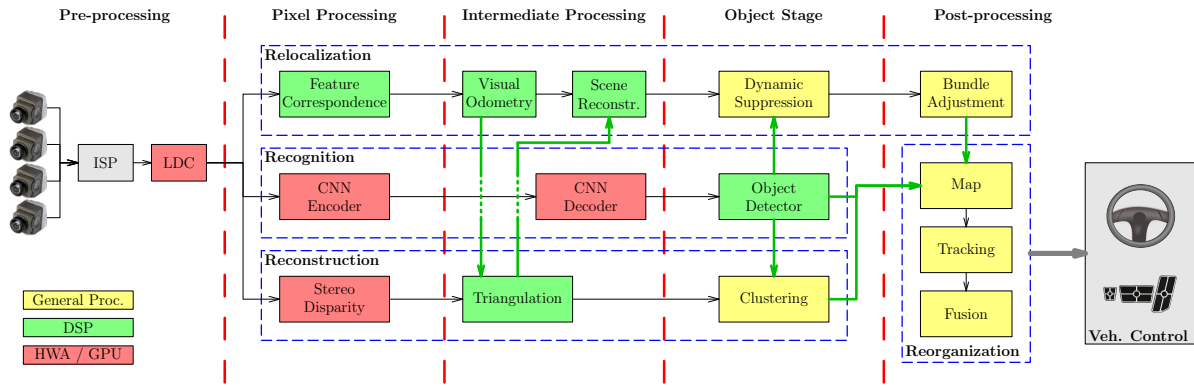


Fig. 12: Example of a high level architecture proposal considering a 4R framework. The 4Rs are shown with each of the processing stages. Useful communication between each of the 4Rs is shown as green arrows. Each block has an example target processing unit, being Hardware Accelerator (HWA) or GPU, Digital Signal Processor (DSP), or general-purpose processing unit, such as ARM. LDC = Lens Distortion Correction, ISP = Image Signal Processor.

Now that we have briefly discussed a simple example of how the 4Rs architecture might work, we will describe in more detail what the 4Rs are.

A. Recognition

The *Recognition* task identifies the semantics of the scene via pattern recognition. In automotive, the first successful application was pedestrian detection which was performed using a combination of hand-designed features like Histogram of Oriented Gradients and a machine learning classifier like Support Vector Machines (for example [60] and many others). Recently, this task has been dominated by CNNs, which have demonstrated remarkable performance leaps for various computer vision tasks in object recognition applications [61]. However, this comes at a cost. Firstly, automotive scenes are very diverse, and the system is expected to work across countries as well as varying weather and lighting conditions, and thus one of the main challenges is to build an effective dataset which covers diverse aspects [62]. Secondly, CNNs are computationally intensive, typically requiring dedicated hardware accelerators or GPUs (in contrast to classical machine learning approaches that are feasible on general purpose computation cores). As such, efficient design techniques are critical to be incorporated in any design [63], [64]. Finally, while CNNs are well studied for rectilinear images, as mentioned previously, the assumption of translation invariance is broken in fisheye images, which pose additional challenges as discussed in [65]. In particular, standard bounding box object detection representation breaks for fisheye images [66], though special consideration should also be given to the design of semantic segmentation approaches for fisheye [67].

In our example recognition pipeline, a multi-task deep learning network for identifying objects based on their appearance patterns is proposed. It comprises of three tasks, namely bounding box objection detection (pedestrians, vehicles, and cyclists), semantic segmentation (road, curbs, and road markings) and lens soiling detection (opaque, semi-transparent, transparent, clear). Object detection and semantic segmentation are standard tasks and for more implementation

details the reader is referred to our FisheyeMultiNet paper [68]. One of the challenges is to balance the three tasks' weights during training phase as one task may converge faster than the others [69]. Additional auxiliary tasks like end-to-end driving which do not have associated annotation costs can aid the training of expensive annotation tasks such as segmentation [70].

Fisheye cameras are mounted relatively low on a vehicle (~ 0.5 to 1.2 m above ground) and are susceptible to lens soiling due to road spray from other vehicles or water from the road. Thus, it is vital to detect soiling on the camera lens to alert the driver to clean the camera or to trigger a cleaning system. The soiling detection task and its usage for cleaning and algorithm degradation is discussed in detail in SoilingNet [71]. A closely related task is desoiling where the soiled areas are restored through inpainting [72], but these desoiling techniques remain in the domain of visualization improvements rather than usage for perception for now. It is an ill-defined problem as it is not possible to predict behind the occlusion, though this can be improved by leveraging temporal information. As the CNN processing capacity is limited on the low power automotive ECU, we make use of multi-task architecture where majority of the computation is shared in the encoder as illustrated in Figure 13.

B. Reconstruction

As mentioned already, *Reconstruction* means inferring scene geometry from a video sequence. This typically means, for example, estimating a point cloud or voxelized representation of a scene. However, we can also consider a temporal aspect of this – if the object is moving, we wish to know its vector of motion.

The first aspect, the reconstruction of static objects, is traditionally done using approaches such as motion stereo [56] or triangulation in multi-view geometry [73]. In the context of designing a depth estimation algorithm, a brief overview of how humans infer depth is given in [74] with useful further referencing. There are four basic approaches to inferring depth:

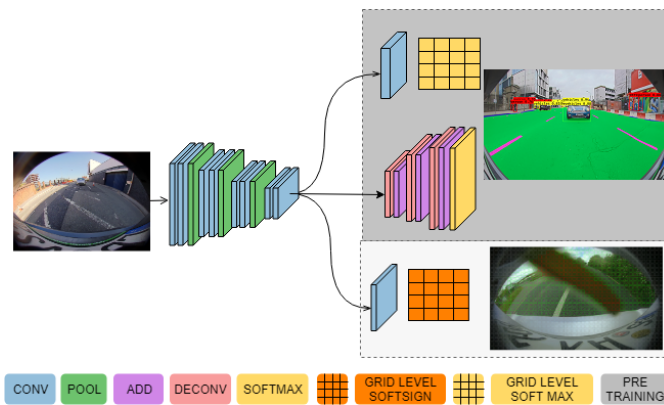


Fig. 13: Illustration of multi-task Recognition architecture comprising of object detection, semantic segmentation and soiling detection tasks.

monocular visual cues, motion-parallax, stereopsis, and depth from focus. Each has its equivalent in computer vision.

Based on earlier theoretical work by Marr & Poggio [75], Grimson provided a computational implementation of stereo vision in the early 1980s [76]. Since then, work has continued on stereo vision (see [77] for an overview of early works). However, stereo systems do not achieve ubiquitous deployment on vehicles, and as such, monocular motion-parallax methods remain popular in automotive research. Computationally, depth from motion parallax is traditionally done through feature triangulation [78], but motion stereo has also proven popular [79]. Neural network approaches to extracting pixel-wise depth from moving cameras have been successful recently [80].

Perhaps one of the more interesting approaches is the use of monocular cues for depth. These are things like texture scale change, occlusions (if A occludes B, then B must be behind A), shading and lighting, object scale (small vs far away), etc. As noted in [74], such monocular cues require contextual interpretation. To work well, they require knowledge of the entire image, and patch-based approaches generally fail. Van Dijk and Croon [81] have shown that, for four publicly available mono-depth networks at least, neural networks learn a correlation between vertical position of the object in the image and the depth to the object. Single-view (one camera, one frame) approaches are extremely desirable, as they remove the need for the camera to move for depth extraction. However, they have proven not to be robust enough [81], though they may be useful still in certain scenarios (when the car is still may be better than no estimate).

Considering fisheye imagery adds significant complexity to the reconstruction task. Most work in multi-view geometry, stereo vision and depth estimation in general assumes a planar perspective image of the scene, i.e., that projective geometry provides a good model of the image. Traditional stereo approaches add a further restriction that the epipolar lines in the image must be horizontal. However, this is rarely the case with real cameras where lens distortions exist thus breaking the planar projection model. It is generally addressed through calibration and rectification of the image. For fisheye

imagery where the lens distortion is extreme, though, it is not feasible to maintain the wide field of view in rectification.

Several approaches have been proposed to address fisheye stereo depth estimation. A common approach is multi-planar rectification, in which the fisheye image is mapped to several perspective planes [82]. However, any planar rectification, even with multiple planes, suffers from significant resampling distortion, as discussed earlier. To minimize this resampling distortion, rectification to non-planar images has been proposed. Some approaches warp to different image geometries that maintain the stereo requirement of epipolar lines being straight and horizontal [83]. Still other approaches bypass the requirement for epipolar lines to be horizontal. For example, the plane sweep method [84], [85] has more recently been applied to fisheye [86]. A related issue with any resampling of the fisheye image is that the noise function is distorted by the resampling process, which is a problem for any method that attempts to minimize a reprojection error (e.g., the widely used optimal triangulation method [87]). Kukulova et al. [73] address this using an iterative technique for standard field of view cameras that minimizes reprojection error while avoiding undistortion. However, this approach depends on a specific camera model, and as such is not directly applicable to fisheye cameras.

