

Computer Vision for Autonomous Vehicles: Problems, Datasets and State of the Art

Joel Janai^{1,*} Fatma Güney^{2,*} Aseem Behl^{1,*} Andreas Geiger¹

¹ Autonomous Vision Group, MPI for Intelligent Systems and
University of Tübingen, Germany

² College of Engineering, Koç University, Turkey

December 18, 2019

* The first three authors contributed equally

Contents

1	Introduction	11
2	History of Autonomous Driving	15
3	Sensors	23
3.1	Camera Models	24
3.1.1	Omnidirectional Cameras	24
3.1.2	Event Cameras	25
3.2	Calibration	26
4	Datasets & Benchmarks	29
4.1	Computer Vision Datasets	31
4.1.1	Object Recognition	32
4.1.2	Object Tracking	33
4.1.3	Stereo and 3D Reconstruction	33
4.1.4	Optical Flow	34
4.2	Autonomous Driving Datasets	36
4.2.1	Object Detection and Semantic Segmentation	37
4.2.2	Tracking	39
4.2.3	Traffic Sign Detection	39
4.2.4	Road and Lane Detection	39
4.2.5	Flow and Stereo	40
4.2.6	Long-Term Autonomy	40
4.3	Synthetic Data Generation using Game Engines	41
5	Object Detection	45
5.1	Problem definition	45
5.2	Methods	46
5.2.1	Classical Pipeline	47
5.2.2	Part-based Approaches	49
5.2.3	Deep Learning for Detection	50

5.2.4	Real-time Pedestrian Detection	52
5.2.5	Human Pose Estimation	52
5.2.6	Traffic Sign Detection	53
5.2.7	3D Object Detection from 2D Images	54
5.2.8	3D Object Detection from 3D Point Clouds	55
5.3	Datasets	56
5.4	Metrics	57
5.5	State of the Art on KITTI	57
5.6	Discussion	63
6	Object Tracking	69
6.1	Problem Definition	69
6.2	Methods	70
6.2.1	Tracking by Detection	71
6.2.2	Tracking with Stereo	73
6.2.3	Pedestrian Tracking	74
6.2.4	Joint Detection and Tracking	74
6.2.5	Deep Learning for Multi-Object Tracking	75
6.3	Datasets	76
6.4	Metrics	78
6.5	State of the Art on MOT & KITTI	78
6.6	Discussion	82
7	Semantic Segmentation	85
7.1	Problem Definition	85
7.2	Methods	85
7.2.1	Deep Learning for Semantic Segmentation	87
7.2.2	Videos	91
7.2.3	Street Side Views	93
7.2.4	3D Data	93
7.2.5	Road Segmentation	97
7.2.6	Free Space Estimation	99
7.2.7	Stixels	100
7.2.8	Aerial Images	102
7.3	Datasets	104
7.4	Metrics	105
7.5	State of the Art on Cityscapes	105
7.6	Discussion	106
8	Semantic Instance Segmentation	107
8.1	Problem Definition	107
8.2	Methods	107
8.2.1	Proposal-based Approaches	107

8.2.2	Proposal-free Approaches	109
8.2.3	Panoptic Segmentation	111
8.3	Datasets	112
8.4	Metrics	113
8.5	State of the Art on Cityscapes	113
8.6	Discussion	113
9	Stereo	115
9.1	Problem Definition	115
9.2	Methods	115
9.2.1	Matching Cost	116
9.2.2	Energy Optimization	117
9.2.3	Higher-Order Models	118
9.2.4	Piecewise Planar Priors	119
9.2.5	Segmentation-based Models	119
9.2.6	Deep Learning for Stereo Matching	120
9.2.7	Variable Baseline	122
9.2.8	Omnidirectional Cameras	122
9.3	Datasets	123
9.4	Metrics	123
9.5	State of the Art on KITTI	123
9.6	Discussion	124
10	Multi-view 3D Reconstruction	127
10.1	Problem Definition	127
10.2	Structure from Motion	128
10.3	Multi-view Stereo	129
10.3.1	Planarity and Primitives	131
10.3.2	Shape Priors	131
10.3.3	Semantics	133
10.3.4	Efficient Reconstruction	135
10.3.5	Deep Learning for Multi-View Stereo	135
10.3.6	Omnidirectional Cameras	136
10.4	Datasets	137
10.5	Metrics	137
10.6	State of the Art on ETH3D & Tanks and Temples	137
10.7	Discussion	139
11	Optical Flow	141
11.1	Problem Definition	141
11.2	Methods	141
11.2.1	Sparse Matches	143
11.2.2	Epipolar Flow	145

11.2.3	Semantic Segmentation	145
11.2.4	Deep Learning for Optical Flow	146
11.2.5	High-Speed Flow	149
11.2.6	Confidences	150
11.3	Datasets	150
11.4	Metrics	150
11.5	State of the Art on KITTI	150
11.6	Discussion	152
12	3D Scene Flow	155
12.1	Problem Definition	155
12.2	Methods	155
12.2.1	Piecewise Rigidity	156
12.2.2	Semantic Segmentation	157
12.2.3	Scene Flow from 3D Point Clouds	158
12.3	Datasets	159
12.4	Metrics	159
12.5	State of the Art on KITTI	160
12.6	Discussion	161
13	Mapping, Localization & Ego-Motion Estimation	163
13.1	Problem Definition	163
13.2	Mapping	164
13.2.1	Metric Maps	164
13.2.2	Semantic Maps	165
13.3	Localization	165
13.3.1	Structure-based Localization	167
13.3.2	Cross-view Localization	171
13.3.3	Semantic Alignment from LiDAR	172
13.4	Ego-Motion Estimation	173
13.4.1	Drift	177
13.4.2	Loop Closure Detection	177
13.4.3	Simultaneous Localization and Mapping (SLAM)	178
13.5	Datasets	181
13.6	Metrics	182
13.7	State of the Art on KITTI	182
13.8	Discussion	186
14	Scene Understanding	189
14.1	Problem Definition	189
14.2	Methods	190
14.2.1	Road Topology and Traffic Participants	190
14.2.2	Physical and Temporal Relationships	191

14.3 Discussion 192

15 End-to-End Learning for Autonomous Driving 195

15.1 Problem Definition 195

15.2 Methods 195

15.2.1 Behavior Cloning 196

15.2.2 Reinforcement Learning 199

15.2.3 Combined Methods 201

15.2.4 Intermediate Representations 201

15.2.5 Transferring from Simulation to the Real World 203

15.3 Datasets 204

15.4 Metrics 206

15.5 Discussion 206

16 Conclusion 209

16.1 Acknowledgement 210

Abstract

Recent years have witnessed enormous progress in AI-related fields such as computer vision, machine learning, and autonomous vehicles. As with any rapidly growing field, it becomes increasingly difficult to stay up-to-date or enter the field as a beginner. While several survey papers on particular sub-problems have appeared, no comprehensive survey on problems, datasets, and methods in computer vision for autonomous vehicles has been published. This book attempts to narrow this gap by providing a survey on the state-of-the-art datasets and techniques. Our survey includes both the historically most relevant literature as well as the current state of the art on several specific topics, including recognition, reconstruction, motion estimation, tracking, scene understanding, and end-to-end learning for autonomous driving. Towards this goal, we analyze the performance of the state of the art on several challenging benchmarking datasets, including KITTI, MOT, and Cityscapes. Besides, we discuss open problems and current research challenges. To ease accessibility and accommodate missing references, we also provide a website that allows navigating topics as well as methods and provides additional information.

Chapter 1

Introduction

Since the first successful demonstrations in the 1980s [165, 167, 649], great progress has been made in the field of autonomous vehicles. However, despite these advances and ambitious commercial goals, fully autonomous navigation in general environments has not been realized to date. The reason for this is two-fold: First, autonomous systems which operate in complex dynamic environments require models which generalize to unpredictable situations and reason in a timely manner. Second, informed decisions require accurate perception, yet most of the existing computer vision models are still inferior to human perception and reasoning.

Existing approaches to self-driving can be roughly categorized into modular pipelines and monolithic end-to-end learning approaches. Both approaches are contrasted at a conceptual level in Figure 1.1. The modular pipeline is the standard approach to autonomous driving, mostly followed in the industry. The key idea is to break down the complex mapping function from high-dimensional inputs to low-dimensional control variables into modules which can be independently developed, trained, and tested. In Figure 1.1 (top), these modules comprise low-level perception, scene parsing, path planning, and vehicle control. However, this is just one particular example of modularizing a self-driving stack and other or more fine-grained modularizations are also possible. Existing approaches typically leverage machine learning (e.g., deep neural networks) to extract low-level features or to parse the scene into individual components. In contrast, path planning and vehicle control are dominated by classical state machines, search algorithms, and control models.

The major advantage of modular pipelines is that they deploy human interpretable intermediate representations such as detected objects or free space information which allow gaining insights into failure modes of the system. Furthermore, the development of modular pipelines can be easily parallelized

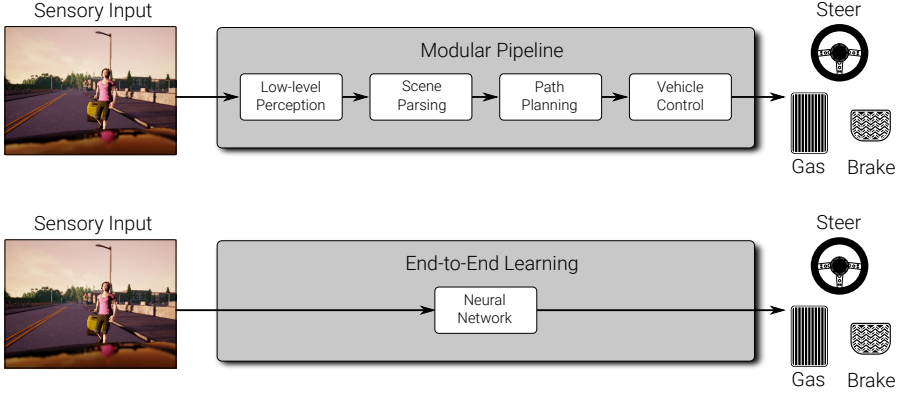


Figure 1.1: **Approaches to Self-Driving.** Classical modular pipeline (top) vs. monolithic end-to-end learning approach (bottom). See text for details.

within companies where typically different teams work on different aspects of the driving problem simultaneously. Furthermore, it is comparably easy to integrate first principles and prior knowledge about the problem into the system. Examples include traffic laws that can be explicitly enforced in the planner or knowledge about the vehicle dynamics, which lead to improved vehicle control. Other aspects that are more difficult to specify by hand, such as the appearance of pedestrians, are learned from large annotated datasets.

A major drawback of modular approaches is the fact that human-designed intermediate representations are not necessarily optimal for the driving task, which typically includes aspects like safety, comfort, and time for reaching the goal. Moreover, most modules are trained and validated independently from each other, making use of auxiliary loss functions. Consider the problem of object detection as an example. Most objects in the scene are not directly relevant for the driving task, yet the learning algorithm is not informed about the relevance of each object and therefore tasks a neural network to detect all objects with equal importance. Thus, the network is wasting capacity on irrelevant objects while not being able to detect the driving relevant objects with the necessary accuracy. This demonstrates the difficulty of defining appropriate intermediate representations and auxiliary loss functions.

An alternative to modular pipelines is end-to-end learning-based models which try to learn a policy, i.e., a function from observations to actions using a generic model such as a deep neural network. This approach is illustrated in Figure 1.1 (bottom) and discussed in detail in Chapter 15. The network parameters can be learned either via imitation learning by replicating the behavior of a teacher or using reinforcement learning by exploring the world and taking actions that are likely to yield a high user-specified reward.

However, reinforcement learning approaches suffer from the credit assignment and reward shaping problems, are typically slow and can only be applied in non-safety-critical simulation environments. Imitation learning, on the other hand, suffers from overfitting and does not easily generalize to novel scenarios. Furthermore, holistic neural network-based approaches are often hard to interpret as they present themselves as “black boxes” to the user which do not reveal *why* a certain error has occurred.

In this survey, we focus on perception for autonomous vehicles. In particular, we discuss the perception-related modules of the modular pipeline as well as end-to-end learning-based approaches. Other aspects of the self-driving problem are discussed in related surveys: For example, Winner et al. [705] put emphasis on driver assistance systems, considering both their structure and their function. Similarly, Klette [349] provides an overview of vision-based driver assistance systems. They describe most aspects of the perception problem at a high level but do not provide an in-depth review of the state of the art in each task as we pursue in this survey. Complementary to our work, Zhu et al. [786] provide an overview of environment perception for intelligent vehicles, focusing on lane detection, traffic sign/light recognition as well as vehicle tracking. In contrast, our goal is to bridge the gap between the robotics, intelligent vehicles, and computer vision communities by providing an extensive overview and comparison, including works from all three fields.

This survey is structured as follows: first, we provide a brief history of autonomous driving, followed by an introduction to camera models and calibration techniques. We then provide an overview of autonomous driving-related datasets with a particular focus on perception before surveying the relevant perception tasks and the state-of-the-art algorithms for solving them. More specifically, we review object detection, tracking, semantic (instance) segmentation, reconstruction, motion estimation, and scene understanding techniques. Each chapter starts with the problem definition, an overview over the most important methods and main design choices, a qualitative and quantitative analysis of the top-performing techniques on the most popular datasets, as well as a discussion of the state of the art in this area. Finally, we provide an overview of state-of-the-art end-to-end models for autonomous driving before concluding this survey. To ease navigation, we also provide an interactive online tool¹ which visualizes the surveyed papers with an interactive graph and additional information in an easily accessible manner. We hope that our survey will become a useful tool for researchers in the field of autonomous vision and lowers the entry barrier for beginners by providing a thorough overview of the field.

¹http://www.cvlibs.net/projects/autonomous_vision_survey

Chapter 2

History of Autonomous Driving

Similar to the invention of the automobile by Carl Benz in 1886, self-driving technology promises to profoundly impact our mobility. In this chapter, we briefly review the history of driverless and self-driving vehicles from 1925 to 2019.

The first demonstration of a driverless vehicle was reported in 1925 when Houdina Radio Control demonstrated the “American Wonder”, a remote-controlled vehicle that traveled along Broadway in New York City trailed by an operator in another vehicle[654]. Several years later, General Motors approached Norman Bel Geddes to sketch his vision about mobility 20 years into the future, culminating in Futurama, the most successful exhibition at the New York World Fair in 1939. Besides multi-lane highways, this vision sketched radio-controlled electric cars that navigated via electromagnetic circuits installed in the roadway. This vision led to several prototypes such as the GM Firebird II [100] in 1956, and RCA Labs’ wire controlled car in 1960 as well as a demonstration of Citroen with its DS 19 and the Cabintaxi¹ of Demag/MBB in 1970. However, the idea of infrastructure-based autonomous navigation is largely restricted to specific use cases such as ground transportation at airports, park shuttles, or automated facilities due to its limited scalability and high cost.

In 1986, the first self-driving car prototypes which did not rely on dedicated infrastructure hit the road. This pioneering effort was led by the Navlab team at CMU in the US as well as Ernst Dickmanns’s team at the Bundeswehr University Munich in Germany. Carnegie Mellon University’s Navlab team [649] (Figure 2.1) achieved another major milestone in 1995, by driving from

¹<https://www.youtube.com/watch?v=ERdFOFK-2io>



Figure 2.1: **The Navlab.** The self-contained laboratory from CMU for navigational vision system research. Figure courtesy of Thorpe et al. [649] © 1988 Springer.

Washington, D.C., to San Diego, CA, 98% autonomously with manual longitudinal control in the 'No hands across America' tour [527, 525]. With ALVINN [526], the Navlab team at CMU demonstrated an imitation learning approach where a relatively small neural network was optimized in an end-to-end fashion to keep the vehicle on the road based on user demonstrations. On the contrary, Dickmanns presented a modular approach in which a vehicle and road model was used for continuously estimating the state and controlling the vehicle [166]. The project was conducted in the context of the European PROMETHEUS project, which involved more than 13 vehicle manufacturers and several research units from governments and universities of 19 European countries. In 1995, the PROMETHEUS team demonstrated the first autonomous long-distance drive from Munich, Germany, to Odense, Denmark, at velocities up to 175 km/h with about 95% autonomous driving [168, 207, 164].

Motivated by the success of the PROMETHEUS projects to drive autonomously on highways, Franke et al. [208] describe a real-time vision system for autonomous driving in complex urban traffic situations. While highway scenarios have been studied intensively, urban scenes have not been addressed



Figure 2.2: **Winner of DARPA Grand Challenge.** The first autonomous vehicle to complete the DARPA Grand Challenge. Figure courtesy of Thrun et al. [651] © 2006 Wiley.

before. Their system included depth-based obstacle detection and tracking from stereo as well as a framework for monocular detection and recognition of relevant objects such as traffic signs. Many approaches to the challenging task of autonomous driving developed during these projects are presented and discussed in [48]. They concluded that sufficient computing power is becoming increasingly available, but difficulties like reflections, wet roads, direct sunshine, tunnels, and shadows still make data interpretation challenging. Thus, they suggested the enhancement of sensor capabilities. They also pointed out that the legal aspects related to the responsibility and impact of automatic driving on human passengers need to be considered carefully. In summary, the automation will likely be restricted to special infrastructures and will be extended gradually.

While full self-driving has remained unsolved to date, driver assistance systems have reached commercial success, enriching driving comfort and safety. In 1995, Mitsubishi presented the first LiDAR-based distance control [136], and in 1999 Mercedes-Benz implemented the radar-assisted adaptive cruise control. In 2000, navigation systems and digital road maps became available. Today, differential GPS in combination with inertial measurement units (IMU) allows for localization at an accuracy of 5cm in good conditions, enabling the use of detailed lane-level road maps (HD maps) and providing redundancy for noisy vision-based localization.

In 2004, the Defense Advanced Research Projects Agency (DARPA) of the US Department of Defense started to organize and sponsor a series of 3 races



Figure 2.3: **Waymo Autonomous Vehicle.** Figure courtesy of www.waymo.com © 2019 Waymo LLC.

to foster the development of self-driving technology [152]. The first race, the Darpa Grand Challenge 2004, was limited to US participants. DARPA offered a prize money of \$1 million for the first team autonomously completing a 240km long dirt route from California to Nevada through the Mojave desert, guided by GPS waypoints. However, none of the robot vehicles completed the route. One year later, in 2005, DARPA announced a second edition of its challenge with 5 vehicles successfully completing the route [89] and Stanford taking the lead (Figure 2.2), arriving 10 minutes before the CMU team which ranked second. In 2007, DARPA organized the last race of this series, the Darpa Urban Challenge [90], where also international participants were allowed. In contrast to the previous challenges, this competition required vehicles to drive a 96 km route through a mock-up town at George Air Force Base while obeying traffic laws, avoiding obstacles, negotiating with other vehicles, and merging into traffic. This time, the CMU team finished first, followed by the Stanford team, which ranked second. Notably, most of the successful teams relied heavily on the emerging multi-beam LiDAR technology developed in a pioneering effort by Velodyne². This spinning multi-beam LiDAR scanner allowed for obtaining precise depth measurements with a 360-degree field-of-view around the vehicle, which turned out crucial for navigating urban environments.

In 2009, Google took the lead and hired a range of star scientists who had participated in the Darpa Challenges (including Sebastian Thrun, Chris Urmson, and Mike Montemerlo). They started their own self-driving car

²<https://www.velodynelidar.com/>

program which included the development of a new driving platform and a custom, affordable multi-beam LiDAR scanner. According to accident reports [473], Google’s self-driving cars were involved in 14 collisions, while 13 were caused by others until 2016.

In 2010, the VisLab team led by Alberto Broggi at the University of Parma in Italy conducted the VisLab Intercontinental Autonomous Challenge (VIAC)[79]. Based on the experience with various prototype vehicles [75, 67, 253], VIAC [47] was an effort to drive semi-autonomously from Parma in Italy to Shanghai in China. In this demonstration, a second vehicle automatically followed a route defined by a manually driven lead vehicle either visually or based on GPS waypoints sent by the lead vehicle. The onboard system allowed for detecting obstacles, lane marking, ditches, berms, and to identify the presence and position of the preceding vehicle.

In the same year, Audi demonstrated a self-driven car ride to the summit of Pikes Peak at 4300 meters above sea level and the Technical University of Braunschweig showcased their Stadtpilot[70] which was able to navigate in a small geofenced innercity area based on LiDAR, cameras, and HD maps. In 2015, the VisLab team conducted the PROUD project [77], a demonstration of inner-city and freeway driving in Parma.

In 2011, TNO organized the Grand Cooperative Driving Challenge [385], a competition focusing on autonomous cooperative driving behavior. It was held in Helmond, Netherlands in 2011 for the first time and in 2016 for the second edition. During the competition, the semi-autonomous vehicles had to negotiate convoys, join convoys, and lead convoys. While longitudinal control was autonomous, lateral control was provided by a human safety driver. The winner (team KIT in 2011 [235] and team Halmstad in 2016) was selected based on a system that assigned points to randomly mixed teams.

In 2012, the KITTI Vision Benchmark³ [238, 237] was released. For the first time, researchers around the globe were able to evaluate their progress on various self-driving perception tasks (including reconstruction, motion estimation, and object recognition) in a fair and objective manner. At the same time, deep learning started to revolutionize many fields, including computer vision and robotics, which laid the foundations for significant improvements in particular in terms of accuracy, robustness, and run-time of the perception components of self-driving vehicles.

In 2013, Mercedes Benz demonstrated the S500 Intelligent Drive, a 103 km autonomous ride on the historic Bertha Benz route from Mannheim to Pforzheim in Germany. The system was developed by Daimler research in collaboration with the Karlsruhe Institute of Technology (KIT) [798]. The Mercedes S500 vehicle was equipped with close-to-production sensor hardware. Object detection and free-space analysis were performed using radar

³<http://www.cvlibs.net/datasets/kitti/>

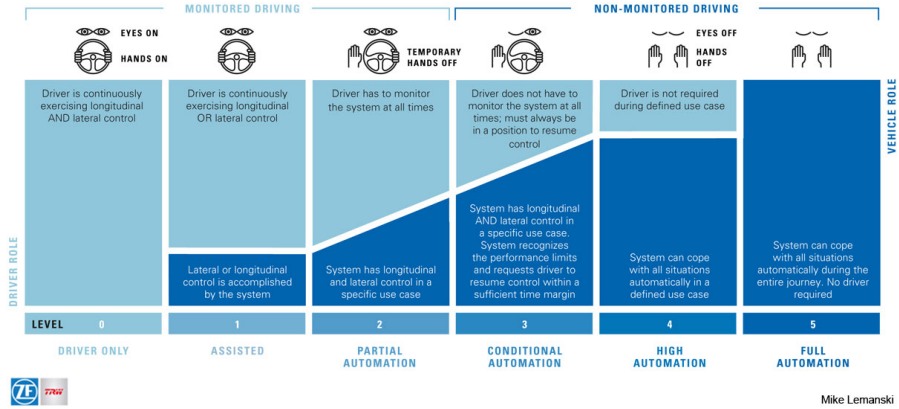


Figure 2.4: **SAE Levels of Autonomy.** Figure courtesy of Hawes [276]
 © 2016 University of Birmingham.

and stereo vision. Monocular vision was used for traffic light detection and object classification. Two complementary vision algorithms, point-feature-based and lane-marking-based, allowed for centimeter-accurate localization relative to manually annotated HD maps. While focusing on a single route, the effort demonstrated that autonomous driving in complex inner-city environments based on close-to-production hardware and HD maps is feasible.

The EU funded collaborative project V-Charge [216] conducted by Volkswagen, Bosch, and several academic partners (ETHZ, Oxford, Parma, Braunschweig) aimed for fully autonomous charging and parking of electric vehicles. In the context of this project, a fully operational system has been demonstrated which included vision-only localization, mapping, navigation and control. The project supported many publications on different problems such as calibration [290, 289], stereo [267], reconstruction [264, 263, 262], SLAM [252] and free space detection [268].

In 2014, the society of automotive engineers released their classification of autonomous driving systems into 6 SAE levels of autonomy, ranging from level 0 (no autonomy) to level 5 (full autonomy), illustrated in Figure 2.4. In the same year, Mercedes released its S Class and Tesla its Autopilot [648] with level 2 autonomy (the driver has to monitor the system at all times), providing autonomous steering, lane keeping, acceleration, and braking on the highway. One year later, ride-hailing company Uber launched its own self-driving effort [663], hiring a large number of robotics researchers from CMU. From October 2016, all vehicles produced by Tesla are equipped with eight cameras, twelve ultrasonic sensors, and a forward-facing radar with the goal of enabling full self-driving in the future. However, both Uber and Tesla

witnessed fatal accidents in which neither the driver was attentive, nor the self-driving system was functioning properly.

In 2016, after completing over 1,5 million miles, Google's self-driving efforts became Waymo, a stand-alone subsidiary of Alphabet Inc. Today, Waymo offers 400 citizens of Phoenix access to its early rider program [695] which features full self-driving in several geo-fenced districts of Phoenix (Figure 2.3) with a safety driver on the back seat.

In the same year, NVIDIA [60] demonstrated a 98% autonomous ride from Holmdel to Atlantic Highlands in Monmouth County NJ using a single convolutional neural network. The network was trained via imitation learning to predict vehicle control directly from input images. In 2018, several last-mile delivery projects were launched, including Nuro [460], a project founded by two former Google self-driving car engineers and Scout [594], a fully-electric delivery system designed to safely get packages to Amazon customers using autonomous delivery devices. In 2019, Bosch and Daimler announced a fleet of autonomous cars, providing customers a shuttle service with automated vehicles on selected routes in California [2].

Chapter 3

Sensors

The navigation of autonomous systems is usually addressed with a sensor suite which comprises various different types of sensors, including cameras, wheel odometry, and range sensors (SONAR, RADAR, and LiDAR). As an example, Tesla uses several cameras, RADAR, and ultrasonics, as illustrated in Figure 3.1 for their advanced driver-assistance system Autopilot. Fusing information from several sensors allows exploiting their complementary characteristics and addressing the limitations of individual sensors, e.g., the loss of structure information in cameras or missing color information in range data.

Wheel odometry measures the rotation of a wheel and can be used to estimate the distance covered by the autonomous vehicle. However, wheel odometry does not provide the full vehicle pose (i.e., all six degrees of freedom) and is thus typically combined with visual odometry or SLAM techniques discussed in Chapter 13. Range sensors, i.e., SONAR, RADAR, LiDAR, provide additional information about the geometry and structure of the scene. Ultrasonic sensors (SONAR) emit high-frequency sound waves and measure the time for sound waves to travel to nearby objects. The distance to objects is computed from the travel time since the speed of sound waves is known. RADAR and LiDAR work with the same principle but use electromagnetic waves and laser light pulses instead of sound waves. Because of the larger wavelength, RADAR sensors benefit from a larger working distance than LiDAR and SONAR but at the price of lower accuracy.

As cameras are cheap, passive, and easy to deploy, they are an attractive sensor choice for self-driving cars, and several existing driver assistance systems rely on cameras for lane keeping or pedestrian detection. We now briefly discuss the most dominant camera types and give a short overview of popular calibration pipelines for estimating intrinsic and extrinsic sensor parameters.

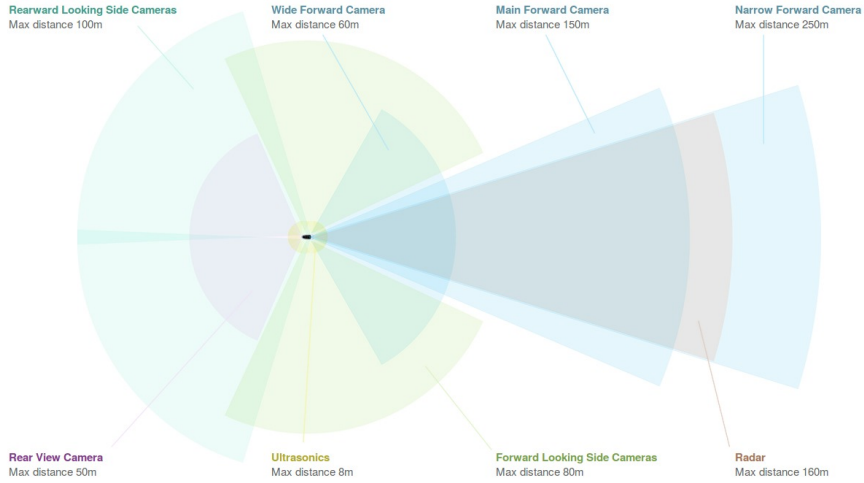


Figure 3.1: **Sensor Suite of Tesla Autopilot.** The sensor suite of Tesla's Autopilot consisting of cameras, RADAR, and ultrasonics. Figure courtesy of <https://www.tesla.com> © 2019 Tesla.

3.1 Camera Models

Most conventional cameras comprise an aperture and one or multiple lenses and can be well approximated by the pinhole camera model (Figure 3.2). Omnidirectional cameras allow to significantly increase the field of view by exploiting mirrors or special lenses. Event cameras enable the acquisition of intensity changes at very high temporal resolutions. In the following, we provide a brief overview of omnidirectional and event cameras. We refer the reader to [640, 274] for an in-depth discussion of the pinhole camera model and projective geometry.

3.1.1 Omnidirectional Cameras

A panoramic field of view is desirable in autonomous driving to gain maximum information about the surrounding area for safe navigation. Omnidirectional cameras with a 360-degree field of view (see Figure 3.3) provides enhanced coverage by eliminating the need for more cameras or mechanically turnable cameras. There are different types of omnidirectional cameras. Catadioptric cameras combine a standard camera with a shaped mirror, such as a parabolic, hyperbolic, or elliptical mirror, while dioptric cameras use purely dioptric fisheye lenses. Polydioptric cameras use multiple cameras with overlapping

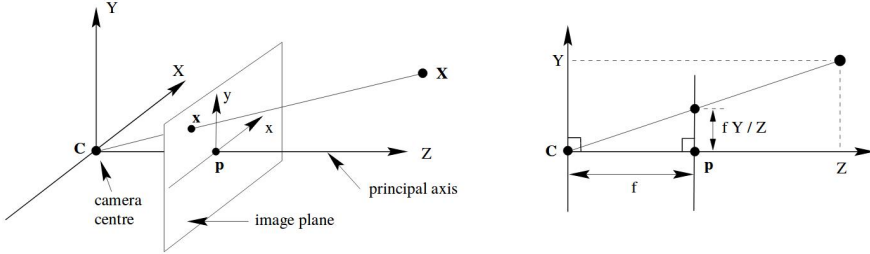


Figure 3.2: **Pinhole Camera Model.** In the pinhole model the three-dimensional world coordinates (X, Y, Z) are mapped to a two-dimensional image plane (x, y) using a perspective projection defined by the principle point (p) and focal length (f) . Figure courtesy of Hartley and Zisserman [274] © 2004 Cambridge University Press.

field of view to provide a full spherical field of view.

Geyer and Daniilidis [243] provide a unifying theory for all central catadioptric systems which is known as unified projection model in the literature and widely used by different calibration toolboxes [453, 290, 289]. Scaramuzza and Martinelli [579] propose to model the imaging function using the Taylor series expansion. Mei and Rives [453] improve upon the unified projection model of [243] to account for real-world errors by modeling distortions. Schönbein et al. [587] propose a fast approximation to computationally expensive non-central camera models.

Omnidirectional cameras are gaining popularity in autonomous driving research. For feature-based applications such as navigation, motion estimation, and mapping, a large field of view enables the extraction and matching of interest points from all around the car. Thus, omnidirectional cameras have been successfully used to improve ego-motion estimation of vehicles Scaramuzza and Siegwart [576] and 3D reconstruction of static scenes Schönbein and Geiger [586] and Häne et al. [267].

3.1.2 Event Cameras

Contrary to conventional frame-based cameras, event cameras produce a stream of asynchronous events of brightness changes surpassing a pre-defined threshold at microsecond resolution, as illustrated in Figure 3.4. An event comprises the location, sign, and timestamp of the change. As events are sparse in both space and time, this representation has the potential to reduce transmission and processing demands. The high temporal resolution enables the development of highly reactive systems.

Dynamic and Active-Pixel Vision Sensors (DAVIS) output both CMOS



Figure 3.3: **Omnidirectional Cameras.** Equirectangular (a) and spherical (b) view of a panorama from Google Street View [468].

images at fixed frame rates as well as asynchronous events, hence combining the benefits of both sensors. Mueggler et al. [477] provide a collection of real and synthetic datasets captured with DAVIS to push research on event-based methods. Binas et al. [53] present the DAVIS Driving Dataset and demonstrate end-to-end learning of steering angles. Recent work exploits DAVIS for feature tracking [232] and SLAM [676], improving accuracy and robustness over using only a single modality.

Several methods have been developed which exploit the high temporal resolution and the asynchronous nature of event sensor for various problems. The majority of these methods focus on the application in unmanned aerial vehicles (UAVs) since very efficient methods are necessary to navigate these systems. In this context, event-based cameras have been used for ego-motion estimation Mueggler et al. [476], simultaneous localization and mapping (SLAM) Rebecq et al. [539] as well as for finding feature correspondences Gallego et al. [222]. More recently, the benefits of event-based sensors have been exploited for autonomous vehicles by learning steering angles end-to-end Maqueda et al. [442].

3.2 Calibration

Geometric calibration is the problem of estimating intrinsic and extrinsic parameters of one or multiple sensors in order to accurately relate 3D world points to 2D measurements. Fiducial markers and checkerboards are often used to facilitate parameter estimation [772, 62, 334, 10, 239].

Various methods for camera calibration can be found since the beginning of the 1970s. Heikkila and Silven [287] were the first to consider the entire calibration pipeline, including control point extraction, model fitting, and image

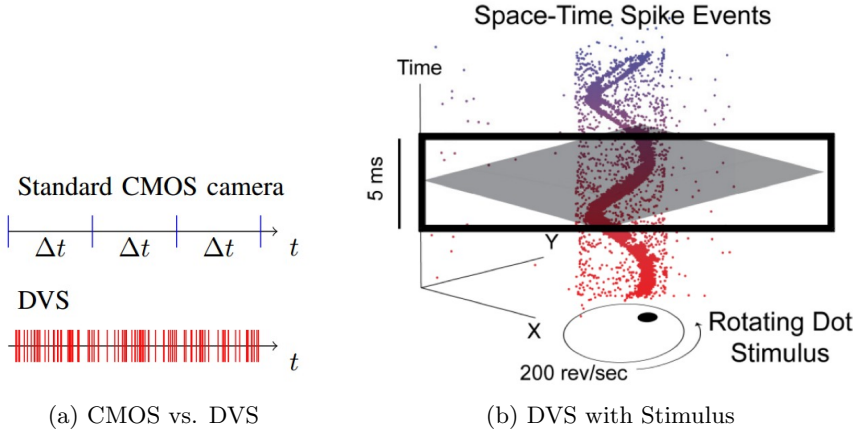


Figure 3.4: **Event Cameras.** (a) A standard CMOS camera sends images at a fixed frame rate (blue) while a Dynamic Vision Sensor (DVS) sends spike events at the time they occur (red). Each event corresponds to a local, pixel-level change of brightness. (b) Visualization of the output of a DVS looking at a rotating dot. Colored dots mark individual events. Events that are not part of the spiral are caused by sensor noise. Figure courtesy of Mueggler et al. [476] © 2015 Wiley-Blackwell.

correction. They proposed a four-step procedure to obtain the parameters of a physical camera model and address the problem of compensating image distortions.

Modern vehicles are typically equipped with multiple different sensors with the goal of increasing robustness and coverage. Several calibration procedures have been proposed to address the needs of such big sensor suites. While early approaches [772, 62] rely on manual extraction of interest points in laser scans, Kassir and Peynot [334] and Andreasson and Lilienthal [10] propose the first complete automatic camera-to-range calibration systems. Geiger et al. [239] demonstrate how to automatically calibrate a setup involving two cameras and a single range sensor such as Kinect or Velodyne laser scanner. Heng et al. [290] tackle the problem of estimating the intrinsic and extrinsic parameters of a multi-camera rig without overlapping field of view. Heng et al. [289] extend this work by removing the requirement to modify the environment by using a map and natural features instead of fiducial markings.

Chapter 4

Datasets & Benchmarks

Datasets have played a key role in the progress of many research fields by providing problem-specific examples with ground truth. Quantitative evaluations of different approaches provide key insights about their capacities and limitations. Landmark examples in the field of computer vision include the Middlebury benchmarks for stereo and optical flow [581] and the PASCAL VOC object recognition challenges [194]. In particular, many of these datasets [581, 194, 29, 238, 92, 389, 133, 357, 591] also provide online evaluation servers that allow for a fair comparison on held-out test sets and provide researchers in the field an up-to-date overview over the state of the art. This way, current progress and remaining challenges can be easily identified by the research community.

In the context of autonomous vehicles, [238, 133, 357, 487, 8, 172, 389] have introduced challenging benchmarks for reconstruction, motion estimation, recognition tasks, and tracking, and contributed to closing the gap between laboratory settings and challenging real-world situations. Kang et al. [332] provide a detailed overview of different datasets and testing environments in the context of autonomous driving.

Only a few years ago, datasets with a few hundred annotated examples were considered sufficient for many problems. The introduction of datasets with many hundred to thousands of labeled examples has led to spectacular breakthroughs in many computer vision disciplines by allowing to train high-capacity deep models in a supervised fashion. However, collecting a large amount of annotated data is not an easy endeavor, in particular for tasks such as optical flow or semantic segmentation where pixel-level annotations are required. For optical flow, Scharstein and Szeliski [581] and Baker et al. [29] acquire dense pixel-level annotations in a controlled lab environment using a time-consuming procedure whereas Geiger et al. [238] and Kondermann et al. [357] are only able to provide sparse pixel-level annotations of real street

Dataset	Realism	Diversity	Autonomous Driving	Evaluation Server	Stereo	Reconstruction	Optical Flow	Object Detection	Traffic Sign Detection	Semantic Segmentation	Road Detection	Lane Detection	Tracking
Middlebury [581]	+	-		✓	XS	XS	XS						
EPFL Multi-View [627]	++	+		✓		XS							
DTU MVS [319]	+	-				S							
ETH3D [591]	++	+		✓	S	S							
Tanks and Temples [351]	++	+		✓		S							
SlowFlow [316]	++	++					S						
HCI Benchmark [357]	++	+	✓	✓			M						
MPI Sintel [92]	O	+		✓	M		M						
Flying Chairs [174]	-	-					L						
Flying Things [450]	-	O	(✓)		L		L						
ImageNet [160]	++	++						XL		XL			
PASCAL VOC [194]	++	++						XL		XL			
Microsoft Coco [420]	++	++						XL		XL			
Cityscapes [133]	++	+	✓	✓				L		L			
EuroCity Persons Dataset [68]	++	++	✓	✓				L					
Mapillary [487]	++	++	✓	✓						L			
ApolloScape [307]	++	+	✓	✓				L		XL		XL	XL
NuScenes [93]	++	+	✓	✓				XL		XL			
Berkeley DeepDrive [755]	++	+	✓	✓				XL		XL	XL	XL	
German Traffic Sign Recognition Benchmark [623]	++	+	✓	✓				XL	L	XL	XL	XL	
German Traffic Sign Detection Benchmark [299]	++	+	✓	✓				XL	M	XL	XL	XL	
Tsinghua-Tencent 100K [793]	++	+	✓	✓				XL	XL	XL	XL	XL	
SYNTHIA [558]	O	+	✓							XL			
Playing for Data [551]	+	+	✓							L			
Playing for Benchmarks [550]	+	+	✓	✓			XL	XL		XL			XL
Caltech Lanes Dataset [8]	++	+	✓									M	
VPGNet Dataset [394]	++	+	✓									L	
MOTChallenge [389]	++	+		✓									M
Caltech Pedestrian Detection [172]	++	+	✓										XL
KITTI [238]	++	+	✓	✓	S	S	S	M		S	S	S	M
VirtualKITTI [221]	O	+	✓	✓	S	S	L	L		L			L

Table 4.1: **Popular Datasets in Computer Vision and Self-Driving.** Overview of popular datasets for Stereo, Reconstruction, Optical Flow, Object Detection, Traffic Sign Detection, Semantic Segmentation, Road Detection, Lane Detection, Tracking. Datasets specific to the autonomous driving scenario are marked with a checkmark in the corresponding column. The size of extra small datasets (XS) are in the order of tens examples/scenes for training, small sized (S) in the order of hundreds, medium sized (M) in the order of thousands, large (L) and extra large (XL) sized datasets in the order of 10 and >100 thousands, respectively. We (subjectively) rate realism and diversity with {--, -, O, +, ++} from low to high.

scenes using a LiDAR laser scanner. Janai et al. [316] pursued a different approach to obtain dense pixel-level annotations in arbitrary real scenes by using a high-speed camera to solve the optical flow problem in a simpler setting. Recently, crowdsourcing with Amazon’s Mechanical Turk platform¹ has been popularized for annotating large scale datasets, e.g., [160, 420, 390, 463, 172]. However, the annotation quality obtained via Mechanical Turk is often not sufficient and significant efforts in post-processing and clean-up are typically required.

An alternative to manual annotation is offered by modern computer graphic techniques which allow generating large-scale synthetic datasets with pixel-level ground truth. However, the creation of photorealistic virtual worlds is time-consuming and expensive. Nevertheless, the popularity of movies and video games has led to an industry creating very realistic 3D content which nourishes the hope to replace real data completely using synthetic datasets. Consequently, several synthetic datasets [92, 174, 450, 221, 558] have been proposed and are being used by AI researchers. It remains an open question, however, whether the realism and variety attained will be sufficient to replace real-world datasets and if models trained on synthetic data will be able to generalize to real-world inputs. Challenges include complex object shape and appearances as well as adversarial environmental conditions such as direct lighting, reflections from specular surfaces, fog, or rain.

Studying the performance of a system over time, e.g., in case of environmental changes or rare situations, is another important aspect for autonomous vehicles. In Section 4.2.6 we discuss several recent datasets for long-term autonomy. While most of these datasets focus on environmental changes, it is more difficult to capture rare situations which can only be captured with a large fleet of vehicles that log these situations in real-world driving. A notable exception is the Tesla Shadow Mode [647] of the Autopilot system which is a dormant logging-only mode that allows validating the Autopilot system running in the background in real and particularly rare situations.

In the following, we will first introduce the most popular computer vision datasets and benchmarks addressing tasks relevant to autonomous vehicles. Thereafter, in Section 4.2, we will focus on datasets particularly dedicated to autonomous vehicles. We also provide a detailed overview of the most popular datasets in computer vision in Table 4.1 and discuss them in the following.

4.1 Computer Vision Datasets

In this section, we introduce the most popular computer vision datasets and benchmarks relevant to autonomous driving tasks. In particular, we discuss

¹<https://www.mturk.com/mturk/welcome>



Figure 4.1: **MS COCO Object Recognition Dataset.** Examples from the MS COCO [420] object detection task. Figure courtesy of www.cocodataset.org © 2015 COCO Consortium.

datasets for object recognition and tracking, stereo and 3D reconstruction, and optical flow estimation.

4.1.1 Object Recognition

The availability of large-scale, publicly available datasets such as ImageNet [160], PASCAL VOC [194] and Microsoft COCO [420] propelled the development of novel computer vision algorithms, in particular, deep learning techniques, for recognition tasks such as object classification, detection, and semantic segmentation.

The EU funded PASCAL Visual Object Classes (VOC) challenge² by Everingham et al. [194] is a benchmark for object classification, object detection, object segmentation, and action recognition. It consists of challenging consumer photographs collected from Flickr with high-quality annotations and contains a large variability in pose, illumination, and occlusion. Since its introduction, the VOC challenge has become one of the most popular testbeds for benchmarking recognition algorithms. It has been regularly adapted to the needs of the community until the end of the PASCAL program in 2012. Over the years, the benchmark grew in size, reaching a total of 11,530 images with 27,450 annotated objects in 2012.

In 2014, Lin et al. [420] introduced the Microsoft COCO dataset³ (Figure 4.1) for object detection, instance segmentation, and contextual reasoning. They provide images of complex everyday scenes containing common objects in their natural context. The dataset comprises 91 object classes, 2.5 million annotated instances, and 328k images in total. Microsoft COCO is significantly larger in the number of instances per class than the PASCAL VOC object segmentation benchmark. All objects have been annotated with per-instance segmentations.