The second aspect of reconstruction is the extraction of moving objects from the video sequence (motion segmentation). 3D-reconstruction of dynamic objects results in position inaccuracy in the global sense, as triangulation assumptions are broken. Typical attempts to reconstruct the geometry of an object under motion requires image motion segmentation, relative fundamental matrix estimation and reconstruction (with scale/projective ambiguity). Of course, significant advances have been made. For example, using Multi-X [88] the first two steps can essentially be combined, as the segmentation can be done based on the fundamental matrix estimation. However, such approaches tend to be either computationally too expensive or not robust enough for embedded automotive applications. Additionally, scale must be resolved for such reconstruction, and deformable objects (such as pedestrians) can have different fundamental matrices for different parts of the body. Therefore, in automotive, the task of dynamic object detection is usually simply motion segmentation.

Clappstein et al. [89] describe a geometric approach to motion segmentation in the automotive context. This work is significantly extended to the surround-view camera case by Mariotti and Hughes [90]. However, in both cases the geometry cannot perfectly distinguish all types of moving feature. That is, there is a class of object motion that makes associated features indistinguishable from static features. Thus, a global or semi-global approach must be taken. In traditional approaches, this is done by grouping optical flow vectors with similar properties to ones that are classed as under motion. CNNs offer globality in a more native way [58], and even offer the potential for instance motion segmentation, though this has yet to be extended to the fisheye case [91]. However, as with the static object reconstruction, the results from [58] seem to indicate that it is performing recognition rather than geometric motion estimation, as still pedestrians are often

classified as under motion. It is therefore likely that a much-improved overall motion segmentation will be obtained by incorporating the geometric constraints of [90] with the more global CNN approach of [58], though this is certainly non-trivial and remains work in progress.

Generally, a key input into motion segmentation is knowledge of the motion of the camera. That is, the essential matrix of the camera (or fundamental matrix in the uncalibrated case) must be known. This is assumed in [89] and [90]. This can be achieved in a couple of ways. Firstly, we can directly use signals on the vehicle network, such as steering angle and wheel velocities, to estimate the motion of the vehicle (e.g., as discussed in [92]), and thus the motion of the cameras. Alternatively, visual approaches to estimate the motion directly from image sequences can be employed [93]. An alternative to the explicit estimation of the motion of the camera is to model the background motion in the image. It has been proposed to use an affine model of background motion [94]. However, this assumes that the background is distant or approximately planar, and that radial distortion is absent or negligible. The latter is clearly not the case with automotive fisheye, but perhaps not as obviously, nor is the former. It is obvious that the background for a typical automotive scene (being the road surface, buildings, etc.) cannot be modelled as a single plane, nor is it very distant.

Figure 14 shows an example of different reconstruction stages, including dense motion stereo, 3D point cloud and a clustering of static obstacle, alongside a dense optical flow based motion segmentation. While the use of fisheye images certainly has an impact on design decisions and may be considered a problem that is not yet fully solved from a theoretical standpoint, it is clear that in practice many of the techniques for fisheye discussed above give acceptable results, depending on specific applications.

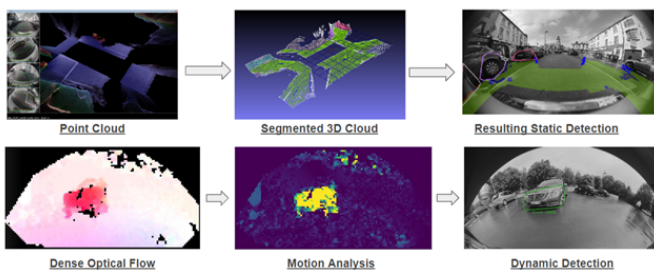


Fig. 14: An example of a reconstruction pipeline with some sample outputs at different processing stages. Top row shows the static pipeline and bottom row shows the dynamic pipeline.

C. Relocalization

Visual Simultaneous Localization and Mapping (VSLAM) is a well-studied problem in robotics and autonomous driving. There are primarily three types of approaches namely (1) Feature based methods, (2) Direct SLAM methods and (3) CNN approaches. Feature based methods make use of descriptive image features for tracking and depth estimation [95] which results in sparse maps. MonoSLAM [96], Parallel

Tracking and Mapping (PTAM) [97] and ORBSLAM [98] are seminal algorithms of this type. Direct SLAM methods work on the entire image instead of sparse features to aid building a dense map. Dense Tracking and Mapping (DTAM) [99] and Large-Scale Semi Dense SLAM (LSD-SLAM) [100] are the popular direct methods which are based on minimization of photometric error. CNN based approaches are relatively less mature for Visual SLAM problems and they are discussed in detail in [101].

Mapping is one of the key pillars of autonomous driving. Many first successful demonstrations of autonomous driving (e.g., by Google) were primarily reliant on localization to pre-mapped areas. HD maps such as TomTom RoadDNA [102] provide a highly dense semantic 3D point cloud map and localization service for majority of European cities with a typical localization accuracy of 10 cm. When there is an accurate localization, HD maps can be treated as a dominant cue, as a strong prior semantic segmentation is already available, and it can be refined by an online segmentation algorithm [103]. However, this service is expensive as it requires regular maintenance and upgrades of various regions in the world. Due to privacy laws and accessibility, such a commercial service cannot be used in many situations and a mapping mechanism must be built within a vehicle’s embedded system. For example, a private residential area cannot be mapped legally in many countries, such as Germany [104].

Visual SLAM (VSLAM), in the automotive context, consists of building a map of the environment surrounding the vehicle while simultaneously estimating the current pose of the car within that map [105]. One of the key tasks of VSLAM is the localization of the vehicle against a previously recorded trajectory [106]. A trained trajectory is typically represented by a group of key poses surrounded by landmarks spanned from the vehicle’s origin to destination positions. These landmarks are represented using robust image features that are unique in the captured images.

A classical feature-based relocalization pipeline is shown in Figure 15. In feature-based SLAM, the first step is the extraction of salient features. A salient feature in an image could be a region of pixels where the intensity changes in a particular way, such as an edge, a corner or a blob [107], [108], [109]. To estimate landmarks in the world, tracking is performed, wherein two or more views of the same features can be matched. Once the vehicle has moved enough, VSLAM takes another image and extracts features. The corresponding features are reconstructed to get their coordinates and poses in real world. These detected, described, and localized landmarks are then stored in persistent memory to describe the relative position of the vehicle for a trajectory. If the vehicle returns to same general location the live feature detections are matched against the stored landmarks to recover the vehicle’s pose relative to the stored trajectory.

D. Reorganization

Reorganization performs three functions - 1) Fusion of Recognition and Reconstruction, 2) Mapping of objects in a centralized world co-ordinate system across cameras and 3)

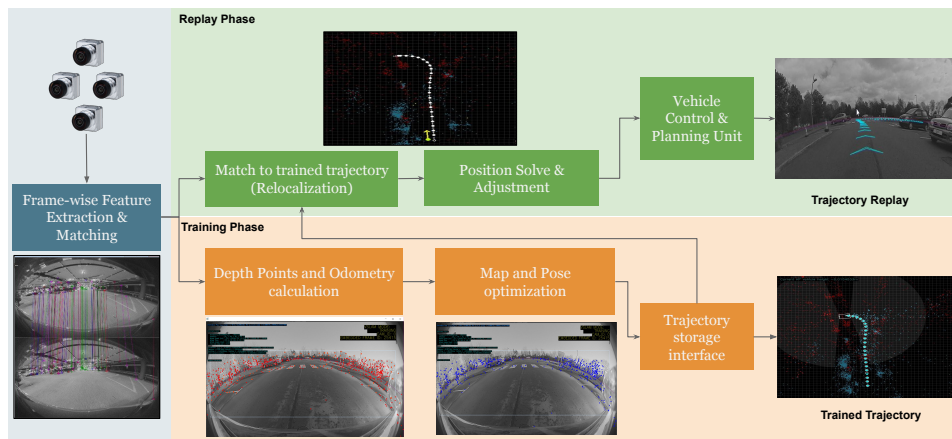


Fig. 15: Relocalization pipeline with intermediate outputs. Reproduced from our paper [110].

Temporal tracking of objects. Although it would be possible for the recognition and reorganization blocks to feed directly into the environment map, we contend there are distinct advantages to implementing some fusion at the vision layer. Let us consider this first with an example. As shown in Figure 16, let us assume we have a system that has a monocular depth estimation, motion segmentation and bounding box vehicle detection. A classical approach to fusing this information is to convert all the data into a world coordinate system, and then associate and fuse the data. This type of approach has advantages. Some automotive sensors, such as laser scanner, provide native Euclidean data, and a fusion system based on such a Euclidean map makes inclusion of these additional sensors easy. However, camera-based detection accuracy will always suffer with the conversion to a Euclidean map. Projections from the image domain to the world domain are known to be error prone, as they are subject to errors from poor calibration, flat ground assumptions, variations in footpoint detection, pixel density and imperfect camera models. Even consider the case in which we have a perfect semantic segmentation, as in Figure 17. If the object does not actually touch the ground at the point of interest, then there will be significant error with the flat-ground assumption for projection to a world coordinate system.