²<http://host.robots.ox.ac.uk/pascal/VOC/>

³<http://mscoco.org/>

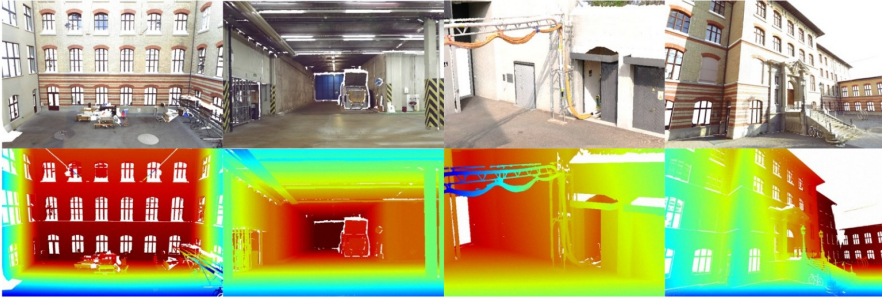


Figure 4.2: **ETH3D Reconstruction Dataset.** Examples from the ETH3D [591] dataset. Colored 3D point cloud renderings in the upper row and depth in the lower row. Figure courtesy of www.eth3d.net.

ImageNet [160], PASCAL VOC [194] and Microsoft COCO [420] are to date the largest and most diverse datasets for object classification, detection, and segmentation (Table 4.1).

4.1.2 Object Tracking

For tracking multiple objects, the first centralized benchmark, MOTChallenge⁴, was introduced by Leal-Taixé et al. [390] and Milan et al. [463]. The benchmark contains 14 challenging video sequences in unconstrained environments filmed with static and moving cameras. MOTChallenge combines several existing multi-object tracking benchmarks such as PETS [203] and KITTI [238]. Public detections provided by the benchmark allow analyzing the performance of tracking systems independent of the detector.

4.1.3 Stereo and 3D Reconstruction

For stereo vision and multi-view reconstruction, there are several publicly available datasets. The Middlebury stereo benchmark⁵ introduced by [581, 582, 580] was proposed with the goal of providing a unified testbed for a fair comparison of stereo matching algorithms. An evaluation server was created, allowing for a direct comparison of the latest approaches. The success of the Middlebury stereo benchmark in fostering research in binocular vision motivated Seitz et al. [596] to create the Middlebury multi-view stereo (MVS) benchmark⁶. The dataset consists of calibrated high-resolution multi-view

⁴<https://motchallenge.net/>

⁵<http://vision.middlebury.edu/stereo/>

⁶<http://vision.middlebury.edu/mview/>

images with registered 3D ground truth models and played a key role in advancing research in MVS.

However, the Middlebury datasets lack in size and diversity in comparison to other datasets for stereo and reconstruction (Table 4.1). The DTU MVS dataset⁷ by Jensen et al. [319] provides 124 different scenes which were recorded in a controlled laboratory environment. Reference data is obtained by combining structured light scans from different camera positions. While the DTU MVS dataset is more diverse than Middlebury in terms of the number of objects used as well as their complexity, neither of these two datasets exhibits the full spectrum of complexities of real-world scenes.

With the goal of moving multi-view stereo out of the laboratory, Strecha et al. [627] presented the EPFL Multi-View dataset⁸, which comprises images and LiDAR scans of 5 different buildings as well as a fountain.

Recently, Schöps et al. [591] published the ETH3D⁹ dataset (Figure 4.2) providing high-resolution DSLR imagery as well as synchronized low-resolution stereo videos for a variety of indoor and outdoor scenes. They used a high-precision laser scanner as [627] and registered all images using a robust optimization technique.

Similarly, Tanks and Temples¹⁰ presented by Knapitsch et al. [351] used a high-precision laser scanner and two high-resolution cameras (one with global and the other with rolling shutter) to create a novel dataset of outdoor and indoor scenes. The dataset consists of 14 scenes comprising sculptures, large vehicles, house-scale buildings as well as large indoor and outdoor scenes.

For large-scale reconstruction, multiple Internet photo collections have been proposed over time. The most popular collections are combined in the BigSFM dataset¹¹ and comprise Vienna [313], Dubrovnik [411], and Rome [139]. While Dubrovnik and Rome were retrieved from Flickr, Vienna was recorded with a calibrated camera. Besides large-scale reconstruction, these datasets are also frequently used for evaluating loop-closure detection (Section 13.4.2) and localization methods (Section 13.3).

4.1.4 Optical Flow

Similar to stereo vision, the Middlebury flow benchmark¹² by Baker et al. [29] provided the first unified test environment and evaluation server for optical flow approaches. The benchmark comprises sequences with non-rigid motion, synthetic sequences, and a subset of the Middlebury stereo benchmark (static

⁷http://roboimagedata.compute.dtu.dk/?page_id=36

⁸<https://www.epfl.ch/labs/cvlab/data/data-strechamvs/>

⁹<https://www.eth3d.net>

¹⁰<https://www.tanksandtemples.org>

¹¹<http://www.cs.cornell.edu/projects/bigsfm/>

¹²<http://vision.middlebury.edu/flow/>

scenes). For all non-rigid sequences, ground truth flow is obtained by tracking hidden fluorescent textures sprayed onto the objects. In comparison to other optical flow datasets (Table 4.1), the Middlebury flow dataset is limited in size and missing real-world challenges like complex structures, lighting variation, and shadows due to the laboratory conditions in which it has been recorded. In addition, Middlebury only contains small motions of up to twelve pixels which do not allow the investigation of challenges related to fast motions.

The acquisition of optical flow ground truth is very difficult since no sensor exists that can capture optical flow ground-truth in general natural scenes. While [238, 357] use a LiDAR laser scanner for this purpose, they only obtain sparse pixel-level annotations and are restricted to static scenes (only camera motion). Janai et al. [316] present a novel approach to obtain accurate reference data from High-Speed video cameras by tracking pixels through densely sampled space-time volumes. This method allows the acquisition of optical flow ground truth in challenging everyday scenes and the data augmentation with realistic effects such as motion blur to compare methods in varying conditions. Janai et al. [316] provide 160 diverse real-world sequences of dynamic scenes with a significantly larger resolution (1280×1024 pixels) than previous optical datasets.

The problem of acquiring optical flow ground truth can also be resolved by creating synthetic datasets. Towards this goal, Butler et al. [92] take advantage of the open-source movie Sintel, a short animated film. They create the MPI Sintel optical flow benchmark¹³ by rendering scenes with optical flow ground truth. Sintel consists of 1,628 frames and provides three different datasets with varying complexity that are obtained using different passes of the rendering pipeline. Similar to Middlebury, they provide an evaluation server for comparison.

The limited size of optical flow datasets hampered the training of deep high-capacity models. Thus, Dosovitskiy et al. [174] introduced a simple synthetic 2D dataset of flying 3D chairs rendered on top of random background images from Flickr to train a convolutional neural network. As the limited realism of this dataset proved insufficient to learn highly accurate models, Mayer et al. [450] presented another large-scale dataset consisting of three synthetic stereo video datasets with optical flow ground truth: FlyingThings3D, Monkaa, Driving. FlyingThings3D provides everyday 3D objects flying along randomized 3D trajectories in a randomly created scene. Inspired by the KITTI dataset, a driving dataset has been created which uses car models from the same pool as FlyingThings3D and additionally highly detailed tree and building models from 3D Warehouse. Monkaa is an animated short movie similar to Sintel used in the MPI Sintel benchmark.

While synthetic optical flow datasets provide numerous examples for train-

¹³<http://sintel.is.tue.mpg.de/>

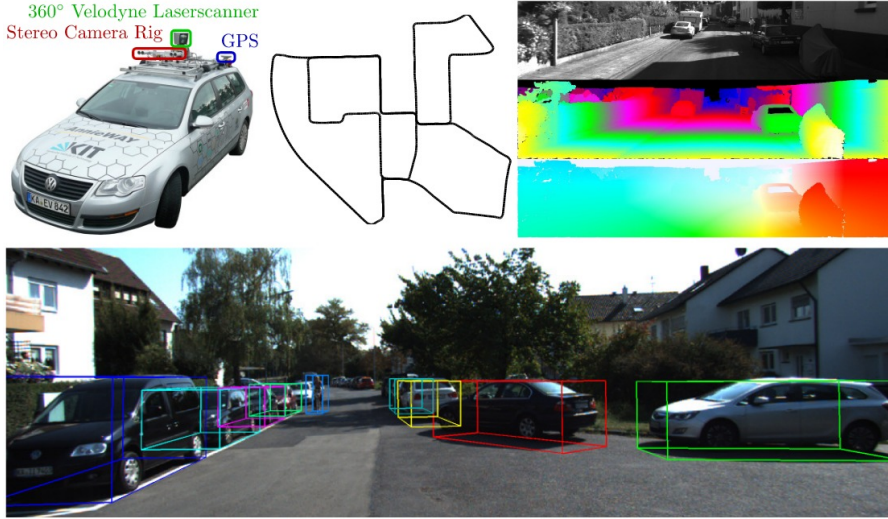


Figure 4.3: **KITTI Dataset.** The recording platform with sensors (top left), trajectory (top center), disparity and optical flow (top right) and 3D object labels (bottom) from the KITTI benchmark proposed by Geiger et al. [238]. Figure courtesy of Geiger et al. [238] © 2012 IEEE.

ing deep neural networks, they lack realism and are limited in diversity, as indicated in Table 4.1. Therefore, large-scale synthetic datasets are typically used for pre-training, and, afterwards, the pre-trained models are fine-tuned on small, more realistic datasets.

4.2 Autonomous Driving Datasets

Several datasets have been proposed to specifically address the problem of autonomous driving. The KITTI Vision Benchmark¹⁴ introduced by Geiger et al. [238, 237] was the first publicly available benchmark for stereo, optical flow, visual odometry/SLAM, and 3D object detection (Figure 4.3) in the autonomous driving context. The dataset has been captured from an autonomous driving platform equipped with high-resolution color and grayscale stereo cameras, a Velodyne 3D laser scanner, and high-precision GPS/IMU inertial navigation system.

Due to the limitations of the rotating laser scanner used as reference sensor, the stereo and optical flow benchmark were restricted to static scenes with camera motion. In the 2015 version of the optical flow and stereo Benchmark,

¹⁴<http://www.cvlibs.net/datasets/kitti/>



Figure 4.4: **Mapillary Vistas Dataset.** Examples colorized according to the class definition of the Mapillary Vistas Dataset proposed by Neuhold et al. [487]. Figure courtesy of Neuhold et al. [487] © 2017 IEEE.

Menze and Geiger [455] provide ground truth for dynamic scenes by fitting 3D CAD models to all vehicles in motion. This new version of KITTI also combined the stereo and flow ground truth to form a novel 3D scene flow benchmark. For the KITTI object detection challenge, a special 3D labeling tool has been developed to annotate all 3D objects with 3D bounding boxes in 7481 training and 7518 test images. The benchmark for object detection was separated into a vehicle, pedestrian and cyclist detection tasks, allowing to focus the analysis on the most important problems in the context of autonomous vehicles. The visual odometry / SLAM challenge consists of 22 stereo sequences, with a total length of 39.2 km. The ground truth pose is obtained by using GPS/IMU localization unit which was fed with RTK correction signals.

The KITTI dataset has established itself as one of the standard benchmarks in all of the aforementioned tasks, in particular in the context of autonomous driving applications. While KITTI provides annotated data and an evaluation server for all problems considered in this work (Table 4.1), it is still comparably limited in size. Therefore, the KITTI dataset is usually used mostly for evaluation and fine-tuning.

4.2.1 Object Detection and Semantic Segmentation

The Cityscapes Dataset¹⁵ by Cordts et al. [133] provides a benchmark and large-scale dataset for pixel-level and instance-level semantic labeling that captures the complexity of real-world urban scenes. High-quality pixel-level annotations are provided for 5,000 images, while 20,000 additional images have been annotated with coarse labels obtained using crowdsourcing. While Cityscapes provides an evaluation server for a fair comparison of methods, the dataset is limited in size and diversity.

For object detection, Braun et al. [68] presented a large-scale dataset recorded in 31 cities of 12 European countries. Similar to Cityscapes, an evaluation server allows a fair comparison of methods. However, they only

¹⁵<https://www.cityscapes-dataset.com/>

provide bounding box, occlusion, and orientation annotations for pedestrians, cyclists, and other riders in urban traffic.

The crowdsourcing company Mapillary¹⁶ has collected 282 million street-level images covering 4.5 million road kilometers around the world. Based on this data, the Mapillary Vistas Dataset¹⁷ [487] has been created and shared with the community, providing 25,000 high-resolution images with dense annotation for 66 object categories and instance-specific labels for 37 classes (Figure 4.4).

The Berkeley DeepDrive dataset¹⁸ [755] for object detection, instance segmentation, road, and lane detection provides 100K partially annotated driving videos from New York, Berkeley, San Francisco, and the Bay Area. The dataset is more diverse in scenes and weather conditions than Cityscapes, but it is still limited in the number of cities used for the recording. In this context, the Mapillary Vistas Dataset is the most diverse autonomous driving-related dataset for semantic segmentation and object recognition (Table 4.1). However, datasets like Mapillary Vistas Dataset, ImageNet, PASCAL VOC, and Microsoft Coco are less suited for training and testing temporal coherence of methods since they provide only single images in contrast to KITTI, Cityscapes, and Berkeley DeepDrive which provide image sequences.

Very recently, major companies working on autonomous driving solutions also started making their annotated data publicly available. The autonomous driving project Apollo from Baidu created the Data Open Platform¹⁹ consisting of simulation, annotation, and demonstration data for autonomous driving. The ApolloScape dataset[307] provides annotated street view images for semantic (144K images) and instance segmentation (90K images), lane detection (160K images), car detection (70K) and tracking of traffic participants (100K images). The dataset allows evaluating the performance of methods in various weather conditions and at different day times.

The company Nutonomy released the NuScenes dataset²⁰ [93], which provides data from an entire sensor suite with annotations for semantic segmentation and object detection. The dataset consists of over 1 million camera images. However, both ApolloScape and NuScenes have been recorded only in one or two cities, respectively, and are therefore still limited in diversity.

So far, datasets for 3D semantic segmentation have been limited in size [480, 40, 260, 773] and the number of classes [237] due to the large labeling effort required for annotating detailed object boundaries. Recently, Behley et al. [38] present a large dataset for 3D semantic segmentation based on the KITTI Visual Odometry Benchmark [237]. In contrast to previous annota-

¹⁶<https://www.mapillary.com/app>

¹⁷<https://www.mapillary.com/dataset/vistas>

¹⁸<https://bdd-data.berkeley.edu>

¹⁹<http://data.apollo.auto>

²⁰<https://www.nuscenes.org>

tions for KITTI, they provide dense point-wise annotations for the complete 360-degree field-of-view of the LiDAR. The dataset comprises over 20,000 scans with 25 different classes.

4.2.2 Tracking

The Caltech Pedestrian Detection Benchmark²¹ proposed by Dollar et al. [172] provides 250,000 frames of sequences recorded while driving through regular traffic in an urban environment. 350,000 bounding boxes and 2,300 unique pedestrians were annotated, including temporal correspondence between bounding boxes and detailed occlusion labels.

4.2.3 Traffic Sign Detection

While all previously discussed detection datasets focus on the detection of generic objects or traffic participants, only a few datasets exist for the recognition and detection of traffic signs. The most popular datasets for this task are the German Traffic Sign Recognition Benchmark (GTSRB²²) [623] and the German Traffic Sign Detection Benchmark (GTSDB²³) [299]. GTSRB considers the task of classifying traffic signs into their corresponding category and consists of 50,000 images. In contrast, GTSDB provides 600 training and 300 test images for the task of detecting traffic signs. Reliable ground truth annotations for 40 different classes were created using a semi-automatic annotation tool. Recently, the limits of both datasets have been reached by state-of-the-art detection systems and Zhu et al. [793] presented Tsinghua-Tencent 100K²⁴, a new traffic sign detection benchmark. In contrast to GTSDB, their benchmark consists of 100,000 images with 30,000 signs. They provide high resolution images with pixel mask annotations and bounding boxes for each traffic sign.

4.2.4 Road and Lane Detection

The KITTI benchmark was extended by Fritsch et al. [211] to the task of road/lane detection. In total, 600 diverse training and test images have been selected for manual annotation of road and lane areas. Mattyus et al. [448] used aerial images to enhance the KITTI dataset with fine-grained segmentation categories such as parking spots and sidewalk as well as the number and location of road lanes.

A larger dataset for lane detection, the Caltech Lane Detection dataset²⁵,

²¹http://www.vision.caltech.edu/Image_Datasets/CaltechPedestrians/

²²<http://benchmark.ini.rub.de/?section=gtsrb&subsection=dataset>

²³<http://benchmark.ini.rub.de/?section=gtsdb&subsection=news>

²⁴<https://cg.cs.tsinghua.edu.cn/traffic-sign/>

²⁵<http://www.mohamedaly.info/datasets/caltech-lanes>

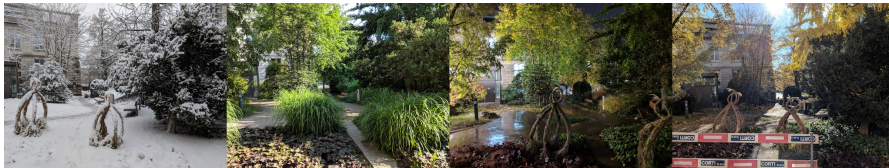


Figure 4.5: **Long-Term Autonomy.** Examples for different weather conditions, seasons and day times for a scene from the Workshop organized by Hammarstrand et al. [265]. Figure courtesy of Hammarstrand et al. [265].

has been proposed by Aly [8]. The dataset was recorded in Pasadena in California at different day times and consists of over 1200 frames. The first large-scale lane detection dataset was presented by [394] and provides over 20,000 images. In contrast to previous datasets, they also consider different weather conditions. The Berkeley DeepDrive dataset²⁶ [755] with 100,000 images is so far the largest and most diverse lane/road detection dataset.

4.2.5 Flow and Stereo

Complementary to the datasets presented in Section 4.1.4 and KITTI, the HCI benchmark²⁷ proposed by Kondermann et al. [357] includes realistic, systematically varied radiometric and geometric challenges for autonomous driving. Overall, a total of 28,504 stereo pairs with stereo and flow ground truth is provided. The major limitation of the HCI Benchmark is that all sequences were recorded in a single street section, and thus the dataset lacks diversity. However, the controlled environment allows for more easily simulating rare events such as accidents which are of great interest for validating autonomous driving systems.

4.2.6 Long-Term Autonomy

Several datasets such as KITTI or Cityscapes focus on the development of algorithmic competences for autonomous driving but do not address challenges of long-term autonomy, as for example environmental changes over time. In order to address this problem, Carlevaris-Bianco et al. [96] presented a new long-term vision and LiDAR dataset comprising 27 sessions. However, the dataset was not recorded from a vehicle but instead using a Segway robot on the campus of the University of Michigan. A novel dataset for long-term autonomous driving has been presented by Maddern et al. [438]. They collected images, LiDAR, and GPS data while traversing 1,000 km in central

²⁶<https://bdd-data.berkeley.edu>

²⁷<http://hci-benchmark.org>

Oxford, UK during an entire year. This allowed them to capture large variations in scene appearance due to illumination, weather and seasonal changes, dynamic objects, and constructions. Such long-term datasets allow for an in-depth investigation of problems that detain the realization of autonomous vehicles such as localization at different times of the year, as illustrated in Figure 4.5.

Several datasets have been proposed which address environmental changes for multi-view reconstruction. The structure-from-motion dataset BigSfM discussed in Section 4.1.3, for instance, consists of Internet photos taken with different cameras at different times. Recently, Sattler et al. [571] presented three datasets for visual localization (Aachen Day-Night, RobotCar Seasons and CMU Seasons) recorded under different weather conditions, seasons and during night and day. While the Aachen Day-Night dataset consists of images recorded using consumer cameras, RobotCar Seasons, and CMU Seasons were obtained using a car-mounted camera. More recently, Scape Technologies²⁸ presented a long-term dataset captured around the Imperial College London campus using a low-end, consumer spherical camera [30]. The dataset was recorded over a period of one year and incorporates different weather conditions, day times, and seasons.

4.3 Synthetic Data Generation using Game Engines

Data from animated movies as used in [92, 450] is very limited since the content is hard to change, and such movies are rarely open-source. Moreover, rendering 3D models into random scenes as in [174, 450] lacks realism and diversity. In contrast, game engines allow for creating an infinite amount of more realistic and diverse data.

One of the first datasets exploring game engines is the Virtual KITTI dataset²⁹ presented by Gaidon et al. [221]. They propose a real-to-virtual world cloning method to create realistic proxy worlds that resemble real scenarios. A cloned virtual world allows varying conditions such as weather or illumination and using different camera settings. This way, the proxy world can be used for virtual data augmentation to train deep networks. Virtual KITTI contains 35 photo-realistic synthetic videos with a total of 17,000 high resolution frames. They provide ground truth for object detection, tracking, scene and instance segmentation, depth, and optical flow.

In concurrent work, Ros et al. [558] created SYNTHIA³⁰, a synthetic collection of imagery and annotations of urban scenarios for semantic segmen-

²⁸<https://scape.io/>

²⁹<https://europe.naverlabs.com/research/computer-vision/proxy-virtual-worlds>

³⁰<http://synthia-dataset.net/>



Figure 4.6: **Synthetic Datasets.** Examples from Virtual KITTI [221], SYNTHIA [558], Carla [175] and Playing for Data [551]. Figure courtesy of Gaidon et al. [221], Ros et al. [558], Dosovitskiy et al. [175], and Richter et al. [551]

tation. They rendered a virtual city using the Unity Engine. The dataset consists of 13,400 randomly taken virtual images from the city and four video sequences with 200,000 frames in total. Pixel-level semantic annotations are provided for 13 classes.

In the Playing for Data project³¹, Richter et al. [551] extracted pixel-accurate semantic label maps for images from the commercial video game Grand Theft Auto V. Towards this goal, they developed a tool that operates between the game and the graphics hardware to obtain pixel-accurate object signatures across time. Their algorithm allows them to produce dense semantic annotations for 25,000 images synthesized by the photorealistic open-world computer game with minimal human supervision. This work was extended in Playing for Benchmarks³² [550] to obtain dense correspondences and semantic

³¹https://download.visinf.tu-darmstadt.de/data/from_games/

³²<https://playing-for-benchmarks.org/>

instances from the game engine. The benchmark consists of about 250,000 images with dense annotations for semantic segmentation, instance segmentation, object detection, tracking, 3D scene layout, visual odometry, and optical flow. They provide an online evaluation server for semantic segmentation, instance segmentation, visual odometry, and optical flow. Similarly, Qiu et al. [531] provide an open-source tool to create virtual worlds by accessing and modifying the internal data structure of Unreal Engine 4. They show how virtual worlds can be used to test deep learning algorithms by linking them with the deep learning framework Caffe [322].

Recently, Carla³³, an open-source simulator for autonomous driving, was introduced by Dosovitskiy et al. [175]. Carla allows generating synthetic data for control and perception of an autonomous driving system in urban environments. Complete access to the engine and digital assets are provided for non-commercial usage. Based on the Unreal Engine 4, extensions for Carla can be easily integrated by the community.

Modern game engines as used in Carla [175] and Playing for Data [551] allow creating impressively realistic data for training large models, as shown in Figure 4.6. While there is still a large gap between real and synthetic data and the creation of 3D content is costly and time-consuming, game engines enable the generation of large datasets and the investigation of dangerous situations that can only be rarely observed in real data.

³³<http://carla.org/>

Chapter 5

Object Detection

5.1 Problem definition

Reliable detection of objects, as shown in Figure 5.1, is a crucial requirement to realize autonomous driving. As the vehicle shares the road with many other traffic participants, particularly in urban areas, the awareness of other traffic participants or obstacles is necessary to avoid accidents that might be life-threatening. The detection in urban areas is hard because of the wide variety of object appearances and occlusions caused by other objects or the object of interest itself. In addition, the resemblance of objects to each other or to the background and physical effects like cast shadows or reflections can make the detection of objects difficult.

Reliable pedestrian detection is particularly difficult because of their complex, highly varying motion and the large variety of appearances due to different clothing and articulated poses. Furthermore, the interaction of pedestrians with each other and the world often cause partial occlusions. This problem has been deeply investigated as for example in advanced driver assistance systems to increase road safety. Pedestrian protection systems (PPS) detect the presence of stationary and moving people around a moving vehicle in order to warn the driver against dangerous situations. Geronimo et al. [242] survey pedestrian detection for Advanced Driver Assistance Systems. While the driver can still handle missed detections of a PPS, an autonomous car needs a flawless pedestrian detection system which is robust against all weather conditions and efficient for real-time detection.

The object detection problem has been approached using a variety of input modalities. Video cameras are the cheapest and most commonly used type of sensors for the detection of objects. The visible spectrum (VS) is typically used for daytime detections, whereas the infrared spectrum offers

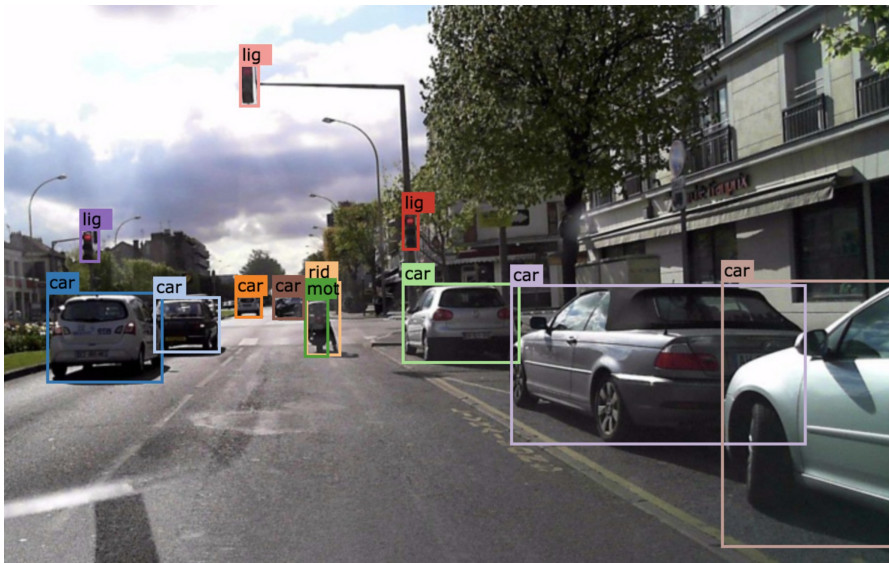


Figure 5.1: **Object Detection.** In object detection, we are interested in finding all objects of certain classes in an image. These detections are usually represented with bounding boxes. Figure courtesy of Berkeley DeepDrive [756].

more visibility for night-time detection[629]. Thermal infrared (TIR) cameras capture relative temperature, which allows distinguishing warm objects like pedestrians from cold objects like vegetation or the road. Active sensors that emit signals and observe their reflection, like laser scanners can provide range information which is helpful for detecting an object and localizing it in 3D. However, laser scanners often have a smaller resolution compared with video cameras. Depending on the weather conditions, time of day, or material properties, it can be problematic to rely on a single type of sensor alone. VS cameras and laser scanners are affected by reflective or transparent surfaces, while hot objects (like engines) or warm temperatures can influence TIR cameras. The combination of information from different sensors via sensor fusion [188, 115, 250] allows for the robust integration of this complementary information.

5.2 Methods

Classical object detection systems usually consist of multiple steps that are applied consecutively to solve the object detection task. With the success of

deep neural networks, most of these steps [602, 247, 283, 246] and even the complete pipeline have been replaced by learned models [601, 544, 540, 424, 541, 419]. We start our discussion with classical pipelines, followed by more modern approaches.

5.2.1 Classical Pipeline

A classical detection pipeline usually comprises the following steps: preprocessing, region of interest extraction (ROI), object classification, and verification or refinement. In the preprocessing step, tasks such as exposure and gain adjustment, as well as camera calibration and image rectification, are usually performed. Some approaches leverage temporal information with a joint detection and tracking system. We discuss tracking approaches in-depth in Chapter 6.

Regions of interest can be extracted using a sliding window approach, which shifts a window over the image at different scales. As exhaustive search is very expensive, several heuristics have been proposed for reducing the search space. Typically, the number of evaluations is reduced by assuming a certain ratio, size, and position of candidate bounding boxes. Apart from that, image features, stereo, or optical flow can be leveraged for focusing the search on relevant regions. Broggi et al. [76], for instance, leverage morphological characteristics (size, ratio, and shape), vertical symmetry of human shape, and distance information obtained from stereo for the extraction of relevant ROIs. Selective Search [665] is an alternative approach to generate regions of interest. Instead of an exhaustive search over the full image domain, selective search exploits a segmentation of the image to extract approximate locations efficiently. For a more detailed discussion, we refer the reader to Dollar et al. [169], presenting an extensive evaluation of pedestrian detection systems from monocular images with a focus on sliding window approaches.

The next step is the processing of candidate image regions from sliding window to verify them and classify objects. The classification of all candidates in an image can be quite costly due to the vast amount of image regions that need to be processed. Therefore, a fast decision is necessary which quickly discards candidates in the background region of the image. Viola et al. [679] combine simple and efficient classifiers, learned using AdaBoost, in a cascade that allows them to quickly discard false candidates while spending more time on promising regions. With the work of Dalal and Triggs [150], linear Support Vector Machines (SVMs) in combination with Histogram of Orientation (HOG) features have become popular tools for classification. Enzweiler and Gavrilu [186] provide an overview of classical approaches for monocular pedestrian detection. They make the observation that SVM with HOG features work well at higher resolutions while having a higher processing time than cascaded approaches that are superior at lower resolutions and achieve

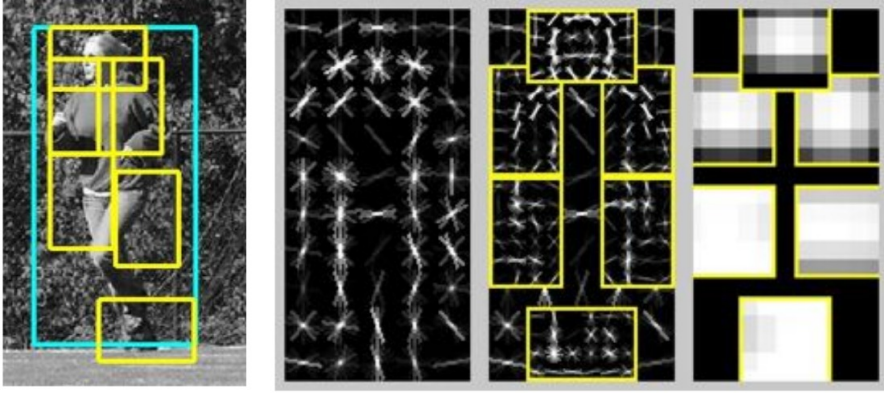


Figure 5.2: **Part-based Approaches.** Illustration of the Deformable Part Model (DPM) proposed by Felzenszwalb et al. [201]. The model consists of a coarse global template (middle-left), several high resolution part templates (middle-right) and the location (right). Figure courtesy of Felzenszwalb et al. [201] © 2008 IEEE.

near real-time performance. In their survey, Benenson et al. [43] found no clear evidence that a certain type of classifier (e.g., SVM or decision forests) is better suited than any other. In particular, Wojek and Schiele [710] show that AdaBoost and linear SVM perform roughly the same if enough features are given. Benenson et al. [43] conclude that the number and diversity of features is clearly an important factor for the performance of classifiers since the classification problem becomes easier with higher dimensional representations. Consequently, today, all state-of-the-art object detection systems use convolutional neural networks to learn expressive features in an end-to-end fashion from large datasets [94, 728, 792, 744, 114, 544, 246].

Multi-Cue Object Detection: While most object detection systems rely on single images as input, there are several approaches which show that using multiple cues such as temporal and structure information can boost performance. Temporal information from video sequences can provide important additional constraints to solve the detection task better. Shashua et al. [609] integrate additional cues measured over time (dynamic gait, motion parallax) and situation-specific features (such as leg positions at certain poses) into a detection system to obtain more reliable detections. Wojek et al. [708] show significant improvement in detection performance by incorporating motion cues and combining different complementary feature types. The dense correspondences between two frames (optical flow) [685] or joint tracking as discussed in Chapter 6 also lead to significant performance gain since more

information about the same object can be aggregated over time. Structure information can be beneficial to generate region of interests and provide additional information about the shape of objects to improve classification. Towards this goal, Keller et al. [335] jointly detect objects and estimate dense depth maps from stereo images.

Generative Models for Augmenting Training Data: As object detection is typically formulated as a supervised learning task, large amounts of annotated training data are required to obtain good performance. Unfortunately, generating examples belonging to the target class is usually time-consuming because of manual labeling, while negative examples can be more easily obtained. Enzweiler and Gavrila [187] address this bottleneck by creating synthesized virtual samples with a learned generative model. The generative model consists of probabilistic shape and texture models for a set of generic poses. As the discriminative model, they consider a neural network [706] and SVMs with Haar features [501] to demonstrate the generality of their approach. The generative model captures prior knowledge about the pedestrian class and allows significant improvement in classification performance.

5.2.2 Part-based Approaches

Learning the appearance of articulated objects is difficult because all possible articulations need to be considered. The idea of part-based approaches is to split the complex appearance of non-rigidly moving objects like humans into simpler parts and to represent articulation using these parts, as illustrated in Figure 5.2. This provides greater flexibility and reduces the number of training examples required for learning the appearance of each part.

The Deformable Part Model (DPM), by Felzenszwalb et al. [201], attempts to break down the complex appearance of objects into easier parts. As a classifier, they train a SVM with latent structure variables which represent the model configuration (part positions) and need to be inferred at training time. They use a coarse global template covering the entire object and higher resolution part templates to model the appearance of each part. An alternative to this representation is the Implicit Shape Model proposed by Leibe et al. [396], which learns a highly flexible representation of object shape. They extract local features around interest points and perform clustering to construct a codebook of local appearances that are characteristic for the particular object class under consideration. Finally, they learn the occurrences of codebook entries for each object. However, Benenson et al. [43] observe in their survey on detection approaches that part-based models like [201, 396] improve results only slightly compared to the much simpler approach of Dalal and Triggs [150].

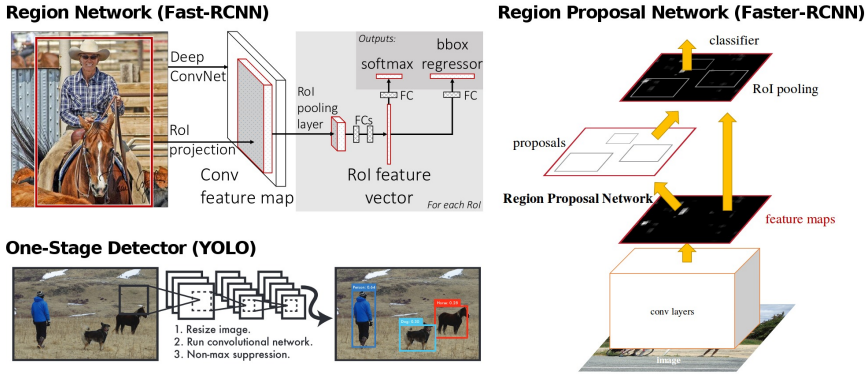


Figure 5.3: **Object Detection Networks.** Illustration of the three popular object detection networks. Upper left: Region-based network Fast-RCNN [246] that works on regions. Right: Region proposal network Faster-RCNN [544] that learn to extract regions. Lower left: One-stage detector YOLO [540] that formulates the detection task as regression problem. Figures courtesy of Girshick [246] © 2015 IEEE, Redmon et al. [540] © 2016 IEEE, and Ren et al. [544] © 2015 NeurIPS

The discussed part-based models can not represent relationships between different objects, their parts, and the scene, which, for instance, is necessary to reason about occlusions. Usually, a separate context model [296, 661, 161, 748] is learned which puts the detected objects in context to the 3D scene. In contrast, Wu et al. [720] propose to learn an And-Or model that embeds a grammar to represent large structural and appearance variations in a reconfigurable hierarchy. The learned model takes into account structural and appearance variations at multi-car, single-car, and part-levels jointly to represent both context and occlusions.

5.2.3 Deep Learning for Detection

All previous methods rely on hand-crafted features that are difficult to design and limited in their representation capabilities. With the renaissance of deep learning [360], convolutional neural networks have been applied to the object detection problem, resulting in significantly increased performance. Examples of the three most popular architectures are illustrated in Figure 5.3.

Sermanet et al. [602] introduced CNNs to the pedestrian detection problem by learning the extraction of expressive features in an unsupervised fashion using convolutional sparse auto-encoders. Eventually, they train a classifier in an end-to-end supervised fashion while extracting the features with a

sliding window scheme and jointly fine-tuning the auto-encoders. However, they use a shallow network with a small receptive field, which allows precise localization of the objects using a sliding window approach. In contrast, deeper networks with larger receptive fields complicate the precise localization because local information is extracted in earlier layers, while high-level information is represented in deeper layers. Therefore, Girshick et al. [247] propose R-CNNs to solve the CNN localization problem via a “recognition using regions” paradigm. They generate many region proposals using selective search [665], extract a fixed-length feature vector for each proposal using a CNN and classify each region with a linear SVM. Region-based CNNs are computationally expensive but several improvements have been proposed to reduce the computational burden [283, 246]. He et al. [283] use spatial pyramid pooling which allows computing a convolutional feature map for the entire image with only one run of the CNN in contrast to R-CNN that needs to be applied on many image regions. Girshick [246] (Fast-RCNN) further improve upon these results by proposing a single-stage training algorithm using a multi-task loss that jointly learns to classify object proposals and refine their spatial locations.

In region-based CNNs, the classical region proposal algorithm remained the primary computational bottleneck and the main factor limiting performance. Therefore, Ren et al. [544] (Faster-RCNN) introduced Region Proposal Networks (RPN), which share full-image convolutional features with the detection network and thus do not incur additional computational costs. RPNs are trained end-to-end to generate high-quality region proposals, which are classified using the Fast R-CNN detector [246].

Eventually, one-stage detectors [601, 540, 424, 541, 419] completely removed the region proposal step by formulating the object detection task as a regression problem. The first one-stage detector by Sermanet et al. [601] was a deep convolutional version of the sliding window approach. They extract features with a CNN and apply a classifier network based on AlexNet [360] on the extracted feature maps in a sliding window fashion. Redmon et al. [540] (YOLO) instead suggest to jointly learn spatially separated bounding boxes and class probabilities from the topmost feature maps of a network based on GoogLeNet [639]. This allows them to achieve real-time performance and eventually YOLO9000 [541] to outperform the Region Proposal Networks. Liu et al. [424] further improve in accuracy and efficiency by incorporating feature maps from different scales and considering a fixed set of bounding boxes. However, one-stage detectors [601, 540, 424, 419] could not compete with region proposal algorithms. One reason for the performance gap is the foreground-background class imbalance [419]. To alleviate this problem and improve training, Lin et al. [419] propose a dynamically scaled cross-entropy loss allowing them to reduce the contribution of easy examples.

All previous one-stage detectors use anchor bounding boxes that are densely

placed over the image, verified, and refined using regression. In contrast, Law and Deng [386] propose to directly predict heatmaps for the top-left and bottom-right corners of all bounding boxes. Finally, they need to identify corners belonging to the same bounding box. Towards this goal, they train a network to predict similar embedding vectors for corners of the same bounding box, which allows them to group the corners according to the distance between the embeddings.

Part-based models have also been introduced to CNN-based approaches. Zhang et al. [771] propose to extract deep convolutional features from bottom-up proposals obtained from a selective search algorithm and learn part appearance models. This allows them to enforce geometric constraints between parts and to outperform previous methods.

5.2.4 Real-time Pedestrian Detection

In case of a potential collision with pedestrians, a fast detection system allows the early intervention of the autonomous system. In classical literature, Benenson et al. [42] provide fast pedestrian detections based on better handling of scales and exploiting depth extracted from stereo. Instead of resizing the images, they scale HOG features similar to Viola and Jones [680]. However, CNN-based approaches recently also reached real-time efficiency due to strong parallelization on the GPU. While Fast R-CNN [246] could only be applied at 0.5 Hz, the faster version with the Region Proposal Network Faster-RCNN [544] already achieves 17 Hz. Finally, YOLO9000 [541] can be applied at up to 90 Hz at 288×288 pixels resolution and achieves 40 Hz at 544×544 pixels resolution.

5.2.5 Human Pose Estimation

The pose of a person provides important information to the autonomous vehicle about the behavior and intention of the person. However, the pose estimation problem is challenging since the pose space is very large, and typically, people can only be observed at low resolutions because of their size and distance to the vehicle. Several approaches have been proposed to jointly estimate the pose and body parts of a person. Traditionally, a two-staged approach was used by first detecting body parts and then estimating the pose as in [520, 248, 634]. This is problematic in cases when people are in proximity to each other because body-parts can be wrongly assigned to different instances.

Pishchulin et al. [519] present DeepCut, visualized in Figure 5.4, a model that jointly estimates the poses of all people in an image. The formulation is based on partitioning and labeling a set of body-part hypotheses obtained from a CNN-based part detector. The model jointly infers the number of

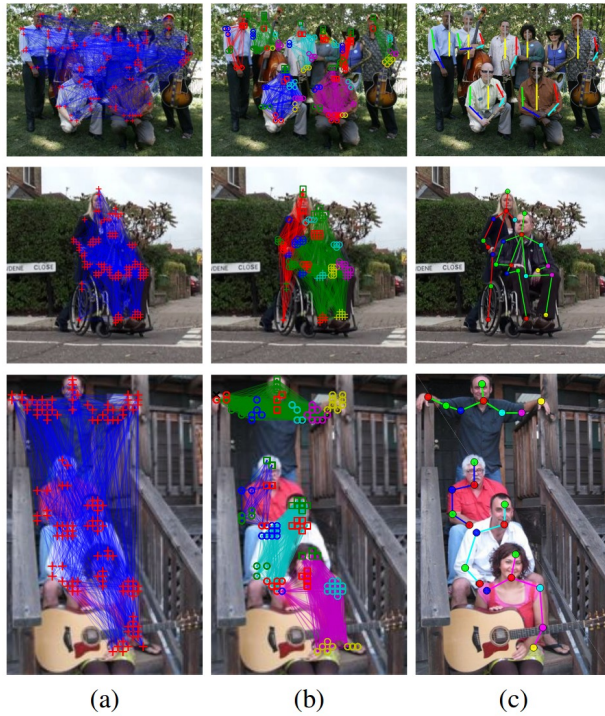


Figure 5.4: **Human Pose Estimation.** Illustration of the DeepCut [519] showing the initial detections and pairwise terms (a), joint clustering of nodes belonging to the same person visualized by colors (b) and the predicted poses (c). Figure courtesy of Pishchulin et al. [519] © 2016 IEEE.

people, their poses, spatial proximity, and part-level occlusions. Bogo et al. [59] use DeepCut to estimate the 3D pose and 3D shape of a human body from a single unconstrained image. Towards this goal, SMPL, a 3D body shape model proposed by Loper et al. [427], is fit to predictions of the 2D body joint locations from DeepCut. SMPL captures correlations in human shape across the population, which allows fitting human poses robustly even in the presence of weak observations.

5.2.6 Traffic Sign Detection

Reliable detection and recognition of traffic signs are essential for autonomous vehicles. The introduction of the German Traffic Sign Recognition Benchmark (GTSRB) by Stallkamp et al. [623] and the German Traffic Sign Detection

Benchmark (GTSDB) by Houben et al. [299] are the most popular datasets for traffic sign detection. However, recent CNNs already reach the limits of GTSRB and GTSDB with a recall and precision of 100%. Therefore, Zhu et al. [793] recently presented Tsinghua-Tencent 100K, a new traffic sign detection benchmark, introducing new challenges to the community.

Several object detectors have been considered for traffic sign detection, i.e., SVMs [439], pattern matching techniques [78], voting schemes such as radial symmetric detectors [35] and integral channel features [171, 445]. However, the recent progress in deep learning also led to better traffic sign classifiers [127, 603, 126, 324]. Ciresan et al. [127] propose a committee consisting of a CNN trained on images and an MLP trained on HOG feature descriptors to classify traffic signs. In contrast, [603] propose a multi-scale CNN to learn meaningful features instead of using handcrafted features such as HOG. For faster training, Jin et al. [324] present a stochastic gradient descent method with a cost function similar to the objective function of the SVM. Similar to [602], Aghdam et al. [4] propose a sliding window detector that extracts features using a CNN. However, they apply the CNN using dilated convolutions on several resolutions to learn the detection of traffic signs at different scales. Finally, they train a convolutional network with fully connected layers to classify the extracted features.

García et al. [230] compare generic object detectors on the popular GTSDB dataset. Region-based networks [247, 283, 246] and one-stage generic detectors [540, 541, 419] have difficulties with traffic signs at small scales. Traffic signs can appear very small in the image depending on the size, distance, and occlusions. Region-Proposal networks [544] give the best performance of generic detectors and in combination with Inception V2 [312] for feature extraction achieve comparable results with [4] on GTSDB. Yang et al. [747] adapt Faster-RCNN [544] to the traffic sign detection task by extracting region proposals in a coarse-to-fine fashion. A novel attention network is proposed to roughly locate and classify RoIs before using the finer Region Proposal Network. This allowed them to improve upon Faster-RCNN on both GTSDB and Tsinghua-Tencent 100K datasets.