However, detections in the image domain, prior to projections to the world, are not subject to such error and therefore association of detections from different vision algorithms in the image domain is more robust. In fact, simple detection overlap measures typically prove robust. Figure 18 shows an implementation of the image-based fusion of a CNN-based vehicle detection and an optical flow-based motion segmentation. Even though significant error exists in the motion segmentation, the fusion successfully classifies the detected objects as both vehicle and dynamic.

One must also consider how the distortion correction impacts measurement noise. Many common algorithms for fusion and tracking, such as the Kalman filter or the particle filter, begin with the assumption of mean-zero, Gaussian noise. For interest point measurement in computer vision (e.g., an image feature or a bounding box footpoint estimate), this is commonly considered to be a valid assumption. However, the fisheye undistortion and ground plane projection process

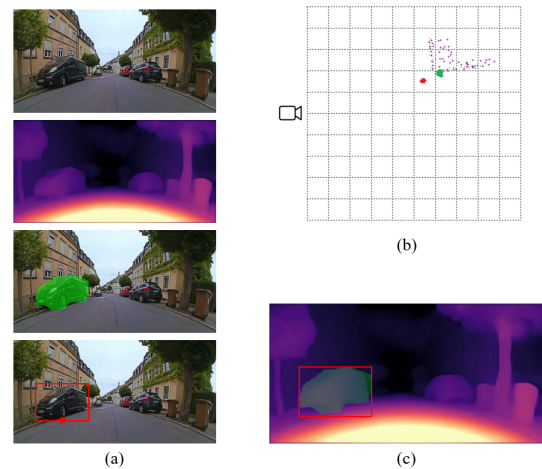


Fig. 16: Two fusion paradigms. (a) shows the inputs, being a monocular depth estimation, motion segmentation and bounding box. (b) shows the world coordinate conversion of these detections, generating a point cloud for the monocular depth and the projection of a reference point for each of the motion segmentation and bounding box detection. It is not clear how we could associate all the detections in a fusion system. In contrast, (c) shows all the detections in a single image plane, where it is intuitive that such fusion is almost trivial (a simple metric such as overlap would suffice).

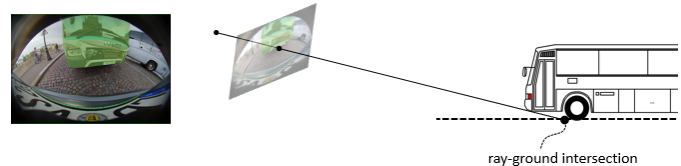


Fig. 17: Even if we had a hypothetical perfect semantic segmentation, if the object does not intersect the ground at the point of interest (e.g., the front of the bus), then the ray projection-based estimation of the world coordinate system object position will necessarily be in error.

distorts this noise model (Figure 19). Addressing this is additionally complicated by the fact that the distortion of the measurement noise is dependent on the location of the interest point in the image and the position of the camera relative to



Fig. 18: Reorganization combines Reconstruction’s dynamic object detection output (red polygon) and Recognition’s object detection (yellow box).

the road surface.

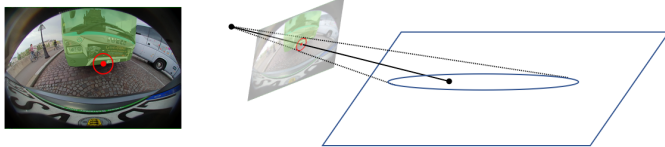


Fig. 19: Any projection from a fisheye image to the world coordinate system will necessarily induce a non-linear distortion of the measurement noise model, due to the non-linearities of the fisheye unprojection and flat ground intersections. Measurement noise is no longer zero-mean, Gaussian, nor even symmetric.

So far in this section, we have predominantly focused on fusion. However, tracking the objects in the scene will suffer from similar artefacts as fusion if we attempt to do this in the world coordinate system. Following detection of an object, one must associate that detection with a previous detection of the same object. If we rely on something like the prediction of the footpoint of an object in world coordinates (per Figure 19), noise will generally mean that the association is very difficult. However, in image coordinates, geometric measures of overlap (e.g., intersection-over-union) are more robust against noise, and the association problem, as with the fusion problem, is significantly simpler. Naturally, we don’t intend to give a complete survey of multi-object tracking. The reader is instead referred to [111], which discusses classical as well as more modern neural network approaches.

The two common mapping paradigms in automotive perception are vector maps [112] and occupancy grids [113], and sensor fusion generally occurs in such a map (that is, fusion between different sensor types, such as camera, laser scanner, radar, etc.). A survey of sensor fusion is beyond the scope of this paper; we would refer the reader to [114] for a very comprehensive review of automotive sensor fusion modalities. Ultimately, we will want the vehicle to make control decisions based on image perception. Given that the vehicle exists in a Euclidean space, and planning and control decisions are generally made on Euclidean data represented by the perception map, map information must be generated from image perception at some point. By fusing detections in the image space, we can, for example, smooth the detections of object footpoints. This will mean a better localization in world

coordinates. More importantly, however, we can more easily fuse object detection with depth information from an image. Using Figure 16 as an example, if we have an accurate depth estimation in our image space, we can attach the semantic labels to the depth prior to generating the map information.

E. Discussion

Overall, we argue that the 4R approach provides a localized semantic-geometric representation of the vehicles environment. By *localized semantic-geometric* representation of the vehicles environment, we mean *localized*: that the 4R processing pipeline provides information about where the vehicle is (can be globally or against pre-learned trajectories), *geometric*: information about the spatial relationship between the vehicle and obstacles in its local environment, and *semantic*: the obstacles will be recognized as belonging to a class of obstacle.

V. SYSTEM SYNERGIES

In this section, we will discuss system synergies. We will look at how Relocalization, Reconstruction and Recognition tasks can support one another, and we will describe the importance of dual sources of detection in providing redundancy in safety critical applications.

A. Recognition and Reconstruction

As already mentioned, depth estimation is important in geometric perception application. In addition to previous material already discussed, the current state-of-the-art are neural network-based methods [115], [116], learnable in a self-supervised manner through reprojection loss [117]. It has been shown that state-of-the-art single frame attempts at monocular depth estimation typically results in recognition tasks [118], and then using cues such as vertical position in the image to infer depth [81]. Moving object detection appears to have a heavy reliance on recognition as well. This is evidenced by the fact that both [48] and [58] show false positives on static objects that are commonly moving (pedestrians, for example - see Figure 20). This does not, in any way, reduce the importance of such attempts. Rather, it points to a very deep connection between recognition and reconstruction, and that from one, you can infer the other.

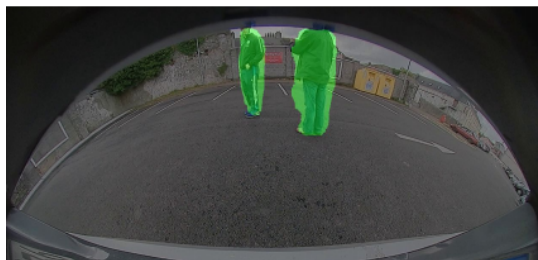


Fig. 20: Learning geometry can have a heavy reliance on recognition - static pedestrians detected as moving [58].

When bounding box pedestrian detection was state of the art, before semantic and instance segmentation, most researchers in automotive pedestrian detection will have considered encoding a depth based on the height of the bounding box, or the vertical position of the pedestrian in the image. This is discussed in detail in [81]. However, it is somewhat intuitive that recognition based on deep neural networks can lead to object depth, especially as the accuracy of neural networks improves. Recent work demonstrates the validity of joint learning of semantic labels and depth [119]. For example, in [120], it is shown that, for monocular depth estimation, adding semantic guidance in each of the distance decoder layers (per Figure 21) improves performance at edges of objects, and even returns reasonable distance estimates for dynamic objects. Table II shows an extract of the results from our work in [120], in comparison to other mono-depth approaches.

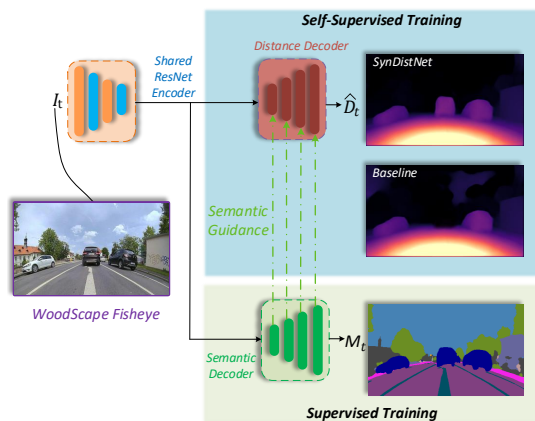


Fig. 21: Overview over the joint prediction of distance and semantic segmentation from a single input image [120].

Thus, we demonstrate the strong link between recognition and reconstruction. This idea is not particularly new. There was research investigating the potential of joint semantic labelling and depth as early as 2010 [125] (building upon even earlier work in geometric/semantic consistency in images [126]). However, it is fair to say that with the advent of neural networks in the last few years, the true potential of this research is beginning to come to fruition.