5.2.7 3D Object Detection from 2D Images

Geometric 3D representations of object classes can recover far more details than just 2D or 3D bounding boxes, however, most of today's object detectors are focused on robust 2D matching. In contrast, Zia et al. [796] exploit the fact that high-quality 3D CAD models are available for many important classes. From these models, they obtain coarse 3D wireframe models using principal components analysis and train detectors for the vertices of the wireframe. At test time, they generate evidence for vertices by densely applying the detectors. Zia et al. [797] extend this work by directly using detailed 3D CAD

models in their formulation, combining them with explicit representations of likely occlusion patterns. Further, a ground plane is jointly estimated to stabilize the pose estimation process. This extension outperforms the pseudo-3D model [796] and shows the benefits of reasoning in true metric 3D space.

While these 3D representations provide more expressive descriptions of objects, they can not yet compete with state-of-the-art detectors using 2D bounding boxes. To overcome this problem, Pepik et al. [506] propose a 3D extension of the deformable parts model [201] that combines the 3D geometric representation with robust matching to real-world images. They further add 3D CAD information of the object class of interest as geometry cue to enrich the appearance model.

Kundu et al. [370] train a CNN to map 2D object proposals to full 3D shape and pose. They add region-wise subnetworks for 3D shape and 3D pose prediction to a Faster-RCNN/Network-on-Convolution [544, 545] architecture. To facilitate the problem, they learn a low dimensional shape-space from CAD models and use it as shape prior. The 3D shape estimation is then formulated as a prediction problem of a set of low dimensional shape parameters. With a differentiable Render-and-Compare loss, they are able to learn 3D shape and pose from 2D supervision (instance segmentation or depth). In contrast, Ku et al. [363] suggest a more flexible approach using LiDAR point clouds as supervision to avoid the dependency on annotated datasets of CAD models. They use 2D detections of MS-CNN [94] and learn a model based on Faster-RCNN [544] to regress amodal, oriented 3D bounding boxes. Manhardt et al. [440] also first extract 2D detection using an architecture based on [544]. They propose a fully-differentiable mapping to lift the 2D detections, orientation, and scale estimation to the 3D space while using monocular depth predictions [513] to guide the distance reasoning.

5.2.8 3D Object Detection from 3D Point Clouds

In contrast to cameras, laser range sensors directly provide accurate 3D information, which simplifies the extraction of object candidates and can be helpful for the classification task as it provides 3D shape information.

Li et al. [406] exploit a fully convolutional neural network for detecting vehicles from range data. They use a 2D representation of the 3D range data analogous to cylindrical images with the channels encoding the 3D location of the points. Given this representation, they simultaneously predict an objectness confidence and bounding box using a single 2D CNN. In contrast, Wang and Posner [687] propose an efficient scheme to apply the common 2D sliding window detection approach to 3D data. More specifically, they discretize the space into a 3D voxel grid and exploit the sparse nature of the problem with a voting scheme on top of a linear classifier, which is shown to be equivalent to convolutions on the full 3D point cloud. Engelcke et

al. [185] extend this feature-centric voting scheme by implementing a novel convolutional layer to apply sparse convolutions across the 3D point cloud. Additionally, they encourage sparsity in the intermediate representation using ReLU non-linearities and L_1 penalty. While [687, 185] extract hand-crafted features from the voxels, VoxelNet from Zhou and Tuzel [785] learns the features in an end-to-end trainable deep network. They propose a voxel feature encoding layer that learns a unified feature representation for the points of the voxels. Eventually, a region proposal network generates detections from these feature representations.

Relying only on laser range data makes the detection task challenging due to the limited density of the laser scans and lack of appearance information. Thus, existing LiDAR-based approaches perform weaker compared to their image-based counterparts on the 2D detection problem of KITTI. However, recently, it has been shown that the fusion of LiDAR and camera information allows reducing the gap and eventually even outperforming state-of-the-art 2D detectors [116, 362, 138, 529, 177]. We will discuss these methods in detail in Section 5.5.

5.3 Datasets

The most popular datasets for object detection are ImageNet [160], PASCAL VOC [194], Microsoft COCO [420], KITTI [238] and Caltech Pedestrian Detection [173]. While ImageNet, PASCAL VOC, and Microsoft COCO consider the general detection problem, KITTI and Caltech Pedestrian Detection benchmark focus on classes that are relevant for the autonomous driving context. KITTI provides separate benchmarks for 2D and 3D detection of cars, pedestrians, and cyclists with 2D and 3D input modalities for both benchmarks. In contrast, the Caltech Detection benchmark focuses on the pedestrian detection problem only.

Recently, EuroCity Persons [68], a new large-scale benchmark for pedestrian detection, was presented. Also, several companies, i.e., ApolloScape [307], NuScenes [93], and Berkeley DeepDrive [755], presented new publicly available datasets for object detection in street scenes. Similarly to KITTI, ApolloScape provides annotation for 3D car detection but is not considering other classes than cars. The Berkeley DeepDrive dataset even provides additional classes (traffic light, traffic sign, train) for the road object detection problem. However, these datasets and benchmarks [68, 307, 93, 755] are not yet established in the field.

In this work, we focus our attention on the KITTI benchmark since it allows us to compare generic object and specific pedestrian detection systems on the same data. We refer the interested reader to the survey papers [43, 774] for an in-depth comparison of pedestrian detection systems on Caltech-USA.

5.4 Metrics

The most popular measures for the performance of object detection systems are the average precision (AP) and average recall (AR) [160, 194, 420, 133, 238]. In addition, the precision-recall curve is usually used to evaluate methods [194, 238]. For calculating precision and recall, the detections are categorized into true positives, false positives, and false negatives. Towards this goal, the intersection-over-union (IOU) between the detected bounding boxes and the ground truth bounding boxes is considered. A popular threshold for true positives is an IOU of at least 50%. The AP with an IOU of 50% is known as the PASCAL VOC [194] metric and used in many different benchmarks [238, 173, 160]. In addition to the standard PASCAL VOC metric, MS COCO [160] considers several additional metrics: the AP with an IOU of at least 75%, for small, medium and large objects, and several AR metrics.

In our discussion here, we consider the metrics reported on the KITTI benchmark [238]. The performance is assessed for three levels of difficulty (easy, moderate, hard) using PASCAL VOC intersection-over-union (IOU) [194]. While the PASCAL VOC metric (IOU of 50%) is used for pedestrians and cyclists on KITTI, the metric used for cars is more strict and requires an overlap of 70%. Easy examples have a minimum bounding box height of 40 px and are fully visible, whereas moderate examples have a minimum height of 25 px and include partial occlusions. Hard examples have the same minimum height but include large levels of occlusion. In Table 5.2, the estimation of the object’s orientation is evaluated using the average orientation similarity (AOS) proposed in [238].

5.5 State of the Art on KITTI

In Tables 5.1 and 5.3, we show the current state of the art on the KITTI benchmark for object, pedestrian, and cyclist detection from images. Note that for all result tables in this book, we list only public methods that have a technical paper associated with them that describes the details of the method.

Region-based networks [247, 283, 246] have proven to be very successful on the PASCAL VOC benchmark. However, they could not achieve similar performance on KITTI benchmark. The main reason is that objects occur at many different scales, and objects are often partially occluded. These objects are hard to detect using generic region-based networks.

In contrast, Region Proposal Networks [544, 744, 94, 728] have been more successful on the KITTI dataset. In the case of small objects, strong activations of convolutional neurons are more likely to occur in earlier layers. Therefore, Yang et al. [744] (SDP+PRN) propose cascaded rejection classifiers that gradually reject negative proposals using stronger features. Combined with a

	Method	Moderate	Easy	Hard	Runtime
1.	RRC [543]	90.23 %	90.61 %	87.44 %	3.6 s / GPU
2.	SJTU-HW [775]	90.08 %	90.81 %	79.98 %	0.85s / GPU
3.	Deep MANTA [102]	90.03 %	97.25 %	80.62 %	0.7 s / GPU
4.	sensekitti [741]	90.00 %	90.76 %	81.83 %	4.5 s / GPU
5.	SINet+ [302]	89.73 %	90.51 %	77.82 %	0.3 s /
11.	SubCNN [727]	88.86 %	90.75 %	79.24 %	2 s / GPU
12.	Deep3DBox [475]	88.86 %	90.47 %	77.60 %	1.5 s / GPU
13.	MS-CNN [94]	88.83 %	90.46 %	74.76 %	0.4 s / GPU
23.	Faster R-CNN [544]	79.11 %	87.90 %	70.19 %	2 s / GPU
41.	YOLOv2 [541]	19.31 %	28.37 %	15.94 %	0.02 s / GPU

(a) KITTI Car Detection Leaderboard

	Method	Moderate	Easy	Hard	Runtime
1.	RRC [543]	75.33 %	84.16 %	70.39 %	3.6 s / GPU
2.	SJTU-HW [775]	74.24 %	85.42 %	69.34 %	0.85s / GPU
3.	MS-CNN [94]	73.62 %	83.70 %	68.28 %	0.4 s / GPU
4.	GN [326]	71.55 %	80.73 %	64.82 %	1 s / GPU
5.	SubCNN [728]	71.34 %	83.17 %	66.36 %	2 s / GPU
10.	sensekitti [741]	67.28 %	80.12 %	62.25 %	4.5 s / GPU
11.	Mono3D [113]	66.66 %	77.30 %	63.44 %	4.2 s / GPU
12.	Faster R-CNN [544]	65.91 %	78.35 %	61.19 %	2 s / GPU
43.	YOLOv2 [541]	16.19 %	20.80 %	15.43 %	0.02 s / GPU

(b) KITTI Pedestrian Detection Leaderboard

	Method	Moderate	Easy	Hard	Runtime
1.	RRC [543]	76.49 %	84.96 %	65.46 %	3.6 s / GPU
2.	MS-CNN [94]	74.45 %	82.34 %	64.91 %	0.4 s / GPU
3.	Deep3DBox [475]	73.48 %	82.65 %	64.11 %	1.5 s / GPU
4.	SDP+RPN [744]	73.08 %	81.05 %	64.88 %	0.4 s / GPU
5.	sensekitti [741]	72.50 %	81.76 %	64.00 %	4.5 s / GPU
6.	SubCNN [728]	70.77 %	77.82 %	62.71 %	2 s / GPU
12.	Mono3D [113]	63.85 %	75.22 %	58.96 %	4.2 s / GPU
13.	Faster R-CNN [544]	62.81 %	71.41 %	55.44 %	2 s / GPU
26.	YOLOv2 [541]	4.55 %	4.55 %	4.55 %	0.02 s / GPU

(c) KITTI Cyclist Detection Leaderboard

Table 5.1: **KITTI Object Detection Leaderboard.** Only image-based methods are shown in these tables, i.e., no laser scan data is used. The numbers represent average precision at different levels of difficulty based on the object size and the level of occlusion/truncation. Higher numbers indicate better performance. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

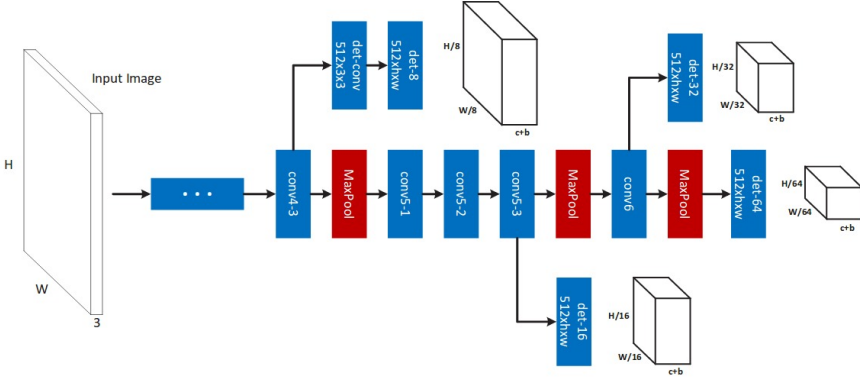


Figure 5.5: **Multi-scale Deep CNN for Object Detection.** The proposal sub-network presented by Cai et al. [94] performs detection at multiple output layers to match objects at different scales. Scale-specific detectors are combined to produce a strong multi-scale object detector. Figure courtesy of Cai et al. [94] © 2016 IEEE.

scale-dependent pooling approach that provides convolutional features from the corresponding scale for each proposal, they achieve competitive results on KITTI cyclist 5.1c). Xiang et al. [728] (SubCNN) improve on the orientation estimation task by guiding the proposal generating and detection network using subcategory information obtained from 3DVP [727]. Object subcategories are defined for objects with similar properties or attributes such as appearance, pose, or shape. This formulation allows them to achieve the best performance in pedestrian orientation estimation (Table 5.2b). The best performing Region-Proposal networks are presented by Cai et al. [94] and Chabot et al. [102]. MS-CNN [94] consists of two subnetworks, i.e., a multi-scale proposal network and a detection network. The proposal network, illustrated in Figure 5.5, has several output layers corresponding to different scales. On each output layer, the detection network is applied that allows detecting objects at different scales. Their multi-scale CNN performs well on pedestrians and cyclists (Tables 5.1b, 5.1c). Deep MANTA [102] leverages a 3D vehicle model dataset for 3D vehicle detection from images. They first extract region proposals from the input image with an iterative refinement of region proposals using a coarse-to-fine CNN (Deep MANTA network). Afterwards, they use a network to choose the closest 3D model from the 3D dataset and perform matching between the 2D regions from the image and 3D models to recover the vehicle orientation and 3D location. This allows them even to detect parts of cars that are occluded and estimate the orientation of cars.

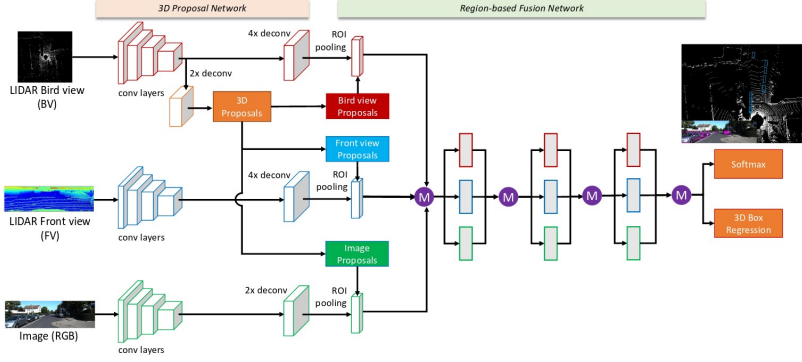


Figure 5.6: **Multi-View 3D Object Detection.** The network proposed by Chen et al. [116] combines region-wise features from the bird’s eye view, the front view of the LiDAR point cloud as well as the RGB image as input for a deep fusion network. Figure courtesy of Chen et al. [116] © 2017 IEEE.

They achieve competitive results on car detection (Table 5.1a) and the best performance on the car orientation estimation task (Table 5.2a).

One-stage detectors [540, 541, 419] have similar difficulties with objects at different scales and occlusions as region-based networks on the KITTI dataset. However, by leveraging feature pyramids as in [424], one-stage detectors [543, 775] can reach state-of-the-art performance. Zhang et al. [775] (SJTU-HW) propose to improve the localization by embedding a localization-quality estimation into the detector. They fuse features from classification and box regression subnetworks to estimate the localization quality. During inference, they combine the localization quality with the classification confidence to obtain more accurate detections. This approach outperforms all Region Proposal Networks on the pedestrian and car detection tasks. However, the best performance on all detection tasks (Tables 5.1) is achieved by Ren et al. [543]. Inspired by feature pyramids used in [424], they propose a Recurrent Rolling Convolution architecture that aggregates contextual information from multiple scales. By providing this rich contextual information to the classifier and box regressors, they achieve state-of-the-art performance on all KITTI detection tasks.

3D Object Detection from 3D Point Clouds: In Table 5.3, we show the LiDAR-based state of the art on the KITTI benchmark for object, pedestrian, and cyclist detection. The performance is assessed similarly to the image-based approaches using the intersection-over-union by projecting the 3D bounding boxes into the image plane.

Chen et al. [116] encode sparse point clouds using a compact multi-view

	Method	Moderate	Easy	Hard	Runtime
1.	Deep MANTA [102]	89.86 %	97.19 %	80.39 %	0.7 s / GPU
2.	Deep3DBox [475]	88.56 %	90.39 %	77.17 %	1.5 s / GPU
3.	SubCNN [728]	88.43 %	90.61 %	78.63 %	2 s / GPU
4.	AVOD (LiDAR) [362]	87.46 %	89.59 %	79.54 %	0.08 s /
5.	AVOD-FPN (LiDAR) [362]	87.13 %	89.95 %	79.74 %	0.1 s /
13.	Pose-RCNN [69]	75.35 %	88.78 %	61.47 %	2 s / $i8$ cores
25.	sensekitti [741]	44.56 %	47.06 %	41.50 %	4.5 s / GPU

(a) KITTI Car Detection and Orientation Estimation Leaderboard

	Method	Moderate	Easy	Hard	Runtime
1.	SubCNN [728]	63.41 %	71.39 %	56.34 %	2 s / GPU
2.	Pose-RCNN [69]	62.25 %	74.85 %	55.09 %	2 s / $i8$ cores
3.	Deep3DBox [475]	59.37 %	68.58 %	51.97 %	1.5 s / GPU
4.	3DOP (Stereo) [114]	58.59 %	71.95 %	52.35 %	3s / GPU
5.	AVOD-FPN (LiDAR) [362]	57.53 %	67.61 %	54.16 %	0.1 s /
8.	AVOD (LiDAR)[362]	54.43 %	64.36 %	47.67 %	0.08 s /
9.	Mono3D [113]	53.11 %	65.74 %	48.87 %	4.2 s / GPU
10.	FRCNN+Or [254]	50.91 %	63.41 %	45.46 %	0.09 s /
11.	sensekitti [741]	42.12 %	46.65 %	36.66 %	4.5 s / GPU

(b) KITTI Pedestrian Detection and Orientation Estimation Leaderboard

	Method	Moderate	Easy	Hard	Runtime
1.	SubCNN [728]	66.28 %	78.33 %	61.37 %	2 s / GPU
2.	Pose-RCNN [69]	59.89 %	74.10 %	54.21 %	2 s / $i8$ cores
3.	3DOP (Stereo) [114]	59.79 %	73.46 %	57.04 %	3s / GPU
4.	DeepStereoOP [512]	59.28 %	73.37 %	56.87 %	3.4 s / GPU
5.	Mono3D [113]	58.12 %	68.58 %	54.94 %	4.2 s / GPU
6.	AVOD-FPN (LiDAR) [362]	44.92 %	53.36 %	43.77 %	0.1 s /
7.	SECOND [740]	43.51 %	51.56 %	38.78 %	38 ms /
8.	DPM-VOC+VP [506]	39.83 %	53.66 %	35.73 %	8 s / 1 core
9.	sensekitti [741]	37.50 %	43.55 %	35.08 %	4.5 s / GPU
10.	AVOD (LiDAR) [362]	36.38 %	44.12 %	31.81 %	0.08 s /

(c) KITTI Cyclist Detection and Orientation Estimation Leaderboard

Table 5.2: KITTI Detection and Orientation Estimation Leaderboard. Only image-based methods are shown in these tables, i.e., no laser scan data is used. The numbers represent average orientation similarity as described in [238]. Higher numbers indicate better detection and orientation estimation. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

	Method	Moderate	Easy	Hard	Runtime
1.	PC-CNN-V2 [177]	90.15 %	90.79 %	87.58 %	0.5 s / GPU
2.	F-PointNet [529]	90.00 %	90.78 %	80.80 %	0.17 s / GPU
3.	MV3D [116]	89.17 %	90.53 %	80.16 %	0.36 s / GPU
4.	MM-MRFC [138]	88.20 %	90.93 %	78.02 %	0.05 s / GPU
5.	AVOD [362]	88.08 %	89.73 %	80.14 %	0.08 s /
18.	CSoR [521]	26.13 %	35.24 %	22.69 %	3.5 s / 4 cores
19.	mBoW [39]	23.76 %	37.63 %	18.44 %	10 s / 1 core

(a) KITTI Car Detection Leaderboard

	Method	Moderate	Easy	Hard	Runtime
1.	F-PointNet [529]	77.25 %	87.81 %	74.46 %	0.17 s / GPU
2.	MM-MRFC [138]	69.96 %	82.37 %	64.76 %	0.05 s / GPU
3.	AVOD-FPN [362]	58.42 %	67.32 %	57.44 %	0.1 s /
4.	MV-RGBD-RF [250]	56.59 %	73.05 %	49.63 %	4 s / 4 cores
5.	Vote3Deep [184]	55.38 %	67.94 %	52.62 %	1.5 s / 4 cores
8.	AVOD [362]	43.49 %	51.64 %	37.79 %	0.08 s /
9.	Vote3D [687]	35.74 %	44.47 %	33.72 %	0.5 s / 4 cores
10.	mBoW [39]	31.37 %	44.36 %	30.62 %	10 s / 1 core

(b) KITTI Pedestrian Detection Leaderboard

	Method	Moderate	Easy	Hard	Runtime
1.	F-PointNet [529]	72.25 %	84.90 %	65.14 %	0.17 s / GPU
2.	Vote3Deep [184]	67.96 %	76.49 %	62.88 %	1.5 s / 4 cores
3.	AVOD-FPN [362]	59.32 %	68.65 %	55.82 %	0.1 s /
4.	AVOD [362]	56.01 %	65.72 %	48.89 %	0.08 s /
5.	BirdNet [41]	49.04 %	64.88 %	46.61 %	0.11 s /
6.	MV-RGBD-RF [250]	42.61 %	51.46 %	37.42 %	4 s / 4 cores
7.	Vote3D [687]	31.24 %	41.45 %	28.60 %	0.5 s / 4 cores
8.	mBoW [39]	21.62 %	28.19 %	20.93 %	10 s / 1 core

(c) KITTI Cyclist Detection Leaderboard

Table 5.3: **KITTI LiDAR Detection Leaderboard.** Methods that focus on LiDAR scans and methods combining LiDAR with RGB images are presented. The numbers represent average precision at different levels of difficulty. Higher numbers indicate better performance. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

representation. While the proposal generation network utilizes the bird’s eye view to generate 3D candidates, they eventually combine region-wise features from multiple views via deep fusion for the final detection and box regression scheme, as illustrated in Figure 5.6. Instead of using an intermediate representation, Ku et al. [362] propose to directly share features extracted from LiDAR point clouds and RGB images with a Region Proposal Network (RPN) and detector network. Costea et al. [138] improve on the pedestrian and car detection task by considering dense optical flow as additional input. They use multi-modal, multi-resolution filtering of intensities, gradient magnitudes and orientations to obtain discriminative features for detection. In contrast to [116, 362, 529, 177], they follow a boosting-based sliding window approach and achieve competitive results while being faster than the deep learning-based approaches.

Qi et al. [529] propose to directly work on the 3D point clouds by reducing the search space using 2D detections in image space. This allows them to use two variants of PointNet [229]; one for 3D object instance segmentation and the other for 3D box regression. With this approach, they outperform all other 3D-based detectors on the categories pedestrian and cyclist (Tables 5.3b, 5.3c) and even all image-based detectors on pedestrians (Table 5.1b). Similar to [529], Du et al. [177] leverage 2D detections to obtain accurate 3D detections. Instead of using PointNets, they propose to fit a generalized 3D car model to the points corresponding to 2D detections. Finally, they use the points matching the model in a two-stage refinement CNN to predict the final 3D box and an objectiveness score. The combination of 2D and 3D detection allows them to outperform all 3D-based detectors on cars (Table 5.3a) while achieving a performance on par with the best performing 2D-based detector [543] on the pedestrian category (Table 5.1a).

5.6 Discussion

Object detection has demonstrated impressive performance in case of high resolution images with little occlusions. For the easy and moderate cases of the car detection task (Table 5.1a), many methods provide accurate detections. The pedestrian and cyclist detection task (Tables 5.1b, 5.1c) is more challenging, as demonstrated by the overall weaker performance of all methods. One reason for this is the limited number of training examples and the possibility of confusing cyclists and pedestrians which differ only via their context and semantics. Remaining major problems across tasks are the detection of small objects and highly occluded objects. In the leaderboards, this manifests in a significant drop in performance when comparing easy, moderate, and hard examples. Qualitatively, this can be observed in Figures 5.7, 5.8, 5.9 where we show typical estimation errors of the best-performing methods on the KITTI

dataset. Major sources of errors are crowds of pedestrians, groups of cyclists, and parked cars that cause occlusions and lead to missing detections for all methods. Furthermore, distant objects still prove to be challenging for modern methods due to the low amount of image evidence provided for these objects.



(a) Images with Largest Number of True Positive Detections



(b) Images with Largest Number of False Positive Detections



(c) Images with Largest Number of False Negative Detections

Figure 5.7: KITTI Vehicle Detection Analysis. Each figure shows images with a large number of true positive (TP) detections, false positive (FP) detections and false negative (FN) detections, respectively. If all detectors agree on TP, FP or FN, the object is marked in red. If only some of the detectors agree, the object is marked in yellow. The ranking has been established by considering the 15 leading methods published on the KITTI evaluation server at time of submission.



(a) Images with Largest Number of True Positive Detections



(b) Images with Largest Number of False Positive Detections



(c) Images with Largest Number of False Negative Detections

Figure 5.8: **KITTI Pedestrian Detection Analysis.** Each figure shows images with a large number of true positive (TP) detections, false positive (FP) detections and false negative (FN) detections, respectively. If all detectors agree on TP, FP or FN, the object is marked in red. If only some of the detectors agree, the object is marked in yellow. The ranking has been established by considering the 15 leading methods published on the KITTI evaluation server at time of submission.



(a) Images with Largest Number of True Positive Detections



(b) Images with Largest Number of False Positive Detections



(c) Images with Largest Number of False Negative Detections

Figure 5.9: **KITTI Cyclist Detection Analysis.** Each figure shows images with a large number of true positive (TP) detections, false positive (FP) detections and false negative (FN) detections, respectively. If all detectors agree on TP, FP or FN, the object is marked in red. If only some of the detectors agree, the object is marked in yellow. The ranking has been established by considering the 15 leading methods published on the KITTI evaluation server at time of submission.

Chapter 6

Object Tracking

6.1 Problem Definition

In tracking, the goal is to estimate the state of one or multiple objects over time given measurements of a sensor. This is in contrast to object detection where each frame is typically processed independently and no associations over time are established. Typically, the state of an object is represented by its location, velocity and acceleration at a certain time. Tracking of other traffic participants is a very important task for autonomous driving. Consider for instance, the braking distance of a vehicle which increases quadratically with its speed. Because of the braking distance it is necessary to detect possible collisions with other traffic participants early on. This is only possible with good predictions of future trajectories. In the case of pedestrians and bicyclists, it is particularly difficult to predict the future behavior because they can abruptly change the direction of their movements. Therefore, humans tend to drive more carefully around pedestrians and bicyclists. Similarly, tracking in combination with the classification of traffic participants allows adapting the speed of the vehicle accordingly. In addition, tracking of other cars can be used for automatic distance control and to anticipate possible driving maneuvers of other traffic participants (such as takeovers) early on.

Tracking systems must cope with a variety of challenges such as cluttered backgrounds, the variety and complexity of motion, and occlusions. The problem of associating instances of the same object over time becomes particularly challenging due to the resemblance of different objects, especially of the same class. In addition to the lack of discriminative information due to similarities with other objects, instances of the same object might not look similar enough for association in different time steps. Often objects are partially or fully occluded by other objects or themselves. The interaction of objects,

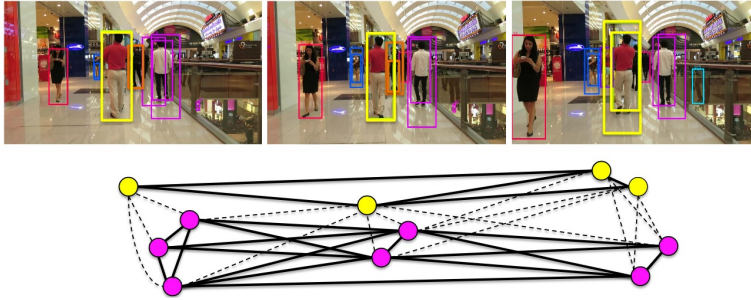


Figure 6.1: **Graph-based Data Association.** Graph-based representation solved with a multi-cut formulation presented by Tang et al. [642]. The graph is created from detections in the upper images and the colorization as well as connections in the graph are obtained by solving the multi-cut problem. Figure courtesy of Tang et al. [642] © 2016 IEEE.

especially in the case of pedestrians, further increases the amount of occlusions and makes it difficult to track each individual object. Difficult lighting conditions and reflections in mirrors or windows pose additional challenges.

6.2 Methods

Historically, tracking has been formulated as a Bayesian inference problem [650] where the goal is to estimate the posterior probability density function of a state given the current observation and the previous state(s). The posterior is usually updated in a recursive manner with a prediction step using a motion model and a correction step using an observation model. In each iteration, the data association problem is solved to assign new observations to the tracked objects. Extended Kalman and particle filtering algorithms [245, 72, 119] are widely used models in this context. Unfortunately, the recursive approach makes it hard to recover from detection errors and to track through occlusions because of missing observations. Therefore, non-recursive approaches [13, 14] that optimize a global energy function with respect to all trajectories in a temporal window, have gained popularity. However, the large number of possible target trajectories per object and the large number of potential objects in a scene lead to a very large search space.

6.2.1 Tracking by Detection

Given the success of static object detectors, a common paradigm often used in tracking is tracking-by-detection. This approach splits the task into two steps: first detect the people and second associate detections of the same person across time. Tracking-by-detection has become very popular since the tracking problem is reduced to a data association problem. However, the tracking system still needs to handle and recover from errors of the detection system, such as false and missing detections.

Tracking on Graphs: Graph representations illustrated in Figure 6.1 are widely adopted for inferring associations in tracking. In the simplest case, bipartite matching between the trajectories and the detections can be considered a graph-based approach with two disjoint sets of nodes. The assignment between the two sets can be performed either greedily [717, 73, 614] or by applying the optimal Hungarian algorithm [508, 303, 732, 542, 530] running in polynomial time.

In network flow approaches [323, 770, 44, 721, 45, 518, 121, 722], a graph is first constructed with nodes as detections and edges representing spatial and temporal links between detections. Then, a simple set of constraints is defined to ensure that produced tracks are valid and continuous between the start and the end nodes. Typically, these constraints are formulated as an integer program which is then relaxed to a linear program in order to avoid the NP-Hardness of the integer program. Various dynamic programming approaches have been proposed to solve the network flow using linear programming [323, 44], k-shortest paths [45, 518, 121] or set cover [721] for optimization.

Another line of work on graphs phrases tracking as clustering problem. Minimum Clique [762, 153] and Minimum Cost Multicut approaches [643, 642, 644, 19] find a decomposition of the graph that has the minimal sum of costs. Maximum-weight independent set formulations [605, 74] first solve the pairwise (two-frame) association problem independently and link the pairwise solutions using a learned distance measure. Graphical models [742, 743, 465, 387] minimize a global energy function defined on the nodes with pairwise and higher-order potentials.

Continuous Optimization: As an alternative to discretization, continuous energy minimization approaches have been proposed. For this highly non-convex problem, Andriyenko and Schindler [13] use a heuristic energy minimization scheme with repeated jump moves to prevent poor local minima and better explore the variable-dimensional search space. The effects of different components of their energy function are illustrated in Figure 6.2. Milan et al. [462] extend the continuous energy function of [13] to take into account physical constraints such as target dynamics, mutual exclusion, and track persistence. Assigning each observation to a certain target in data asso-

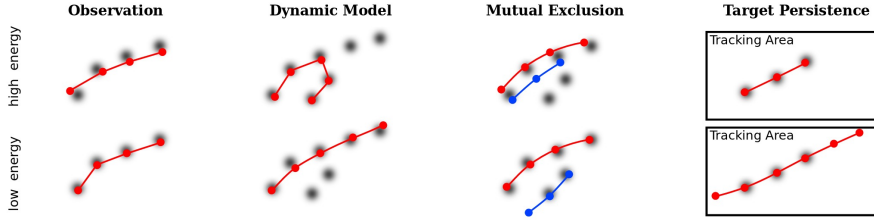


Figure 6.2: **Continuous Energy Formulation.** Components of the energy function proposed by Andriyenko and Schindler [13]. The upper and lower row show a configuration with a higher and smaller energy. The darker grey-values correspond to higher target likelihoods. Figure courtesy of Andriyenko and Schindler [13] © 2011 IEEE.

ciation is intrinsically a discrete optimization problem. Therefore, Andriyenko et al. [14] argue that a joint discrete and continuous formulation describes the tracking problem more naturally. Their method alternates between solving the data association problem using discrete optimization with label costs and analytically fitting continuous trajectories while disregarding the label costs. Milan et al. [465] propose a mixed discrete-continuous conditional random field model that specifically addresses mutual exclusion in the data association and the trajectory estimation. During data association, each observation should be assigned to at most one target while in the trajectory estimation, two trajectories should always remain spatially separated.

Multiple Cues: For data association, various complementary cues can be used in combination in order to improve the robustness of tracking systems. Giebel et al. [245] learn a spatio-temporal shape representation based on distinct linear subspace models. They handle appearance changes by combining shape, texture, and depth from stereo in the observation model of a particle filter. Gavrilu and Munder [231] employ the same set of cues with a cascade of modules in a detection and tracking system, namely region of interest generation, shape-based detection, texture-based classification, and stereo-based verification. Their system can focus on relevant image regions inferred by a stereo-based region of interest approach. They propose a novel mixture-of-experts architecture by weighting texture-based component classifiers according to the outcome of the shape matching. In their appearance-based approach, Choi et al. [119] use a combination of detection systems, each specialized in a different task such as pedestrian and upper body, face, skin color, depth-based shape, and motion. The response of all detection systems is combined in the observation likelihood to improve matching.

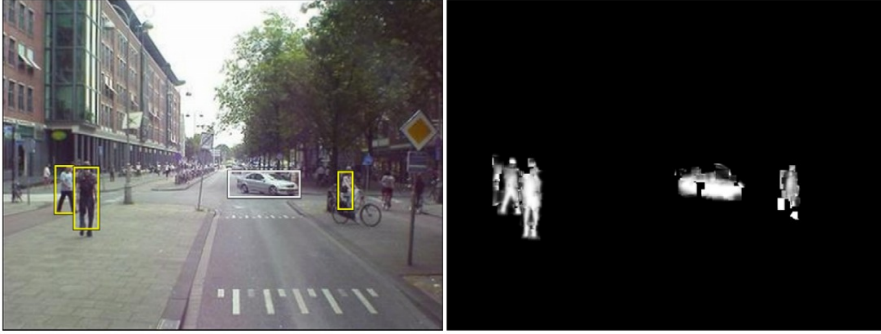


Figure 6.3: **Object Detections and Segmentations for Tracking.** The detections (left) and corresponding top-down segmentations (right) used by Leibe et al. [397] to learn an object-specific color model for tracking. Figure courtesy of Leibe et al. [397] © 2011 IEEE.

6.2.2 Tracking with Stereo

Some works have investigated a joint formulation for object tracking and stereo depth estimation to obtain the structure of the scene while estimating the trajectories of the objects in the scene. The structure of the scene allows the tracking system to focus on more plausible solutions. Leibe et al. [395, 397] propose an approach integrating scene geometry estimation, 2D object detection, 3D localization, trajectory estimation, and tracking. They learn object-specific color models using the detection and top-down segmentation of objects, as illustrated in Figure 6.3. The structure of the scene guides the extraction of physically plausible space-time trajectories and a final global optimization criterion takes object-object interactions into account to refine the 3D localization and trajectory estimation results. Ess et al. [191] jointly estimate the camera position, stereo depth, object detection, and the pose of all objects over time using a graphical model. Thereby, the graphical model represents the interplay between the different components and incorporates object-to-object interactions.

Tracking-Before-Detection: In addition to facilitating the tracking problem, depth also allows segmenting the scene into different objects, independently of their class. In tracking-before-detection, these segmented class agnostic objects are directly considered as observations in the tracking formulation. This way, the tracking system is independent of a classifier and, thus, is able to track unknown objects which have not been seen before or for which only a little amount of training data exists. Furthermore, motion information from the object’s estimated trajectory can be used as another cue to

detect a certain class of objects. Mitzel and Leibe [469] extract observations of objects by segmenting the scene using depth from stereo. With a compact 3D representation, they can robustly track known and unknown object categories. This representation also allows them to detect anomalous shapes such as carried items.

6.2.3 Pedestrian Tracking

Tracking of pedestrians is of particular importance for autonomous driving as mentioned before. However, the identification of pedestrians remains difficult, especially because of false positives of detection systems. Andriluka et al. [11] address this problem with a joint detection and articulated human pose tracking formulation. They extend an existing person detector to a limb-based structure model and model the dynamics of the detected limbs with a hierarchical Gaussian process latent variable model (hGPLVM). This allows them to detect people more reliably than approaches considering only one frame. In [12], they extend this idea towards 3D pose estimation from monocular images. In the first stage, they estimate 2D articulation and the viewpoint of people and associate them across a small number of frames. This accumulated 2D image evidence is then used to estimate the 3D pose with a hGPLVM. This approach allows them to accurately estimate the 3D poses of multiple people from monocular images. In combination with a Hidden Markov Model (HMM), these approaches can track people over very long sequences.

6.2.4 Joint Detection and Tracking

While the typical tracking-by-detection approach assumes detections to be available, Dehghan et al. [154] and Tian et al. [653] propose to solve detection and association jointly with a network flow approach by learning a model for each target and modifying the graph to encode the assignment probabilities between the targets and nodes.

Kang et al. [331, 330] introduce a tubelet proposal module that combines object detection and tracking for video object detection. A tubelet represents detections of the same object over consecutive frames. The performance is improved by first generating static object proposals as spatial anchors (e.g., from a Region Proposal Network) and then predicting the relative movements to adjust the anchors. Instead of propagating bounding boxes, Tang et al. [641] link objects in the same frame and propagate box scores across frames. In addition, per-frame proposals in [331, 330] are replaced by spatio-temporal proposals that are directly generated for video segments.

Another line of work uses optical flow for feature aggregation in videos [791, 790, 789]. The feature maps of nearby frames are warped according to

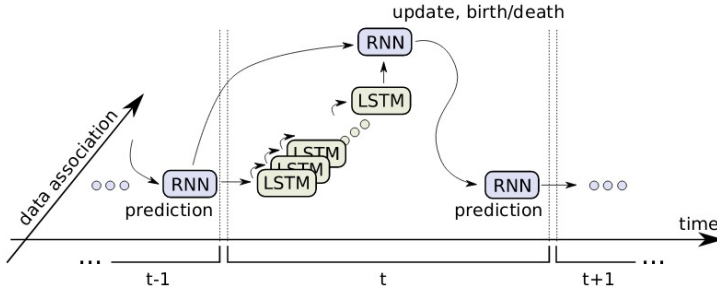


Figure 6.4: **Deep Learning for Tracking.** The end-to-end learning method by Milan et al. [464] uses RNNs [563] for state estimation and LSTMs [295] for data association. Figure courtesy of Milan et al. [464] © 2017 AAAI.

the estimated optical flow and aggregated by learning an adaptive weighting. The motivation is to improve the detection of fast-moving objects which are hard to detect on some frames due to motion blur. A more efficient version is proposed by Zhu et al. [789] with key-frame selection, i.e., selecting frames or parts of the frames to aggregate. Wang et al. [691] also use flow at different levels, namely at the pixel-level by per-pixel warping and at the instance-level by predicting instance movements. Then, the two levels are combined according to the motion pattern observed, e.g., by relying on pixel-level more in the case of non-rigid motion. To avoid expensive optical flow computation, Bertasius et al. [46] propose a spatio-temporal sampling mechanism based on deformable convolutional layers.

6.2.5 Deep Learning for Multi-Object Tracking

Tracking has strongly benefited from the success of deep learning in object detection discussed in Chapter 5.2.3. Moreover, deep learning has been used for representation learning to verify detections belonging to the same person [388, 642, 644] or more recently, for learning track representations using sequential models [565, 343, 19]. Learned sequential models combined with a traditional model for the association have shown to improve performance in comparison to their predecessors. Examples include the combination of appearance, motion, and interaction LSTM networks [565] in contrast to Markov decision process tracking with hand-crafted features [726], a modified bi-linear LSTM [343] in contrast to the multiple hypothesis tracking model with CNN features [342], and a hierarchical clustering method based on tracklet similarity using an RNN [19] in contrast to the lifted multi-cut approach with Siamese networks [644]. In these examples, the common approach is to learn a good

track representation and then use an established method for the association.

Recently, several approaches [593, 464, 212] proposed end-to-end learning of multi-object tracking. The challenges are mainly the scarcity of labeled data, the structured nature of the problem both in the input and the output space, and the combinatorial search space. Schuster et al. [593] propose a network layer to learn the network flow cost functions based on hand-designed representations of bounding boxes. The first end-to-end learning method for tracking presented by Milan et al. [464] illustrated in Figure 6.4 uses a RNN to estimate the states of targets and LSTMs for the association. The model, however, is trained on synthetic data and lacks an appearance model, which makes it unable to match the performance of previous approaches. Frossard and Urtasun [212] propose an end-to-end learning method for detection and tracking of vehicles in 3D using a deep structured loss to backpropagate through the linear program which solves the association problem. In contrast, Feichtenhofer et al. [200] present a more general approach for end-to-end learning of detection and tracking by extending the convolutional object detector proposed in [148] with a tracking loss that regresses object coordinates across frames. However, they only evaluate on ImageNet VID challenge [564] which mostly consists of sequences with one or a few objects at the center of the video.

6.3 Datasets

Early datasets for multi-object tracking include independent sequences such as PETS [203], TUD [11], and ETHZ [192]. The separate evaluation of these sequences led to tracking algorithms over-fitting to some of these sequences while performing worse on others. The MOT Challenge [390, 463] combines most of these sequences into one framework by providing a centralized evaluation and comparison. While some sequences, e.g., PETS and TUD, are captured from a static observer, other sequences that are more relevant for autonomous driving are acquired from a mobile platform. The KITTI dataset [238, 237] provides tracking data specific to autonomous driving with separate evaluations for tracking car and pedestrian classes.

The MOT Benchmark published over three consecutive years, MOT15 [390], MOT16, and MOT17 [463], consists of sets of sequences with tracking labels and provides an official evaluation protocol based on CLEAR metrics [624]. The earliest MOT15 uses a classical object detector based on aggregated channel features (ACF) [170]. In MOT16, detections are obtained using Deformable Parts Model (DPM) [201] while in MOT17, three different sets of object detections are provided using DPM [201], Faster R-CNN [544], and Scale Dependent Pooling (SDP) [744]. Providing sets of detections allows comparing approaches based on their ability to track objects independent

	Method	MOTA	IDF1	MT	ML	IDS	FRAG	Hz
1.	HCC [435]	49.3%	50.7%	17.8 %	39.9%	391	535	0.8
2.	eTC [688]	49.2%	56.1%	17.3 %	40.3 %	606	882	0.7
3.	AFN [610]	49.0%	48.2%	19.1 %	35.7 %	899	1,383	0.6
4.	KCF [123]	48.8%	47.2%	15.8 %	38.1 %	906	1,116	0.1
5.	LMP [644]	48.8%	51.3%	18.2 %	40.1 %	481	595	0.5
16.	NOMT [120]	46.4%	53.3%	18.3%	41.4%	359	504	2.6
17.	JMC [642]	46.3%	46.3%	15.5%	39.7%	657	1,114	0.8
18.	STAM [125]	46.0%	50.0%	14.6%	43.6%	473	1,422	0.2
22.	MHT_DAM [342]	42.9%	46.1%	13.6%	46.9%	499	659	0.8
38.	GMMCP [154]	38.1%	35.5%	8.6 %	50.9 %	937	1,669	0.5
44.	CEM [462]	33.2%	0.0%	7.8%	54.4%	642	731	0.3
49.	DP_NMS [518]	32.2%	31.2%	5.4%	62.1%	972	944	5.9

(a) MOT16 Leaderboard using Public DPM detections

	Method	MOTA	IDF1	MT	ML	IDS	FRAG	Hz
1.	LMP [644]	71.0%	70.1%	46.9 %	21.9 %	434	587	0.5
2.	KDNT [753]	68.2%	60.0%	41.0%	9.0%	933	1,093	0.7
3.	POI [753]	66.1%	65.1%	34.0%	20.8%	805	3,093	9.9
8.	NOMTwSDP16 [120]	62.2%	62.6%	32.5%	31.1%	406	642	3.1
9.	DeepSORT.2 [713]	61.4 %	62.2%	32.8 %	18.2%	781	2,008	17.4
10.	SORTwHPD16 [50]	59.8%	79.6%	25.4%	22.7%	1,423	1,835	59.5
11.	IOU [57]	57.1%	46.9%	23.6%	32.9%	2,167	3,028	3,004.6

(b) MOT16 Leaderboard using a Private Detector

Table 6.1: MOT16 Multi Object Tracking Leaderboard. We report the Multiple Object Tracking Accuracy (MOTA), F1 score on identified detections (IDF1), the ratio of mostly tracked (MT) and mostly lost trajectories (ML), number of ID switches (IDS) and track segmentations (FRAG), and run time. The metrics are detailed in [463]. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

of errors caused by different detectors. For MOT16, the leaderboard with methods using the public (DPM) detections is provided in Table 6.1a and the methods using a private detector are shown in Table 6.1b. For MOT17, Table 6.2 shows the average results over three provided detectors on the same set of sequences.