B. Relocalization and Recognition

Relocalization is the process of a vehicle recognizing a previously learned position or path, as discussed. However, in the real automotive world, many things can disturb this.

For example, the scene can change due to movable objects - for example, parked vehicles can move between the time the scene is learned and when relocalization is requested. In such a case, semantic segmentation approaches, (e.g. [127], [16]), can be used to identify objects that may potentially move (vehicles, bicycles, pedestrians), and remove mapped features associated with such objects. Further opportunities exist for the support of traditional Visual-SLAM pipelines with deep learning techniques (Figure 22), as described in detail in [101].

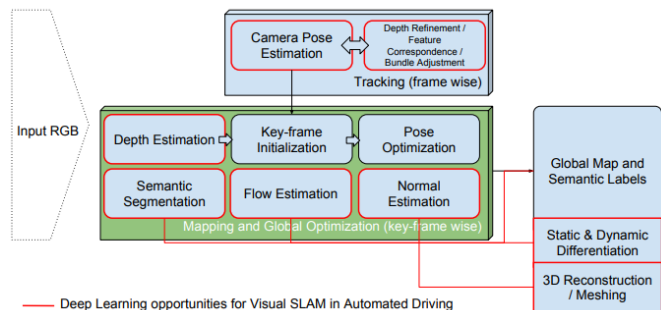


Fig. 22: The Fundamental pipeline of Visual SLAM is composed of multiple geometric vision tasks including depth estimation, optical flow and pose estimation. Those tasks have well known solutions based on CNNs in their individual domain. In contrast, the overall Visual-SLAM is not dominated by Deep Learning. [101]

Place recognition in Visual-SLAM has a couple of applications. Firstly, it allows loop closure to correct for accumulated drift, and secondly it allows for building and maintaining maps from multiple passes through the same scene. Classical approaches using Bag of Words (e.g. [128]) proved reasonably successful, if perhaps lacking in terms of robustness. CNN-based approaches are proving to be more robust, with appearance-invariant approaches showing promising, if initial, results [129]. The recognition of places when significant time has passed is an important topic. Table III shows a small set of results for a Visual-SLAM pipeline, and demonstrates that errors increase significantly with a six-month time difference between training and relocalization.

Finally, view invariant localization can be considered. This is important when the camera viewpoint at the relocalization time is significantly different to the camera viewpoint at training, for example due to a rotation of the vehicle caused by approaching the trained trajectory at a large angle. Traditional Visual-SLAM methods based on feature descriptors fail, as the same surfaces of the landmarks may not even be visible. It has been shown that attaching semantic labels to scene landmarks (via bounding box classification) can significantly improve the performance of viewpoint invariance [130].

C. Relocalization and Reconstruction

This is perhaps the most straightforward synergy to discuss. Relocalization, and Visual-SLAM in general, can be considered as the storage of scene reconstruction (i.e., building a map) along with iterative refinement of said map through bundle adjustment (refer to Figure 15). In this way, reconstruction and visual odometry become a seed for the traditional Visual-SLAM approaches. There are direct methods that bypass

TABLE II: Extract of our results presented in [120]. Comparison of the SynDistNet network, which combines semantic segmentation and depth estimation in a single network, with other monocular methods. KITTI dataset [121] is used in all cases.

Method	Resolution	Abs Rel	Sq Rel	RMSE	RMSE _{log}	$\delta < 1.25$ $\delta < 1.25^2$ $\delta < 1.25^3$		
						higher is better		
EPC++ [122]	640 x 192	0.141	1.029	5.350	0.216	0.816	0.941	0.976
Monodepth2 [123]	640 x 192	0.115	0.903	4.863	0.193	0.877	0.959	0.981
PackNet-SfM [124]	640 x 192	0.111	0.829	4.788	0.199	0.864	0.954	0.980
FisheyeDistanceNet [95]	640 x 192	0.117	0.867	4.739	0.190	0.869	0.960	0.982
SynDistNet [120]	640 x 192	0.109	0.843	4.594	0.186	0.878	0.968	0.986
Monodepth2 [123]	1024 x 320	0.115	0.882	4.701	0.190	0.879	0.961	0.982
FisheyeDistanceNet [95]	1024 x 320	0.109	0.788	4.669	0.185	0.889	0.964	0.982
SynDistNet [120]	1024 x 320	0.103	0.705	4.386	0.164	0.897	0.980	0.989

TABLE III: Quantitative results of our relocalization algorithm on selected WoodScape [20] dataset scenes. The time difference is the number of days between training and relocalization. The distance is the starting distance from the trained trajectory. The average offset is given in terms of position and angle.

Scene		Difference		Average Offset	
Training Date	Replay Date	Time (days)	Dist. (m)	Pos. (m)	Angle (deg)
20161208	20161208	0.003	4.723	0.468	4.704
20161208	20161208	0.005	2.483	0.355	5.366
20161208	20161208	0.006	2.692	0.3	5.149
20161208	20170607	181.156	2.49	1.085	8.162
20161208	20170607	181.155	0.066	0.903	9.498
20161208	20170607	181.154	4.96	0.896	10.751

this seeded approach, for example LSD-SLAM [49] (and its Omnidirectional camera extension [100]), where photometric error is minimized as opposed to reprojection error. However, if one considers a time-slicing of a bundle-adjusted map, it can also be seen that Visual-SLAM can be used to refine the reconstruction (both scene structure and visual odometry). In addition, it is well known moving objects (as distinct from *moveable* objects discussed in the previous section) can cause significant degradation in the performance of any Visual-SLAM pipeline [131]. Dynamic object detection (e.g. [90], [48], [58]) can therefore be used as an input into a Visual-SLAM pipeline to suppress outliers caused by said moving objects.

D. Synergies in next generation

Table IV compares the current 4R architecture with previous and next generation architectures. Previous generation has simplistic features due to limited compute availability. For Recognition, it had only pedestrian and park-slot detection using classical machine learning. Reconstruction was performed using sparse optical flow in software without any hardware accelerators. There was no image level reorganization or relocalization. CNN models have progressed rapidly to provide state-of-the-art results for geometric tasks like reconstruction [17] and relocalization [133]. For the next generation, a unified CNN model with high synergies would be the likely path. We have recently published an initial prototype Omnidet [132] showing joint modelling of reconstruction and recognition. Figure 23 illustrates its high level architecture with cross links shown across the different tasks. In Table V, an extract of the results from that work are presented. While this is not

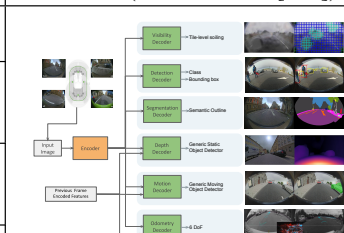
a complete implementation of what we would consider next generation, it does indicate that there is significant potential in jointly considering vision tasks.

While we believe that the proposal to use a complete CNN pipeline will offer optimal performance, we openly acknowledge that the approach has shortcomings. There are open challenges to improve the interpretability and trustworthiness of the perception model [134]. Uncertainty estimation plays a critical role in providing the confidence of the prediction to effectively fuse with other sensor perception and it will also enable the determination of out of distribution input samples for safe handover to a manual driver. The ability to debug failures during development and deployment phase is another challenge which has recently received increased attention through explainable AI techniques [135]. However, until interpretability and debuggability approaches for CNN-based processing is significantly more mature, we would propose that some level of redundancy is required by supplementing complete CNN-based approaches with classical computer vision and machine learning. We discuss this in more detail in the next section. In addition, it should be noted that some prominent scholars argue that improved performance is obtained by designing principled algorithms for the geometric estimations, and using deep neural networks for the extraction of robust visual features [136], where they argue that taking this approach results in a system capable of redeployment into new scenes without fine-tuning or retraining. That said, in future work, we plan to explore inclusion of relocalization and reorganization in neural network frameworks [137].

E. Dual-sources of detection

We have thus far discussed possible synergies between Reconstruction, Recognition and Relocalization. There is, however, another overarching synergistic consideration: that of redundancy. In automated vehicles, redundancy plays a significant role in the safety of the application. When a system component fails, then another must be available to ensure that the vehicle remains in a safe state. For example, FuseModNet [138] illustrates a synergistic fusion of cameras which provide dense information and lidar which performs well at low light. In terms of sensing, this would traditionally be achieved using multiple sensor types, such as computer vision systems, radar and laser scanner (Figure 24). For near-field sensing, an array of ultrasonic sensors is a mature low-cost sensor which provides robust safety around the vehicle [139].

TABLE IV: Features provided by current generation 4R architecture and comparison with previous and next generation.

Module	Previous Gen	Current Gen 4R framework	Next Gen (Unified CNN [132])
Recognition	Pedestrian (PD) Park-slot detection (PSD)	Bounding Box - PD, Cyclist, Vehicles Segmentation - Road, curb, road markings	
Reconstruction	Sparse 3D Reconstruction Sparse Flow Visual Odometry(VO)	Dense Depth Dense Flow clustering Multicam VO	
Reorganization	-	Static Obj Fusion - Freespace, Curb Dynamic Obj Fusion - Vehicles, PD, Cyclists Lane Handler - Multicamera 3D lane fit	
Relocalization	-	Sparse feature geometric map	

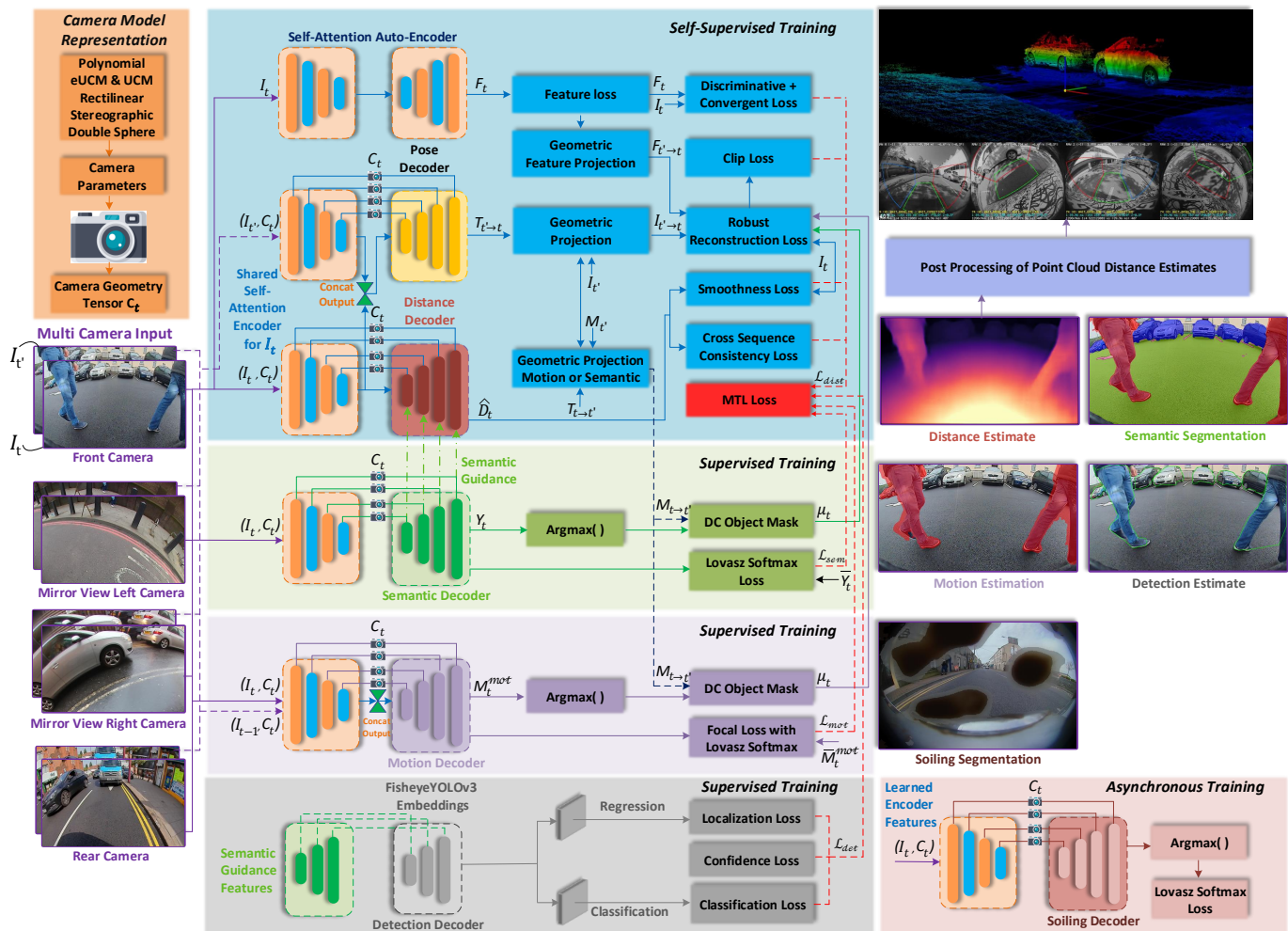


Fig. 23: Overview of our next generation unified multi-task visual perception framework. Refer to our OmniDet paper [132] for more details.

It is our contention that added safety is achieved through the parallel usage of different computer vision algorithm types. That is, a computer vision system architecture can be configured to maximize redundancy. This is particularly true as the sources are completely different types of processing – for example, statistical processing from the recognition pipeline and geometric processing from the reconstruction pipeline (Figure 12). In addition, such processing will typically run on different silicon components within an SoC. However, one must be aware that if you maximize the other synergies, the potential for redundancy is reduced. For example, if you use a CNN-based depth as a seed for a Visual-SLAM algorithm, you cannot claim the CNN as a redundancy for the Visual-SLAM,

as Visual-SLAM is now dependent on the CNN processing. One must also be aware that the two processing elements will likely use the same video feed – and so the safety of the camera itself and associated hardware/software, may also be a limiting factor. However, one should consider the potential for added safety in a system design following the 4R principles.

VI. CONCLUSIONS

In this paper, we provided a high-level survey of visual perception on surround-view cameras targeting commercial grade automated driving systems. We structure our survey into modular components namely Recognition, Reconstruction, Relocalization and Reorganization, jointly called 4Rs,

TABLE V: Extract of the results from our previous OmniDet work [132]. It can be seen that jointly learning the tasks outperforms treating each task separately (\downarrow means lower is better, \uparrow means higher is better, PA denotes pixel accuracy). VarNorm task weighting is used.

	Distance Estimation		Semantic Segmentation		Motion Segmentation		Object Detection
	Sq. Rel \downarrow	Abs Rel \downarrow	mIoU \uparrow	PA \uparrow	mIoU \uparrow	PA \uparrow	mAP \uparrow
Single Task	0.060	0.304	72.5	94.8	68.1	94.1	63.5
OmniDet	0.046	0.276	76.6	96.4	75.3	96.1	68.4

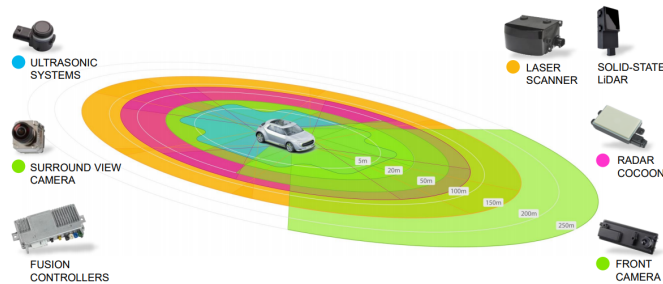


Fig. 24: Perception cocoon of redundant vehicle sensors. Safety is achieved through at least one sensor being available if others become unavailable.

and we argue that designing a vision architecture for vehicle automation along the lines of a 4R architecture can lead to system efficiencies. We discussed each component in detail and then we discussed how they are synergized to provide a more accurate system. We also provide a system and application context helping understand an industrial system. We have presented several architectures and frameworks, augmented with results predominantly from our previous publications, that support our argument in this direction.

The first three of the 4Rs (Recognition, Reconstruction, Relocalization) provide the means for the detection of objects and the extraction of their geometry and location in reference to the autonomous vehicle. However, that is not a complete description of the scene, and the fourth R (Reorganization) provides a higher-level scene understanding that can include the contextual spatial and temporal relationships between objects in the scene and the autonomous vehicle. Though massive advances have been made in the last decade in computer vision, we cannot yet claim to have achieved this complete scene understanding. It is likely that full vehicle autonomy will not be feasible until we have such a high level of visual reasoning deployed on vehicles. However, we propose that the 4R architecture can encapsulate, and provide a framework for, this level of vehicular cognition.

ACKNOWLEDGMENT

We would like to thank our employer Valeo for encouraging advanced research. Many thanks to Edward Jones (NUI Galway) and Matthieu Cord (Sorbonne University and Valeo.ai) for providing a detailed review prior to submission. We would also like to thank our colleagues Fabian Burger, Nagarajan Balmukundan, Pantelis Ermilios and Nivedita Tripathi for supporting the paper.