For autonomous driving application specifically, KITTI [238] provides two benchmarks, one for tracking of cars (KITTI car) in Table 6.3a and the other for tracking of pedestrians in Table 6.3b. Methods marked with an asterisk use Regionlet detections [692] for an independent comparison of the tracking performance. The separate challenges for cars and pedestrians allow focusing on each class separately and investigating the problems specific to a class deeply.

	Method	MOTA	IDF1	MT	ML	IDS	FRAG	Hz
1.	JBNO [292]	52.6%	50.8%	19.7 %	35.8 %	3,050	3,792	5.4
2.	FAMNet [124]	52.0%	48.7%	19.1%	33.4 %	3,072	5,318	-
3.	eTC [688]	51.9%	58.1%	23.1%	35.5 %	2,288	3,071	0.7
4.	eHAF17 [611]	51.8%	54.7%	23.4%	37.9%	1,834	2,739	0.7
5.	AFN [610]	51.5%	46.9%	20.6%	35.5%	2,593	4,308	1.8
6.	FWT [291]	51.3%	47.6%	21.4%	35.2%	2,648	4,279	0.2
7.	jCC [341]	51.2%	54.5%	20.9%	37.0%	1,802	2,984	1.8
9.	MHT_DAM [342]	50.7%	47.2%	20.8%	36.9%	2,314	2,865	0.9
14.	DMAN [787]	48.2%	55.7%	19.3%	38.3%	2,194	5,378	0.3
17.	MHT_bLSTM [343]	47.5%	51.9%	18.2%	41.7%	2,069	3,124	1.9

Table 6.2: **MOT17 Multi Object Tracking Leaderboard using Provided Detections.** We report the Multiple Object Tracking Accuracy (MOTA), F1 score on identified detections (IDF1), the ratio of mostly tracked (MT) and mostly lost trajectories (ML), number of ID switches (IDS) and track segmentations (FRAG), and run time. The metrics are detailed in [463]. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

6.4 Metrics

In Tables 6.1a, 6.1b, 6.2, 6.3a, 6.3b, we consider the commonly used tracking measures, Multiple Object Tracking Accuracy (MOTA) and Multiple Object Tracking Precision (MOTP) introduced by [624], the ratio of mostly tracked (MT) and mostly lost trajectories (ML), number of ID switches (IDS) and track segmentations (FRAG). For the MOT leaderboards, following the benchmark page, we show the IDF1 score introduced by [555] instead of MOTP. The IDF1 score is the F1 score of the identification precision and recall, i.e., the ratio of correctly identified detections over the average number of ground truth and computed detections. Mostly tracked and mostly lost trajectories show the percentage of trajectories that are covered by a hypothesis at least 80% or at most 20% of the time, respectively. For descriptions of the metrics and the detailed tables with additional metrics such as False Negatives, False Positives, ID recall, and ID precision, please check the KITTI [238] and MOT benchmarks [390, 463] as well as [555].

6.5 State of the Art on MOT & KITTI

MOT16 Benchmark: Classical approaches such as near-online multi-target tracking approach [120], multiple hypothesis tracking approach [342] and tracking based on Markov decision processes [726] still perform consistently well on MOT benchmarks in comparison to newly proposed methods.

	Method	MOTA	MOTP	MT	ML	IDS	FRAG	Runtime
1.	MOTBeyondPixels [608]	84.24 %	85.73 %	73.23 %	2.77 %	468	944	0.3 s / 1 core
2.	IMMDP [726]	83.04 %	82.74 %	60.62 %	11.38 %	172	365	0.19 s / 4 cores
3.	JCSTD [652]	80.57 %	81.81 %	56.77 %	7.38 %	61	643	0.07 s / 1 core
4.	3D-CNN/PMBM [584]	80.39 %	81.26 %	62.77 %	6.15 %	121	613	0.01 s / 1 core
7.	NOMT* [120]	78.15 %	79.46 %	57.23 %	13.23 %	31	207	0.09 s / 16 cores
10.	DSM [212]	76.15 %	83.42 %	60.00 %	8.31 %	296	868	0.1 s / GPU
11.	SCEA* [751]	75.58 %	79.39 %	53.08 %	11.54 %	104	448	0.06 s / 1 core
12.	CIWT* [497]	75.39 %	79.25 %	49.85 %	10.31 %	165	660	0.28 s / 1 core
14.	SSP* [399]	72.72 %	78.55 %	53.85 %	8.00 %	185	932	0.6 s / 1 core
18.	RMOT* [752]	65.83 %	75.42 %	40.15 %	9.69 %	209	727	0.02 s / 1 core

(a) KITTI Car Tracking Leaderboard

	Method	MOTA	MOTP	MT	ML	IDS	FRAG	Runtime
1.	IMMDP [726]	47.22 %	70.36 %	24.05 %	27.84 %	87	825	0.9 s / 8 cores
2.	NOMT* [120]	46.62 %	71.45 %	26.12 %	34.02 %	63	666	0.09 s / 16 cores
4.	JCSTD [652]	44.20 %	72.09 %	16.49 %	33.68 %	53	917	0.07 s / 1 core
5.	SCEA* [751]	43.91 %	71.86 %	16.15 %	43.30 %	56	641	0.06 s / 1 core
6.	RMOT* [752]	43.77 %	71.02 %	19.59 %	41.24 %	153	748	0.02 s / 1 core
8.	CIWT* [497]	43.37 %	71.44 %	13.75 %	34.71 %	112	901	0.28 s / 1 core
10.	NOMT [120]	36.93 %	67.75 %	17.87 %	42.61 %	34	789	0.09 s / 16 core
11.	RMOT [752]	34.54 %	68.06 %	14.43 %	47.42 %	81	685	0.01 s / 1 core
13.	SCEA [751]	33.13 %	68.45 %	9.62 %	46.74 %	16	717	0.05 s / 1 core

(b) KITTI Pedestrian Tracking Leaderboard

Table 6.3: KITTI Tracking Leaderboard. We report the Multiple Object Tracking Accuracy (MOTA), Multiple Object Tracking Precision (MOTP), the ratio of mostly tracked (MT) and mostly lost trajectories (ML), number of ID switches (IDS) and track segmentations (FRAG), and run time. The metrics are detailed in [238]. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

Their deep learning counterparts with better appearance models that are learned as explained in Section 6.2.5 perform even better [565, 343]. The comparison of [565] to [726] can be found on MOT15¹ due to the date of the first publication preceding MOT16 and MOT17. The learning-based method proposed in [343] is trained on ground truth detections and performs worse compared to [342] on MOT17 (Table 6.2) due to noisy DPM detections, which affect overall performance as explained in [343].

The success of single object trackers (SOT) triggered methods that combine several single object detectors for MOT (Tables 6.1a, 6.2) by learning a tracker for each object [125, 123]. These approaches initialize a new single object tracker whenever a new object is detected with high confidence. They assign detections of known objects to each single object tracker by restricting the search space according to a motion model and choosing the best detection candidate using a binary classifier. In Chu et al. [125], a spatial-temporal at-

¹https://motchallenge.net/results/2D_MOT_2015/

tention mechanism is proposed to handle the drift caused by occlusion and interactions among targets. Chu et al. [123] encode awareness both within and between object models and proposes an adaptive model refreshment strategy to eliminate noise in model initialization.

The customized tracker proposed by Ma et al. [435] is the best-published method on MOT16 using the provided detections (Table 6.1a). The method is based on the formulation of tracking as a minimum cost lifted multi-cut problem similar to [643, 642, 644]. Despite being offline and therefore not directly applicable to autonomous driving, this type of graph-based clustering formulation performs very well on MOT. In contrast to previous methods, Ma et al. [435] learn a sequence-specific tracker by fine-tuning a re-identification network using the test sequences. They use the assumption that non-overlapping tracklets represent different individuals to adapt a generic re-identification CNN on test sequences.

Comparing public and private detectors on MOT16 (Tables 6.1a, 6.1b) shows the importance of good object detectors. Tracking algorithms perform much better using private (usually better) object detections compared to public detections, e.g., LMP [644] 71.0% versus 48.8% and NOMT [120] 62.2% versus 46.4% in MOTA. Recent object detectors combined with simple tracking algorithms perform significantly better than any tracker with public detections such as the simple tracker IOU [57] or SORT [50], which are based on the Hungarian method in combination with Kalman filtering. Wojke et al. [713] improve the performance of SORT further by incorporating deep features into the pipeline for appearance matching. Similarly, Yu et al. [753] also use a tracking algorithm based on the Hungarian algorithm and Kalman filtering in combination with deep features for appearance matching. However, their detector is trained on additional data including a self-collected surveillance dataset, which is not public.

MOT17 Benchmark: Top-performing methods on MOT17 (Table 6.2) follow a graph clustering scheme by associating tracklets, i.e., a short sequence of detections, which can be easily and reliably associated, instead of detections. Wang et al. [688] first create tracklets based on IOU and epipolar geometry in the case of a moving camera. Tracklets represent nodes on a graph which are then clustered based on a greedy search-based clustering method. Shen et al. [610] incorporate the score of the tracklets into the learning-based network flow approach proposed in [593]. While these methods have a preprocessing step to generate tracklets, a more recent approach called FAMNet [124] combines feature extraction, affinity estimation and the assignment problem in a single network. Furthermore, single object tracking is incorporated into the tracking system in order to recover from missing detections.

Recently, several approaches [341, 291, 292] propose to use additional cues such as head detections and motion segmentation to improve tracking. Two of

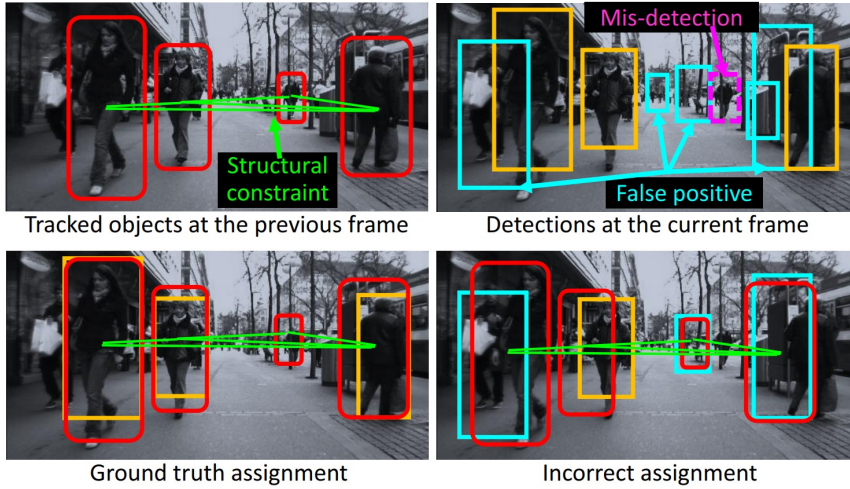


Figure 6.5: **Tracking with Structural Motion Constraints.** Structural motion constraints introduced by Yoon et al. [751] to resolve errors caused by false positives. The correct detections are marked with red and yellow boxes. Figure courtesy of Yoon et al. [751] © 2016 IEEE.

the best-performing methods on MOT17 (Table 6.2) fuse head, body, and joint detectors into a tracking system [291, 292]. Another method, proposed by Keuper et al. [341], addresses multi-object tracking with top-down clustering of bounding boxes and bottom-up motion segmentation by grouping point trajectories.

KITTI Benchmark: In contrast to the MOT Challenge, the KITTI benchmark focuses on the challenging scenario of tracking pedestrians (Table 6.3b) and cars (Table 6.3a) in traffic scenes. Similarly to MOT, classical approaches perform reasonably well such as tracking based on Markov decision process (IMMDP) [726], improved min-cost network flow [399], or the near-online multi-target tracking algorithm (NOMT) [120]. In IMMDP, a policy is learned using reinforcement learning, which corresponds to learning a similarity function for data association. An improved version with Region Proposal Network [544] is the best performing method on the car tracking task. Lenz et al. [399] propose a computational and memory bounded version of the min-cost network flow formulation presented in [770]. This approach achieves good accuracy and precision while being amongst the fastest approaches on KITTI car. NOMT [120] proposes Aggregated Local Flow Descriptor (ALFD) which encodes relative motion patterns. Thanks to these features, distant detec-

tions can be robustly matched. Using multiple feature cues, their method outperforms all the online tracking approaches on KITTI car.

Recent approaches leverage domain-specific information such as the motion of the car or the structure of the scene. Yoon et al. [752] factor out the camera motion by constructing a network to describe the relative motion between objects. They further improve in [751] by exploiting structural motion constraints defined by the location and velocity difference between two objects as illustrated in Figure 6.5. Jointly reasoning about the structure allows them to alleviate problems that are common to 2D trackers (e.g., occlusions) and outperform them, especially in the car tracking task. Frossard and Urtasun [212] propose to learn tracking in a network flow approach based on 3D detections. The structured hinge loss is adapted to backpropagate through the Integer Program. Other top-performing 3D algorithms are [497] coupling image and world-space estimations using a novel 2D-3D Kalman filter and [584] proposing a Poisson multi-Bernoulli mixture (PMBM) tracker. Sharma et al. [608] exploit the geometry of urban road scenes to infer 3D cues for tracking such as 3D pose and shape based on single view reconstruction of objects. This approach outperforms all others in accuracy (MOTA) and precision (MOTP) in the KITTI car leaderboard.

6.6 Discussion

Reliable tracking-by-detection can only be achieved by using very accurate object detections. The impact of the detection system can be observed when comparing the methods marked with and without asterisks in the KITTI leaderboards (Tables 6.3a, 6.3b). In the MOT16 leaderboards this can be observed when comparing the tables for methods using public detections in Table 6.1a and private object detectors in Table 6.1b. However, we discuss the problem of object detection in detail in Section 5.6 and focus our attention in this section on the tracking problem. Similar to the detection, tracking pedestrians is typically more challenging than cars. The reason is the complex motion of pedestrians which is hard to predict, in contrast to the rigid motion of cars which are bound by the road region and follow a less erratic behavior due to their large mass and dynamical constraints. 3D reasoning can help to improve tracking performance, especially for cars, by identifying plausible solutions according to geometric relationships.

In traffic scenes, detectors frequently fail for partially or fully occluded objects. In these cases, the tracking system needs to re-identify the tracked objects later in time but this can be difficult due to changes in lighting conditions or similarity to other objects in the proximity. These problems cause a reinitialization of trajectories, which can be observed in the high number of fragmentations (FRAG) and ID switches (IDS) in the MOT and KITTI

benchmarks. Furthermore, we note that most tracking systems comprise complex pipelines and very few end-to-end multiple target tracking algorithms have been proposed in the literature. Bridging this gap from detection to tracking with the goal of a generic and end-to-end trainable model will be an important direction for future research in this area.

Chapter 7

Semantic Segmentation

7.1 Problem Definition

Semantic segmentation is a fundamental problem in computer vision and an intermediate goal towards solving higher-level tasks such as scene understanding or sensorimotor control. The goal of semantic segmentation is to assign each pixel in the image a label from a predefined set of categories. The task is illustrated in Figure 7.1 using an example from the Cityscapes dataset¹ by Cordts et al. [133]. Segmentation of images into semantic regions that are typically found in street scenes, such as cars, pedestrians, or road allows for a comprehensive understanding of the surrounding which is essential to autonomous navigation. The task is difficult due to the complexity of the scene, complicated object boundaries, small objects and the large size of the label space.

7.2 Methods

The goal of semantic segmentation is to assign a semantically meaningful class label (e.g., road, sidewalk, pedestrian, sky) to each pixel of an image. Traditionally, the problem was posed as maximum-a-posteriori (MAP) inference in a conditional random field (CRF), defined over pixels [284, 674, 612] or superpixels [285, 352]. Hierarchical [284, 367, 377, 379] and long-range connectivity as well as higher-order potentials defined on image regions [285, 352] have been exploited to compensate for limitations of CRFs with local connections and to model long-range interactions within the image. Krähenbühl and Koltun [358] propose a tractable inference algorithm for fully connected

¹<https://www.cityscapes-dataset.com/>

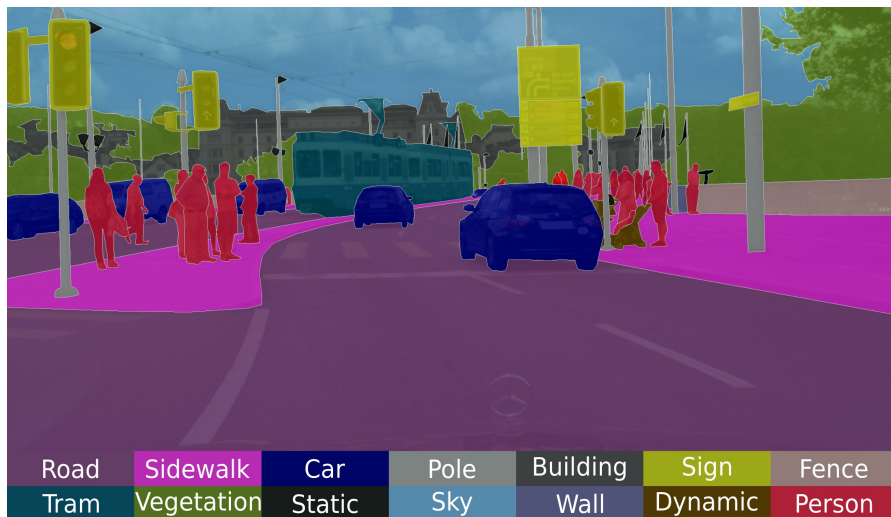


Figure 7.1: **Semantic Segmentation.** In semantic segmentation, the goal is to assign a semantic class label to each pixel in the image. Example from the Cityscapes dataset by Cordts et al. [133].

CRF models which model pairwise potentials between all pairs of pixels in the image. While previous methods using fully connected CRFs [534, 223, 352] could only be applied to smaller image regions due to the computational and memory complexity of these algorithms, [358] allows deploying fully connected CRF models at pixel-level. Figure 7.2 illustrates the results of [358] and compares them to pixel-wise classification and inference over superpixels [352].

An alternative to inference in graphical models for the task of semantic segmentation is presented by Munoz et al. [479]. They train a sequence of inference models in a hierarchical procedure that captures context over larger image regions. This allows them to bypass the difficulties of training structured prediction models when exact inference is intractable and yields a very efficient and accurate scene labeling algorithm.

While most previous approaches rely on very simple features such as color, edge and texture information, Shotton et al. [612] observed that more powerful features have the potential to significantly boost performance. They propose an approach based on a novel feature type called texture-layout filter that exploits the textural appearance of objects, its layout as well as textural context. They combine texture-layout filters with lower-level image features in a CRF to obtain pixel-level segmentations.

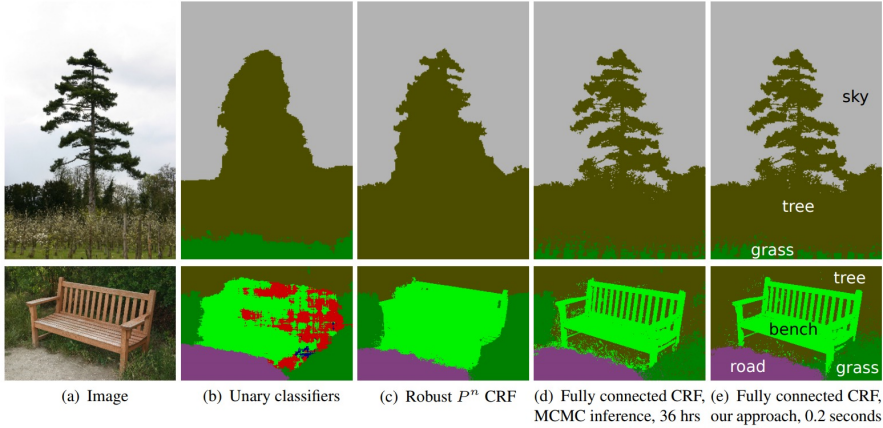


Figure 7.2: **Fully Connected Conditional Random Field.** Semantic segmentation results of a per-pixel classifier [358], a superpixel-based CRF [352] and a fully connected CRF [358]. Figure courtesy of Krähenbühl and Koltun [358] © 2016 NeurIPS.

Co-occurrence of Object Classes: The methods so far consider each object class independently. However, the co-occurrence of object classes is typically not random and can thus be an important cue for semantic segmentation, e.g., cars are more likely to occur in a street scene than in an office scene and co-occur with other street scene objects such as traffic signs. Ladicky et al. [378] propose to explicitly incorporate object class co-occurrence as global features into a CRF. They optimize the CRF using graph cuts and demonstrate better performance compared to pairwise models. Zhang and Chen [776] extend this idea by encoding spatial arrangements of different object categories. Myeong et al. [483] propose a retrieval-based approach which extracts contextual relationships from annotated region pairs.

7.2.1 Deep Learning for Semantic Segmentation

The success of deep convolutional neural networks for image classification and object detection has sparked interest in leveraging their potential for solving pixel-level tasks, in particular semantic segmentation. The fully convolutional neural network [425] illustrated in Figure 7.3 is one of the earliest works which applies CNNs to the image segmentation problem. Modern convolutional neural networks for image classification combine multi-scale contextual information by consecutive pooling and sub-sampling layers that lower the resolution. However, semantic segmentation requires multi-scale contextual reasoning together with full-resolution predictions, i.e., dense predictions.

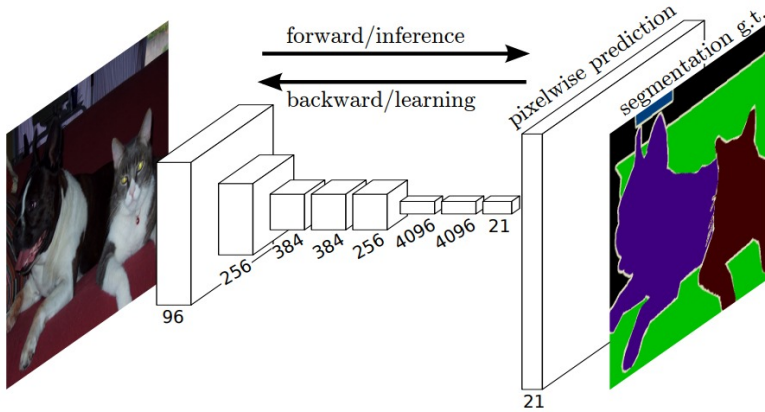


Figure 7.3: **Convolutional Neural Network.** Fully convolutional neural network for semantic segmentation proposed by Long et al. [425]. Figure courtesy of Long et al. [425] © 2015 IEEE.

Several methods [110, 754, 244, 25] have therefore been proposed to tackle the opposing needs of multi-scale inference and full-resolution outputs. Dilated convolutions [110, 754] enlarge the receptive field of neural networks without loss of resolution. The dilated convolution operation corresponds to a regular convolution that skips pixels while applying the filter. This allows for efficient multi-scale reasoning without increasing the number of model parameters. Chen et al. [109] extend this idea by using multiple dilated convolutions with different sampling rates in parallel.

In contrast, Badrinarayanan et al. [25] propose an encoder-decoder network with skip connections. Each decoder layer maps a low resolution feature map of an encoder (max-pooling) layer to a higher resolution feature map. In particular, the decoder in their model takes advantage of the pooling indices computed in the max-pooling (i.e., downsampling) step of the corresponding encoder to implement the upsampling process. This eliminates the need to learn the upsampling and thus results in a smaller number of parameters. Furthermore, sharper segmentation boundaries can be obtained using this approach.

While activation maps at lower-levels of the CNN hierarchy lack information specific to object categories, they provide information of higher spatial resolution. Ghiasi and Fowlkes [244] leverage this assumption and propose to construct a Laplacian pyramid based on a fully convolutional network. Aggregating information at multiple scales allows them to successively refine the boundary reconstructed from lower-resolution layers. They achieve this by us-

ing skip connections from higher resolution feature maps and multiplicative confidence gating, penalizing noisy high-resolution outputs in regions where low-resolution predictions have high confidence.

Combining CNNs and CRFs: A different way to address the needs of multi-scale inference and full resolution prediction is the combination of CNNs with CRF models. Chen et al. [110] and Chen et al. [109] propose to refine the label map obtained using a convolutional neural network using a fully connected CRF model [358]. The CRF allows them to capture fine details based on the raw RGB input which are missing in the CNN output due to the limited spatial accuracy of the CNN model. In a similar spirit, Jampani et al. [314] generalize bilateral filters and backpropagate through the CRF inference [413] which allows for end-to-end training of the (generalized) filter parameters from data. This effectively allows for reasoning over larger spatial regions within one convolutional layer by leveraging input features as a guiding signal.

Inspired by higher-order CRFs for semantic segmentation, Gadde et al. [218] propose a new Bilateral Inception module for CNN architectures as an alternative to structured CNNs and CRF techniques. They use the assumption that pixels which are spatially and photometrically similar are more likely to have the same label. This allows them to directly learn long-range interactions, thereby removing the need for post-processing using CRF models. Specifically, the proposed modules propagate edge-aware information between distant pixels based on their spatial and color similarity, incorporating the spatial layout of superpixels. Propagation of information is achieved by applying bilateral filters with Gaussian kernels at various scales.

Deeper CNNs: Simonyan and Zisserman [616] and Szegedy et al. [639] have shown that the depth of a CNN is crucial to represent rich features. However, increasing the depth of a network leads to an increase in complexity as well as to saturation and degradation in accuracy. He et al. [281] proposed the deep residual learning framework (ResNet) illustrated in Figure 7.4 to address this problem. In deep residual networks, each stacked layer learns a residual mapping instead of the original mapping. This facilitates the backpropagation of gradients and thus training and results in higher accuracy in comparison to regular deep networks. Pohlen et al. [522] present a ResNet-like architecture which preserves high-resolution information throughout the entire network by combining two different processing streams. One stream passes through a sequence of convolution and pooling layers, whereas the other stream processes feature maps at full image resolution by adding successive residuals from the other stream. Both processing streams are connected using full resolution residual units.

Wu et al. [723] propose a more efficient ResNet architecture by analyzing the effective depth of residual units. They point out that ResNets behave

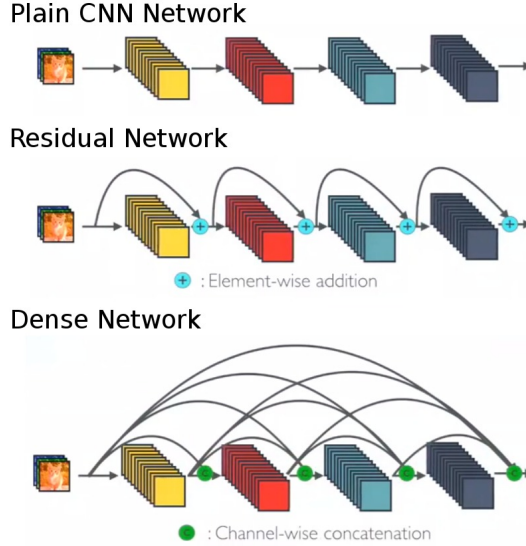


Figure 7.4: **Deep Convolutional Neural Networks.** Comparison of plain, Residual[281] and Dense[304] convolutional neural networks. Figure courtesy of Tsang [660].

as linear ensembles of shallow networks. Based on this understanding, they design a group of relatively shallow convolutional networks for the task of semantic image segmentation, which performs better. To better incorporate global context information into the pixel-level prediction task, Zhao et al. [779] propose a pyramid scene parsing network (PSPNet), illustrated in Figure 7.5. They apply a pyramid parsing module to the last convolutional layer of a CNN which fuses features of several pyramid scales to combine local and global context information. The resulting representation is fed into a convolution layer to obtain the final per-pixel predictions. Inspired by this work, [111] revisited the Atrous Spatial Pyramid Pooling (ASPP) [109] by experimenting with cascading and parallel application of dilated convolutions. This allows them to improve upon their previous work [109] while achieving comparable results to PSPNet [779].

Motivated by deeper architectures like ResNet, Huang et al. [304] propose dense convolutional networks that connect a layer with all preceding layers by concatenation. This allows maximal information throughput from lower to higher levels. In Figure 7.4, plain, residual, and dense architectures are illustrated. Jégou et al. [318] extend dense CNNs to the semantic segmentation problem by constructing a downsampling and upsampling path using dense modules and connecting them with skip connections [557].

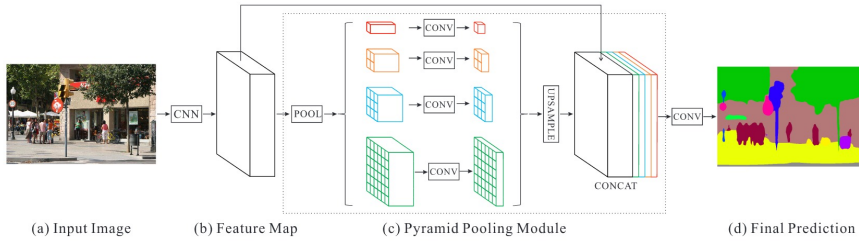


Figure 7.5: **Pyramid Pooling Module.** Overview of the method proposed by Zhao et al. [779]. The pyramid parsing module (c) is applied to the output of a CNN feature map (b) and fed into a convolutional layer for per-pixel estimation of semantic class labels (d). Figure courtesy of Zhao et al. [779] © 2017 IEEE.

7.2.2 Videos

In robotic applications such as autonomous driving we usually have access to videos rather than single image frames. The temporal correlation between adjacent frames can be exploited to improve segmentation accuracy, efficiency and robustness. The scene usually changes only slightly between two adjacent frames. Thus, given correspondences between two frames, semantic labels can be propagated in time or corrected using temporal information.

Floros and Leibe [204] propose a graphical model for semantic segmentation operating on video sequences in order to enforce temporal consistency between frames. Specifically, they present a CRF where temporal consistency between consecutive video frames is ensured by linking corresponding image pixels to the inferred 3D scene points obtained by Structure-from-Motion (SfM). Compared to an image-only baseline, they achieve improved segmentation performance and observe better generalization to varying image conditions. 3D reconstruction works relatively well for static scenes but is still an open problem in dynamic scenes. The presence of both camera and object motion makes temporal association in videos a challenging task. In case of significant motion, Euclidean distance in the space-time volume is not a good measure for finding correspondences. In order to tackle this problem, Kundu et al. [371] propose a method for optimizing the feature space of a dense CRF for spatio-temporal regularization. Specifically, the feature space is optimized such that distances between features associated with corresponding points are minimized using correspondences from optical flow. The resulting mapping is exploited by the CRF to achieve long-range regularization over the entire space-time volume.

Label Propagation: Another way to explore temporal correlations in

videos for semantic segmentation is label propagation. Creating large scale image datasets with highly accurate pixel-level annotations is labor-intensive, and thus obtaining the desired degree of quality is very expensive. Semi-supervised methods for annotating video sequences can help to reduce this cost. Compared to annotating individual images, video sequences offer the advantage of temporal consistency between consecutive frames. Label propagation techniques take advantage of this fact by propagating annotations from a small set of annotated keyframes to all unlabeled frames of the video by exploiting color and motion information.

Towards this goal, Badrinarayanan et al. [24] propose a coupled Bayesian network which employs a propagation scheme based on correspondences obtained from patch-based similarities and semantically consistent regions. This allows them to transfer label information to unlabeled frames between annotated keyframes. Budvytis et al. [88] extend this approach by proposing a hybrid model of the generative propagation introduced in [24] as well as a discriminative classification stage which tackles occlusions and disocclusions, and allows to propagate over larger time intervals. To correct erroneously propagated labels, Badrinarayanan et al. [23] propose a superpixel based mixture-of-tree model for temporal correlation where each component of the mixture contains a tree-structured temporal linkage between superpixels of different frames. Vijayanarasimhan and Grauman [677] tackle the problem of selecting the most promising keyframes for manual labeling such that the expected propagation error is minimized.

While the aforementioned methods transfer annotations in 2D, Chen et al. [107] and Xie et al. [731] propose to annotate directly in 3D and then transfer these annotations into the image domain. Given 3D information (e.g., from stereo or LiDAR), these approaches are able to produce time coherent semantic labels with limited annotation costs. Towards this goal, Chen et al. [107] use annotations from KITTI [237] and leverage 3D CAD models of cars to infer separate figure-ground segmentations for all cars in the image. In contrast, Xie et al. [731] reason jointly about all objects in the scene by also handling categories for which CAD models or 3D point measurements are not available. To this end, they propose a non-local CRF model which reasons jointly about semantic and instance labels of all 3D points and pixels in the image.

Scene Understanding: Scene understanding approaches such as [193, 236] discussed in Chapter 14 exploit semantic segmentation as a cue for reasoning about road topologies and traffic participants. While Ess et al. [193] use semantic information to classify a scene into different road topologies based on a short video sequence, Geiger et al. [236] formulate a probabilistic model which explains semantic segmentation together with vehicle trajectories, vanishing points, scene flow and occupancy information. However, both approaches

do not leverage temporal correlations to improve the semantic segmentation itself.

7.2.3 Street Side Views

One specific application scenario of semantic segmentation which has important applications for autonomous vehicles is the segmentation of street-side images (i.e., building facades) into their components (wall, door, window, vegetation, balcony, store, mailbox, etc.). Such semantic segmentations are useful for accurate 3D reconstruction [263, 262, 118], memory-efficient 3D mapping, robust localization [590] as well as path planning. As an example, in 3D reconstruction applications such side information allows for ignoring vegetation that is difficult to model and will change over time.

Xiao and Quan [730] propose a multi-view semantic segmentation framework for images captured by a camera mounted on a car driving along the street. Specifically, they define a pairwise MRF across superpixels in multiple views, where the unary terms are based on 2D and 3D features. Furthermore, they minimize color differences for spatial smoothness and use dense correspondences to enforce smoothness across different views. Xiao et al. [729] go one step further and generate photo-realistic 3D models from images captured at ground level. In particular, they segment each image into semantically meaningful areas, such as building, sky, ground, vegetation or car. Then, they partition buildings into independent blocks exploiting architectural priors for inference. This allows them to cope with noisy and missing 3D data and produces visually compelling results. While Xiao and Quan [730] and Xiao et al. [729] represent facades with planes or simple geometric primitives, Mathias et al. [444] propose a more flexible 3-layered method for segmentation of building facades. First, the facade is segmented into semantic classes which are combined with the output of detectors for architectural elements such as windows and door. Finally, weak architectural priors such as alignment, symmetry and co-occurrence are exploited to encourage the reconstruction to be architecturally consistent. The complete pipeline is illustrated in Figure 7.6.

7.2.4 3D Data

While the problem of semantic object labeling has been studied extensively, most of these algorithms work in the 2D image domain where each pixel in the image is labeled with a semantic category such as car, road or pavement. However, 2D images lack important information such as the 3D shape and scale of objects, which are strong cues for object class segmentation and facilitate the detection and separation of individual object instances. Furthermore, semantic segmentation of 3D data enables autonomous systems to recognize

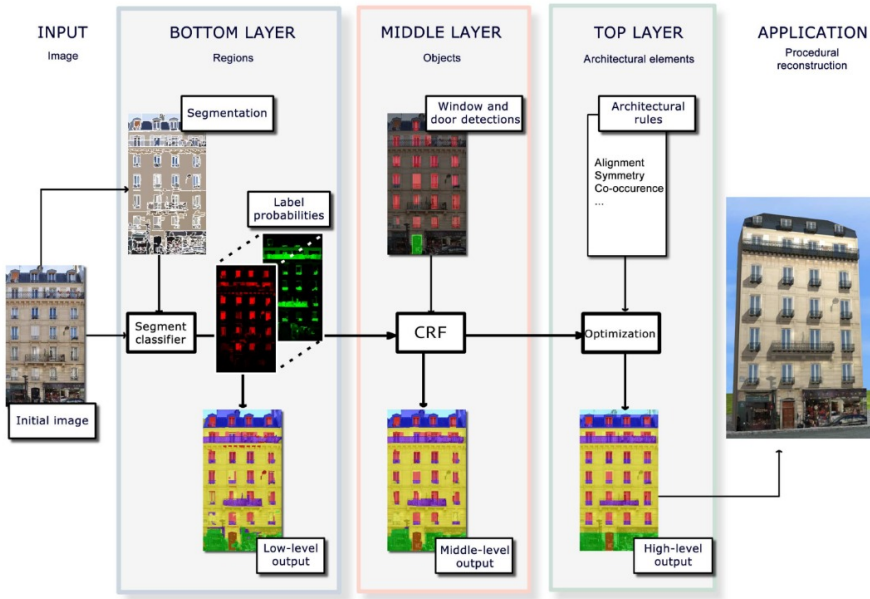


Figure 7.6: **Facade Parsing.** The three-layered approach proposed by Mathias et al. [444] for facade parsing. They first segment the facade and assign probability distributions to semantic classes considering extracted visual features. In the next layer they use detectors for specific objects such as doors and windows to improve the classifier output. Finally, they incorporate weak architectural priors and search for the optimal facade labeling using a sampling-based approach. Figure courtesy of Mathias et al. [444] © 2016 Springer.

their surroundings, identify and interact with objects of interest in physical 3D space.

The problem of 3D semantic segmentation has been addressed using different input modalities, i.e., monocular image sequences [443], stereo image sequences [672, 599] or 3D point clouds [733, 300, 261, 219]. While [443, 672] use multi-view reconstruction approaches which we discuss in Chapter 10 for estimating the 3D structure of the scene from monocular image sequences, [219, 261] directly work with 3D point clouds, e.g., from LiDAR. Sengupta et al. [599] propose to project 2D semantic segmentation into a 3D model obtained from depth map fusion using ego-motion estimation from visual odometry as illustrated in Figure 7.7. In parallel, the input images are semantically labeled using a CRF model. The results of this segmentation are

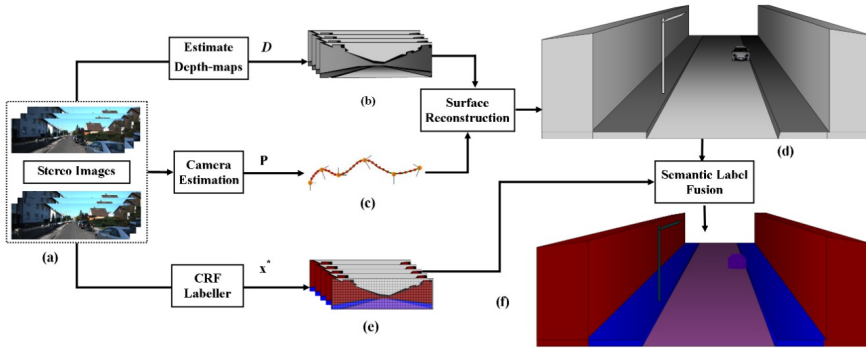


Figure 7.7: **Semantic Segmentation of 3D Data.** From a stereo image pair (a) Sengupta et al. [599] compute the disparity map (b) and track the camera motion (c). They use both outputs to obtain a volumetric representation (d) and fuse the semantic segmentation of street images (e) into a 3D semantic model of the scene (f). Figure courtesy of Sengupta et al. [599] © 2013 IEEE.

then aggregated across the sequence to generate the final 3D semantic model.

Several approaches [300, 672, 443, 219, 261] tackle the problem of semantic scene reconstruction directly in 3D space as shown in Figure 7.8. Valentin et al. [672] apply a cascaded classifier to learn geometric cues from the mesh and appearance cues from images. In contrast, Martinović et al. [443] avoid time-consuming conversions between 2D and 3D representations by training Random Forest classifiers on 3D features. Afterwards, they separate individual facades based on their semantic structure and impose weak architectural priors. Instead of imposing architectural priors, Gadde et al. [219] implement a sequence of boosted decision tree classifiers, stacked using auto-context features. They demonstrate that the system is fast at inference time and easily adapts to new datasets. Hackel et al. [261] propose a fast semantic segmentation approach for large 3D point clouds, which can also handle strongly varying densities. They construct approximate multi-scale neighborhoods by down-sampling the point cloud in order to generate a pyramid with decreasing density. This scheme allows extracting a rich feature representation that captures the geometry in a point’s local neighborhood such as roughness, surface orientation, and height over ground. A random forest classifier finally predicts class-conditional probabilities.

Image-based 3D semantic segmentation approaches like [599] lead to redundant computations due to the overlap of images used for reconstruction of the 3D model. Therefore, approaches directly working in the 3D space are usually more efficient. Riemenschneider et al. [554] exploit the inherent redundancy in the labeling of all overlapping images to further increase the efficiency

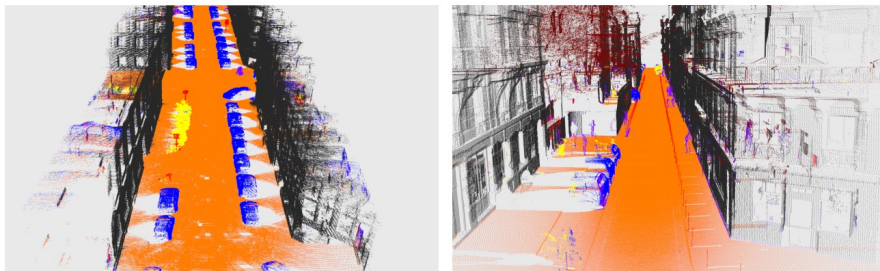


Figure 7.8: **3D Semantic Segmentation.** Semantic segmentation of two 3D scenes using the method of Hackel et al. [261] with facades (gray), ground (orange), cars (blue), motorcycles (yellow), traffic signs (red), pedestrians (violet) and vegetation (bordeaux). Figure courtesy of Hackel et al. [261] © 2016 ISPRS.

of image-based 3D semantic segmentation. They propose an approach that exploits the geometry of a 3D mesh model obtained from multi-view stereo to predict the best view for each face of the mesh before inferring the semantic class label. This allows them to accelerate their pipeline by two orders of magnitude, however, with lower accuracy than Martinović et al. [443].

Online Methods: While all aforementioned methods work in batch mode, i.e., they process all data at once, online methods allow the flexible incorporation of new measurements. This is particularly useful in the context of autonomous driving where new data arrives continuously. Towards online 3D semantic segmentation, Xiong et al. [733] train a sequence of classifiers to make predictions on different scales in a coarse-to-fine fashion (from regions to points). Predictions from the preceding scale are used as additional information for the current scale. They extend this work in [300] with a hierarchical representation of the 3D data and an improved inference procedure. Vineet et al. [678] propose an end-to-end system which processes data incrementally while performing real-time dense stereo reconstruction and semantic segmentation of outdoor environments. They achieve this using voxel hashing [490], a hash-table-driven 3D volumetric representation that ignores unoccupied space in the target environment. Furthermore, they employ an online volumetric mean-field inference technique that incrementally refines the voxel labeling and achieve real-time rates by harnessing the processing power of modern GPUs. McCormac et al. [451] present a pipeline for dense 3D semantic mapping designed to work online by fusing semantic predictions of a CNN with the geometric information from a SLAM system (ElasticFusion by Whelan et al. [704]). Specifically, ElasticFusion provides correspondences between 2D frames and a globally consistent map of surface elements or “sur-

fels”. Furthermore, they use a Bayesian update scheme which computes the class probabilities for each surfel based on the CNN’s predictions.

3D CNN: While convolutional networks have proven very successful in segmenting 2D images semantically, there exists comparably little work on labeling 3D data using convolutional networks. Maturana and Scherer [449] were one of the first to apply 3D Convolutional Neural Network (3D-CNN) for object recognition of volumetric 3D data. Their VoxNet approach classifies 32^3 voxel volumes using a convolutional neural network. In contrast, Huang and You [305] propose a framework to directly label 3D point cloud data using a 3D-CNN. Specifically, they compute 3D occupancy grids of size 20^3 centered at a set of randomly generated keypoints. The occupancy and the labels form the input to a 3D CNN which is composed of convolutional layers, max-pooling layers, a fully connected layer and a logistic regression layer. Towards processing larger volumes, Riegler et al. [553] propose OctNets, a 3D convolutional network, that allows for training deep architectures at significantly higher resolutions. They build on the observation that 3D data (e.g., point clouds, meshes) is often sparse in nature. OctNet exploits this sparsity property by hierarchically partitioning the 3D space into a set of octrees and applying pooling in a data-adaptive fashion. This leads to a reduction in computation and memory requirements as the convolutional network operations are defined on the structure of these trees. Thus, resources can be allocated dynamically depending on the structure of the input.

7.2.5 Road Segmentation

Segmentation of road scenes is a crucial problem in computer vision for autonomous driving. For instance, in order to navigate, an autonomous vehicle needs to determine the drivable area ahead and determine its own position on the road with respect to the lane markings. However, the problem is challenging due to the presence of a variety of differently shaped objects such as cars and people, different road types and varying illumination and weather conditions. Traditionally, the problem of autonomous driving has been tackled by detecting lane markings [61, 719, 394, 408]. However, as lane marking features are often not reliable (bad weather, construction sites, missing lane markings), more holistic approaches which consider the entire road area have been explored lately.

Alvarez et al. [5] propose a Bayesian framework to classify road sequences by combining low-level appearance cues with contextual 3D road cues such as the horizon, vanishing points, the 3D scene layout and 3D road models. In addition, they extract temporal cues for temporally smoothing the results. In follow-up work, Álvarez and López [7] convert the image into an illumination invariant feature space to make their method robust to shadows. Mansinghka

et al. [441] propose an inverse-graphics inspired method by employing generative probabilistic graphics programs (GPGP) to infer roads in images taken from vehicle-mounted cameras. GPGPs consist of a stochastic scene generator for generating random samples from a road scene prior, a graphics renderer for rendering the image segmentation of each sample and a stochastic likelihood model linking the renderer’s output and the data. Kuehnl et al. [364] present a method to improve appearance-based classification by incorporating the spatial layout of the scene. Specifically, they suggest a two-stage approach for road segmentation. First, they represent the road surface and delimiting elements such as curbstones and lane-markings using confidence maps based on local visual features. From these confidence maps, they extract SPatial RAY (SPRAY) features that incorporate global properties of the scene and train a classifier on those features. Their evaluation shows that spatial layout helps especially in cases where there is a clear structural correspondence between properties at different spatial locations.

Deep Learning: Recently, the problem of road segmentation has been addressed using convolutional neural networks [470, 494]. Mohan [470] propose a scene parsing system by using deconvolutional networks [764] in combination with traditional CNNs for feature learning. Deconvolutional networks learn features that capture mid-level cues such as edge intersections, parallelism and symmetry in image data and thus obtain a more robust representation. Oliveira et al. [494] investigate the trade-off between segmentation quality and runtime using U-Nets [557]. Specifically, they introduce a new mapping between classes and filters at the up-convolutional part of the network for reducing runtime. They further segment the entire image with a single forward pass, resulting in a more efficient approach compared to patch-based ones [470]. However, as road segmentation is a subproblem of semantic segmentation, today most state-of-the-art results on road segmentation are achieved using generic off-the-shelf semantic segmentation networks.

Data Acquisition: All existing algorithms for labeling road scenes are based on machine learning where the parameters of the respective model must be estimated from large annotated datasets. To alleviate the burden of annotating large datasets manually, Álvarez et al. [6] propose a method for road segmentation where noisy training labels for road images are generated using a convolutional neural network trained on a general image database. Laddha et al. [376] follow a different approach and obtain ground truth labels by exploiting OpenStreetMap information projected into the image domain using the vehicle pose provided by the GPS sensor.

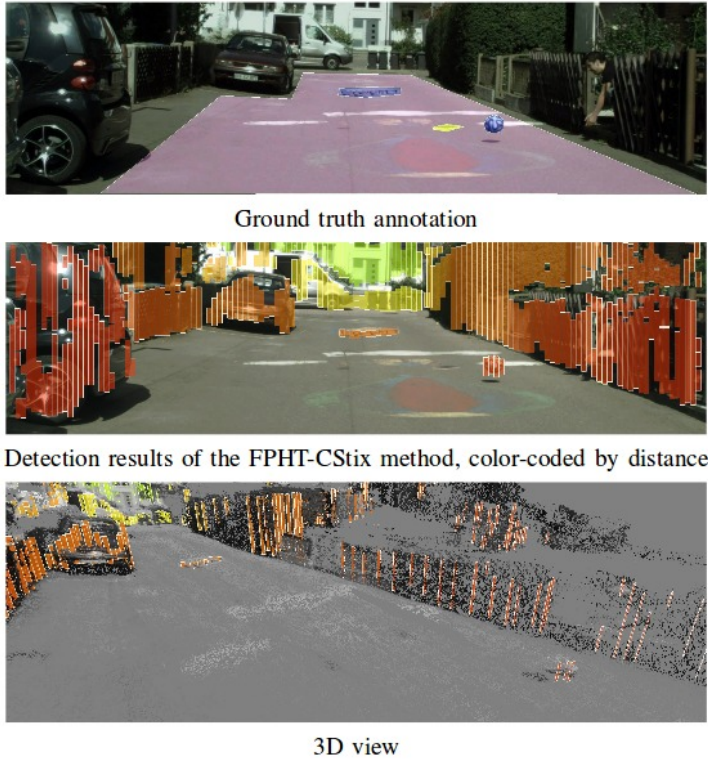


Figure 7.9: **Free Space Estimation.** Free space and detected obstacles on the Lost and Found dataset[515]. Figure courtesy of Pinggera et al. [515] © 2016 IEEE

7.2.6 Free Space Estimation

Accurate and reliable estimation of free space and the detection of obstacles are core problems that need to be solved for enabling autonomous driving. Free space is defined as the available space on the ground surface where navigation of vehicle is guaranteed without collision. Obstacles refer to structures that block the path of the vehicle by sticking out of the ground surface. In contrast to road segmentation approaches, methods estimating free-space in front of a vehicle often rely on geometric features which can be derived from a depth map computed from stereo sensors. However, both complementary approaches can be advantageously combined.

Badino et al. [20] propose a method for free space estimation by computing stochastic occupancy grids based on stereo information, where cells in

a stochastic occupancy grid carry information about the likelihood of occupancy. Stereo information is integrated over time in order to reduce depth uncertainty. The boundary between free space and occupied space is robustly obtained using dynamic programming on the occupancy grid. This work laid the foundations for the Stixel representation, see Section 7.2.7 for an in-depth discussion. While the original method of Badino et al. [20] makes the assumption of a planar road surface, this assumption is often violated in practice. In order to tackle more complicated road surfaces, Wedel et al. [697] propose an algorithm which models non-planar road surfaces using B-splines. The surface parameters are estimated from stereo measurements and tracked over time using a Kalman filter. In contrast, Suleymanov et al. [630] propose a complete pipeline to detect and drive on collision-free traversable paths, based on stereo information using a variational approach. In addition to free space detection, their approach also establishes a semantic segmentation of the scene, where labels include ground, sky, obstacles and vegetation.

Fisheye cameras discussed in Section 3.1.1 provide a wider field of view compared to regular cameras and allow for the detection of obstacles closer to the car. Häne et al. [268] propose a method for obstacle detection using monocular fisheye cameras. In order to reduce runtime, they avoid using visual odometry for accurate vehicle poses and instead, rely on less accurate pose estimates from wheel odometry. While they show good accuracy in the estimation of distances between objects, their experiments are limited to objects in close proximity to the sensor.

Long Range Obstacle Detection: The accuracy of obstacle detection at long-range is crucial for timely obstacle localization when the observer (i.e., the ego-vehicle) moves at high speed, e.g., in highways. Unfortunately, the error of stereo vision systems increases quadratically with depth in contrast to laser range sensors or radar sensors. In order to tackle this problem, Pinggera et al. [514, 515] propose long range obstacle detection algorithms using stereo vision. They formulate obstacle detection as a statistical hypothesis test, exploiting geometric constraints on camera motion and planarity. Independent hypothesis tests are performed on small local patches distributed across the input images. Detection results for a scene from their dataset are illustrated in Figure 7.9.

7.2.7 Stixels

Stixels are a compact mid-level representation of 3D traffic scenes with the goal to bridge the gap between pixels and objects [21]. The so-called “Stixel World” representation originates from the observation that free space in front of the vehicle is mostly limited by vertical surfaces. Stixels are represented by a set of rectangular sticks standing vertically on the ground to approximate



Figure 7.10: **Multi-layer Stixel World.** The multi-layer Stixel World representation of Pfeiffer and Franke [511]. The scene is segmented into planar segments termed “Stixels”. In contrast to the Stixel World of [21], objects are allowed to be located at multiple depths within a single image column. The color represents the distance to the obstacle with red being close and green far away. Figure courtesy of Pfeiffer and Franke [511] © 2016 BMVA.

these surfaces. Assuming a constant width, each stixel is defined by its height and its 3D position relative to the camera. The main goal of stixels is to gain efficiency through a compact, complete, stable, and robust representation. In addition, the stixel representation provides an encoding of the free space and the obstacles in the scene.

Using depth maps from SGM [294] as input, Badino et al. [21] use dynamic programming based on occupancy grids to compute free space, determining the stixels’ lower positions. Pfeiffer and Franke [511] extend [21] to a unified probabilistic scheme. They furthermore lift the constraint on stixels to touch the ground and allow multiple stixels for each image column, leading to a more flexible representation as illustrated in Figure 7.10.

In the dynamic stixel world representation introduced by Pfeiffer and Franke [510] the stixel representation was extended to dynamic scenes by tracking stixels using 6D Kalman filters based on optical flow. In contrast, Günyel et al. [258] show that motion estimation for stixels can be reduced to a 1D problem and can be solved efficiently via 2D dynamic programming, avoiding costly dense optical flow computation. Based on the dynamic stixel world representation, Erbs et al. [190, 189] present a CRF framework for semantically segmenting traffic scenes.

Several approaches proposed to leverage high-level information for inferring stixel representations more robustly. Cordts et al. [134] incorporate top-down object-level cues into the bottom-up stixel representation using a probabilistic approach. With the success of deep learning, Schneider et al. [585] present a semantic stixel representation to jointly infer the semantic and geometric layout of the scene from a dense disparity map and pixel-level semantic

scene labeling. Towards this goal, they used a deep learning-based scene labeling approach. In contrast, Levi et al. [401] propose StixelNet to directly infer the foot point of each stixel from the input image.

7.2.8 Aerial Images

The aim of aerial image parsing is the automated extraction of urban objects from data acquired by airborne sensors. The need for accurate and detailed information for urban objects such as roads is rapidly increasing because of its applications in the navigation of autonomous driving systems. For example, the output of aerial image parsing can be used to automatically build road maps (even in remote areas) and keep them up-to-date. Furthermore, information from aerial images can be used for localization. However, the problem is challenging because of the heterogeneous appearance of objects like buildings, streets, trees and cars which results in high intra-class variance but low inter-class variance. Furthermore, the complex structure of road networks and the difficulty of representing their geometry and topology accurately makes this problem hard. Roads must form a connected network of thin segments with slowly changing curvatures which meet at junctions. This type of prior knowledge is more challenging to formalize and integrate into a structured prediction formulation than standard smoothness assumptions.

Graphical Models: Graphical models have been a very popular way of addressing the problem of semantic segmentation in aerial images [700, 701, 472, 675, 447, 448, 699]. Wegner et al. [700] propose a CRF formulation for road labeling in which the prior is represented by cliques that connect sets of superpixels along straight line segments. Specifically, they formulate the constraints as high-order cliques with asymmetric P^N -potentials which express a preference to assign all rather than just a few of their constituent superpixels to the road class. This allows the road likelihood to be amplified for thin chains while still being amenable to efficient inference using graph cuts. Wegner et al. [701] also model the road network using a CRF with long-range, higher-order cliques. However, unlike [700], they allow for arbitrarily shaped segments which adapt to more complex road shapes by searching for putative roads with minimum cost paths based on local features. Montoya et al. [472] extend this formulation to multi-label classification of aerial images with class-specific priors for buildings and roads. In addition to the road network prior of [701], they introduce a second higher-order potential for cliques specific to buildings. In contrast, Verdie and Lafarge [675] propose the application of Markov point processes for recovering specific structures from images, including road networks. Markov point processes are a generalization of traditional MRFs which can address object recognition problems by directly manipulating parametric entities such as line segments. Importantly, they

implicitly solve the model-selection problem, i.e., they allow for an arbitrary number of variables in the MRF which can be associated with the parameters of the objects of interest.

Aerial Image Parsing using Maps: Instead of framing the problem of detecting topologically correct road networks as a semantic segmentation problem, Mattyus et al. [447] exploit map information from the free and community-driven mapping project OpenStreetMap (OSM)². Given a road map from OSM, Mattyus et al. [447] propose an MRF which reasons about the location of the road centerline and its width for each road segment in OSM. In addition, they incorporate smoothness between consecutive line segments by encouraging their widths to be similar. This formulation has the advantage of being efficient at inference time due to the restriction of the road topology to the input maps. However, it cannot recover from errors or missing information in the original map. Very recently, Facebook has announced a new set of tools³ that leverage AI to help the OSM community to build maps more efficiently.

Fine-grained Image Parsing: While aerial images provide full coverage of a significant portion of the world, they are of much lower resolution than ground images. In aerial imagery, the resolution relates to the ground area covered by one pixel. Whereas 1 meter resolution is already a high resolution for satellite imagery, the standard resolution for most publicly accessible image databases (e.g., Google Earth⁴) is 0.30 meter. Resolutions of 0.15 to 0.03 meter are considered high resolutions for aerial imagery and are usually not made publicly available. This makes fine-grained segmentation from aerial images a challenging problem. In contrast, ground images provide additional information which enables fine-grained semantic segmentation. Motivated by the complementary nature of these cues, several methods [448, 699] for fine-grained segmentation have been recently proposed which jointly reason about co-located aerial and ground image pairs.

Mattyus et al. [448] extend the approach of Mattyus et al. [447] by introducing a formulation that reasons about fine-grained road semantics such as lanes and sidewalks. To infer this information, they jointly consider monocular aerial images and high resolution stereo images captured from ground vehicles. Specifically, they formulate the problem as energy minimization in a MRF, inferring the number and location of the lanes for each road segment, all parking spots and sidewalks as well as the alignment between the ground and aerial images. Towards this goal, they exploit deep learning to estimate semantics from aerial and ground images and define potentials exploiting both cues. Wegner et al. [699] build a map of trees for urban planning applications from aerial images, street view images and semantic map data. They train

²<https://www.openstreetmap.org/>

³<https://mapwith.ai/>

⁴<https://www.google.com/earth/>

	Method	Coarse	Depth	IoU class	iIoU class	IoU category	iIoU category
1.	DRN_CRL_Coarse [795]	✓		82.8	61.1	91.8	80.7
2.	DPC [106]	✓		82.7	63.3	92.0	82.5
3.	RelationNet_Coarse [794]	✓		82.4	61.9	91.8	81.4
4.	SSMA [670]	✓	✓	82.3	62.3	91.5	81.7
5.	GFF-Net [410]			82.3	62.1	92.0	81.4
10.	DeepLabv3 [111]	✓		81.3	62.1	91.6	81.7
11.	AdapNet++ [670]	✓		81.3	59.5	91.0	80.1
12.	PSPNet [779]	✓		81.2	59.6	91.2	79.2
14.	ResNet-38 [723]	✓		80.6	57.8	91.0	79.1

Table 7.1: **CITYSCAPES Semantic Segmentation Leaderboard.** Segmentation performance is measured by class intersection-over-union and instance-level intersection-over-union. All methods are trained on the dense dataset consisting of 5000 frames and methods trained on the coarse dataset consisting of additional 20000 frames are marked in the corresponding column. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

CNN-based object detection algorithms on human-annotated data.

7.3 Datasets

There exist many large-scale realistic datasets for semantic segmentation as discussed in Chapter 4. The most popular datasets are PASCAL VOC [194], Microsoft COCO [420] and Cityscapes [133]. Recently, several companies also created new datasets which focus on the autonomous driving scenario such as Mapillary [487], ApolloScape [307] and Berkeley DeepDrive [755]. In addition, there exist several synthetic datasets for semantic segmentation, e.g., SYNTHIA [558] and Playing for data [551]. Here, we focus on the comparison of different semantic segmentation approaches on the popular Cityscapes dataset⁵ by Cordts et al. [133] as it is most relevant to the autonomous driving scenario. Cityscapes provides 5,000 images with high-quality dense annotations and 20,000 additional images with coarse labels obtained using a novel crowdsourcing platform.

In contrast to 2D semantic segmentation, there are only a few datasets that address the 3D semantic segmentation problem. Furthermore, these datasets are either very limited in size [480, 40, 260, 773] or in the number of classes [237]. Recently, Behley et al. [38] presented a large-scale dataset for 3D semantic segmentation.

⁵<https://www.cityscapes-dataset.com/>

7.4 Metrics

The performance of methods for semantic segmentation is usually evaluated using the intersection-over-union metric (IoU) which is defined as the number of true positive pixels divided by the sum over true positive, false positive and false negative pixels. As classes with larger segments will have a larger effect on the IoU score, Cityscapes [133] also report the instance-level intersection-over-union (iIoU) metric which weights the contribution of each true positive and false negative pixel by the ratio of the average instance size of the respective class with respect to the respective ground truth instance size. Cityscapes [133] report the IoU and iIoU metrics for two semantic granularities, i.e., classes and categories.

7.5 State of the Art on Cityscapes

Table 7.1 shows the leaderboard of Cityscapes for the pixel-level semantic labeling task. All methods are trained on the dense dataset comprising 5,000 densely annotated frames. Methods that are additionally trained on the coarse dataset with additional 20,000 frames are marked in the table. The state of the art in semantic segmentation shows very similar accuracy in terms of IoU and iIoU. Li et al. [410] extend the pyramid scene parsing network (PSPNet) of [779] with an advanced fusion mechanism. They propose Gated Fully Fusion modules which enable for every pixel to fuse only the relevant information from different feature maps. This allows better accuracy on fine-level details than PSPNet[779]. In contrast, Valada et al. [670] present a multi-modal fusion approach that fuses features extracted from images and depth. The feature extraction network is based on the full pre-activation ResNet-50 [282] with multi-scale residual units proposed in [671] as well as an efficient variant of Atrous Spatial Pyramid Pooling (ASPP) [111]. Zhuang et al. [794] follow a different approach and introduce a Relation Module that correlates features with their spatial neighborhood by shifting the features in four directions (left-right, top-down) using pre-defined offsets while passing them through Gated Recurrent Units. The features are extracted using a ResNet-like architecture [723] modified with dilated and deformable convolutions [149]. Zhuang et al. [795] extend this idea by exploiting additional offsets to correlate features over larger neighborhoods. This allows them to outperform all other methods on Cityscapes (Table 7.1).

Most existing network architectures for semantic segmentation are designed by the developer. Recently, a new line of work proposes to search for novel architectures in a properly defined search space. Chen et al. [106] address three dense prediction problems, i.e., street scene parsing, person-part segmentation and semantic image segmentation, with an architecture search.

They use Xception [122, 149, 112] as the backbone network and build a recursive search space from three popular operators, i.e., 1x1 convolution, 3x3 atrous convolution and average spatial pyramid pooling. Finally, they adapt a random search algorithm to explore the recursive search space. By evaluating 28K architectures on 370 GPUs, they find an architecture that achieves state-of-the-art performance on Cityscapes.

7.6 Discussion

The focus on multi-scale inference has led to impressive results in pixel-level semantic segmentation on Cityscapes. Today, the top methods on Cityscapes (Table 8.1) reach an impressive IoU of almost 83% over classes and 92% over categories. In contrast, the instance-weighted IoU still ranges around 63% over classes and 82% over categories. This indicates that semantic segmentation works well for instances covering large image areas but is still challenging for instances covering smaller regions which provide less information about the semantic label and require context reasoning. Furthermore, segmenting small, and possibly occluded objects is a challenging task which might benefit from accurate depth estimation. Recently, multi-modal fusion approaches leveraging depth data have shown great performance for indoor [279, 144] and outdoor [670] semantic segmentation, the latter achieving state-of-the-art performance on Cityscapes as discussed in the previous section. Furthermore, exploiting temporal correlations as in [371] has the promise to further improve semantic segmentation accuracy and temporal consistency.

Chapter 8

Semantic Instance Segmentation

8.1 Problem Definition

The goal of semantic instance segmentation is to simultaneously detect, segment and classify every individual object in an image. Unlike semantic segmentation, a solution to this task provides information about the position, semantics, shape, and count of individual objects, and therefore has many applications in autonomous driving.

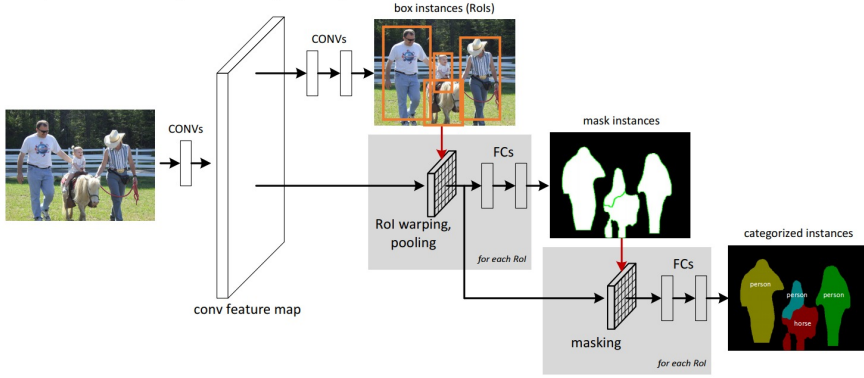
8.2 Methods

There exist two major lines of research for the task of semantic instance segmentation: Proposal-based and proposal-free instance segmentation. While proposal-based approaches usually consist of two steps, i.e., proposal extraction and proposal classification, proposal-free methods predict pixel labels directly from the image.

8.2.1 Proposal-based Approaches

Proposal-based instance segmentation methods extract class-agnostic proposals which are classified as an instance of a semantic class in order to obtain pixel-level instances. There exist several region proposal methods like Constrained Parametric Min-Cut (CMPC) [98], Multiscale Combinatorial Grouping (MCG) [17], DeepMask [516], and SharpMask [517] returning generic class-agnostic region proposals which can be directly used as instance segments. Several object detection classifiers were proposed which simultane-

Sequential Pipeline (MNC)



Joint Pipeline (Mask R-CNN)

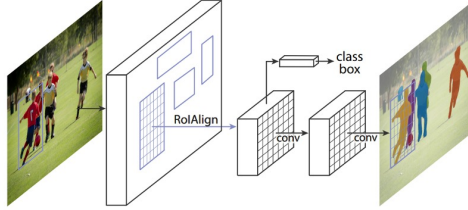


Figure 8.1: **Proposal-based Instance Segmentation Networks.** Architectures for proposal-based instance segmentation. Sequential pipelines as MNC [147] (upper row) use a detection and segmentation network sequentially. In contrast, joint formulations as Mask R-CNN [280] (bottom row) usually have the two networks in parallel, i.e., an object mask prediction and bounding box recognition network. Figure courtesy of Dai et al. [147] and He et al. [280] © 2016,2017 IEEE.

ously address object detection and semantic segmentation by leveraging region features from instance segments to improve the detection accuracy, i.e., O²P [97], Simultaneous Detection, and Segmentation (SDS) [271], Convolutional Feature Masking (CFM) [146], HyperColumn [270].

Proposal-based algorithms are slow at inference time due to the computationally expensive proposal generation step. To avoid this bottleneck, Dai et al. [147] propose Multi-task Network Cascade (MNC) a fully convolutional network with three stages illustrated in Figure 8.1. They extract box proposals, use shared features to refine these to segments, and finally classify them into semantic categories. The causal relations between the outputs of the stages complicate training of the multi-task cascade. In order to overcome

these difficulties, a fully differentiable mask prediction layer is presented to train the whole model in an end-to-end fashion. Box proposals can also induce errors into the proposal-based instance segmentation method due to wrongly scaled or shifted bounding boxes. In order to tackle this problem, Hayder et al. [277] present a shape aware object mask network that predicts a binary mask for each bounding box proposal, potentially extending beyond the box itself. They integrate the object mask network into the Multi-task Network Cascade framework of Dai et al. [147] by replacing the original mask prediction stage.

While earlier methods address the detection and segmentation problem with two sub-networks sequentially, recent work [412, 280, 108] propose to jointly address these problems. We illustrate an example of a sequential and joint formulation in Figure 8.1. All joint formulations use ResNet-like architectures [281] for feature extraction. Li et al. [412] propose FCIS, the first fully convolutional neural network for end-to-end instance semantic segmentation. They extend the fully convolutional mask proposal network [145] by sharing the convolutional representation of the proposals with a detection and segmentation sub-network. In contrast to FCIS, Mask R-CNN [280] and MaskLab [108] both build on Faster R-CNN [544]. He et al. [280] extend Faster R-CNN [544] by an additional branch for predicting segmentation masks. Chen et al. [108] combine box predictions from Faster R-CNN with semantic segmentation logits for pixel-wise classification and direction prediction logits estimating the direction towards instance centers. The direction towards instance centers allows them eventually to separate instances from the same class.

8.2.2 Proposal-free Approaches

Due to the problem of proposal-based approaches to inherit errors of the proposal generation, a number of alternative methods have been proposed recently. These methods jointly infer the segmentation and the semantic category of individual instances by casting instance segmentation directly as a pixel labeling task.

Several approaches [778, 777, 664] show how depth information can be used to identify different object instances. Zhang et al. [778, 777] train a fully convolutional neural network (FCN) to directly predict pixel-level instance segmentations of densely sampled image patches while the instance ID encodes a depth ordering. They improve the predictions and enforce consistency with a subsequent Markov Random Field. Uhrig et al. [664] propose an FCN to jointly predict semantic segmentation as well as depth and an instance-based direction relative to the centroid of each instance. This relative direction cue is then used for clustering pixels into individual instances. The instance segmentation pipeline is illustrated in Figure 8.2. However, all [778, 777, 664]

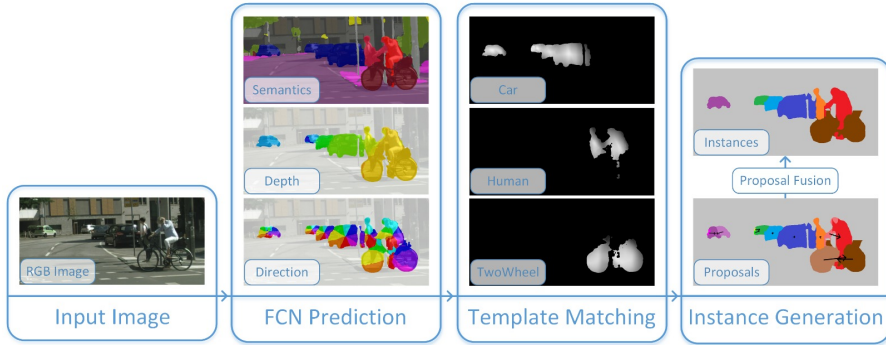


Figure 8.2: **Proposal-free Instance Segmentation Pipeline.** Uhrig et al. [664] predict semantics, depth, and instance center direction from the input image to compute template matching scores for all semantic maps. They fuse them after generating instance proposals to obtain an instance segmentation. Figure courtesy of Uhrig et al. [664] © 2016 Springer.

require ground-truth depth data for training their model.

Instead of relying on depth information, concurrent work [346, 27, 18] present proposal-free approaches based on an initial semantic segmentation. Kirillov et al. [346] combine semantic segmentation and object boundary detection via global reasoning in a multi-cut formulation to infer semantic instance segmentation. Bai and Urtasun [27] combine ideas from classical watershed transform with deep learning to create an energy map from an initial semantic segmentation and the input image where the basins correspond to object instances. This allows them to cut at a single energy level for obtaining a pixel-level instance segmentation. Arnab and Torr [18] propose to refine an initial semantic segmentation using an instance subnetwork. The initial category-level segmentation is used along cues from the output of an object detector within an unrolled Conditional Random Field [780] to predict pixel-level instances.

A new line of work is presented by Liu et al. [422]. They follow a sequential strategy with increasing semantic complexity. Several neural networks are applied in sequential order, each grouping pixels with different strategies starting by finding vertical and horizontal breakpoints, then connecting them to vertical and horizontal lines, grouping pixels in between these lines, and finally, extracting instances from the grouped pixels.

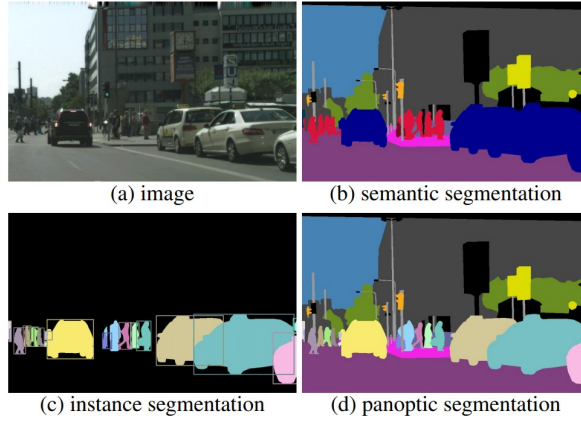


Figure 8.3: **Panoptic Segmentation.** Difference between semantic (b), instance (c) and panoptic segmentation (d). Figure courtesy of Kirillov et al. [345].

8.2.3 Panoptic Segmentation

Instance segmentation focuses on instances of objects and usually ignores classes that are not amenable to this task like sky or road as illustrated in Figure 8.3. In contrast, panoptic segmentation, firstly introduced by Kirillov et al. [345], addresses the dense estimation of a semantic label and instance id. Several approaches [345, 409, 137, 734, 344] have been proposed to address the problem.

Proposal-free instance segmentation approaches like [18, 27, 346] can be used directly to learn panoptic segmentation. However, ground truth for training them on this problem is very limited. Thus, Li et al. [409] use [18] in a semi-supervised fashion to learn panoptic segmentation. They use interactive foreground extraction (GrapCut) [560], proposal segmentation [528] and gradient-based localization of classes [598] to train the network in a semi-supervised fashion.

In contrast, several approaches [137, 734, 344] address panoptic segmentation with a joint semantic and instance segmentation formulation based on Mask R-CNN [280]. Costea et al. [137] propose to fuse object detections, semantic, and instance segmentation. They use semantic segmentation to distinguish between fore- and background regions. While the semantic class of background regions is directly obtained from the semantic segmentation, they use object detection, instance, and semantic segmentation to determine the class of foreground regions. In contrast, concurrent work [734, 344] ex-

Method	Coarse	COCO[420]	Depth	AP	AP 50%	AP 100m	AP 50m
1. PANet [423]		✓		36.4	63.1	49.2	51.8
2. UPSNet [734]				33.0	59.6	46.8	50.7
3. Mask R-CNN [280]		✓		32.0	58.1	45.8	49.5
4. PANet [423]				31.8	57.1	44.2	46.0
5. Mask R-CNN [280]				26.2	49.9	37.6	40.1
6. PolygonRNN++ [1]				25.5	45.5	39.3	43.4
7. SGN [422]	✓			25.0	44.9	38.9	44.5
8. Pixelwise Inst. Seg. with a DIN [18]	✓			23.4	45.2	36.8	40.9
9. Multitask Learning [336]				21.6	39.0	35.0	37.0
10. Deep Watershed Transformation [27]				19.4	35.3	31.4	36.8
11. Sem. Inst. Seg. with a DLF [65]				17.5	35.9	27.8	31.0
12. Boundary-aware Inst. Seg. [277]				17.4	36.7	29.3	34.0
13. InstanceCut [346]	✓			13.0	27.9	22.1	26.1
14. Foveal Vis. for Inst. Seg. of Road Images [496]			✓	12.5	25.2	20.4	22.1
15. Joint Graph Decomp. & Node Labeling [403]				9.8	23.2	16.8	20.3
16. Pixel-level Encoding for Inst. Seg. [664]			✓	8.9	21.1	15.3	16.7
17. R-CNN + MCG convex hull [133]				4.6	12.9	7.7	10.3

Table 8.1: **CITYSCAPES Instance Segmentation Leaderboard.** Instance detection performance is measured in terms of several average precision variants. The coarse annotations only provide rough class-level labels and are thus only used by a few methods. Since the dense annotations are quite limited, Microsoft COCO [420] is also sometimes used for training. More details in [133]. Accessed on: June 2019.

tends Mask R-CNN by an additional semantic segmentation branch removing the requirement of a heuristic fusion. Xiong et al. [734] combine a semantic segmentation network based on deformable convolutions with Mask R-CNN. They predict dense class labels by applying a softmax on concatenated channels of the semantic and instance segmentation networks. In contrast, Kirillov et al. [344] train the semantic and instance segmentation networks simultaneously without concatenating the channels. They apply non-maximum suppression [345] to avoid overlapping instances.

8.3 Datasets

Only a few datasets for instance segmentation exist. Microsoft COCO [420], consisting of 328K dense annotations, and Cityscapes [133], consisting of 5K dense and 25K sparse annotations, are the most popular datasets. While the KITTI [238] dataset also provides instance-level semantic annotations, the dataset consists of only 200 training and test scenes. The original PASCAL VOC [194] does not provide instance-aware annotations but Hariharan et al. [272] extended the dataset by semantic contours which are instance-aware. Still, the extension of PASCAL VOC is rarely used. The new datasets Mapillary [487], ApolloScape [307] and Berkeley DeepDrive [755] also provide instance-level annotations for 25K, 90K, and 10K images, respectively, but still need to prevail in the community.

Similar to semantic segmentation, we compare different methods on the

Cityscapes dataset¹ by Cordts et al. [133] because of the autonomous driving context, the online leaderboard and its acceptance in the community.

8.4 Metrics

The performance of instance segmentation methods is typically assessed by measuring average precision (AP) on instance regions that reach a certain overlap with ground truth regions. Usually, different thresholds are considered for the overlap and comparisons are performed according to the average over these thresholds as well as all classes. Cityscapes [133] use the same metric reported in Microsoft COCO [420] which considers 10 thresholds between 50% and 95%. In addition, the AP for an overlap value of 50 % (AP 50%) and for objects within 100 m and 50 m (AP 100m, AP 50m) are considered separately.

8.5 State of the Art on Cityscapes

In Table 8.1, we show the leaderboard of semantic instance segmentation methods on the Cityscapes dataset. The state of the art in instance segmentation is dominated by proposal-based approaches [1, 280, 734, 423]. However, they are closely followed by proposal-free approaches [27, 18, 422] with the sequential approach from Liu et al. [422] being the best performing proposal-free approach. While the proposal-based approaches [1, 280, 734] are built on Faster R-CNN [544], Liu et al. [423] propose a new feature hierarchy to propagate features from all levels (i.a. accurate localization signals from lower layers) to proposal sub-networks. They outperform all other methods with additional training on the Microsoft COCO dataset [420] since the dense annotations of Cityscapes are rather limited. The panoptic segmentation approach presented by Xiong et al. [734] is the best performing method when training is restricted to the dense annotations of Cityscapes.

8.6 Discussion

The instance segmentation task is much harder than the semantic segmentation task. Each instance needs to be carefully labeled separately whereas in semantic segmentation groups of one semantic class can be labeled together when they occur next to each other. In addition, the number and size of instances vary greatly between different images. In the autonomous driving context, often a wide view is present. Therefore, a large number of instances appear rather small in the image, making them challenging to detect. In contrast to bounding box detections discussed in Section 5.6, the exact shape of

¹<https://www.cityscapes-dataset.com/>

each object instance needs to be inferred in this task. Thus, the state of the art is still struggling on the Cityscapes dataset (Table 8.1) reaching an average precision of 36% or less. Proposal-based approaches which jointly address detection and segmentation with parallel sub-networks are currently the most promising direction. The joint formulation allows improving the generation of small instance proposals, which is important for segmenting instances in the context of autonomous driving.

Chapter 9

Stereo

9.1 Problem Definition

Stereo estimation is the process of extracting 3D information from passive 2D images captured by stereo cameras without the need for dedicated active light-emitting range measurement devices. In particular, stereo algorithms estimate depth information by finding correspondences between two images taken at the same point in time, typically by two cameras mounted next to each other on a fixed rig. These correspondences are projections of the same physical surface in the 3D world. Depth information is crucial for applications in autonomous driving or driver assistance systems. Accurate estimation of dense depth maps is a necessary step for 3D reconstruction, and many other problems such as obstacle detection, free space analysis, and tracking benefit from the availability of accurate depth estimates.

9.2 Methods

In stereo matching, the images from two cameras are usually projected onto a common parallel during rectification. This reduces the matching problem to a 1D search along the epipolar line, as illustrated in Figure 9.1, and the distance on this line is usually referred to as disparity.

The stereo literature can be separated into two groups. Feature-based methods [640, 583] provide only sparse depth maps, while dense methods generate dense outputs at the expense of computation time. In our survey, we focus on dense methods since they are more popular, and with the introduction of deep learning, they also became much more efficient. We further distinguish between local and global methods. Local methods compute the disparity by simply selecting the lowest matching cost, which is known as



Figure 9.1: **Stereo Matching Problem.** Visualization of the stereo matching problem. Given two rectified images (from KITTI training [238]), stereo matching reduces to a 1D search problem along the epipolar line (blue rectangle).

the winner takes all (WTA) solution [293, 640]. However, they usually result in very noisy estimates caused by ambiguities. In contrast, global methods formulate disparity computation as an energy-minimization problem integrating smoothness assumptions between neighboring pixels or regions [294, 240, 227, 264, 71, 372, 537]. Optimization can be carried out using variational approaches in the continuous domain and discrete approaches such as graph cuts or belief propagation for discrete label spaces.

9.2.1 Matching Cost

Stereo matching is a correspondence estimation problem where the goal is to identify the matching points between the left and right image based on a cost function. The algorithms usually assume rectified images, and the search space is reduced to a horizontal line (Figure 9.1). The matching cost computation is the process of computing a cost function at each pixel for all possible disparities, which is minimal at the true disparity. However, it is hard to design such a cost function in practice. Therefore stereo algorithms typically use the assumption of constant appearance between matching points. This assumption is often violated in real-world situations, such as cameras with slightly different settings causing exposure changes, vignetting, image noise, non-Lambertian surfaces, illumination changes, etc. Hirschmüller and Scharstein [293] systematically investigate the effect of these radiometric changes on commonly used matching cost functions, namely absolute differences, filter-based costs (Laplacian of Gaussian, Rank and Mean), hierarchical mutual information (HMI), and normalized cross-correlation. They found that the performance of a cost function depends on the stereo method that uses it. On images with simulated and real radiometric differences, rank filter performed best for correlation-based methods. For global methods, in tests

with global radiometric changes or noise, HMI performed best, while in the presence of local radiometric variations, Rank and Laplacian of Gaussian filters performed better than HMI. Qualitative results show that filter-based costs cause blurred object boundaries when used with global methods. None of the matching costs under consideration could succeed in handling strong lighting changes.

9.2.2 Energy Optimization

The inherent ambiguity in appearance-based matching costs can be overcome by regularization, i.e., introducing prior knowledge about the expected disparity map into the stereo estimation process. Therefore, an energy consisting of the matching cost and a smoothness constraint is usually optimized in contrast to WTA over the matching costs. The simplest prior favors neighboring pixels to take on the same disparity value (local smoothness).

Discrete Optimization: Discrete optimization methods optimize an energy with respect to a discrete set of disparities. While the resulting minimization problem is NP-hard, good approximations can be obtained using belief propagation [202] and graph cuts [64].

Semi-Global Matching (SGM) proposed by Hirschmüller [294] is the most prominent discrete optimization method for stereo matching. They hierarchically compute the matching cost by considering Mutual Information. A global smooth energy is approximated with cost aggregation by summing costs along 1D paths from multiple directions towards each pixel using dynamic programming. SGM became an influential stereo matching technique for autonomous driving due to its speed and high accuracy, as evidenced in various benchmarks such as Middlebury [581] and KITTI [238].

There are a few follow-up works investigating the practical and theoretical sides of SGM. Gehrig et al. [233] propose a real-time, low-power implementation of SGM with algorithmic extensions for automotive applications on a reconfigurable hardware platform. Drory et al. [176] offer a principled explanation for the success of SGM by clarifying its relation to belief propagation and tree-reweighted message passing. They show that SGM is equivalent to early stopping for a particular variant of belief propagation, effectively approximating the solution.

The performance of SGM can be further improved by incorporating a confidence measure. Seki and Pollefeys [597] leverage CNNs to predict confidence for stereo estimations. Taking into account ideas from conventional methods, they design a two-channel disparity patch which is used as input to a CNN. The first channel uses local smoothness, and the second enforces left-right consistency (disparity estimation using the other image should yield corresponding results). The confidences are incorporated into SGM by weighting

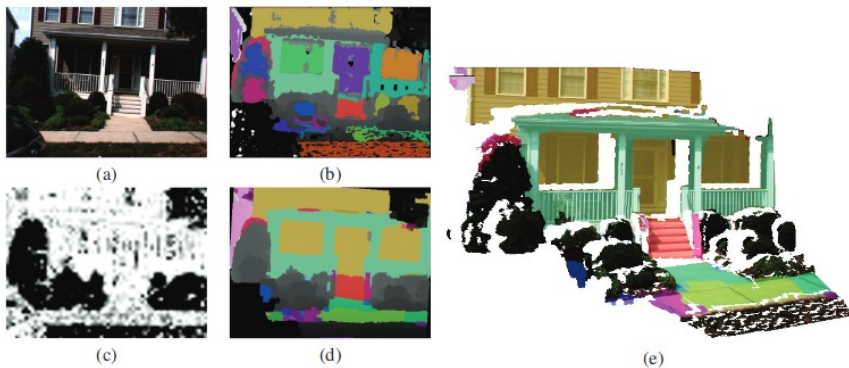


Figure 9.2: **Piecewise Planarity.** Gallup et al. [227] enforces piecewise planarity on planar structures of the scene found with RANSAC and plane classifiers. Plane candidates obtained with RANSAC (b), planar class probabilities (c), final plane assignment (d) and the 3D model with highlighted planes (e) are shown. Figure courtesy of Gallup et al. [227] © 2010 IEEE.

each pixel according to the estimated confidence.

Continuous Optimization: Variational approaches optimize the energy function with respect to continuous disparities. Data costs using the image intensities are usually non-convex and, thus, the global optimum can only be approximated. Coarse-to-fine approaches are used to handle large disparities by going from a low to a high resolution solution of the matching problem. For each resolution, the previous lower resolution solution is used as initialization. Coarse-to-fine approaches are typically used for optical flow estimation and will be discussed in detail in Chapter 11.

A commonly used smoothness prior is Total Variation (TV) [562] that penalizes the absolute difference between neighboring disparities. In the presence of weak and ambiguous observations, TV does not produce convincing results since it encourages piecewise constant disparities leading to stair-casing artifacts.

9.2.3 Higher-Order Models

Pairwise smoothness priors fail to reconstruct poorly-textured and slanted surfaces, as they favor fronto-parallel planes. A more generic approach to handle arbitrary smoothness priors is to exploit high-order correlations between pixels. Higher-order priors are able to express more realistic assumptions about depth images, but usually at additional computational costs.

Woodford et al. [715] introduce second-order priors for a graph cut stereo

formulation. While incorporating higher-order priors in discrete optimization has long been considered computationally infeasible, they propose an efficient optimization strategy for inference with triple cliques. In addition, they present an asymmetrical occlusion model that is combined with the second-order prior.

For continuous TV formulations, Haene et al. [264] introduce patch-based priors in the form of small, piecewise planar dictionaries. Total Generalized Variation (TGV) [71] is argued to be a better prior than TV, since it does not penalize piecewise affine solutions. However, it is restricted to convex data terms in contrast to TV, where global solutions can be computed even in the presence of non-convex data terms. Coarse-to-fine approaches often end up with a loss of details. In order to preserve fine details, Kuschik and Cremers [372] integrate an adaptive regularization weight into the TGV framework by using edge detection and report improved results compared to coarse-to-fine approaches. Ranftl et al. [537] obtain even better results by proposing a decomposition of the non-convex functional into two subproblems.

9.2.4 Piecewise Planar Priors

One common way to deal with slanted surfaces in the literature is to assume piecewise planarity. Geiger et al. [240] build a prior over the disparity space by forming a triangulation on a set of robustly matched correspondences, called support points. This reduces matching ambiguities and results in an efficient algorithm by restricting the search to plausible regions. Gallup et al. [227], illustrated in Figure 9.2, first train a classifier to segment an image into piecewise planar and non-planar regions. Afterwards, they enforce a piecewise planarity prior only on planar regions using plane hypotheses obtained from RANSAC. Non-planar regions are modeled by the output of a standard multi-view stereo algorithm.