REFERENCES

- [1] E. Yurtsever, J. Lambert, A. Carballo, and K. Takeda, "A survey of autonomous driving: Common practices and emerging technologies," *IEEE Access*, vol. 8, pp. 58 443–58 469, 2020.
- [2] K. Stricker, T. Wendt, W. Stark, M. Gottfredson, R. Tsang, and M. Schallehn, "Electric and autonomous vehicles: The future is now," *Bain & Company Brief*, 2020, (accessed 30/03/2021). [Online]. Available: <https://www.bain.com/insights/electric-and-autonomous-vehicles-the-future-is-now/>
- [3] CB Insights, "40+ corporations working on autonomous vehicles," Available at <https://www.cbinsights.com/research/autonomous-driverless-vehicles-corporations-list/> (accessed 13/07/2020), 2020.
- [4] C. Badue, R. Guidolini, R. V. Carneiro, P. Azevedo, V. B. Cardoso, A. Forechi, L. Jesus, R. Berriel, T. Paixão, F. Mutz, L. Veronese, T. Oliveira-Santos, and A. F. D. Souza, "Self-driving cars: A survey," *Expert Systems with Applications*, vol. 165, p. 113816, 2021.
- [5] J. Malik, P. Arbeláez, J. Carreira, K. Fragkiadaki, R. Girshick, G. Gkioxari, S. Gupta, B. Hariharan, A. Kar, and S. Tulsiani, "The three R's of computer vision: Recognition, reconstruction and reorganization," *Pattern Recognition Letters*, vol. 72, pp. 4–14, 2016.
- [6] M. Heimberger, J. Horgan, C. Hughes, J. McDonald, and S. Yogamani, "Computer vision in automated parking systems: Design, implementation and challenges," *Image and Vision Computing*, vol. 68, pp. 88–101, 2017.
- [7] G. Velez and O. Otaegui, "Embedding vision-based advanced driver assistance systems: a survey," *IET Intelligent Transportation Systems*, vol. 11, no. 3, 2016.
- [8] R. P. Loce, E. A. Bernal, W. Wu, and R. Bala, "Computer vision in roadway transportation systems: a survey," *Journal of Imaging*, vol. 22, no. 4, 2013.
- [9] C. Hughes, M. Glavin, E. Jones, and P. Denny, "Wide-angle camera technology for automotive applications: a review," *IET Transactions on Intelligent Transport Systems*, vol. 3, pp. 19–31, 2009.
- [10] C. Badue, R. Guidolini, R. V. Carneiro, P. Azevedo, V. B. Cardoso, A. Forechi, L. Jesus, R. Berriel, T. M. Paixão, F. Mutz, L. de Paula Veronese, T. Oliveira-Santos, and A. F. De Souza, "Self-driving cars: A survey," *Expert Systems with Applications*, vol. 165, p. 113816, 2021.
- [11] F. Li, P. Bonnifait, and J. Ibanez-Guzman, "Estimating localization uncertainty using multi-hypothesis map-matching on high-definition road maps," in *Proceedings of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2017.
- [12] J. Martinez-Carranza, A. Calway, and W. Mayol-Cuevas, "Enhancing 6D visual relocalisation with depth cameras," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2013.
- [13] F. Gustafsson, "Automotive safety systems: Replacing costly sensors with software algorithms," *IEEE Signal Processing Magazine*, vol. 26, no. 4, pp. 32–47, 2009.
- [14] D. Alvarez-Coello, B. Klotz, D. Wilms, S. Fejji, J. M. Gómez, and R. Troncy, "Modeling dangerous driving events based on in-vehicle data using random forest and recurrent neural network," in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2019.
- [15] G. Borghi, R. Gasparini, R. Vezzani, and R. Cucchiara, "Embedded recurrent network for head pose estimation in car," in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2017.
- [16] M. Siam, M. Gamal, M. Abdel-Razek, S. Yogamani, and M. Jagersand, "RTSeg: Real-time semantic segmentation comparative study," in *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, 2018, pp. 1603–1607.
- [17] S. Qiao, Y. Zhu, H. Adam, A. Yuille, and L. Chen, "ViP-DeepLab: Learning visual perception with depth-aware video panoptic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

- [18] V. R. Kumar, S. Milz, C. Witt, M. Simon, K. Amende, J. Petzold, S. Yogamani, and T. Pech, "Near-field depth estimation using monocular fisheye camera: A semi-supervised learning approach using sparse LiDAR data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2018.
- [19] R. Li, S. Wang, and D. Gu, "DeepSLAM: a robust monocular SLAM system with unsupervised deep learning," *IEEE Transactions on Industrial Electronics*, vol. 68, no. 4, pp. 3577–3587, 2021.
- [20] S. Yogamani, C. Hughes, J. Horgan, G. Sistu, P. Varley, D. O’Dea, M. Uříčář, S. Milz, M. Simon, K. Amende *et al.*, "WoodScape: A multi-task, multi-camera fisheye dataset for autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 9308–9318.
- [21] P. Walzer and H.-W. Grove, "IRVW Futura - The Volkswagen Research Car," in *Passenger Car Conference and Exposition*. SAE International, 1990.
- [22] T. Wu, P. Tsai, N. Hu, and J. Chen, "Research and implementation of auto parking system based on ultrasonic sensors," in *Proceedings of the International Conference on Advanced Materials for Science and Engineering (AMSE)*, 2016, pp. 643–645.
- [23] Y. Song and C. Liao, "Analysis and review of state-of-the-art automatic parking assist system," in *Proceedings of the IEEE International Conference on Vehicular Electronics and Safety (ICVES)*, 2016.
- [24] C. Wang, H. Zhang, M. Yang, X. Wang, L. Ye, and C. Guo, "Automatic parking based on a bird’s eye view vision system," *Advances in Mechanical Engineering*, vol. 6, p. 847406, 2014.
- [25] U. Schwesinger, M. Bürki, J. Timpner, S. Rottmann, L. Wolf, L. M. Paz, H. Grimmert, I. Posner, P. Newman, C. Häne, L. Heng, G. H. Lee, T. Sattler, M. Pollefeys, M. Allodi, F. Valenti, K. Mimura, B. Goebelsmann, W. Derendarz, P. Mühlfellner, S. Wonneberger, R. Waldmann, S. Grysczyk, C. Last, S. Brüning, S. Horstmann, M. Bartholomäus, C. Brummer, M. Stellmacher, F. Pucks, M. Nicklas, and R. Siegwart, "Automated valet parking and charging for e-mobility," in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2016, pp. 157–164.
- [26] S. Ma, W. Fang, H. Jiang, M. Han, and C. Li, "Parking space recognition method based on parking space feature construction in the scene of autonomous valet parking," *Applied Sciences*, vol. 11, no. 6, 2021.
- [27] I. Isaksson-Hellman and M. Lindman, "The effect of a low-speed automatic brake system estimated from real life data," *Annals of Advances in Automotive Medicine*, vol. 56, pp. 231–240, 2012.
- [28] S. Lüke, O. Fochler, T. Schaller, and U. Regensburger, *Traffic-Jam Assistance and Automation*. Springer International Publishing, 2014, pp. 1–13.
- [29] L. Ulrich, "Bmw’s autonomous "extended traffic jam assistant" needs supervision," *IEEE Spectrum*, 2019.
- [30] S. J. Rao and G. J. Forkenbrock, "Test procedures traffic jam assist test development considerations," National Highway Traffic Safety Administration, Tech. Rep. DOT HS 812 757, July 2019.
- [31] T. Nothdurft, P. Hecker, S. Ohl, F. Saust, M. Maurer, A. Reschka, and J. R. Böhmer, "StadtPilot: First fully autonomous test drives in urban traffic," in *Proceedings of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2011, pp. 919–924.
- [32] C. R. Sunstein, "Rear visibility and some unresolved problems for economic analysis (with notes on experience goods)," Harvard Law School, Harvard University, Tech. Rep., July 2019. [Online]. Available: <https://dx.doi.org/10.2139/ssrn.3411057>
- [33] J. B. Cicchino, "Real-world effects of rear automatic braking and other backing assistance systems," *Journal of Safety Research*, vol. 68, pp. 41–47, 2019.
- [34] V. Usenko, N. Demmel, and D. Cremers, "The double sphere camera model," in *2018 International Conference on 3D Vision (3DV)*, 2018, pp. 552–560.
- [35] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *International Journal of Robotics Research (IJRR)*, 2013.
- [36] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, "1 Year, 1000km: The Oxford RobotCar Dataset," *The International Journal of Robotics Research (IJRR)*, vol. 36, no. 1, pp. 3–15, 2017.
- [37] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 3213–3223.
- [38] C. Hughes, R. McFeely, P. Denny, M. Glavin, and E. Jones, "Equidistant fish-eye perspective with application in distortion centre estimation," *Image and Vision Computing*, vol. 28, no. 3, pp. 538–551, 2010.
- [39] E. Plaut, E. Ben Yaacov, and B. El Shlomo, "3d object detection from a single fisheye image without a single fisheye training image," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2021, pp. 3659–3667.
- [40] T. S. Cohen, M. Geiger, J. Koehler, and M. Welling, "Spherical cnns," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [41] B. Coors, A. Paul Condurache, and A. Geiger, "SphereNet: Learning spherical representations for detection and classification in omnidirectional images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 518–533.
- [42] M. Eder and J.-M. Frahm, "Convolutions on spherical images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPR-W)*, 2019.
- [43] T. Ho and M. Budagavi, "Dual-fisheye lens stitching for 360-degree imaging," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 2172–2176.
- [44] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, "ENet: A deep neural network architecture for real-time semantic segmentation," *arXiv preprint arXiv:1606.02147*, 2016.
- [45] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 6, pp. 1137–1149, 2017.
- [46] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [47] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2015, pp. 2650–2658.
- [48] M. Siam, H. Mahgoub, M. Zahran, S. Yogamani, M. Jagersand, and A. El-Sallab, "MODNet: Motion and appearance based moving object detection network for autonomous driving," in *Proceedings of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 2859–2864.
- [49] J. Engel, T. Schöps, and D. Cremers, "LSD-SLAM: Large-scale direct monocular SLAM," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 834–849.
- [50] S. Wang, R. Clark, H. Wen, and N. Trigoni, "DeepVO: Towards end-to-end visual odometry with deep recurrent convolutional neural networks," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017.
- [51] S. Mittal, "A survey on optimized implementation of deep learning models on the nvidia jetson platform," *Journal of Systems Architecture*, vol. 97, pp. 428–442, 2019.
- [52] L. Yahiaoui, J. Horgan, B. Deegan, S. Yogamani, C. Hughes, and P. Denny, "Overview and empirical analysis of ISP parameter tuning for visual perception in autonomous driving," *Journal of Imaging*, vol. 5, no. 10, 2019.
- [53] A. Mosleh, A. Sharma, E. Onzon, F. Mannan, N. Robidoux, and F. Heide, "Hardware-in-the-loop end-to-end optimization of camera image processing pipelines," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 7529–7538.
- [54] M. Buckler, S. Jayasuriya, and A. Sampson, "Reconfiguring the imaging pipeline for computer vision," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017.
- [55] R. Szeliski, *Computer Vision: Algorithms and Applications*, 1st ed. Springer, 2010.
- [56] —, "A multi-view approach to motion and stereo," Microsoft Research, Tech. Rep. MSR-TR-99-19, 1999.
- [57] M. R. U. Saputra, A. Markham, and N. Trigoni, "Visual SLAM and structure from motion in dynamic environments: A survey," *ACM Computing Surveys*, vol. 51, no. 2, 2018.
- [58] M. Yahiaoui, H. Rashed, L. Mariotti, G. Sistu, I. Clancy, L. Yahiaoui, and S. Yogamani, "FisheyeMODNet: Moving object detection on surround-view cameras for autonomous driving," in *Proceedings of the Irish Machine Vision and Image Processing (IMVIP)*, 2019.
- [59] G. Sistu, I. Leang, S. Chennupati, S. Yogamani, C. Hughes, S. Milz, and S. Rawashdeh, "NeurAll: Towards a unified visual perception model for automated driving," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, 2019, pp. 796–803.
- [60] S. Bauer, S. Köhler, K. Doll, and U. Brunsman, "Fpga-gpu architecture for kernel svm pedestrian detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition - Workshops (CVPR-W)*, 2010, pp. 61–68.

- [61] M. Hassaballah and A. I. Awad, *Deep learning in computer vision: principles and applications*. CRC Press, 2020.
- [62] M. Uříčář, D. Hurych, P. Křížek, and S. Yogamani, “Challenges in designing datasets and validation for autonomous driving,” in *Proceedings of the International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISAPP)*, 2019, pp. 653–659.
- [63] A. Briot, P. Viswanath, and S. Yogamani, “Analysis of efficient cnn design techniques for semantic segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2018, pp. 663–672.
- [64] A. Das., S. Kandan., S. Yogamani., and P. Křížek, “Design of real-time semantic segmentation decoder for automated driving,” in *Proceedings of the International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISAPP)*, 2019, pp. 393–400.
- [65] V. Ravi Kumar, S. Yogamani, M. Bach, C. Witt, S. Milz, and P. Mader, “UnRectDepthNet: Self-supervised monocular depth estimation using a generic framework for handling common camera distortion models,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- [66] H. Rashed, E. Mohamed, G. Sistu, V. R. Kumar, C. Eising, A. El-Sallab, and S. Yogamani, “Generalized object detection on fisheye cameras for autonomous driving: Dataset, representations and baseline,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021, pp. 2272–2280.
- [67] L. Deng, M. Yang, Y. Qian, C. Wang, and B. Wang, “Cnn based semantic segmentation for urban traffic scenes using fisheye camera,” in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2017, pp. 231–236.
- [68] P. Maddu, W. Doherty, G. Sistu, I. Leang, M. Uříčář, S. Chennupati, H. Rashed, J. Horgan, C. Hughes, and S. Yogamani, “Fisheyemultinet: Real-time multi-task learning architecture for surround-view automated parking system,” in *Proceedings of the Irish Machine Vision and Image Processing Conference (IMVIP)*, 2019.
- [69] I. Leang, G. Sistu, F. Bürger, A. Bursuc, and S. Yogamani, “Dynamic task weighting methods for multi-task networks in autonomous driving systems,” in *Proceedings of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2020.
- [70] S. Chennupati, G. Sistu, S. Yogamani, and S. Rawashdeh, “Auxnet: Auxiliary tasks enhanced semantic segmentation for automated driving,” in *Proceedings of the International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISAPP)*, 2019, pp. 645–652.
- [71] M. Uříčář, P. Křížek, G. Sistu, and S. Yogamani, “Soilingnet: Soiling detection on automotive surround-view cameras,” in *Proceedings of the IEEE Intelligent Transportation Systems Conference (ITSC)*, 2019, pp. 67–72.
- [72] M. Uříčář, J. Ulicny, G. Sistu, H. Rashed, P. Křížek, D. Hurych, A. Vobecky, and S. Yogamani, “Desoiling dataset: Restoring soiled areas on automotive fisheye cameras,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2019.
- [73] Z. Kukulova and V. Larsson, “Radial distortion triangulation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 9681–9686.
- [74] A. Saxena, S. H. Chung, and A. Y. Ng, “3-D depth reconstruction from a single still image,” *International Journal of Computer Vision*, vol. 76, pp. 53–69, 2007.
- [75] D. Marr and T. Poggio, “A computational theory of human stereo vision,” *Proceedings of the Royal Society of London B*, vol. 292, pp. 301–328, 1979.
- [76] W. E. L. Grimson, “A computer implementation of a theory of human stereo vision,” *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, vol. 292, no. 1058, pp. 217–253, 1981.
- [77] E. Davies, “The three-dimensional world,” in *Machine Vision (Third Edition) – Theory, Algorithms, Practicalities*. Springer, 2005, ch. 16, pp. 445–485.
- [78] Z. Kukulova, J. Heller, M. Bujnak, and T. Pajdla, “Radial distortion homography,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 639–647.
- [79] A. Hirata, R. Ishikawa, M. Roxas, and T. Oishi, “Real-time dense depth estimation using semantically-guided lidar data propagation and motion stereo,” *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3806–3811, 2019.
- [80] V. R. Kumar, S. Milz, C. Witt, M. Simon, K. Amende, J. Petzold, S. Yogamani, and T. Pech, “Monocular fisheye camera depth estimation using sparse lidar supervision,” in *Proceedings of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 2853–2858.
- [81] T. van Dijk and G. de Croon, “How do neural networks see depth in single images?” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [82] W. Gao and S. Shen, “Dual-fisheye omnidirectional stereo,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 6715–6722.
- [83] J. Esparza, M. Helmle, and B. Jähne, “Wide base stereo with fisheye optics: A robust approach for 3d reconstruction in driving assistance,” in *Pattern Recognition*, X. Jiang, J. Hornegger, and R. Koch, Eds. Springer International Publishing, 2014, pp. 342–353.
- [84] R. T. Collins, “A space-sweep approach to true multi-image matching,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 1996, pp. 358–363.
- [85] D. Gallup, J.-M. Frahm, P. Mordohai, Q. Yang, and M. Pollefeys, “Real-time plane-sweeping stereo with multiple sweeping directions,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007.
- [86] C. Häne, L. Heng, G. H. Lee, A. Sizov, and M. Pollefeys, “Real-time direct dense matching on fisheye images using plane-sweeping stereo,” in *Proceedings of the IEEE International Conference on 3D Vision (3DV)*, 2014, pp. 57–64.
- [87] R. Hartley and P. Sturm, “Triangulation,” *Computer Vision and Image Understanding*, vol. 68, no. 2, pp. 146–157, 1997.
- [88] D. Barath and J. Matas, “Multi-class model fitting by energy minimization and mode-seeking,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [89] J. Klappstein, F. Stein and U. Franke, “Detectability of moving objects using correspondences over two and three frames,” in *Proceedings of the Joint Pattern Recognition Symposium*, 2007, pp. 112–121.
- [90] L. Mariotti and C. Eising, “Spherical formulation of geometric motion segmentation constraints in fisheye cameras,” *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [91] E. Mohamed, M. Ewaisha, M. Siam, H. Rashed, S. Yogamani, W. Hamdy, M. Helmi, and A. El-Sallab, “Monocular instance motion segmentation for autonomous driving: Kitti instancemotseg dataset and multi-task baseline,” in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2021.
- [92] A. Brunker, T. Wohlgenuth, M. Frey, and F. Gauterin, “Odometry 2.0: A slip-adaptive eif-based four-wheel-odometry model for parking,” *IEEE Transactions on Intelligent Vehicles*, vol. 4, no. 1, pp. 114–126, 2019.
- [93] M. O. A. Aqel, M. H. Marhaban, M. I. Saripan, and N. B. Ismail, “Review of visual odometry: types, approaches, challenges, and applications,” *SpringerPlus*, vol. 5, 2016.
- [94] A. Bugeau and P. Pérez, “Detection and segmentation of moving objects in complex scenes,” *Computer Vision and Image Understanding*, vol. 113, no. 4, pp. 459–476, 2009.
- [95] V. R. Kumar, S. A. Hiremath, M. Bach, S. Milz, C. Witt, C. Pinard, S. Yogamani, and P. Mäder, “FisheyeDistanceNet: Self-supervised scale-aware distance estimation using monocular fisheye camera for autonomous driving,” in *Proceedings of the IEEE/RSJ International Conference on Robotics and Automation (ICRA)*, 2020, pp. 574–581.
- [96] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, “Monoslam: Real-time single camera slam,” *Transaction Pattern Analysis Machine Intelligence*, vol. 29, no. 6, pp. 1052–1067, 2007.
- [97] G. Klein and D. Murray, “Parallel tracking and mapping for small AR workspaces,” in *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2007.
- [98] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, “ORB-SLAM: a versatile and accurate monocular SLAM system,” *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [99] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, “DTAM: Dense tracking and mapping in real-time,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2011, pp. 2320–2327.
- [100] D. Caruso, J. Engel, and D. Cremers, “Large-scale direct slam for omnidirectional cameras,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2015, pp. 141–148.
- [101] S. Milz, G. Arbeiter, C. Witt, B. Abdallah, and S. Yogamani, “Visual slam for automated driving: Exploring the applications of deep learn-