9.2.5 Segmentation-based Models

An alternative way of modeling piecewise planarity is to explicitly partition the image into superpixels (groups of pixels) and modeling the surface at each superpixel as a slanted plane [738, 257]. However, care must be taken to ensure that the superpixelization is indeed an oversegmentation of the image with respect to planarity, i.e., no superpixel contains two surfaces that are not co-planar. Yamaguchi et al. [738] jointly reason about occlusion boundaries and depth in a hybrid MRF composed of both continuous and discrete random variables. Güney and Geiger [257] use a similar framework to incorporate object-category specific 3D shape proposals that regularize over larger distances. By leveraging semantic segmentation and 3D CAD models,

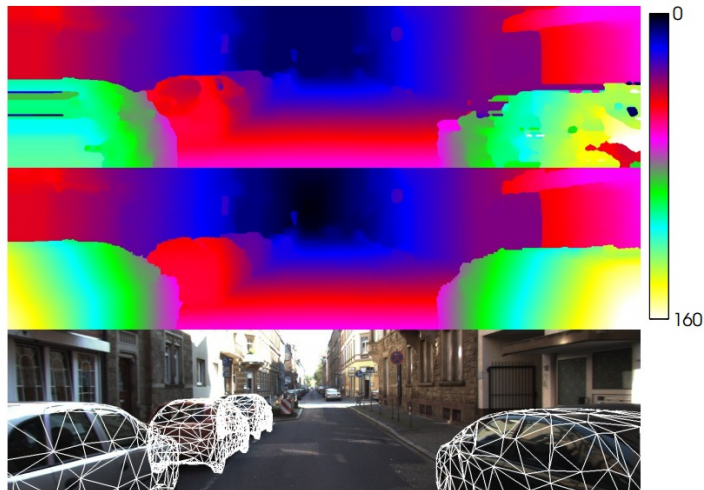


Figure 9.3: **Stereo Matching using Object Knowledge.** Stereo methods often fail at reflecting, textureless or semi-transparent surfaces (top [763]). By using object knowledge, Güney and Geiger [257] encourage disparities to agree with plausible surfaces (center). This improves results both quantitatively and qualitatively while simultaneously recovering the 3D geometry of the objects in the scene (bottom). The disparity is illustrated with a color coding shown on the right side. Figure courtesy of Güney and Geiger [257] © 2015 IEEE.

they resolve ambiguities in reflective and textureless regions originating from highly specular surfaces of cars in the scene, as shown in Figure 9.3.

9.2.6 Deep Learning for Stereo Matching

In the last years, deep learning approaches gained popularity in stereo estimation. While some methods try to learn richer feature representations [763, 432], others learn to directly predict a disparity map from the input stereo image pair [450, 339, 103].

For richer feature representations, Žbontar and LeCun [763] and Luo et al. [432] use a Siamese network that consists of two sub-networks with shared weights and a final score computation layer. The idea is to train the network for computing the matching cost by learning a similarity measure on small image patches. Žbontar and LeCun [763] define positive/negative examples as matching and non-matching patches and use a margin loss to train either a fast architecture with a simple dot-product layer in the end or a slow but more accurate architecture which learns score computation with a set of fully

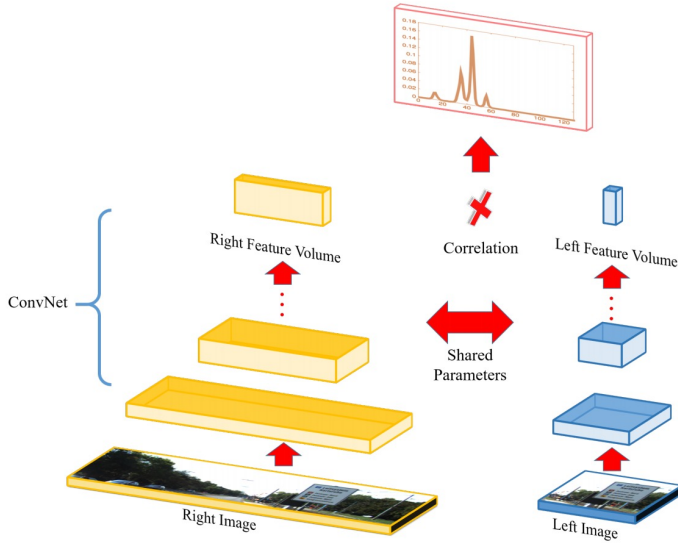


Figure 9.4: **Deep Learning for Stereo Matching.** A Siamese network is trained to extract marginal distributions over all possible disparities for each pixel. Figure courtesy of Luo et al. [432] © 2016 IEEE.

connected layers. Luo et al. [432] use a similar architecture, but formulate the problem as multi-class classification over all possible disparities to capture correlations between different disparities implicitly, as visualized in Figure 9.4. Both approaches rely on SGM [294] as a post-processing step to propagate information to neighboring pixels and estimate the final dense disparity map.

In contrast, Mayer et al. [450] adapt the encoder-decoder architecture proposed by Dosovitskiy et al. [174] for optical flow estimation (Chapter 11) to directly learn a model which predicts the entire disparity map at once without need for additional post-processing. The encoder computes abstract features while the decoder reestablishes the original resolution with additional cross-links between the contracting and expanding network parts for preserving the details. Post-processing and regularization are not necessary since the encoder-decoder architecture implicitly learns the entire mapping end-to-end. However, this architecture has to learn the concept of matching from scratch. Thus, inspired by [174], they also propose an alternative network (DispNetC) that first process each image independently and finally correlates extracted features from both images. Kendall et al. [339] combine ideas from previous methods, i.e., Siamese feature extraction, and cost volume formation. More specifically, they propose to extract deep feature representations

using a Siamese network and correlate these features to create a cost volume. After the cost volume, they use an encoder-decoder architecture to enlarge the receptive field and apply 3D convolutions on each encoder level. A differentiable soft argmin operation allows them to train the network end-to-end. Chang and Chen [103] introduce a spatial pyramid pooling and 3D CNN module to exploit more context information. Spatial pyramid pooling allows the extraction of richer features by taking larger regions into account. The 3D CNN module with multiple stacked encoder-decoder networks enables them to leverage global context information and achieve state-of-the-art performance.

Recently, it has been demonstrated that semantic information can also be exploited in the context of deep learning-based stereo estimation. Yang et al. [745] jointly formulate semantic segmentation and stereo estimation in one framework. This allows them to learn semantic cues and incorporate them into the disparity estimate by introducing a semantic softmax loss that regularizes the disparity with semantic cues. They show the benefit of their joint formulation in the unsupervised and in the supervised setting.

9.2.7 Variable Baseline

Stereo estimates can be fused to yield a more complete reconstruction of the static parts of the three-dimensional scene. However, assuming a fixed baseline, focal length, and field of view might not always be the best strategy. Gallup et al. [225] point out two problems with traditional stereo methods: dropping accuracy in the far range and unnecessary computation time spent in the near range. They, therefore, propose to use a multi-camera rig and to dynamically select the best cameras with the appropriate baseline for accurate estimation. In addition, they reduce the resolution to speed up the computation in the near range. In contrast to traditional fixed-baseline stereo, the proposed variable baseline stereo algorithm achieves constant accuracy over the reconstructed volume by evenly spreading the computation throughout the volume.

9.2.8 Omnidirectional Cameras

Omnidirectional sensors discussed in 3.1.1 allow to significantly increase the field of view for stereo matching. However, only limited work on stereo estimation using omnidirectional sensors exist. Häne et al. [267] extend the plane-sweeping stereo matching for fisheye cameras by incorporating the unified projection model for fisheye cameras [243] directly into the plane-sweeping stereo matching algorithm [226]. This kind of approach allows for producing dense depth maps directly from fisheye images in real-time using GPUs. Schönbein and Geiger [586] consider the stereo matching problem for catadioptric omnidirectional cameras.

9.3 Datasets

The most popular datasets for stereo estimation are the Middlebury [581, 582, 580] and KITTI [238] datasets. The ETH3D [591] also provides a two-view benchmark but is relatively new and does not focus on the autonomous driving scenario. Since only the KITTI dataset considers the autonomous driving context, we focus our attention on the KITTI benchmark.

Larger datasets are necessary for training deep models. In this case, the community relies on synthetic datasets such as SYNTHIA [558], Virtual KITTI [221], Flying Things [450] and Sintel [92]. However, the models trained on synthetic datasets are usually not generalizing to real datasets and need further fine-tuning on real datasets.

9.4 Metrics

Multiple metrics have been proposed to measure the performance of stereo approaches. The most popular measures are the root-mean-squared error (RMS) and outlier ratio, i.e., percentage of bad pixels (pixel with an error larger than a threshold). Typically, the average RMS is reported while the outlier ratio is often evaluated using several thresholds. Middlebury reports results for 0.5, 1, 2, and 4 pixels thresholds. In contrast, the KITTI benchmark reports the percentage of pixels with an error larger than 3 pixels or 5%. In addition, they separately evaluate the percentage of bad pixels over background and foreground regions.

9.5 State of the Art on KITTI

In Table 9.1 we show the ranking of stereo methods on the KITTI stereo 2015 benchmark. Tulyakov et al. [662] combine similar to [174, 339, 103] learning of the feature extraction, correlation, and regularization in an end-to-end trainable model. They extract deep features with a bottleneck architecture in contrast to a regular encoder-decoder network and propose a novel sub-pixel maximum a posteriori (MAP) approximation for inference based on the weighted mean around the disparity with maximum posterior probability. While the bottleneck architecture allows reducing the memory footprint, the sub-pixel MAP approximation enables to handle different disparity ranges than used for training. They achieve competitive results on KITTI. However, the spatial pyramid pooling and 3D CNN proposed by Chang and Chen [103], as discussed in Section 9.2.6, is computationally more efficient and improves significantly in the background regions. Yang et al. [745] jointly address the semantic segmentation problem to incorporate more contextual information as discussed in Section 9.2.3. While they reach similar performance on the

	Method	D1-bg	D1-fg	D1-all	Runtime
1.	EdgeStereo-V2 [622]	1.84 %	3.30 %	2.08 %	0.32s /
2.	Stereo-fusion-SJTU [622]	1.87 %	3.61 %	2.16 %	0.7 s /
3.	SegStereo [745]	1.88 %	4.07 %	2.25 %	0.6 s /
4.	PSMNet [103]	1.86 %	4.62 %	2.32 %	0.41 s /
5.	PDSNet [662]	2.29 %	4.05 %	2.58 %	0.5 s / 1 core
6.	SCV [431]	2.22 %	4.53 %	2.61 %	0.36 s /
7.	CRL [500]	2.48 %	3.59 %	2.67 %	0.47 s /
8.	GC-NET [339]	2.21 %	6.16 %	2.87 %	0.9 s /
15.	Displets v2 [257]	3.00 %	5.56 %	3.43 %	265 s / 8 cores
19.	MC-CNN-acrt [763]	2.89 %	8.88 %	3.89 %	67 s /
20.	PRSM [682]	3.02 %	10.52 %	4.27 %	300 s / 1 core
21.	DispNetC [450]	4.32 %	4.41 %	4.34 %	0.06 s /
41.	SGM_ROB [294]	5.06 %	13.00 %	6.38 %	0.11 s /

Table 9.1: **KITTI 2015 Stereo Leaderboard.** Numbers correspond to percentages of bad pixels according to the 3px/5% criterion defined in [455] in background (bg), foreground (fg) or all regions. The methods below the horizontal line are older entries, serving as reference. Accessed on: June 2019.

background regions, the joint formulation enables to improve also on the foreground regions. The best performance in foreground and background regions is achieved by Song et al. [622]. Similar to [745], they use a joint formulation, but instead of semantic segmentation, they jointly learn image edges using an edge-aware smoothness loss. In combination with a context pyramid to extract multi-scale features and one-stage residual pyramid returning a full-size disparity map, they outperform all other methods, as shown in Table 9.1. However, DispNetC presented by Mayer et al. [450] remains one of the fastest approaches while achieving competitive results on the foreground.

9.6 Discussion

Stereo estimation has shown great progress in the last years both in terms of accuracy and efficiency. However, some inherent problems prevent it from being considered solved. Stereo matching is equivalent to searching for correspondences in two images based on the assumption of constant appearance. However, appearance frequently changes due to non-rigidity or illumination changes. Furthermore, saturated pixels, occluded regions, or pixels leaving the frame cannot be matched. Therefore, failure in those cases is inevitable for methods that solely rely on appearance matching without any other prior assumptions about the geometry. We show the accumulated errors of the top 15 methods on the KITTI stereo benchmark [238] in Figure 9.5. The most

common examples of failure cases in the autonomous driving context are car surfaces that cause appearance changes due to their shiny and reflective nature. This problem can be addressed by leveraging more context information, e.g., using joint formulations [745, 622]. Similarly, windows that are reflective and transparent cannot be matched reliably. Occlusions are another source of error and require geometric reasoning beyond matching. Other examples of problematic regions include thin structures like traffic signs or repetitions as caused by fences. In these cases, continuous disparity estimation and the incorporation of more context information could be promising future directions.



Figure 9.5: **KITTI 2015 Stereo Analysis.** The averaged errors of the 15 best-performing stereo methods published on the KITTI 2015 Stereo benchmark. Red colors correspond to regions where the majority of methods fail according to the 3px/5% criterion defined in [455]. Yellow colors correspond to regions where some of the methods fail. Regions that are correctly estimated by all methods are transparent.

Chapter 10

Multi-view 3D Reconstruction

10.1 Problem Definition

The goal of multi-view 3D reconstruction is to infer 3D geometry from a set of 2D images by inverting the image formation process using appropriate prior assumptions. In contrast to two-view stereo, multi-view reconstruction algorithms recover the complete 3D shape of an object by inferring shape from many viewpoints.

In this survey, we focus on multi-view reconstruction from an autonomous driving perspective which mainly concerns the reconstruction of urban areas. The goal of urban reconstruction algorithms is to produce fully automatic, high-quality, dense reconstructions of urban areas by addressing inherent challenges such as lighting conditions, occlusions, appearance changes, high-resolution inputs, and large scale outputs. In the context of autonomous driving, 3D reconstructions can be used for static obstacle detection (traffic lights, road signs, etc.) and avoidance or precise localization as discussed in Section 13.3.

Musialski et al. [482] provide a survey of urban reconstruction approaches by following an output-based ordering which considers buildings and semantics, facades and images, and finally, city blocks and cities. They list ground, aerial, and satellite imagery, as well as Light Detection and Ranging (LiDAR) scans as the most commonly used sensor modality for urban reconstruction. Ground-level imagery is the most prevalent one due to its ease of acquisition, storage, and exchange. However, more and more aerial and satellite images become available today as well. In contrast to aerial or multi-view imagery, satellite imagery provides worldwide coverage at low costs, but also



Figure 10.1: **Large-Scale 3D Reconstruction.** Upper row: Coarse reconstruction and camera pose estimation (red) from the SfM [588] pipeline of COLMAP. Bottom row: Fine reconstruction with the MVS [589] pipeline of COLMAP. Figure courtesy of Schönberger and Frahm [588] and Schönberger et al. [589] © 2016 IEEE.

with low resolution. LiDAR delivers semi-dense 3D point clouds at high precision, both ground-level and aerial, but the sensor is expensive and the data is sparse. Some approaches [214, 58] also combine these data types in order to leverage their complementary strengths. Several methods [683, 259] leverage additional information, like Digital Surface Models (DSMs), which capture the Earth’s surface, to deal with challenging outdoor conditions. DSMs are 2.5D representations of an urban scene that provide a height for each surface point.

10.2 Structure from Motion

In Structure-from-Motion (SfM), the camera parameters (intrinsic and extrinsic) need to be estimated jointly with the 3D structure while in Multi-View Stereo (MVS), the camera parameters are assumed to be known. Furthermore, while MVS approaches create a dense 3D model of the object or scene of interest, SfM approaches typically recover a sparse 3D point cloud of the scene. Solving for the camera parameters and 3D geometry of the scene is equivalent to solving the correspondence problem based on a photo-consistency function that measures the agreement between different viewpoints. Typically, 3D reconstruction pipelines consist of an SfM method to estimate a coarse 3D reconstruction while recovering the camera parameters followed by an MVS

method to obtain a finer reconstruction, as illustrated in Figure 10.1 using COLMAP[588, 589].

Classical SfM pipelines [620, 619, 135, 474, 3, 718, 215, 638] first extract and match sparse features. Usually, an initial transformation between pairs of cameras (essential matrix) is estimated with RANSAC. Given the initial camera transformations, a geometric verification stage evaluates photometric consistency between re-projected sparse features and excludes outliers. Starting from an initial two-view reconstruction, an incremental reconstruction is performed based on best view selection, triangulation, and bundle adjustment. Due to this incremental approach, SfM pipelines are usually not very efficient and need to be applied offline. Simultaneous Localization and Mapping (SLAM) methods discussed in Section 13.4.3 also address the problem of joint camera estimation (ego-motion) and 3D scene reconstruction. However, SLAM techniques focus primarily on accurate ego-motion estimation and real-time performance, typically sacrificing geometric accuracy for these goals.

The web provides large amounts of publicly available imagery from cities taken by tourists that can be used to reconstruct popular buildings or even entire cities. This task requires a different approach than the ones mentioned earlier because of the large amount of images and the unknown geometric properties of the cameras the images have been taken with. Agarwal et al. [3] address this problem considering Flickr images of Rome. They use SIFT feature matching in combination with an efficient image retrieval approach to reduce the number of comparisons. Afterwards, a fast bundle adjustment method on minimal subsets of images captures the geometry of a scene. Finally, they optimize the whole pipeline in parallel, which allows them to reconstruct cities from 150K images in less than a day using 500 computing nodes. Frahm et al. [205, 206] present a highly efficient system for city-scale reconstruction from millions of images on a single computer by leveraging the high parallelization capabilities of graphics hardware. Recently, Schönberger and Frahm [588] proposed a structure-from-motion pipeline with better completeness and accuracy while better reducing drift in comparison to previous methods [620, 619, 3, 205, 206, 718, 215]. They further propose a more robust best view selection and triangulation method, producing more complete structures. Finally, a novel iterative Bundle Adjustment, re-triangulation, and outlier filtering step lead to significantly more complete and accurate 3D models.

10.3 Multi-view Stereo

Multi-view stereo approaches can be classified according to their scene representation into depth map-, point cloud-, mesh-, and volumetric-based meth-



Figure 10.2: **Point Cloud and Surface Representation.** The different steps of Patch-based Multi-View Stereo (PMVS) [217]. The input image, extracted features, reconstructed patches from the initial matching, reconstruction after expansion and filtering, and the final polygonal surface representation. Figure courtesy of Furukawa and Ponce [217] © 2010 IEEE.

ods. We first introduce and discuss classical approaches by grouping them based on their scene representation (Depth Maps, Point Clouds, Volumetric) and the final representation of the reconstruction (Mesh or Surfaces).

Depth Map: The depth map representation summarizes a 3D scene using one 2.5D depth map for each input view. These depth maps can later be fused into a single coherent 3D reconstruction using 3D fusion techniques [761, 141, 552]. One strategy which is particularly effective for recovering depth maps from urban scenes is the Plane Sweeping Stereo algorithm [131]. This algorithm “sweeps” a family of parallel hypothetical planes through the scene, projects images into each other via the homography induced by these planes, and evaluates photo-consistency. In very large scenes, one of the primary challenges is to handle large amounts of data efficiently. Pollefeys [524] propose a large scale, real-time MVS system based on the depth map representation by exploiting the parallel processing capabilities of modern GPUs.

Point Cloud: The reconstruction problem can also be addressed with a 3D point cloud representation Furukawa and Ponce [217] and Schönberger and Frahm [588]. Patch-based Multi-View Stereo (PMVS) Furukawa and Ponce [217] starts with a feature matching step to generate a sparse set of patches and then iterates between a greedy expansion step and a filtering step to make patches dense and remove erroneous matches. The steps of PMVS are visualized in Figure 10.2.

Volumetric Representation: Volumetric approaches represent geometry using a regularly sampled 3D grid, i.e., volume, either as a discrete occupancy function [373] or a function encoding distance to the closest surface (level-

set) [199]. More recent approaches use a probability map defined at regular voxel locations to encode the probability of occupancy [52, 523, 667]. The amount of memory required is the main limitation of volumetric approaches. There exists a variety of proposals for dealing with this problem, such as voxel hashing [490], data-adaptive discretization of the space in the form of Delaunay triangulation [375], or using octrees [263, 553].

Mesh or Surfaces: The final representation of a 3D reconstruction algorithm is typically a triangular mesh-based surface (right image in Figure 10.2). Volumetric surface extraction techniques can fuse multiple 2.5D measurements (MVS depth maps or laser scans) into a single, coherent 3D mesh model. Seminal work by Curless and Levoy [141] proposes an algorithm to accumulate surface evidence into a voxel grid using signed distance functions. The surface is implicitly represented as the zero crossing of the aggregated signed distance functions and can be extracted using the Marching Cube algorithm Lorensen and Cline [428] to label each voxel as either interior or exterior. Other approaches directly start from images [132, 157, 156] and refine a mesh model using an energy function composed of a data and a regularization term.

10.3.1 Planarity and Primitives

Man-made environments usually consist of regular structures. The introduction of appropriate priors, therefore, allows for more accurate and dense reconstructions. Micusik and Kosecka [461] present a method exploiting image segmentation cues as well as the presence of dominant scene orientations and piecewise planar structures. In particular, they adopt a super-pixel-based dense stereo reconstruction method exploiting the Manhattan world assumption in their MRF formulation. Another way of exploiting piecewise planar structures and repetitive shapes is to detect primitives such as planes, spheres, cylinders, cones, and tori [380, 382, 381]. Primitive-based approaches lead to compact and memory-efficient representations. However, their representations are often simplistic and fail to model fine details and irregular shapes. Therefore, Lafarge et al. [381] propose a hybrid approach that is both compact and detailed. Starting from an initial mesh-based reconstruction, they use primitives for regular structures such as columns and walls, while irregular elements are described using triangular meshes for preserving architectural details (Figure 10.3).

10.3.2 Shape Priors

Advances in sensors to acquire 3D shapes and the performance of object detection algorithms have encouraged the use of 3D shape priors in multi-

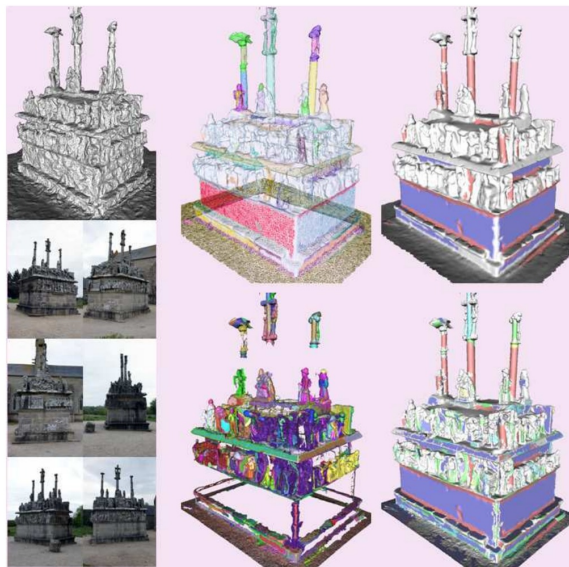


Figure 10.3: **Primitive-based Reconstruction.** The hybrid reconstruction approach of Lafarge et al. [381] uses primitives for regular structures (top right) and meshes for irregular structures (bottom left) to compactly represent a coarse initial mesh (top left). Figure courtesy of Lafarge et al. [381] © 2013 IEEE.

view stereo approaches. Dimensionality reduction is an effective and popular way of representing shape knowledge. Early approaches [659] use linear dimensionality reduction such as Principal Component Analysis (PCA) to capture shape variance in low dimensional latent shape spaces. More recent approaches, like Dame et al. [151], who investigate the importance of shape priors in a monocular SLAM approach, use non-linear dimensionality reduction techniques such as Gaussian Process Latent Variable Models (GP-LVM). In parallel with depth estimation, they refine an object’s pose, shape, and scale to match an initial segmentation and depth cues. Their experiments show improvements on transparent and specular surfaces, and even in unobserved parts of the scene.

In addition to the mean shape, Bao et al. [34] propose to learn a set of anchor points to represent object shape across several instances. They first perform an initial alignment of the mean shape to the point cloud from SfM using 2D object detectors. Finally, they warp and refine the mean shape to approximate the actual shape. Their evaluation demonstrates that the model is general enough to learn semantic priors for different object categories by

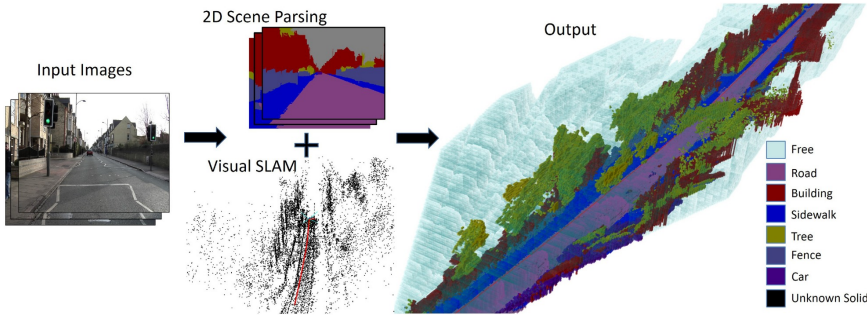


Figure 10.4: **Joint Reconstruction and Semantic Segmentation.** Joint 3D scene reconstruction and segmentation by Kundu et al. [369]. Figure courtesy of Kundu et al. [369] © 2014 IEEE.

handling large shape variations across instances.

An alternative to using latent space representations is to directly leverage 3D CAD models provided by free 3D model repositories. Güney and Geiger [257] propose a model for jointly inferring disparity maps and the geometry, pose, and type of 3D car models in urban scenes. Ulusoy et al. [666] extend this approach to the volumetric multi-view case. While the approaches mentioned earlier [151, 34] fit a parametric shape model to input data, Haene et al. [263, 262] model the local distribution of normals for an object. They also propose an object class-specific shape prior in the form of spatially varying anisotropic smoothness terms.

Zhou et al. [784] propose to jointly learn volumetric shape models for 3D reconstruction of street scenes from a sequence of fisheye cameras. Motivated by recurring objects of similar 3D shapes in outdoor scenes, they first localize buildings and vehicles using 3D object detectors and then jointly reconstruct them while learning a volumetric model of their shape. This allows the reduction of noise while completing missing surfaces as objects of similar shape benefit from all observations of the respective category. Instead of modeling a semantic prior for each object explicitly, Wei et al. [702] propose a data-driven regularization to transfer shape information from semantically matched patches in the training database using the SIFT flow algorithm.

10.3.3 Semantics

Similar to stereo, semantic information allows multi-view stereo approaches to recover from potential failures of photo-consistency in case of imperfect and ambiguous image information, e.g., specularities, lack of texture, repetitive

structures, or strong lighting changes. Semantic labels provide geometric cues about likely surface orientations at a certain location and help to resolve inherent ambiguities as illustrated in Figure 10.4 by the joint reconstruction and semantic segmentation approach of Kundu et al. [369].

Volumetric scene reconstruction typically segments the volume into occupied and free-space regions. Haene et al. [263] present the mathematical framework to extend this approach to multi-label volumetric segmentation, assigning object classes or a free-space label to voxels. They first learn appearance likelihoods and class-specific geometry priors for surface orientations from the training data. Afterwards, these data-driven priors are used to define unary and pairwise potentials in a continuous formulation for volumetric segmentation.

Haene et al. [263] require dense depth measurements, which can be difficult to obtain because of textureless regions and low parallax. Thus, Kundu et al. [369] propose another approach working on sparse 3D point clouds. They model the problem using a higher-order Conditional Random Field in 3D, which allows them to impose realistic scene constraints and priors such as 3D object support. In addition, they explicitly model free space, which provides cues to reduce ambiguities, especially along weakly supported surfaces. Their evaluation on the CamVid and Leuven datasets shows improved 3D structure compared to traditional SfM and state-of-the-art MVS pipelines as well as better segmentation quality over video segmentation methods.

Previous works on semantic reconstruction [263, 369] are limited to small scenes and low resolutions due to their large memory footprint and computational cost. In order to scale to larger scenes, Blaha et al. [55] note that high resolution is not required for large regions such as free space, parts under the ground, or inside the building. They propose an extension of Haene et al. [263] and employ an adaptive octree data structure with coarse-to-fine optimization to generate 3D city models from terrestrial and aerial images. Starting from a coarse voxel grid, they solve a sequence of problems in which the solution is gradually refined near the predicted surfaces. The adaptive refinement saves memory and runs much faster while still being as accurate as the fixed voxel discretization at the highest target resolution, both in terms of geometric reconstruction and semantic labeling. Besides the spatial extent, the number of different semantic labels is also problematic for scalability due to the increasing memory requirements. Cherabier et al. [117] propose to divide the scene into blocks in which only a set of relevant labels is active. Thus, the absence of semantic classes from a specific block can be determined early on. Accordingly, they deactivate labels from the beginning of the optimization, which leads to more efficient processing.

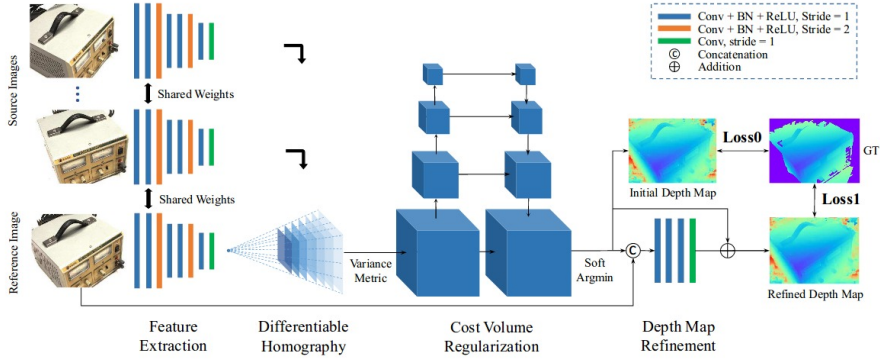


Figure 10.5: **Deep Learning for Multi-View Stereo.** MVSNet by Yao et al. [749] comprising feature extraction networks and a differentiable homography warping stage for constructing cost volumes. The final depth map is obtained using a refinement network. Figure courtesy of Yao et al. [749] © 2018 IEEE.

10.3.4 Efficient Reconstruction

The extraction of detailed 3D information from video streams leads to high computational costs for multi-view stereo algorithms. Cornelis et al. [135] focus on creating compact, memory-efficient 3D city models from a stereo pair at high frame rates based on simplified geometry assumptions such as ruled surfaces for building facades. Since objects such as cars violate these assumptions, they integrate the detection and localization of cars into the reconstruction. In contrast, Geiger et al. [241] propose an efficient stereo matching algorithm to generate accurate piece-wise planar 3D reconstructions in real-time.

10.3.5 Deep Learning for Multi-View Stereo

Several learning-based approaches have been proposed to address the Multi-View Stereo problem. In early works, learning was mainly used to obtain more robust feature representations for establishing better correspondences [266, 763, 432]. Pairwise similarities obtained from these approaches are usually averaged in order to match features from multiple images. In contrast, Hartmann et al. [275] propose to directly learn a matching function using multiple images as input.

Recently, several pipelines for end-to-end learning of multi-view reconstruction [333, 321, 502, 306] have been presented combining learned high-level information with classical constraints. Kar et al. [333] and Ji et al. [321]

propose to unproject features along the viewing rays onto a 3D feature grid for matching. Afterwards, both approaches use 3D convolutional networks to smooth the 3D feature grid. In contrast, Huang et al. [306] propose to learn disparities by combining a plane-sweep approach with an end-to-end trained CNN for feature extraction. Considering a reference view, they create plane-sweep volumes consisting of neighboring views warped according to the hypothetical depth. Afterwards, they extract features using a network on patch pairs (patches from the reference view and plane-sweep volumes). An encoder-decoder network with skip connections is used to combine features over larger regions and, eventually, the disparity map is estimated with a max-pooling layer. While previous approaches leverage physical constraints by projecting features according to the camera transformation, they do not model occlusion relationships. Paschalidou et al. [502] combine learning-based feature extraction with a Markov Random Field that employs high-order ray-potentials [667] to model the image formation process and occlusions.

Inference with 3D feature grids [333, 321], using a voxel grid [502] or plane-sweep volumes [306], is computationally expensive. A more efficient MVS reconstruction approach was presented by Yao et al. [749], see Figure 10.5. Similar to Huang et al. [306], they decouple the MVS reconstruction problem into a depth prediction problem for each view. A differentiable homography warping operation allows them to encode the camera geometry and to build a 3D cost volume. A 3D convolutional network predicts the depth from the 3D cost volume, and the final reconstruction is obtained with a depth map fusion approach [459] which minimizes depth occlusions and differences between viewpoints.

10.3.6 Omnidirectional Cameras

While omnidirectional cameras, as discussed in Section 3.1.1, provide a larger field of view compared to traditional perspective cameras, their special geometric properties need to be addressed during 3D reconstruction. The epipolar geometry of central catadioptric systems was explored by Svoboda and Pajdla [637], who showed that correspondences lie on epipolar conics and who proposed a rectification procedure for this setup. In contrast, Bunschoten et al. [91] and González-Barbosa and Lacroix [249] propose to project the omnidirectional image to a panoramic view and use standard stereo matching methods to search for correspondences. While Bunschoten et al. [91] search on sinusoidal shaped epipolar curves, González-Barbosa and Lacroix [249] rectify the panoramic view to obtain straight epipolar lines. Schönbein and Geiger [586] propose a method for 3D reconstruction through joint optimization of disparity estimates from two temporally and two spatially adjacent omnidirectional views in a unified omnidirectional space using plane-based priors.

10.4 Datasets

Several datasets have been proposed to evaluate multi-view stereo algorithms. Popular datasets include Middlebury [581] and DTU MVS [319]. However, these datasets provide only a few or no examples for urban reconstruction. While EPFL Multi-View [627], Restrepo et al. [547], ETH3D [591] and Tanks and Temples [351] provide urban scenes, they do not focus on the autonomous driving task.

Large-scale reconstruction methods [3, 205, 206, 588] typically use the BigSFM dataset ¹, a collection of smaller datasets from Cornell University which consists of Vienna [313], Dubrovnik [411], Rome and Quad datasets [139]. However, these datasets do not have ground truth data and, therefore, a quantitative evaluation of methods is not possible.

As ETH3D and Tanks and Temples are the MVS datasets closest to the autonomous driving scenario and also provide an online evaluation server, we focus our discussion on these two datasets. As opposed to ETH3D [591], Tank and Temples [351] does not provide camera poses and thus an additional structure-from-motion pipeline [619, 718, 474, 215, 638, 588] is necessary to estimate camera poses.

10.5 Metrics

In MVS, the accuracy and completeness of the output reconstruction are standard measures for evaluation. Accuracy is defined as the percentage of estimated points with a distance smaller than a predefined threshold to the closest ground truth points. Completeness is defined as the percentage of ground truth points with a distance smaller than a predefined threshold to the closest estimated points. Some benchmarks also report the mean (Chamfer) or harmonic mean (F1-measure) of accuracy and completeness.

10.6 State of the Art on ETH3D & Tanks and Temples

In Table 10.1, Table 10.2, we show the leaderboards for the intermediate as well as advanced scenes of ETH3D [591] and Tanks and Temples [351], respectively. Both benchmarks use the F1-measure for comparison.

COLMAP [589, 588] jointly models pixel-level view selection and depth estimation using a graphical model. They incorporate geometric as well as temporal priors for improved view selection and a geometric consistency for simultaneous depth/normal estimation with a PatchMatch sampling scheme.

¹<http://www.cs.cornell.edu/projects/bigsfm/>

Method	All	Low-Res Many-View			High-Res Multi-View		
		All	Indoor	Outdoor	All	Indoor	Outdoor
1. ACMM [736]					80.78	79.84	83.58
2. OpenMVS [101]	72.83	56.18	45.66	63.19	79.77	78.33	84.09
3. LTVRE_ROB [366]	69.57	53.52	45.46	58.89	76.25	74.54	81.41
4. ACMH [736]	67.68	47.97	38.24	54.45	75.89	73.93	81.77
5. COLMAP_ROB [588, 589]	66.92	52.32	42.45	58.89	73.01	70.41	80.81
6. OpenMVS_ROB [101]	64.09	48.56	38.68	55.15	70.56	68.19	77.65
7. CMP-MVS [317]	51.72	7.38	0.03	12.27	70.19	68.16	76.28
8. Gipuma [224]					45.18	41.86	55.16
9. PMVS [217]	37.38	21.09	11.49	27.48	44.16	40.28	55.82
10. MVE [215]	26.22	16.26	16.97	15.79	30.37	25.89	43.81

Table 10.1: **ETH3D Leaderboard.** Evaluation results on two ETH3D [591] challenges: low-resolution multi-view stereo from video data (many-view) and high-resolution multi-view stereo on few images recorded with a DSLR. The average F-measure is reported. Accessed on: May 2019.

COLMAP achieves competitive results on both benchmarks and is considered one of the leading open MVS methods today, serving as the backbone for several other techniques. In contrast to COLMAP, MVSNet [749] discussed in Section 10.3.5 learns depth map inference with an end-to-end deep learning architecture. For unstructured image sequences like Tanks and Temples, they obtain the depth range and camera trajectory using OpenMVG [474]. Yao et al. [750] extend this work with a recurrent version called R-MVSNet which replaces 3D convolutions applied on the cost volume for regularization. While both methods were not evaluated on ETH3D, R-MVSNet improves on the intermediate and advanced scenes from Tanks and Temples. The best performing reconstruction method, ACMM [736], proposes an adaptive checkerboard sampling scheme and a multi-hypothesis joint view selection approach (ACMH) for improved propagation of hypotheses and pixel-wise view selection. In addition, they propose a multi-scale geometric consistency guidance scheme (ACMM) for improved depth estimation in low textured regions. In contrast to Yao et al. [749, 750], they use COLMAP’s structure-from-motion method [588] to obtain the camera trajectory.

The runtime of methods improved significantly with the introduction of learning-based methods. MVSNet [749] is currently the fastest approach, with 230 seconds per scan on the DTU MVS evaluation set [319]. The authors used the DTU dataset to compare the runtime and report a large speedup in comparison to COLMAP [589, 588]. ACMM [736] compares their runtime to COLMAP on Tanks and Temples and achieves a 3-fold speed up.

Method	Intermediate		Advanced	
	Rank	F-measure	Rank	F-measure
ACMM [736]	1.	57.27	1.	34.02
ACMH [736]	2.	54.82	2.	33.73
Dense R-MVSNet [750]	3.	50.55	3.	29.55
R-MVSNet [750]	4.	48.40	5.	24.91
MVSNet [749]	5.	43.48		
COLMAP [588, 589]	6.	42.14	4.	27.24
VisualSfM [718] + PMVS [217]	14.	27.80	15.	10.22
VisualSfM [718] + CMP-MVS [317]	18.	22.40	17.	7.57
Bundler [620, 619] + PMVS [217]	19.	12.86	18.	5.61

Table 10.2: **Tanks and Temples Leaderboard.** Evaluation results for intermediate and advanced scenes from Tanks and Temples [351]. The rank and average F-measure are reported. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

10.7 Discussion

In the last decade, great advances have been made in multi-view reconstruction, as can be observed from Tables 10.1 and Table 10.2. The current state of the art significantly improves upon classical approaches like PMVS [217] and CMP-MVS [317] on all benchmarks. However, the performance on low-resolution images of ETH3D and all scenes from Tank and Temples is still far from perfect. While great advances have also been made for large-scale reconstruction, a unified benchmark that considers the autonomous driving/mapping task is still missing.

An open question that remains for the autonomous driving problem is what kind of accuracy and completeness are necessary to realize safe mapping, localization, and navigation. For localization (Section 13.3) and loop-closure detection (Section 13.4.2) high accuracy is required. In contrast, for obstacle avoidance, high completeness is necessary in order not to miss any obstacle.

Chapter 11

Optical Flow

11.1 Problem Definition

Optical flow is defined as the two-dimensional motion of brightness patterns between two images. This definition only represents the motion of intensity patterns in the image plane but not the 3D motion of the objects in the scene. Recovering the 3D motion itself is the goal in Scene Flow discussed in Chapter 12.1. Figure 11.1 shows the synthetic Yosemite sequence with the optical flow ground truth generated by texture mapping aerial images of Yosemite valley onto an approximate mesh model.

Optical flow provides essential information about the scene and serves as input for several tasks such as ego-motion estimation (Chapter 13), structure from motion (Chapter 13), and tracking (Chapter 6). Research on this problem started several decades ago with the variational formulation by Horn and Schunck [298], assuming the brightness of a pixel to be constant over time. Despite the long history of the optical flow problem, occlusions, large displacement, and fine details are still challenging for modern methods. A fundamental problem with the optical flow definition is that besides the actual motion of interest, illumination changes, reflections, and transparency can also cause intensity changes. In contrast to stereo, the search space for finding correspondences is two-dimensional in the case of optical flow.

11.2 Methods

Traditionally, the optical flow problem has been approached with a variational formulation. Variational methods minimize an energy comprising a data term, assuming little appearance change over time, and a smoothness term, encouraging similarity between spatial neighbors. Horn and Schunck

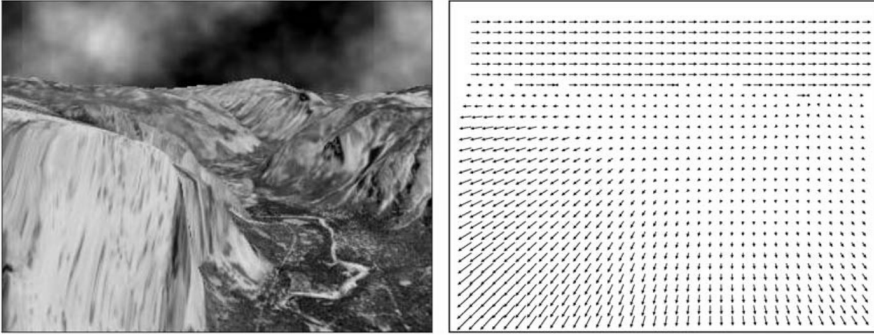


Figure 11.1: **Optical Flow Problem.** The Yosemite sequence generated by Quam [532] and the corresponding ground truth flow created by Heeger [286]. The sequence was later incorporated into the Middlebury dataset of Baker et al. [29]. Figure courtesy of Heeger [286] © 1988 Springer.

[298] introduced the brightness constancy assumption which models the intensity value of a pixel as constant over time. Considering a single pixel in isolation, this assumption yields one equation with two unknowns, which does not result in a unique solution (known as the aperture problem). Additional constraints must, therefore, be introduced in order to solve the aperture problem and estimate optical flow. A common way of regularizing variational optical flow estimation is to encourage similarity of spatially neighboring flow vectors. This prior assumption is motivated by the fact that flow fields are often smooth and discontinuities typically occur only at object boundaries.

The original formulation [298] uses a quadratic penalty function in the data and smoothness term. However, a quadratic penalty cannot handle frequent violations of brightness constancy assumption, e.g., due to varying illumination conditions. One way to alleviate this problem is to use a robust penalty function, as proposed by Black and Anandan [54]. In addition, several different data terms have been proposed that are less affected by illumination changes. Vogel et al. [681] systematically evaluate pixel- and patch-based data costs on the KITTI dataset [238]. On real data, they found patch-based terms to perform better than pixel-based terms.

Flow discontinuities frequently occur near motion boundaries caused by objects moving in front of each other. The original formulation by Horn and Schunck [298], cannot handle these discontinuities due to a homogeneous, non-robust smoothness term. Total Variation regularization used in Zach et al. [759] replaces the quadratic penalization by the L_1 norm to preserve discontinuities in the flow field. However, like the original formulation by Horn and Schunck, this model also biases the solution towards fronto-parallel surfaces

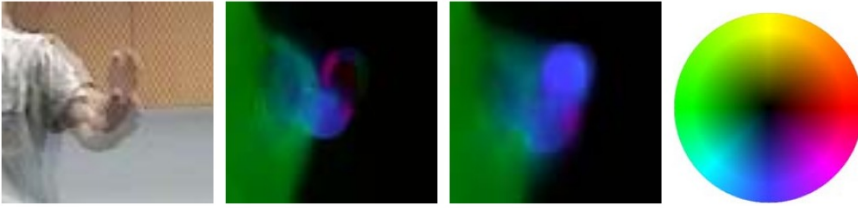


Figure 11.2: **Sparse Matching Guidance.** The fast hand motion (left) is an example where classical warping methods fail (center left), but sparse matches introduced by Brox and Malik [81] help to estimate the flow (center right). The color encoding with the hue and intensity representing the orientation and magnitude of the flow, accordingly, is visualized in the right image. Figure courtesy of Brox and Malik [81] © 2011 IEEE.

leading to artifacts in the estimation results, in particular in the presence of strongly slanted planes (e.g., the road surface). Thus, higher-order regularizations like the Total Generalized Variation (TGV) model have been proposed [71]. TGV priors can better represent real data as they leverage a piecewise affine motion model. The non-local Total Generalized Variation [536] is an extension of this model that enforces the piecewise affine assumption in a local neighborhood. This allows them to improve performance in regions where the data term is ambiguous in comparison to TGV which considers only direct neighbors. Zimmer et al. [799] provide a detailed assessment of image- and flow-driven regularizers for the variational formulation and discuss the qualities of different data terms.