- ing,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2018.
- [102] TomTom N.V., “TomTomRoadDNA,” <https://www.tomtom.com/products/hd-map/>, [Online: 12/2020].
- [103] B. Ravi Kiran, L. Roldao, B. Irastorza, R. Verastegui, S. Suss, S. Yogamani, V. Talpaert, A. Lepoutre, and G. Trehard, “Real-time dynamic object detection for autonomous driving using prior 3d-maps,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [104] Q. Luo, Y. Cao, J. Liu, and A. Benslimane, “Localization and navigation in autonomous driving: Threats and countermeasures,” *IEEE Wireless Communications*, vol. 26, no. 4, pp. 38–45, 2019.
- [105] A. Singandhupe and H. M. La, “A review of slam techniques and security in autonomous driving,” in *Proceedings of the IEEE International Conference on Robotic Computing (IRC)*, 2019, pp. 602–607.
- [106] A. Kasyanov, F. Engelmann, J. Stückler, and B. Leibe, “Keyframe-based visual-inertial online slam with relocalization,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- [107] D. G. Lowe, “Object recognition from local scale-invariant features,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 1999, pp. 1150–1157.
- [108] E. Rublee, V. Rabaud, K. Konolige, and G. R. Bradski, “ORB: An efficient alternative to SIFT or SURF,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2011, pp. 2564–2571.
- [109] P. F. Alcantarilla, J. Nuevo, and A. Bartoli, “Fast explicit diffusion for accelerated features in nonlinear scale spaces,” *IEEE Transactions on Pattern Analysis Machine Intelligence*, vol. 34, no. 7, p. 1281–1298, 2011.
- [110] N. Tripathi and S. Yogamani, “Trained trajectory based automated parking system using Visual SLAM,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2021.
- [111] G. Ciaparrone, F. Luque Sánchez, S. Tabik, L. Troiano, R. Tagliaferri, and F. Herrera, “Deep learning in video multi-object tracking: A survey,” *Neurocomputing*, vol. 381, pp. 61–88, 2020.
- [112] K. Baizid, G. Lozenguez, L. Fabresse, and N. Bouraqadi, “Vector maps: A lightweight and accurate map format for multi-robot systems,” in *Proceedings of the International Conference on Intelligent Robotics and Applications (ICIRA)*, N. Kubota, K. Kiguchi, H. Liu, and T. Obo, Eds., 2016, pp. 418–429.
- [113] S. Thrun, W. Burgard, and D. Fox, *Probabilistic Robotics*. MIT Press, 2005.
- [114] Z. Wang, Y. Wu, and Q. Niu, “Multi-sensor fusion in automated driving: A survey,” *IEEE Access*, vol. 8, pp. 2847–2868, 2020.
- [115] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, “Deep ordinal regression network for monocular depth estimation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [116] V. P. Feihu Zhang, R. Yang, and P. H. S. Torr, “GA-Net: Guided aggregation net for end-to-end stereo matching,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [117] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, “Unsupervised learning of depth and ego-motion from video,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [118] M. Tatarchenko, S. R. Richter, R. Ranftl, Z. Li, V. Koltun, and T. Brox, “What do single-view 3d reconstruction networks learn?” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [119] X. Lin, D. Sánchez-Escobedo, J. R. Casas, and M. Pardás, “Depth estimation and semantic segmentation from a single RGB image using a hybrid convolutional neural network,” *Sensors*, vol. 19, 2019.
- [120] V. R. Kumar, M. Klingner, S. Yogamani, S. Milz, T. Fingscheidt, and P. Mader, “Syndistnet: Self-supervised monocular fisheye camera distance estimation synergized with semantic segmentation for autonomous driving,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021, pp. 61–71.
- [121] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The KITTI dataset,” *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [122] C. Luo, Z. Yang, P. Wang, Y. Wang, W. Xu, R. Nevatia, and A. Yuille, “Every Pixel Counts ++: Joint Learning of Geometry and Motion with 3D Holistic Understanding,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [123] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, “Digging into self-supervised monocular depth estimation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 3828–3838.
- [124] V. Guizilini, R. Ambrus, S. Pillai, and A. Gaidon, “3D packing for self-supervised monocular depth estimation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2485–2494.
- [125] B. Liu, S. Gould, and D. Koller, “Single image depth estimation from predicted semantic labels,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [126] S. Gould, R. Fulton, and D. Koller, “Decomposing a scene into geometric and semantically consistent regions,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2009.
- [127] M. Siam, S. Elkerdawy, M. Jagersand, and S. Yogamani, “Deep semantic segmentation for automated driving: Taxonomy, roadmap and challenges,” in *Proceedings of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2017, pp. 1–8.
- [128] D. Galvez-López and J. D. Tardos, “Bags of binary words for fast place recognition in image sequences,” *IEEE Transactions on Robotics*, vol. 28, no. 5, pp. 1188–1197, 2012.
- [129] S. Garg, N. Suenderhauf, and M. Milford, “Lost? appearance-invariant place recognition for opposite viewpoints using visual semantics,” in *Proceedings of Robotics: Science and Systems (RSS)*, 2018.
- [130] J. Li, D. Meger, and G. Dudek, “Semantic mapping for view-invariant relocalization,” in *Proceedings of the IEEE/RSJ International Conference on Robotics and Automation (ICRA)*, 2019.
- [131] S. Wangsiripitak and D. W. Murray, “Avoiding moving outliers in visual SLAM by tracking moving objects,” in *Proceedings of the International Conference on Robotics and Automation (ICRA)*, 2009.
- [132] V. Ravikumar, S. Yogamani, H. Rashed, G. Sistu, C. Witt, I. Leang, S. Milz, and P. Mader, “Omnidet: Surround view cameras based multi-task visual perception network for autonomous driving,” *IEEE Robotics and Automation Letters*, 2021.
- [133] G. L. Oliveira, N. Radwan, W. Burgard, and T. Brox, “Topometric localization with deep learning,” in *Robotics Research*, N. M. Amato, G. Hager, S. Thomas, and M. Torres-Torriti, Eds. Cham: Springer International Publishing, 2020, pp. 505–520.
- [134] G. Schwalbe and M. Schels, “A survey on methods for the safety assurance of machine learning based systems,” in *Proceedings of the 10th European Congress on Embedded Real Time Software and Systems (ERTS 2020)*, 2020.
- [135] É. Zablocki, H. Ben-Younes, P. Pérez, and M. Cord, “Explainability of vision-based autonomous driving systems: Review and challenges,” *arXiv preprint arXiv:2101.05307*, 2021.
- [136] P.-E. Sarlin, A. Unagar, M. Larsson, H. Germain, C. Toft, V. Larsson, M. Pollefeys, V. Lepetit, L. Hammarstrand, F. Kahl, and T. Sattler, “Back to the feature: Learning robust camera localization from pixels to pose,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [137] Y. Cui, R. Chen, W. Chu, L. Chen, D. Tian, Y. Li, and D. Cao, “Deep learning for image and point cloud fusion in autonomous driving: A review,” *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–18, 2021.
- [138] H. Rashed, M. Ramzy, V. Vaquero, A. El Sallab, G. Sistu, and S. Yogamani, “Fusmodnet: Real-time camera and lidar based moving object detection for robust low-light autonomous driving,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2019, pp. 0–0.
- [139] M. Pöpperli, R. Gulagundi, S. Yogamani, and S. Milz, “Capsule neural network based height classification using low-cost automotive ultrasonic sensors,” in *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, 2019, pp. 661–666.