Besides the model specifications, the choice of the optimization method and its implementation are additional factors that influence the performance of variational optical flow estimation algorithms. A detailed study of optical flow methods is provided by Sun et al. [631]. They investigate the most critical factors for the success of optical flow methods and propose an approach optimizing a classical formulation with modern techniques.

11.2.1 Sparse Matches

Linear approximations that are used to obtain the optical flow equation hold only for pixel motion. Therefore, variational methods cannot handle large displacements without an additional strategy. In variational formulations, this problem is typically addressed with a coarse-to-fine strategy, estimating the flow on a coarser resolution to initialize the estimation on a finer resolution. While this strategy works for large structures of little complexity by capturing

the dominant motion in the scene, fine geometric details are often lost in the process. Besides, textural details important for correspondence estimation are lost at coarse resolutions, hence leading the optimizer to a local minimum. One example of the loss of fine details is illustrated in Figure 11.2, which shows the optical flow field of a fast-moving hand. These problems can be alleviated by integrating sparse feature correspondences into the variational formulation, as proposed by Brox and Malik [81]. The feature matches, obtained from a nearest neighbor search on a coarse grid, are used as a soft constraint in a coarse-to-fine optimization. In Figure 11.2, the classical formulation fails to recover the optical flow for the hand, while integrating feature matches guides the optimizer to a better solution.

Another solution for large displacements is proposed by Revaud et al. [548]. They replace the coarse-to-fine strategy with an interpolation of sparse matches to initialize a dense optimization at full resolution. Sparse matches are obtained using DeepMatching, a deep neural network matching approach introduced by Weinzaepfel et al. [703]. In contrast to DeepMatching, Menze et al. [456] use approximate nearest neighbor search to generate a set of proposals as candidates to be used in a discrete optimization framework. The inference is made feasible by restricting the number of matches to the most likely ones and by exploiting the truncated form of the pairwise potentials. Motivated by the success of Siamese networks in stereo [763] (Chapter 9), Güney and Geiger [256] extend this work to learning features for 2D patch matching. They further investigate the importance of the receptive field size exploiting dilated convolutions as proposed by Yu and Koltun [754] for semantic segmentation. Chen and Koltun [113] argue that the heuristic pruning used to make inference feasible destroys the highly regular structure of the space of mappings and propose a discrete optimization over the full space. Min-convolutions are used to reduce the complexity and to effectively optimize the large label space using a modified version of Tree-Reweighted Message Passing by [353].

Wulff and Black [724] present a different approach to obtain dense optical flow from sparse matches. In their approach, the optical flow field is represented as a weighted sum of basis flow fields learned from reference flow fields which have been estimated from Hollywood movies. They estimate the optical flow by finding the weights which minimize the error with respect to the detected sparse feature correspondences. While this results in overly smooth flow fields, the so-called PCA Flow approach is very fast compared to variational and discrete optimization methods. A slower but more accurate version is also proposed to better handle flow discontinuities by using a layered approach.

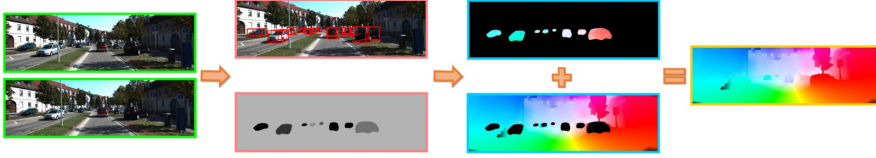


Figure 11.3: **Epipolar Flow.** The full pipeline of Bai et al. [26] first segments the scene into dynamic objects (cars) and the static background. Afterwards the motion is estimated for each object and background, independently, and finally combined to one flow field. Figure courtesy of Bai et al. [26] © 2016 IEEE.

11.2.2 Epipolar Flow

In the context of autonomous driving, application-specific assumptions can be made to simplify the optical flow estimation. The assumption of a static scene or the decomposition of a scene into rigidly moving objects allows for treating optical flow as a matching problem along epipolar lines radiating from the focus of expansion. Yamaguchi et al. [737] propose a slanted-plane Markov random field that represents the epipolar flow of each segment with slanted planes. This formulation needs a time-consuming optimization, which can be avoided with the joint stereo and flow formulation of Yamaguchi et al. [739]. They assume the scene to be static and present a new semi-global block matching algorithm using the joint evidence of stereo and video.

11.2.3 Semantic Segmentation

Scenes in the context of autonomous driving are usually composed of a static background and dynamically moving traffic participants. This observation can be exploited by splitting the scene into independently moving objects. Bai et al. [26] extract traffic participants using instance-level segmentation and estimate the optical flow independently for different instances. Similar to [737, 739], they use the slanted plane model but only for background flow estimation. For each moving object, an independent epipolar flow estimation is performed, as illustrated in Figure 11.3. Sevilla-Lara et al. [604] use semantic segmentation for optical flow estimation. First, semantics provide information on object boundaries and spatial relationships between objects that can be exploited to reason about depth ordering, which in turn determines occlusion relationships in optical flow. Second, the division of the scene into semantic units allows them to exploit different motion models according to the respective object type, similar to [26]. The motion of planar regions is modeled with homographies, whereas independently moving objects, e.g.,

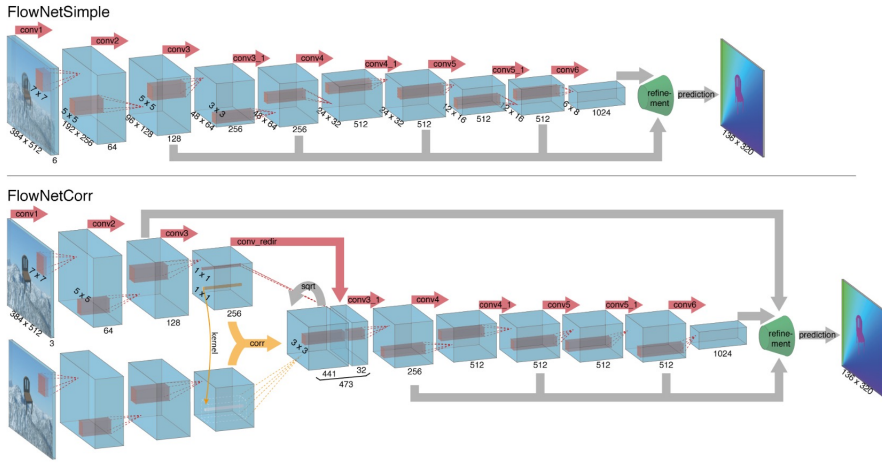


Figure 11.4: **FlowNet Optical Flow Network.** The encoder of the FlowNetSimple and FlowNetCorr architecture proposed by Dosovitskiy et al. [174]. In FlowNetSimple the images are stacked before the first convolutional layer. In contrast, FlowNetCorr processes both images separately and correlates the extracted feature maps. The refinement model is a decoder consisting of deconvolutional layers that get informed by the encoder using skip connections. Figure courtesy of Dosovitskiy et al. [174] © 2015 IEEE.

cars, are modeled by affine motions. Complex objects like vegetation are modeled with a classical spatially varying dense flow field. Finally, the constancy of object identities over time is used to encourage the temporal consistency of the optical flow.

11.2.4 Deep Learning for Optical Flow

Most optical flow approaches do not incorporate high-level information making it hard to overcome ambiguities that require reasoning about larger image regions. The recent success of convolutional neural networks has led to an attempt to use them for the optical flow problem.

Dosovitskiy et al. [174] presented FlowNet to learn optical flow end-to-end using a CNN. FlowNet consists of a contracting part that extracts important features and an expanding part that produces the high-resolution optical flow field as output. They propose two different architectures illustrated in Figure 11.4: a simple network (FlowNetSimple) stacking the images and a complex network (FlowNetCorr) correlating features of the separately processed images. One problem in learning optical flow is the limited amount of training

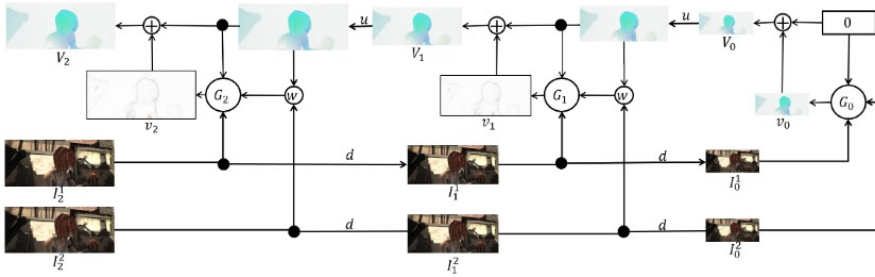


Figure 11.5: **SpyNet Optical Flow Network.** The SpyNet architecture by Ranjan and Black [538] inspired from classical coarse-to-fine approaches. An image pyramid is created and for each resolution a network is trained to predict the residual flow with respect to the previous layer. Figure courtesy of Ranjan and Black [538] © 2017 IEEE.

data. KITTI 2012 [238] and KITTI 2015 [455] only provide around 200 training examples each while Sintel [92] has 1041 training image pairs. Since these datasets are too small to train large CNNs, Dosovitskiy et al. [174] created the Flying Chairs dataset by rendering 3D chair models on top of images from Flickr. This first attempt to end-to-end optical flow learning demonstrated that it was possible to learn optical flow estimation from data, despite not yet reaching the performance of state-of-the-art traditional methods on KITTI or Sintel. However, due to the parallel GPU implementation, FlowNet was able to run in real-time as opposed to most of the classical algorithms implemented on the CPU.

In contrast to the contracting and expanding networks of Dosovitskiy et al. [174], Ranjan and Black [538] present the SpyNet architecture which is inspired by the coarse-to-fine matching strategy leveraged in traditional optical flow estimation techniques. As shown in Figure 11.5, each layer of the network represents a different scale and only estimates the residual flow with respect to the image warped according to the flow of the previous layer. This formulation allowed them to achieve similar performance as FlowNet while being faster and 96 % smaller in terms of network weights, making it attractive for embedded systems with limited compute capabilities. Ilg et al. [311] present FlowNet2, an improved version of FlowNet, by stacking the architectures and fusing the stacked network with a subnetwork specialized in small motions. Similar to SpyNet, they also input the warped image into the stacked networks. Each stacked network estimates the flow between the original frames instead of the residual flow, as in SpyNet. In contrast to FlowNet and SpyNet, they use the FlyingThings3D dataset [450] consisting of 22k renderings of static 3D scenes with moving 3D models from the ShapeNet

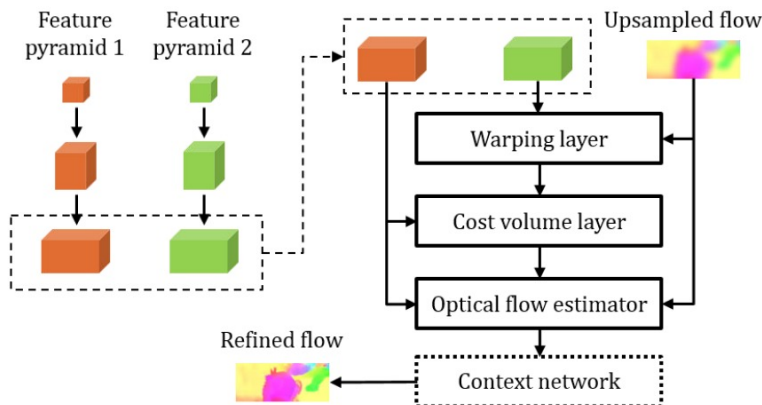


Figure 11.6: **PWCNet Optical Flow Network.** Sun et al. [633] combine coarse-to-fine estimation with cost-volume filtering using a Siamese network for feature extraction. This figure shows one level of the architecture that uses one level of the feature pyramid and the upsampled flow from the previous level for residual flow estimation. Figure courtesy of Sun et al. [633] © 2018 IEEE.

dataset [574]. Recently, PWC-Net [633] illustrated in Figure 11.6 was proposed that combines the classical ideas of coarse-to-fine warping [538] and cost volume filtering [174] with a Siamese network that proved to learn rich feature representations [763].

Unsupervised Learning: Because large annotated datasets for supervised learning of optical flow are rather limited, several recent works [758, 454, 693, 315] address the problem of unsupervised learning of optical flow. Typically, these approaches train one of the standard networks with a photometric loss and a smoothness loss. The photometric loss compares the first image with the second image warped according to the predicted flow. The smoothness loss encourages a similar motion between neighboring pixels. Recently, several approaches [454, 693, 315] noticed that occluded regions introduce errors in the photometric loss that cause misleading gradients during training. They propose to mask out occluded regions in order to avoid this problem. While Meister et al. [454] and Wang et al. [693] both rely on heuristics for estimating occlusions, Janai et al. [315] use a three frame formulation to jointly learn occlusions and optical flow in an unsupervised fashion. Even though occlusion handling results in large improvements and performance comparable to the first fully supervised approaches [174, 538], they are not yet able to compete with the state-of-the-art supervised approaches [311, 633] that dominate the

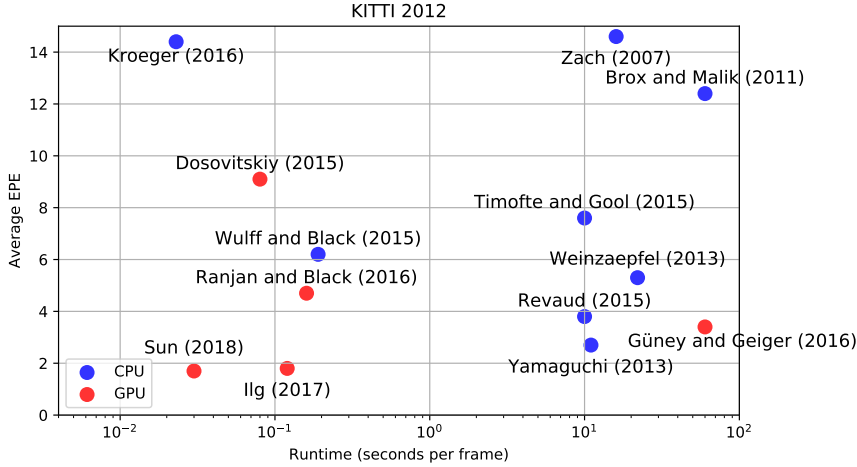


Figure 11.7: **Accuracy vs Efficiency.** The trade-off between performance and speed on KITTI 2012 [238].

leaderboards today.

11.2.5 High-Speed Flow

With some exceptions (Wulff and Black [724], Timofte and Gool [655], Weinzaepfel et al. [703], Farneback [198], and Zach et al. [759]), most of the classical optical flow approaches are very inefficient and cannot be applied in real-time which is necessary for applications in autonomous driving. The trade-off between accuracy and speed for different algorithms on the KITTI 2012 benchmark [238] is illustrated in Figure 11.7. While variational approaches yielded a good precision, they belonged to the slowest set of methods for motion estimation. The duality-based approach for total variation optical flow proposed by Zach et al. [759] allows an efficient GPU implementation that performs in real-time (30 Hz) on a resolution of 320×240 . Sparse matching approaches are usually more efficient than variational formulations but often need variational refinement as a post-processing step to achieve subpixel precision.

The recent introduction of deep learning to the optical flow problem yielded several near real-time approaches (Dosovitskiy et al. [174] and Ranjan and Black [538]) including (Ilg et al. [311] and Sun et al. [633]) which achieve state-of-the-art performance on popular datasets. The approach proposed by Kroeger et al. [361] allows to trade-off accuracy and runtime. They obtain fast patch correspondences with inverse search resulting in a dense flow field when aggregating patches across multiple scales. This allows them to estimate optical flow at up to 600 Hz, but at the cost of accuracy.

11.2.6 Confidences

Considering the remaining challenges in optical flow, a confidence measure to assess the quality of the estimated flow is desirable. Several measures based on spatial and temporal gradients have been proposed [668, 9, 615] to quantify the uncertainty in the optical flow estimate. In contrast, algorithm-specific measures propose confidence estimates for a specific group of methods, i.e., variational methods [83] and methods for pixel-based minimization problems [374]. Learning-based measures [355, 356] learn a model that relates flow algorithm success to spatio-temporal image data or the computed flow field. A detailed evaluation of different confidence measures is given by Mac Aodha et al. [437]. In addition, they present a learning-based approach that uses multiple feature types, such as temporal features, texture, or distance from image edges, to estimate confidences.

11.3 Datasets

Sintel [92] and KITTI [238, 237] discussed in Chapter 4 are the most popular datasets for the evaluation of optical flow algorithms. However, in this survey, we focus on the autonomous driving application. Therefore, we will only refer to the KITTI leaderboard when comparing methods.

11.4 Metrics

The performance of methods is usually assessed considering the endpoint error (Euclidean distance) between the estimated flow vectors and the ground truth. While Sintel reports the average endpoint error for different velocities, occluded and non-occluded regions, the KITTI dataset uses outliers which are computed as the percentage of flow vectors with the absolute endpoint error (EPE) exceeding 3 pixels and 5% of its true values. The percentage of outliers is averaged over background (Fl-bg), foreground (Fl-fg), and all regions (Fl-all), resulting in three different evaluation metrics.

11.5 State of the Art on KITTI

In Table 11.1, we show the leaderboard for the KITTI 2015 benchmark. In addition to the estimation error, the density of the output flow field and the runtime are also provided.

Bai et al. [26] achieve great accuracy in background regions by leveraging semantic segmentation and epipolar geometry. However, their performance drops on foreground regions with dynamic objects that do not follow their

	Method	Fl-bg	Fl-fg	Fl-all	Runtime
1.	PWC-Net+ [632]	7.69 %	7.88 %	7.72 %	0.03 s / GPU
2.	LiteFlowNet [309]	9.66 %	7.99 %	9.38 %	0.0885 s / GPU
3.	PWC-Net [633]	9.66 %	9.31 %	9.60 %	0.03 s / GPU
4.	ContinualFlow_ROB (MF) [485]	8.54 %	17.48 %	10.03 %	0.15 s /
5.	MirrorFlow [310]	8.93 %	17.07 %	10.29 %	11 min / 4 core
6.	FlowNet2 [311]	10.75 %	8.75 %	10.41 %	0.1 s / GPU
7.	SDF [26]	8.61 %	23.01 %	11.01 %	TBA / 1 core
8.	UnFlow [454]	10.15 %	15.93 %	11.11 %	0.12 s / GPU
26.	RicFlow [301]	18.73 %	19.09 %	18.79 %	5 s / 1 core
27.	FlowFields+ [28]	19.51 %	21.26 %	19.80 %	28s / 1 core
28.	PatchBatch [220]	19.98 %	26.50 %	21.07 %	50 s / GPU
29.	DDF [256]	20.36 %	25.19 %	21.17 %	1 min / GPU
36.	Back2FutureFlow (MF) [315]	22.67 %	24.27 %	22.94 %	0.12 s / GPU
37.	MotionSLIC [737]	14.86 %	64.44 %	23.11 %	30 s / 4 cores
39.	FullFlow [113]	23.09 %	24.79 %	23.37 %	4 min / 4 cores
45.	EpicFlow [548]	25.81 %	28.69 %	26.29 %	15 s / 1 core
50.	SPyNet [538]	33.36 %	43.62 %	35.07 %	0.16 s / 1 core
51.	HS [631]	39.90 %	51.39 %	41.81 %	2.6 min / 1 core
52.	DB-TV-L1 [760]	47.52 %	48.27 %	47.64 %	16 s / 1 core

Table 11.1: **KITTI 2015 Optical Flow Leaderboard.** Numbers correspond to percentages of bad pixels according to the 3px/5% criterion defined in [455] averaged over background (bg), foreground (fg), or all regions. Methods followed by (MF) use multiple frames as input. Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

assumptions. Hur and Roth [310] formulate a symmetric optimization problem to jointly reason about optical flow and occlusions. Using an alternating optimization of forward-backward flow and occlusions, they obtain similar results in background regions while improving on foreground regions.

The best performing methods learn optical flow end-to-end [633, 309, 485, 311, 454, 632]. FlowNet2 [311] provides different network variants for the spectrum between 8fps and 140fps, allowing the trade-off between accuracy and computation. The most accurate network achieves comparable results to the state of the art. Neoral et al. [485] use a three frame formulation to jointly learn optical flow and occlusions. This allows them to perform well on background regions that are often occluded by foreground objects. Hui et al. [309] follow a similar approach to PWC-Net by using a Siamese network, coarse-to-fine warping, and computing correlations. Sun et al. [632] propose a new learning rate schedule consisting of several disruptions (a strong increase of the learning rate) and show how this improves the training of the original PWC-Net. PWC-Net [633] with the adapted training protocol [632] outperforms all methods on KITTI 2015 (Table 11.1) and Sintel in both background and foreground objects while being one of the fastest methods on

both benchmarks.

11.6 Discussion

Robust optical flow methods need to handle intensity changes not caused by the actual motion of interest but by illumination changes, reflections, and transparency. In real-world scenes, repetitive patterns, textureless surfaces, saturated image regions, and occlusions are frequent sources of errors. While illumination changes have been addressed with novel data terms [54, 681], the problems caused by reflection, transparency, ambiguities, and occlusions remain mostly unsolved. In Figure 11.8, we show the accumulated error of the 15 best-performing methods on KITTI 2015 [455]. The highest error can be observed for regions moving outside the image domain for which the optical flow has to be guessed, as observations are not available. Untextured, reflective, and transparent regions also result in large errors in many cases. A better understanding of the world is necessary to tackle these problems. Semantics [26] and learned high-capacity models [633, 309, 485, 311, 454, 632] have already proven to improve optical flow estimation by resolving ambiguities in the data. Joint optical flow and occlusion formulations have also shown great potential to alleviate these problems for optimization-based [310] as well as learning-based [454, 693, 315] methods.

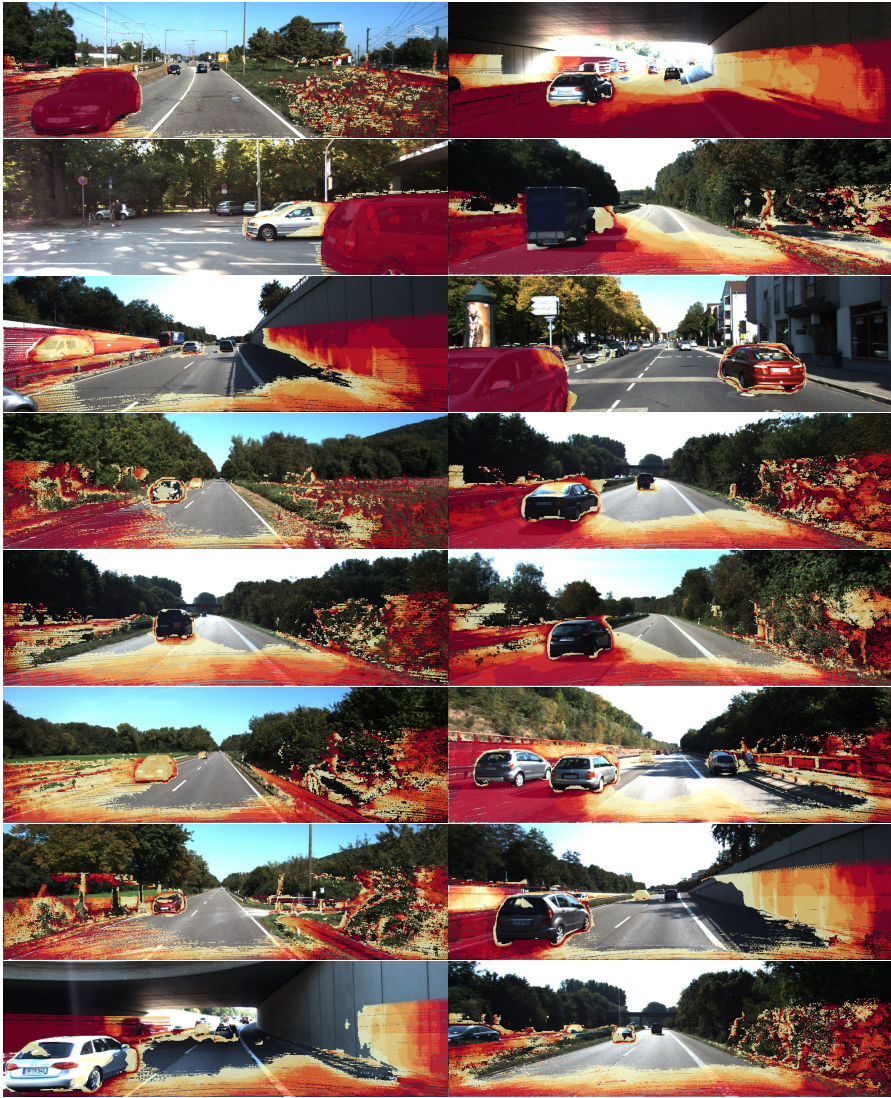


Figure 11.8: **KITTI 2015 Optical Flow Analysis.** The averaged errors of the 15 best-performing optical flow methods published on the KITTI 2015 Flow benchmark. Red colors correspond to regions where the majority of methods fail according to the 3px/5% criterion defined in [455]. Yellow colors correspond to regions where some of the methods fail. Regions that are correctly estimated by all methods are transparent.

Chapter 12

3D Scene Flow

12.1 Problem Definition

Humans are able to effortlessly integrate depth and motion cues from observations over time. That kind of reasoning is essential for many tasks in autonomous driving, such as segmentation of moving objects in the 3D world. Scene flow generalizes optical flow to 3D, or equally, dense stereo to dynamic scenes. Given stereo image sequences, the goal is to estimate the three-dimensional motion field that is a 3D motion vector for every point on every visible surface in the scene. The minimal setup for image-based scene flow estimation is given by two consecutive stereo image pairs, as visualized in Figure 12.1. Establishing correspondences between the four images results in the 3D location of the surface point in both frames and hence fully describes the 3D motion of that surface point. A dense output is preferred, although some early works focused on establishing sparse correspondences [209]. Scene flow shares some of the challenges with stereo and optical flow, such as matching ambiguities in weakly textured regions and the aperture problem, but integrating observations from four images and solving both tasks jointly leads to a better-constrained problem.

12.2 Methods

Following the seminal work by Vedula et al. [673], the problem is traditionally formulated in a variational setting where optimization proceeds in a coarse-to-fine manner, and local regularizers are leveraged to encourage spatial smoothness of depth and motion. Wedel *et al.* [698, 696] propose a variational framework by decoupling the motion estimation from the disparity estimation while maintaining stereo constraints. Starting from a precomputed disparity map at

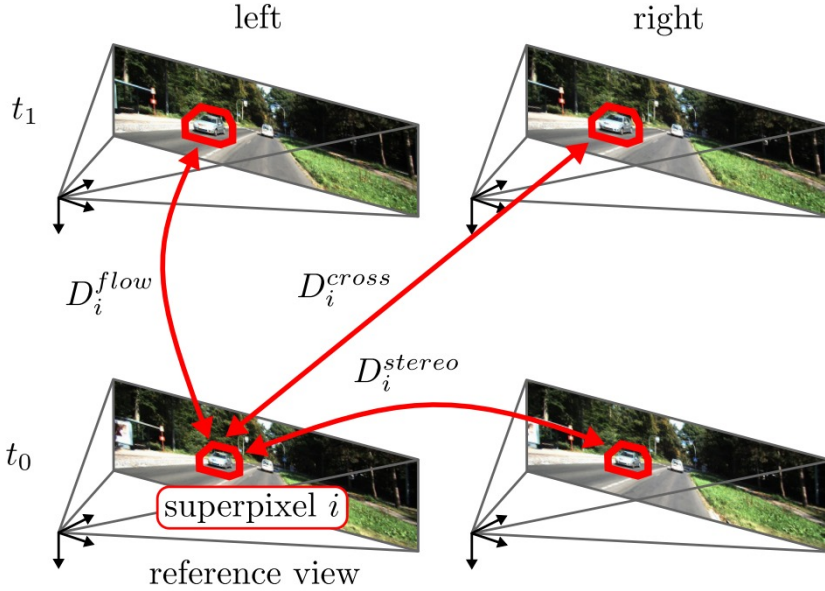


Figure 12.1: **Scene Flow**. The minimal setup for image-based scene flow estimation is given by two consecutive stereo image pairs. Figure courtesy of Menze and Geiger [455] © 2015 IEEE.

each time step, optical flow for the reference frame and disparity for the other view are estimated. The motivation for this decoupling is mainly computational efficiency by choosing the optimal technique for each task. In addition, Wedel et al. [696] propose a solution for varying lighting conditions based on residual images and provide an uncertainty measure which they showed to be useful for object segmentation. Rabe et al. [533] integrate a Kalman filter to the decoupling approach for temporal smoothness and robustness.

12.2.1 Piecewise Rigidity

Similar to stereo and optical flow, prior assumptions about the geometry and motion can be exploited to better handle the challenges of the scene flow problem. Vogel et al. [682] and Lv et al. [433] represent the dynamic scene as a collection of rigidly moving planar regions, as shown in Figure 12.2. Vogel et al. [682] jointly recover this segmentation while inferring the shape and motion parameters of each region. They use a discrete optimization framework and incorporate occlusion reasoning as well as other scene priors in the form of spatial regularization of geometry, motion, and segmentation. In ad-

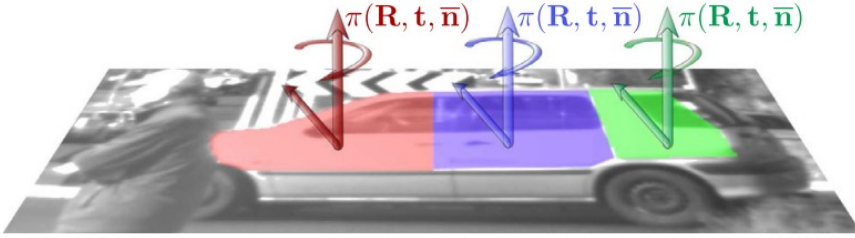


Figure 12.2: **Piecewise Rigidity.** In Vogel et al. [682] the scene is modeled as a collection of rigidly moving planar segments. Figure courtesy of Vogel et al. [682] © 2015 Springer.

dition, they reason over multiple frames by constraining the segmentation to remain stable over a temporal window. Their experiments show that their view-consistent multi-frame approach significantly improves accuracy for challenging scenarios. Using the same representation, Lv et al. [433] focus on an efficient solution to the problem. They assume a fixed superpixel segmentation and perform optimization in the continuous domain for faster inference. Starting from an initialization based on Deep Matching [703], they independently refine the geometry and motion of the scene, and finally perform a global non-linear refinement using the Levenberg-Marquardt algorithm.

Piecewise Rigidity at the Object Level: Menze and Geiger [455] and Behl et al. [36] also follow a slanted plane approach, but in addition to previous methods [682, 433], they model the decomposition of the scene into a small number of independently moving objects and the background. By conditioning on a superpixelization, they jointly estimate this decomposition as well as the rigid motion of the objects and the plane parameters of each superpixel in a discrete-continuous Conditional Random Field (CRF). Compared to [682, 433], they leverage a more compact representation, by implicitly regularizing over larger distances. They also present a new scene flow dataset by annotating dynamic scenes from the KITTI raw data collection using detailed 3D CAD models. Menze et al. [457] propose an extension of this model where the pose and 3D shape of the objects are inferred in addition to the rigid motion and segmentation. In particular, they incorporate a deformable 3D active shape model of vehicles into the scene flow approach.

12.2.2 Semantic Segmentation

Semantic information allows constraining the space of possible rigid body motions. For instance, in an autonomous driving scenario, pixels which are

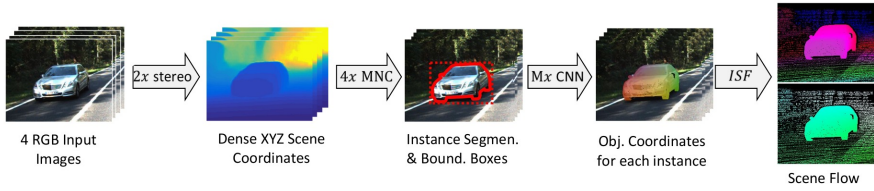


Figure 12.3: **Semantic Segmentation for Scene Flow.** Behl et al. [36] leverage instance segmentation and bounding boxes from an instance segmentation pipeline to independently model the motion of each object and the background. Figure courtesy of Behl et al. [36] © 2017 IEEE.

grouped together in the segmentation are likely to move as a single rigid object in the case of vehicles. Furthermore, a pixel on a vehicle instance in one frame should be mapped to a vehicle instance in the other frame. Behl et al. [36] investigate the impact of bounding box detection, instance segmentation, and 3D object coordinates on scene flow estimations and show which one is most beneficial for scene flow. They obtain the bounding boxes and instance segmentation from the proposal-based instance segmentation method MNC [147] discussed in Chapter 8. 3D object coordinates are predicted with a CNN trained on the 2D instance segmentations as illustrated in Figure 12.3. Using a CRF based on [455], they show that semantic cues lead to significant improvements. However, the benefit of 3D object coordinates over instance segmentations is negligible. Recently, Ma et al. [436] leverage multiple cues consisting of CNNs for instance segmentation (Mask R-CNN [280]), optical flow (PWC-Net [633]) and stereo (PSM-Net [103]) to address the scene flow problem. They formulate an energy combining all the cues with a photometric term. By unrolling the optimization as a recurrent network, they are able to train the whole pipeline end-to-end.

12.2.3 Scene Flow from 3D Point Clouds

The image-based methods discussed earlier estimate scene flow based on two consecutive image pairs of a calibrated stereo camera rig. However, stereo-based scene flow methods suffer from the “curse of two-view geometry”, i.e., the depth error grows quadratically with the distance to the observer. Furthermore, most modern self-driving car platforms rely on LiDAR technology for 3D geometry perception. In contrast to cameras, laser scanners do not suffer from the quadratic error behavior of stereo cameras. In addition, laser scanners provide a 360-degree field of view with just one sensor and are generally unaffected by lighting conditions. Therefore, there have been several

methods proposed recently for estimating 3D scene flow from pairs of unstructured 3D point clouds. Dewan et al. [163] propose a 3D scene flow approach where local SHOT descriptors [656] are associated via a CRF that incorporates local smoothness and rigidity assumptions. However, local shape representations such as SHOT often fail in the presence of noisy or ambiguous inputs. In contrast, Behl et al. [37] address the scene flow problem using a generic end-to-end trainable model that is able to learn local and global statistical relationships directly from data. In order to apply the standard 3D convolution operations, they discretize the point cloud to a grid of voxels. However, because of the sparse nature of the LiDAR data, most of the space is empty, which makes the approach computationally and memory inefficient. To alleviate this problem, Wang et al. [690] propose a novel continuous convolution that operates over non-grid structured data.

12.3 Datasets

Only a few datasets exist for scene flow [238, 450, 92]. Similar to flow and stereo, the KITTI scene flow benchmark [238] is the most popular dataset allowing the comparison of methods on an online evaluation server. For deep learning, the Flying Things datasets [450] is often used for pre-training since KITTI is too small. Recently, MPI Sintel [92] published stereo sequences for the training dataset¹ and is since used to show the generalization of scene flow approaches to other scenes than street scenes from KITTI.

12.4 Metrics

Scene flow methods are usually evaluated by jointly measuring the accuracy of the stereo (Section 9.4) and optical flow estimates (Section 11.4). The KITTI benchmark considers the percentage of erroneous pixels. A pixel is erroneous if the Euclidean distance to the ground truth exceeds a 3 pixels or 5% threshold. The percentage of stereo disparity outliers in the first frame (D1), the percentage of stereo disparity outliers in the second frame (D2), the percentage of optical flow outliers (F1), and the percentage of scene flow outliers (SF), i.e., outliers in either D0, D1 or F1 are reported. The outlier ratio for foreground/background regions can be found separately on the website of the benchmark², but it is omitted here for space reasons.

	Method	D1	D2	F1	SF	Runtime
1.	UberATG-DRISF [436]	2.55 %	4.04 %	4.73 %	6.31 %	0.75 s / CPU+GPU
2.	ISF [36]	4.46 %	5.95 %	6.22 %	8.08 %	10 min / 1 core
3.	PRSM [682] (MF)	4.27 %	6.79 %	6.68 %	8.97 %	300 s / 1 core
4.	OSF+TC [484] (MF)	5.03 %	6.84 %	7.02 %	9.23 %	50 min / 1 core
5.	OSF 2018 [458]	5.28 %	7.06 %	7.41 %	9.66 %	390 s / 1 core
6.	SSF [546]	4.42 %	7.02 %	7.14 %	10.07 %	5 min / 1 core
7.	OSF [455]	5.79 %	7.77 %	7.83 %	10.23 %	50 min / 1 core
8.	FSF+MS [646] (MF)	6.74 %	9.85 %	11.30 %	14.96 %	2.7 s / 4 cores
9.	PWOC-3D [575]	5.13 %	8.46 %	12.96 %	15.69 %	0.13 s / GPU
18.	SGM+C+NL [294]	6.84 %	28.25 %	35.61 %	40.33 %	4.5 min / 1 core
19.	SGM+LDOF [294]	6.84 %	28.56 %	39.33 %	43.67 %	86 s / 1 core
20.	DWBSF [549]	20.12 %	34.46 %	39.14 %	45.48 %	7 min / 4 cores
21.	GCSF [99]	14.21 %	33.41 %	46.40 %	53.54 %	2.4 s / 1 core
22.	VSF [308]	26.38 %	57.08 %	49.28 %	66.90 %	125 min / 1 core

Table 12.1: **KITTI 2015 Scene Flow Leaderboard.** Numbers correspond to percentages of bad pixels according to the 3px/5% criterion defined in [455] for disparity in the first frame (D1), disparity in the second frame (D2), optical flow between both frames (F1) as well as the combination of all criteria yielding the final scene flow metric (SF). Approaches using more than 2 frame pairs are marked by (MF). Methods below the horizontal line show older entries for reference. Accessed on: June 2019.

12.5 State of the Art on KITTI

Table 12.1 shows the ranking of methods on the KITTI Scene Flow 2015 benchmark [455].

All top-performing methods use either semantic cues [36, 436, 546] or the assumption of rigidly moving segments [682, 455, 484, 458]. Modeling the motion of objects using a rigid transformation [455, 457, 458, 484] achieves impressive results on the KITTI dataset. However, this is a very strong assumption even in street scenes as the non-rigidity of pedestrians cannot be handled. In contrast, the segmentation of the scene into superpixels [682] alleviates this problem and allows better performance since non-rigid objects can be modeled by multiple superpixels. However, the best performance is achieved by integrating semantic information [36, 436]. While most scene flow approaches are very inefficient, the method of Ma et al. [436] is a notable exception, requiring only 0.75 seconds. They achieve this efficiency by combining CNNs for instance segmentation (Mask R-CNN [280]), optical flow (PWC-Net [633]) and stereo (PSM-Net [103]) to address the scene flow problem and exploiting a GPU for inference (in contrast to classical scene flow approaches which typically run on the CPU).

¹<http://sintel.is.tue.mpg.de/stereo>

²http://www.cvlibs.net/datasets/kitti/eval_scene_flow.php

12.6 Discussion

The scene flow problem shares many challenges with stereo and optical flow while integrating more information than each task alone and consequently leading to better results. Ideally, methods should exploit depth and motion cues together to reason about dynamic 3D scenes. However, considering the optical flow (Table 11.1) and stereo matching leaderboards (Table 9.1), the joint formulation is more advantageous for the optical flow problem leading to significant improvements as for instance UberATG-DRISF [436] reaches an outlier ratio of 4.73% in comparison to PWC-Net+ [632] reaching 7.72 %. In contrast, the stereo matching performance is comparable with an outlier ratio of 2.55% (UberATG-DRISF [436]) in comparison to 2.08 % (EdgeStereo-V2 [622]).

We show the accumulated errors of the top 5 methods on the KITTI scene flow benchmark in Figure 12.4. Car surfaces are the most problematic regions due to matching problems and the independent motion of cars. Pixels close to the image boundary are another common source of error, especially on the road surfaces in front of the car, where large scale changes occur. Although local planarity and rigidity assumptions alleviate the problem, they are often violated due to complex geometric objects like vegetation, pedestrians, or bicycles. Superpixels grouping different surfaces due to wrong estimation of planes cause additional problems, especially at the boundaries of objects. Semantic image understanding seems a promising direction [546, 36, 436], especially at the object level, by segmenting car instances. However, an additional network has to be trained for obtaining this information, and prediction errors can lead to irreversible errors in the final scene flow estimation. Leveraging temporal information [682, 484] also leads to improvements and should be exploited whenever possible. Especially, long-term temporal interactions could allow to alleviate ambiguities and improve. However, obtaining a robust, accurate, and real-time multi-frame scene flow estimate remains an open problem that requires further work.



Figure 12.4: **KITTI 2015 Scene Flow Analysis.** The averaged errors of the 15 best-performing scene flow methods published on the KITTI 2015 Scene Flow benchmark. Red colors correspond to regions where the majority of methods yield bad pixels according to the 3px/5% criterion defined in [455]. Yellow colors correspond to regions where some of the methods fail. Regions that are correctly estimated by all methods are transparent.

Chapter 13

Mapping, Localization & Ego-Motion Estimation

13.1 Problem Definition

Navigating a vehicle requires a precise understanding of the position and orientation of the car. Localization is a well-studied problem in both robotics and vision, covering a broad range of techniques from indoor localization using noisy sensory measurements to locating where a picture was taken. From an autonomous driving perspective, the main task is to localize the vehicle on a map in order to exploit static features provided by the map. The task of generating a map of the world is defined as the mapping problem. In this chapter, we discuss approaches for generating both metric as well as semantic maps. While metric maps allow for accurate localization, semantic maps provide problem-specific information such as the location of parking areas. Maps for localization can be generated offline, exploiting accurate, but computationally expensive optimization techniques.

In contrast to localization, the ego-motion estimation problem considers the change in position and orientation of the vehicle. While this problem can be addressed more efficiently than the localization problem (the previous position is assumed to be known), small inaccuracies quickly accumulate to larger drifts. Approaches for Simultaneous Localization and Mapping (SLAM) address this problem by detecting loop closures to correct for drift.

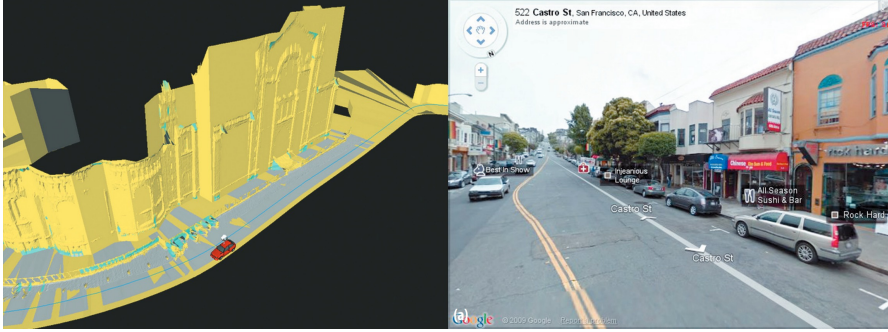


Figure 13.1: **Google Street View.** Dominant scene surfaces reconstructed from images and laser range data (left) and a scene from the Google Street View project [15]. Figure courtesy of Anguelov et al. [15] © 2010 IEEE.

13.2 Mapping

Street, aerial, and satellite imagery enable the generation of precise metric and semantic maps. Depending on the required level of detail, various computer vision techniques, i.e., multi-view reconstruction, scene understanding, or semantic segmentation, are typically employed for generation maps.

13.2.1 Metric Maps

For autonomous driving, 2D metric maps (i.e., representing information in bird's-eye view) are usually sufficient for localization. Methods for scene understanding, such as [193, 709, 707, 236, 657, 595, 766] discussed in Chapter 14 can also be used to extract road features. 3D information can be obtained using multi-view reconstruction techniques operating on street-level [3, 205, 206] (Section 10.2) or aerial [214, 58, 178] images.

The Google Street View project [15] is a prominent example of a large collection of panoramic images that are registered with respect to each other to form a world map, see Figure 13.1. For registering the dataset, Anguelov et al. [15] estimate the pose of the vehicle using a Kalman filter, fusing data from GPS, wheel encoder, and inertial navigation [350]. The pose estimates are refined with a probabilistic graphical model, and the 3D scene geometry is recovered by robustly fitting coarse meshes to the 3D measurements.

Levinson et al. [404] propose to construct a map based on aggregated reflectance measurements from a LiDAR scanner. They exploit these maps for centimeter-accurate LiDAR-based localization during the DARPA Urban Challenge. In contrast, Geiger [234] presents an approach for road mosaicing in dynamic environments with the goal of creating obstacle-free bird's-eye



Figure 13.2: **Appearance Changes in Localization.** Examples for different weather conditions, seasons, and day times for a scene from the Workshop organized by Hammarstrand et al. [265]. Figure courtesy of Hammarstrand et al. [265].

views. The road surface is extracted using optical flow on Harris corners and approximated by a plane. Afterwards, multiple road reconstructions are combined using multi-band blending.

13.2.2 Semantic Maps

Metric maps ignore semantic information, which is important for some tasks such as automated parking. Semantic maps are necessary to address this problem. Several approaches address the creation of semantic maps [700, 701, 472, 675, 447, 448, 699, 448]. Scene understanding approaches like [193, 236] also estimate semantic classes to extract road topologies but do not create a semantic map.

Sengupta et al. [600] present an approach to generate a semantic overhead map of an urban scene from street-level images. They formulate the problem using two CRFs. The first is used for semantic image segmentation of the street view images treating each image independently. Each street view image is then projected into an overhead view. These views are then aggregated over many images to form the input for a second CRF producing a semantic labeling of the ground plane.

In contrast, Grimmett et al. [252] fuse semantic and metric maps for vision-only automated parking. They update the map with static and dynamic labels and use active learning for lane, parking space, and pedestrian crossings detection.

13.3 Localization

Localization can be performed using either a sensor like GPS or visual information based on images. Using GPS alone typically provides an accuracy of around 5 meters. Although centimeter-level precision is possible in unobstructed environments using correction signals and a combination of several

sensors as in the KITTI car [238], it is often rendered infeasible in traffic scenes with several disturbing effects such as occlusions by vegetation and buildings or multi-path effects due to reflections. Therefore, image-based localization independent of satellite systems remains highly relevant.

Visual localization techniques are commonly classified into metric and topological methods. Metric localization [158, 492] is achieved by computing the 3D pose with respect to a map. Topological localization approaches [414, 781, 278] provide a coarse estimate from a finite set of possible locations that are represented as nodes in a graph and connected by edges that link them according to some distance or appearance criteria. Metric localization can be very accurate, but is usually not suitable for very long sequences, while topological localization may be more reliable, but only provides rough estimates.

Metric Localization: The problem of metric map localization has been traditionally addressed using Monte Carlo methods which recover the probability distribution over the agent’s pose by drawing a set of samples. Dellaert et al. [158] define indoor localization in two steps, global position estimation and tracking of the local position over time. Instead of modeling the probability density function itself, they represent uncertainty by a set of samples and update the representation over time using Monte Carlo methods. This allows them to model arbitrary multi-modal distributions in a memory-efficient way.

Outdoor localization is, in general, more challenging compared to the indoor localization task due to its scale and often unreliable sensor information, e.g., GPS failures. Oh et al. [492] use semantic information available in maps to compensate for the failure cases of GPS sensors. By exploiting knowledge about the environment, they assign low probabilities to implausible map locations, e.g., inside buildings. They incorporate these map-based priors into their particle filter formulation to bias the motion model towards areas of higher probability.

Topological Localization: Early image-based techniques [414, 781] approach the problem of localizing in topological maps as classification into one of a predefined set of places which are often referred to as “landmarks”. Others [278, 140, 504, 658, 16] create a database of images with known locations and formulate localization as an image retrieval problem. These methods require a similarity measure to compare images based on local or global appearance cues. The larger the database, the more difficult the localization task becomes. Challenges include appearance changes, similar-looking places, and changes due to viewpoint or position. In Figure 13.2, we show an example for the appearance change of a scene over different seasons from the Workshop organized by Hammarstrand et al. [265].

Lowry et al. [430] provide a comprehensive review of visual place recognition techniques. Given a map of the environment, the goal of place recognition

is to decide whether the current observation is a place already included in the map, and if so, which one.

Topometric Localization: In contrast to purely topological methods, the graph of a topometric localization model is more fine-grained: each node corresponds to a metric location without semantic meaning. Towards this goal, Badino et al. [22] propose to construct a graph using the vehicle’s position from GPS at fixed distance intervals while associating visual or 3D features to the corresponding graph node. At runtime, real-time localization is performed using a Bayes filter to estimate the probability distribution of the vehicle position along the route by matching features extracted from the sensor data to the map’s feature database. Brubaker et al. [82] leverage a graph-based representation. In contrast to traditional localization approaches, however, they do not require a visual feature database of the environment, but instead, directly build this graph from road networks extracted from OpenStreetMap. They further propose a probabilistic model that allows inferring a distribution over the vehicle location along the edges of this road graph using visual odometry measurements. For tractability in very large environments, they leverage several analytic approximations for efficient inference yielding higher stability compared to particle-based filtering techniques.

Scale and Accuracy: The scale of the target area is a distinctive property to compare different approaches and is related to the accuracy achieved. Both scale and accuracy depend on the methodology used, such as map-based approaches [82] which cover a large area but might suffer from the errors on the map compared to descriptor-based approaches [22, 592] on a smaller area. While the descriptor-based method of Badino et al. [22] achieves an average localization accuracy of 1 m over an 8 km route, the localization approach of Brubaker et al. [82] which requires only road networks as input attains an accuracy of 4 m on a 18 km² map containing 2,150 km of drivable roads.

Schreiber et al. [592] point out that the required precision for autonomous driving and future driver assistance systems is in the range of a few centimeters and present a feature- and map-matching-based localization algorithm which can achieve centimeter-level accuracy on approximately 50 km of rural roads. They approach the problem from the perspective of lane recognition. First, they create a highly accurate map that contains road markings and curbs. Then while driving, they detect and match them to the map in order to determine the position of the vehicle relative to the markings.

13.3.1 Structure-based Localization

While the output of the aforementioned localization approaches is either a rough camera position or a distribution over positions, another line of work which is known as “structure-based localization”, aims to estimate all camera

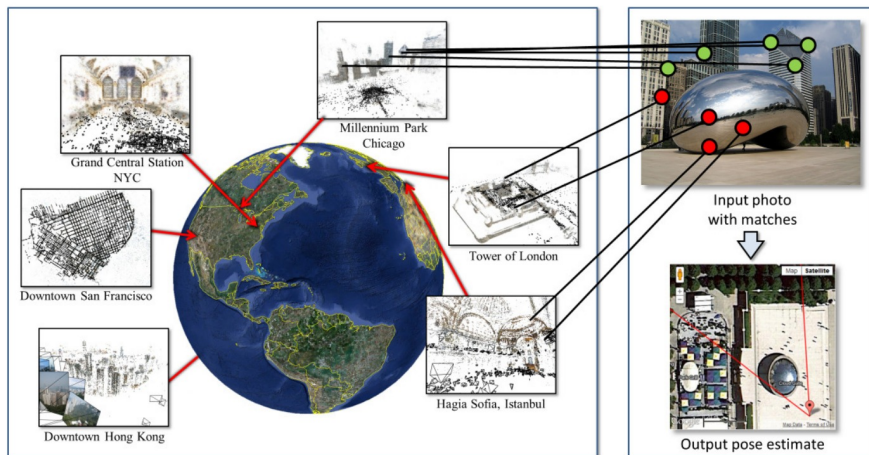


Figure 13.3: **Structure-based Localization.** A query image is matched to a database of geo-referenced structure-from-motion point clouds assembled from photos of places around the world (left). In structure-based approaches, the goal is to compute the geo-referenced pose of new query images by matching to a large database of feature descriptors (right). Figure courtesy of Li et al. [415] © 2012 IEEE.

matrix parameters, including position, orientation, and sometimes also camera intrinsics. Estimating the intrinsics usually enables more accurate results. Localization is realized as a 2D-to-3D matching problem where the 2D points on the images are matched to a large, geo-registered 3D point cloud, and the pose is estimated with respect to correspondences as shown in Figure 13.3.

Direct matching by approximate nearest neighbor search using SIFT features usually results in many incorrect matches. Therefore, many approaches rely on the SIFT ratio test [429] to detect and reject ambiguous matches. This works well on small to medium scale scenes. However, with growing model size, the discriminative power of the descriptors decreases, and many matches will be rejected by the ratio test. On the other hand, relaxing the ratio test leads to many ambiguous and wrong matches.

Several approaches [313, 568, 415, 636, 765] address this problem by restricting the search space. Irschara et al. [313] and Sarlin et al. [568] use image retrieval techniques to identify parts of the scene which likely include the query image. Afterwards, 2D-3D matching is performed to 3D points visible in the retrieved images. In contrast, Li et al. [415] find statistical co-occurrences of 3D model points in images and then use them as a sampling prior for RANSAC to exploit co-visibility relations. In addition, they employ

a bidirectional matching scheme, forward from features in the image to points in the database and inverse from points to image features. They show that the bidirectional approach performs better than forward or inverse matching alone. Svärm et al. [636] and Zeisl et al. [765] propose to use geometric cues to obtain matches that are likely to be inliers. They also exploit the gravity direction obtained from gravitational sensors and an approximation of the camera height to reduce the search space.

Besides ambiguities, the efficiency of the matching stage and memory requirements to store the large number of descriptors contained in the model are also problems related to large scale. Therefore, several approaches use only a subset of the 3D points [434] or present compression schemes for the descriptors [434, 570, 569, 421, 95] for more efficient matching or memory reduction. Sattler et al. [570] and Sattler et al. [569] use quantization into a fine vocabulary to accelerate the matching stage where each descriptor is represented by its word ID. Sattler et al. [570] separate the difficult problem of finding a unique 2D-3D matching into two simpler ones. They first establish locally unique 2D-3D matches using a fine visual vocabulary and a visibility graph which encodes the visibility relation between 3D points and cameras. Then, they disambiguate these matches by using a simple voting scheme to enforce the co-visibility of the selected 3D points. Their experiments show that matching based on a visual vocabulary leads to state-of-the-art results. Sattler et al. [569] propose a prioritized matching scheme based on quantization, focusing on efficiency. They significantly accelerate 2D-to-3D matching by considering more likely features first and terminating the correspondence search as soon as enough matches are found. A hybrid approach combining the idea of working on a subset of 3D points and the compression of the descriptors is presented by Camposeco et al. [95]. For a small subset of 3D points, they keep the full appearance information, while for a larger set of points, they store a compressed descriptor. This enables them to obtain a more complete representation of the scene with a memory consumption similar to the previous approaches.

Deep Learning:

The motivation for using CNNs for structure-based localization is to learn high-level information which might help to handle problems like textureless areas, motion blur, and illumination changes. In contrast to classical localization approaches whose runtime depends on several factors such as the number of features found in a query image or the number of 3D points in the model, the runtime of CNN-based approaches only depends on the size of the network.

Kendall et al. [337] and Walch et al. [684] use a convolutional neural network to regress the camera pose from a single RGB image in an end-to-end manner. Kendall et al. [337] modify GoogLeNet [639] by replacing the softmax

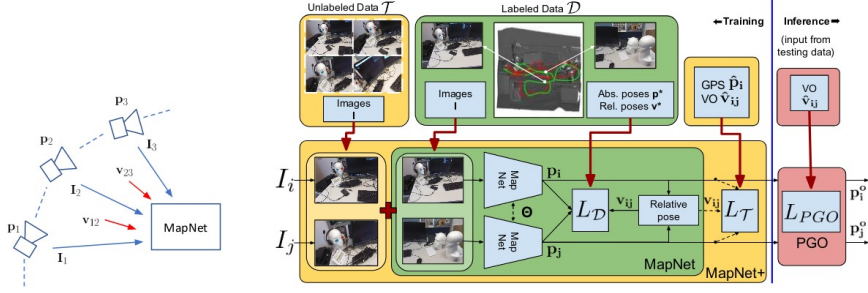


Figure 13.4: **Learning Structure-based Localization.** MapNet proposed by Brahmabhatt et al. [66] learns a map representation from images, visual odometry, and GPS (left). During inference (right) visual odometry is used to update the maps in a self-supervised fashion and pose graph optimization (GPO) allows for further refinement. Figure courtesy of Brahmabhatt et al. [66] © 2018 IEEE.

classifiers with affine regressors and inserting another fully connected layer before the final regressor, which can be used as a localization feature vector for further analysis. The final architecture, dubbed PoseNet, is initialized by using the weights of classification networks trained on giant datasets such as ImageNet [160] and Places [782]. The network is further fine-tuned on a new pose dataset which was automatically created using SfM to generate camera poses from a video of the scene. Walch et al. [684] use a similar approach, but in addition, they spatially correlate each element of the output of the CNN using Long Short-Term Memory (LSTM) units. This way, the network is able to capture more contextual information and outperform PoseNet in different localization tasks, including large-scale outdoor, small-scale indoor, and a newly proposed large-scale indoor localization benchmark.

Recently, Brahmabhatt et al. [66] proposed MapNet for representing maps as deep neural networks. They exploit visual odometry and GPS in addition to images for image-based localization and formulate geometric constraints as additional loss terms. Thus, the model can be updated in a self-supervised fashion using unlabeled data. This allows them to significantly improve in comparison to PoseNet-based approaches. The model is illustrated in Figure 13.4.

While previous methods [337, 684, 66] regress the absolute pose in a given scene, another line of work [566, 31] proposes to learn the relative pose with respect to an image retrieved from a database. Eventually, the absolute pose is obtained from the known pose of the retrieved image and the relative pose.

Sattler et al. [572] notice that PoseNet-based approaches [337, 684] are

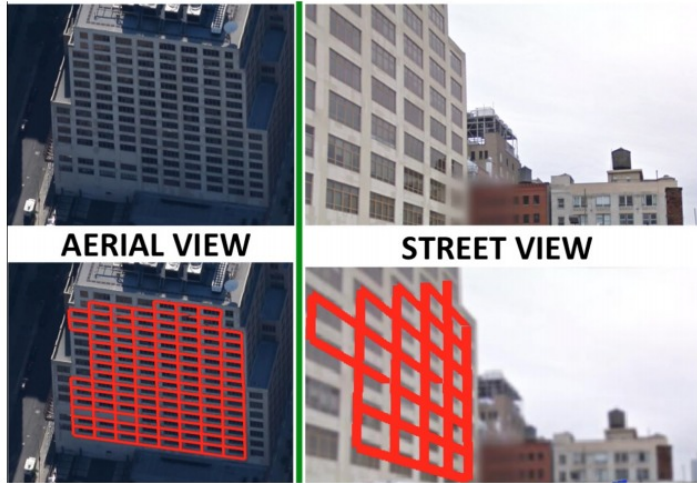


Figure 13.5: **Aerial to Street-View Matching.** Repeating patterns of buildings are exploited by regularity-driven approaches for aerial to street-view matching. Figure courtesy of Wolff et al. [714] © 2016 IEEE.

not able to outperform simple image retrieval approaches [658] and learning-based approaches are in general still inferior to structure-based approaches such as [635].

13.3.2 Cross-view Localization

It is a very difficult endeavor to keep an up-to-date repository of ground-level imagery around the world. In contrast, establishing live maps from aerial and satellite images is comparably easier. This motivated the development of geo-localization approaches that try to register ground-level images to aerial imagery. The underlying idea is to learn a mapping between ground-level and aerial image viewpoints to localize a ground-level query in an aerial reference image database.

Lin et al. [417] match ground-level queries to other ground-level reference photos as in traditional geo-localization, but then use the overhead appearance and land cover attributes of those ground-level matches to build sliding-window classifiers in the aerial and land cover domain. In contrast to previous methods, they are able to localize a query even if it has no corresponding ground-level image in the database by learning the co-occurrence of features in different views. Inspired by the success of face verification algorithms using deep learning, Lin et al. [418] train a Siamese network to match cross-view pairs of the same location. Towards this goal, they collect a cross-view patch

dataset using range data and camera parameters from Google Street View. Finally, they warp the dominant building surface plane to appear approximately as a 45% aerial view. In contrast, Workman et al. [716] use CNNs for extracting ground-level image features and predict these features from aerial images of the same location. This way, the CNN is able to extract semantically meaningful features from aerial images without manually specifying semantic labels. They conclude that the cross-view localization approach can obtain a precise estimate of the geographic locations which are distinctive from above. Otherwise, it can be used as a pre-processing step to a more expensive matching process.

Buildings Facades: Several methods have been developed which specialize in building facades from cross-view matching. The repeating patterns yield valuable matching cues, as illustrated in Figure 13.5. By combining satellite and oblique bird’s-eye views, Bansal et al. [33] first extract building outlines as well as facades and then match the ground image to oblique aerial images based on a statistical description of the facade pattern. Wolff et al. [714] define a matching cost function to compare street-view motifs to aerial view motifs based on the similarity of color, texture, and edge-based context features.

Geo-Referenced Reconstruction: Another line of work addresses the problem of geo-referencing a reconstruction by automatic alignment with a satellite image, floor plan, map, or other overhead views. Kaminsky et al. [329] compute the optimal alignment between SfM reconstructions and overhead images using an objective function that matches 3D points to image edges and imposes free space constraints based on the visibility of points in each camera. Matching ground and aerial images directly is a difficult endeavor due to the large differences in camera viewpoints, occlusions, and imaging conditions. Instead of seeking invariant feature detections, Shan et al. [606] propose a viewpoint-dependent matching technique by exploiting approximate alignment information and the underlying 3D geometry.

13.3.3 Semantic Alignment from LiDAR

Several companies acquire LiDAR data from scanners mounted on cars driving through cities to acquire 3D models of real-world urban environments. However, the accuracy of the 3D point positions acquired by the 3D scanners depends on the scanner poses predicted by GPS, inertial sensors, and structure-from-motion, which often fail in urban environments. These misalignments cause problems for point cloud registration methods. Yu et al. [757] propose to align semantic features that can be matched robustly at different scales. By following a coarse-to-fine approach, they first successively align roads, facades, and poles which can be matched robustly. Afterwards, they match cars and other small objects which require better initial align-

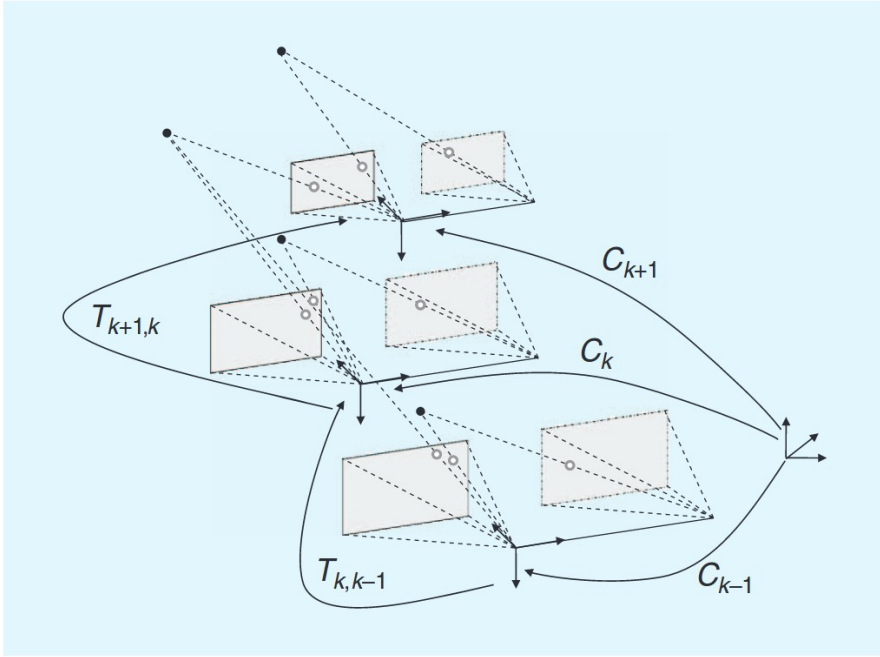


Figure 13.6: **Visual Odometry.** Illustration of the incremental visual odometry approach by Scaramuzza and Fraundorfer [577]. The transformation $T_{k,k-1}$ between two adjacent camera systems is obtained using visual features. The accumulation of all transformations yields the absolute pose C_k with respect to the initial coordinate frame $k = 0$. Figure courtesy of Scaramuzza and Fraundorfer [577] © 2011 IEEE.

ments to find correct correspondences. The use of semantic features provides a globally consistent alignment of LiDAR scans, and their evaluation shows improvement over the initial alignments.

13.4 Ego-Motion Estimation

One of the simplest ways of estimating the ego-motion of a vehicle is to use the wheel angle in combination with the output of wheel encoders which measure the rotation of the wheel. These methods suffer from wheel slip in uneven terrain or adverse conditions and can not recover from errors in the measurements. Visual odometry and LiDAR-based odometry techniques that estimate ego-motion from visual observations (images or laser range measurements) are more robust in many situations and can correct for drift by

loop closure detection, i.e., by recognizing re-visited places (Section 13.4.2). In this section, we provide a summary of the most relevant visual odometry techniques for autonomous driving. For a more detailed survey on visual odometry techniques, we refer the reader to Scaramuzza and Fraundorfer [577] and Fraundorfer and Scaramuzza [210].

In visual odometry, the goal is to recover a trajectory (i.e., a sequence of poses) of one camera or a camera system comprising multiple cameras from images. Most approaches incrementally estimate the relative transformation between two frames and integrate this information over time to recover the full trajectory. The incremental approach is illustrated in Figure 13.6. Methods on visual odometry can be roughly divided into two main categories: feature-based methods [426, 491, 578, 391, 347, 481] that extract features from key points to optimize a geometric error, and direct formulations [489, 340, 180, 182, 195, 197, 788, 746, 181] which directly operate on raw measurements by optimizing the photometric error.

Feature-based methods typically detect corners in the image and match the corresponding feature descriptors across different images. While these approaches are very efficient, they discard valuable information, e.g., straight or curved edges, that are very common in man-made environments. In contrast, direct methods leverage structural information in the entire image. Therefore, these methods usually achieve higher accuracy and robustness in environments with fewer key points. In addition, they allow to simultaneously estimate semi-dense [180, 182] and even dense depth maps [628, 489], as illustrated in Figure 13.7. However, direct methods suffer more from local minima in the optimization problem compared to feature-based methods, in particular when the pose initialization is far from the true solution. Initially, the field was dominated by feature-based methods since they are typically more efficient, but direct formulations have recently grown in popularity due to their increased accuracy [489, 628, 340, 180, 182, 195, 197, 788, 746, 181].

Feature-based Methods

2D-to-2D Matching: Depending on how corresponding points between two time steps are represented (2D or 3D), different methods must be used to obtain the camera transformation. The essential matrix (or fundamental matrix), which represents the epipolar geometry between the two cameras and contains relative pose information, can be recovered from 2D feature matches (2D-to-2D). One of the most popular algorithms for estimating the essential or fundamental matrix is the eight-point algorithm [273]. The five-point algorithm [491] is a minimal solution that only applies to the scenario of calibrated cameras. Scaramuzza et al. [578] estimate the essential matrix from monocular images with only one 2D feature correspondence using non-holonomic

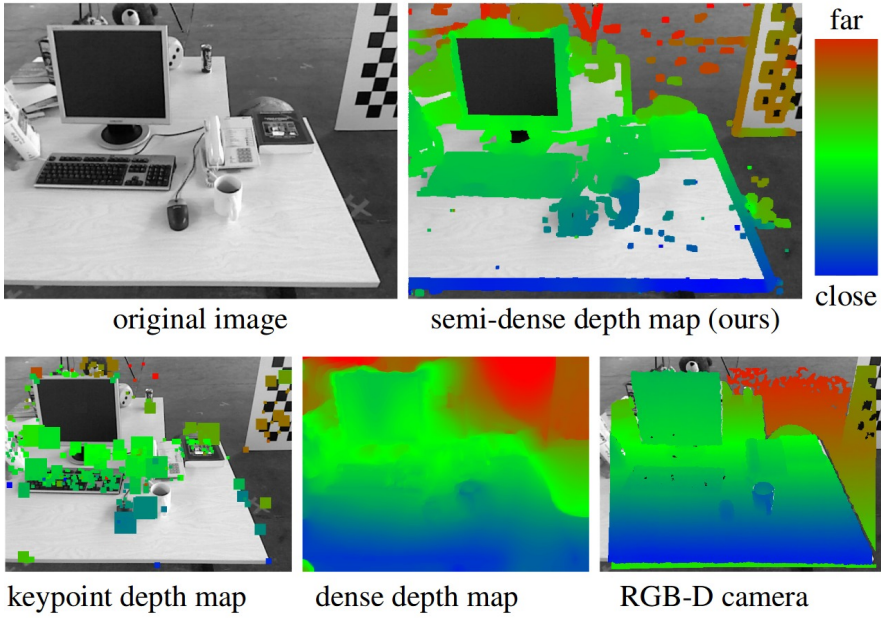


Figure 13.7: **Semi-Dense Depth Maps.** The semi-dense depth map representation of Engel et al. [180] (top right) in comparison to key points [348] (bottom left), a dense depth [446] (bottom middle), and the output of a dedicated RGB-D camera [628] (bottom right). Figure courtesy of Engel et al. [180] © 2013 IEEE.

constraints of wheeled vehicles imposing a restrictive motion model.

In general, visual odometry with monocular images cannot recover the metric scale due to the inherent scale ambiguity. Lee et al. [391] extend [578] to a novel two-point minimal solution that is able to obtain the metric scale using a multi-camera system. In contrast to the non-holonomic constraints, Lee et al. [393] assume the vertical directions to be known (from an Inertial Measurement Unit) and propose a minimal four-point and linear eight-point algorithm for a multi-camera system. Kitt et al. [347] estimate the ego-motion using trifocal geometry, which relates features between three images. Most algorithms employ RANSAC for robust estimation. The number of iterations necessary to guarantee that a correct solution is found with RANSAC depends on the number of points from which the model can be instantiated. Minimal solvers allow to the reduction of the number correspondences leading to a reduced number of iterations and runtime of the approach.

Omnidirectional cameras discussed in Section 3.1.1 enable feature-based

approaches that extract and match interest points from all around the car. The increased field of view makes the visual odometry problem more constrained and consequently allows for more accurate visual odometry. Scaramuzza and Siegwart [576] exploit this observation and estimate the ego-motion of the vehicle relative to the road from a single, central omnidirectional camera using a homography-based tracker for the ground plane and an appearance-based tracker for the rotation of the vehicle.

3D-to-2D Matching: If stereo or RGB-D information is available, a simple solution to the visual odometry problem is to project 3D features from one image into the other view and optimize for the pose by minimizing reprojection errors. Following this idea, Geiger et al. [241] present a real-time visual odometry and sparse 3D reconstruction method. They detect sparse features in stereo images using blob and corner detectors and estimate the vehicle's ego-motion by minimizing the reprojection error of the projected 3D features. In addition, they propose a real-time stereo reconstruction algorithm [240] and fuse disparity maps over time into a coherent city-scale 3D reconstruction.

3D-to-3D Matching: When dealing with 3D correspondences (3D-to-3D), the relative transformation between two time steps can be obtained by aligning the two sets of 3D features, for instance, using the iterative closest point (ICP) algorithm [49]. In visual odometry, the features extracted from images are projected into 3D using depth, whereas LiDAR-based approaches such as Zhang and Singh [768, 769] directly obtain the 3D points from the sensor. However, the triangulated 3D points from stereo will exhibit a large anisotropic uncertainty due to the small baseline and the quadratic increase of errors with respect to distance. Thus it is more natural to minimize reprojection errors in the images where error statistics can be approximated more easily. Laser-based approaches do not suffer from this problem and thus typically optimize in 3D space.

Direct Methods

In contrast to feature-based methods that optimize reprojection errors, direct approaches optimize the photometric error for estimating motion. Engel et al. [180] estimate a semi-dense inverse depth map for whole-image alignment of monocular images. Depth is estimated using multi-view stereo for pixels with non-negligible gradients and is represented by a Gaussian probability distribution. They propagate depth information from frame to frame and obtain camera poses by minimizing the photometric error. With this semi-dense formulation, they achieve comparable performance to fully dense methods [489] while not requiring a depth sensor [340]. Engel et al. [181] present a direct sparse approach for monocular visual odometry. They use a

probabilistic model and jointly optimize all model parameters (camera poses, camera intrinsics, and inverse depth) in real-time.

13.4.1 Drift

The incremental approach to ego-motion estimation greatly suffers from drift caused by the accumulation of estimation errors of the individual transformations. One way of alleviating the drift problem is to use an iterative refinement over several images that are observed most recently. In feature-based approaches, this is done by reprojecting image points into 3D by triangulation and minimizing the sum of squared reprojection errors (sliding window bundle adjustment or windowed bundle adjustment). However, simpler techniques such as a proper selection of the extracted features can also reduce drift. Kitt et al. [347] use bucketing to obtain well distributed corner-like feature matches, whereas Deigmoeller and Eggert [155] use various heuristics on flow and depth estimation to reject non-stable features.

The drift problem can also be addressed with simultaneous localization and mapping (SLAM) discussed in Section 13.4.3, which jointly estimates the location and a map of the environment to recognize places that have been visited before. The detection of already mapped places is also known as “loop closure detection”. If a loop has been detected, additional constraints can be added to the bundle adjustment problem, which leads to globally consistent maps and vehicle poses. However, poses are only corrected in hindsight, and thus, the drift problem persists during longer periods in which no loop closure can be detected. Furthermore, as loop closure detection is computationally expensive and computation increases with the length of the trajectory, such techniques are often only executed sporadically and not with every new incoming frame.

13.4.2 Loop Closure Detection

The relocalization in already mapped areas is an important subproblem of SLAM, known as loop closure detection. Relocalization is used to correct drift in the trajectory and inaccuracies in the map caused by drift.

Cummins and Newman [140] present a probabilistic approach for the recognition of places based on their appearance. They learn a generative model of appearances using a bag-of-words model as distinctive combinations of visual words will often arise from common objects. The generative model is robust and works even in visually repetitive environments. The performance of the approach is demonstrated on a self-recorded dataset and visualized in Figure 13.8. Paul and Newman [504] extend this idea by incorporating pairwise distances between words coupled to the observation of visual words using a random graph. The random graph models the pairwise distance between



Figure 13.8: **Loop Closure Detection.** Loop closure with appearance-based matching overlaid on an aerial image by Cummins and Newman [140]. Images that are matched with a probability larger than 99% are marked in red. Figure courtesy of Cummins and Newman [140] © 2008 SAGE.

words besides their distribution of occurrences. In contrast, Lee et al. [392] consider a pose graph with vertices representing camera poses and edges representing constraints between the poses. They show that the relative pose with metric scale between two loop-closing vertices can be obtained from the epipolar geometry of a multi-camera system with overlapping views.

Image-based loop closure detection can become unreliable in case of strong illumination or viewpoint changes. In contrast, LiDAR-based localization is not affected by changes in illumination and does not suffer as much from changes in viewpoint due to the captured 3D geometry and the large field of view. Dubé et al. [179] propose a loop closure detection algorithm based on matching 3D segments. Segments from the point cloud are extracted and described using a combination of descriptors. Matching of segments is performed by obtaining candidates with k-d tree search in feature space and estimating matching scores using a random forest.

13.4.3 Simultaneous Localization and Mapping (SLAM)

A detailed map of the environment simplifies planning and navigation in autonomous vehicles. However, in places for which no map is provided or the map is outdated or incomplete, the autonomous car must locate itself while generating the map. Further, the map needs to be updated continuously to

reflect environmental changes over time. In this context, SLAM refers to the task of simultaneous estimation of the location of an agent while continuously constructing a map of the environment. While SLAM addresses a similar problem as structure-from-motion techniques discussed in Section 10.2, SLAM approaches focus particularly on large-scale environments, loop-closure detection, and real-time performance.

Traditionally, a map is represented by a set of landmarks that may correspond to semantically meaningful parts or detected image features. Early approaches to SLAM have addressed the problem with Bayesian formulations using extended Kalman filters [618] or particle filters [471]. Given the last state and new observations, the current state, represented by pose, velocity, and the locations of the landmarks is recursively updated. However, this formulation is not applicable to large environments since the belief state and time complexity of the filter update grow quadratically with the number of landmarks in the map (n).

One solution for reducing complexity is to leverage filtering techniques that maintain a tractable approximation of the belief state as proposed by Paskin [503]. However, filtering may lead to inconsistent maps when applied to non-linear SLAM problems [325]. In contrast, full SLAM approaches, such as graph-based or least-squares formulations, provide more accurate solutions as they consider all poses at once. Kaess et al. [328] propose an incremental smoothing and mapping approach based on fast incremental matrix factorization. They extend their earlier work [159] on factorizing the matrix of a non-linear least-squares problem to an incremental approach that only recalculates entries which change in the matrix. Kaess et al. [327] introduce the Bayes tree, a novel data structure, which allows for a better understanding of the connection between inference in graphical models and sparse matrix factorization. Factored probability densities are encoded in the Bayes tree which naturally maps to a sparse matrix. Recently, Lenac et al. [398] proposed a filtering-based SLAM method that is able to compete with graph-based optimization techniques.

Stereo SLAM: Stereo cameras are a popular choice for tackling the SLAM problem since they allow to estimate the depth while simultaneously providing detailed information of an objects' appearance (in contrast to LiDAR sensors). Lategahn et al. [384] propose a dense stereo visual SLAM method that estimates a dense 3D map. Using a sparse visual SLAM system, they obtain the pose and a sparse map. For the dense 3D map, they compute a dense representation from stereo in a local coordinate system and continuously update the map by tracking the local coordinate systems with the sparse SLAM system. Engel et al. [183] propose LSD-SLAM, a real-time large-scale direct SLAM algorithm that couples static stereo from a camera setup with temporal multi-view stereo (Figure 13.9). This allows them to estimate the depth

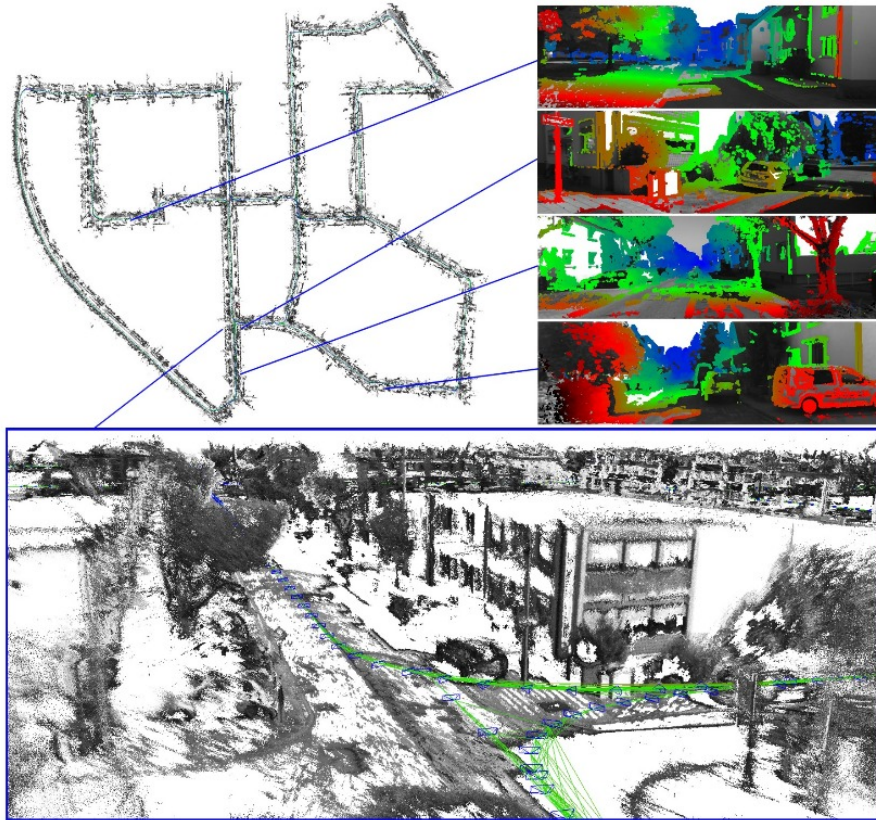


Figure 13.9: **Stereo LSD-SLAM.** Engel et al. [183] compute accurate camera movement as well as semi-dense probabilistic depth maps in real-time. The depth visualization uses blue for far away scene points and red for close objects. Figure courtesy of Engel et al. [183] © 2015 IEEE.

of pixels that are under-constrained in static stereo while avoiding scale-drift that occurs using multi-view stereo. The images are directly aligned based on the photoconsistency of high contrast pixels. Mur-Artal et al. [481] use the ORB features proposed by Rublee et al. [561] for tracking, mapping, relocalization, and loop closure. They combine methods from loop detection [228], loop closing [626, 625], and pose graph optimization [368] into a single system which they call ORB-SLAM and which became one of the most widely used SLAM systems today.

A fusion approach is proposed by Leutenegger et al. [400] in order to take advantage of the complementary nature of visual and inertial cues. They use

a non-linear optimization approach and integrate IMU measurements with reprojection errors into a joint cost function. Similarly, Usenko et al. [669] also propose a joint visual-inertial SLAM method. However, they present a fully direct method based on [183] that estimates geometry from semi-dense depth maps in contrast to sparse key points.

Environmental Changes: Changes in the environment that might not be represented in the map are a major challenge in SLAM. Levinson et al. [404] alleviate this problem by creating a map comprising of features that are very likely to be static over time. Using 3D LiDAR, they retain only flat surfaces and obtain an infrared reflectivity map of overhead views of the road surface. The map is then used to locate a vehicle with a particle filter in real-time. Levinson and Thrun [405] extend this work considering maps as probability distributions over environment properties instead of a fixed representation. Specifically, every cell of the probabilistic map is represented as its own Gaussian distribution. This allows them to represent the world more accurately and localize with fewer errors. In addition, they use offline SLAM to align multiple passes of the same environment at different times to establish an increasingly robust understanding of the world.

13.5 Datasets

Several datasets have been considered in the localization and ego-motion estimation literature. The popular dataset 7 Scenes from Shotton et al. [613] focuses only on indoor scenes. Large-scale reconstruction datasets such as Vienna [313], Dubrovnik [411], Rome [139] are very popular in particular for structure-based localization methods. With the introduction of deep learning to structure-based localization, [337] presented a new outdoor localization dataset (Cambridge Landmarks dataset), which became popular for CNN-based approaches. Most of the aforementioned datasets are limited in their variety in terms of weather conditions and seasons, which are important factors for evaluating the robustness of localization systems. To address this issue, Carlevaris-Bianco et al. [96] proposed a new long-term vision and LiDAR dataset created on the campus of the University of Michigan comprising 27 sessions. Recently, Sattler et al. [571] presented three datasets for the same problem: Aachen Day-Night, RobotCar Seasons, and CMU Seasons. [22] also recorded a dataset in different weather conditions, seasons, and during the night as well as day.

Only few datasets exist which particularly address the visual odometry problem. Most of these datasets are either small [617, 499, 56], provide only low-quality images [255], or are not yet established [438, 307]. A notable exception is the KITTI benchmark [238] discussed in Chapter 4, which provides a large dataset of challenging sequences and evaluation metrics as well as an

online evaluation server. We list the current leading monocular, stereo, and LiDAR methods on the KITTI benchmark in Table 13.1, Table 13.2, and Table 13.3, respectively.

13.6 Metrics

For image-retrieval approaches, a popular metric is the percentage of recognized queries (Recall at N). A place is considered recognized if at least one of the top N retrievals are within 25 meters from the query. For autonomous driving, this precision is not satisfactory since a higher accuracy is necessary to navigate through the environment. Consequently, localization approaches for loop closure detection [140, 504, 392] strive for higher accuracy and typically consider the precision-recall metric.

Structure-based localization approaches consider the position error (Euclidean distance between the estimated pose and the ground truth pose) as well as the orientation error. Sattler et al. [571] report the percentage of localized query images that differ from the ground truth pose using high (0.25m, 2deg), medium (0.5, 5deg), and low (5m, 10deg) accuracy thresholds.

The performance of methods for visual odometry is often measured using the Absolute Trajectory Error (ATE) or Relative Pose Error (RPE). The APE estimates the absolute distance between the estimated and ground truth trajectory. The RPE considers a fixed time interval and measures the local accuracy of the translational and rotational component. The KITTI dataset reports the average translational and rotational error measured for all possible subsequences of length (100, 200, \dots , 800) meters.

13.7 State of the Art on KITTI

Localization: A unified and established benchmark for localization methods is still missing which makes the comparison of different approaches difficult. However, several newly introduced datasets Carlevaris-Bianco et al. [96] and Sattler et al. [571], reveal open challenges to the community. Sattler et al. [571] compare two structure-based methods [569, 635] and three image retrieval approaches [658, 16, 140] on their dataset. While the structure-based methods significantly outperform the image retrieval approaches and show better robustness, all methods fail in more challenging conditions, particularly at night, when foliage changes as well as in suburban and park regions.

Monocular Visual Odometry: Monocular visual odometry methods are able to recover motion only up to a scale factor. The absolute scale can be determined by computing the size of objects in the scene, from motion constraints, or by integrating other sensors.

	Method	Translation	Rotation	Runtime
1.	DVSO [746]	0.90 %	0.0021 [deg/m]	0.1 s / GPU
2.	BVO [507]	1.76 %	0.0036 [deg/m]	0.1 s / 1 core
3.	PMO / PbT-M2 [197]	2.05 %	0.0051 [deg/m]	1 s / 1 core
4.	FTMVO [466]	2.24 %	0.0049 [deg/m]	0.11 s / 1 core
5.	PbT-M1 [195, 196]	2.38 %	0.0053 [deg/m]	1 s / 1 core
6.	MLM-SFM [621]	2.54 %	0.0057 [deg/m]	0.03 s / 5 cores
7.	RMCPE+GP [467]	2.55 %	0.0086 [deg/m]	0.39 s / 1 core
8.	EB3DTE+RJMCM [63]	5.45 %	0.0274 [deg/m]	1 s / 1 core
9.	VISO2-M + GP [621]	7.46 %	0.0245 [deg/m]	0.15 s / 1 core
10.	VISO2-M [241]	11.94 %	0.0234 [deg/m]	0.1 s / 1 core
11.	OABA [213]	20.95 %	0.0135 [deg/m]	0.5 s / 1 core

Table 13.1: **KITTI Monocular Odometry Leaderboard.** The numbers show relative translational errors and relative rotational errors, averaged over all subsequences of length 100 meters to 800 meters. Accessed on: April 2019.

Fanani et al. [195] follow a direct approach and propagate 3D key points into the next frame using relative pose predictions. Combined with the scale estimation method proposed in [196] which uses dense and sparse ground plane estimates for scale correction, they achieve competitive results in Table 13.1. However, their approach is not applicable in real-time. In contrast, Mirabdollah and Mertsching [466] follow a robust feature-based monocular visual odometry approach capable of real-time estimation using the iterative five-point method. They obtain the location of landmarks using a probabilistic triangulation method and estimate the scale of the motion from sparse low-quality features on the ground plane. Fanani et al. [197] improve the scale correction of [196] by utilizing street pixels detected with a convolutional neural network for ground plane pose estimation. Furthermore, they extend the keypoint propagation method presented in [195] which allows them to improve on previous work.

In contrast to other approaches, Pereira et al. [507] consider backward motion with a backward-facing camera or by processing the images of a forward facing-camera in reverse order. They argue that initial depth estimation of sparse feature matching approaches is not very accurate since usually, new features are initialized the first time they have been observed in the far distance. By considering the reverse order for a forward-facing camera, new features will be detected in the nearest frame, which allows more accurate depth estimates in case of forward motion.

Recently, Yang et al. [746] propose to use deep monocular depth predictions for monocular visual odometry by incorporating depth predictions into a windowed direct bundle adjustment. With this direct approach, they outperform all monocular visual odometry methods in Table 13.1. However, we

	Method	Translation	Rotation	Runtime
1.	SOFT2 [142]	0.65 %	0.0014 [deg/m]	0.1 s / 2 cores
2.	LG-SLAM [398]	0.82 %	0.0020 [deg/m]	0.2 s / 4 cores
3.	RotRocc+ [84, 87]	0.83 %	0.0026 [deg/m]	0.25 s / 2 cores
4.	GDVO [788]	0.86 %	0.0031 [deg/m]	0.09 s / 1 core
5.	SOFT [143]	0.88 %	0.0022 [deg/m]	0.1 s / 2 cores
6.	RotRocc [84]	0.88 %	0.0025 [deg/m]	0.3 s / 2 cores
7.	Stereo DSO [689]	0.93 %	0.0020 [deg/m]	0.1 s / 1 core
8.	ROCC [85]	0.98 %	0.0028 [deg/m]	0.3 s / 2 cores
9.	cv4xv1-sc [509]	1.09 %	0.0029 [deg/m]	0.145 s / GPU
10.	MonoROCC [86]	1.11 %	0.0028 [deg/m]	1 s / 2 cores
31.	VISO2-S [241]	2.44 %	0.0114 [deg/m]	0.05 s / 1 core

Table 13.2: **KITTI Odometry Stereo Leaderboard.** The numbers show relative translational errors and relative rotational errors, averaged over all subsequences of length 100 meters to 800 meters. Methods below the horizontal line show older entries for reference. Accessed on: April 2019.

remark that the KITTI dataset requires metric output, thus scale drift and scale estimation have a strong impact on the performance of the approaches.

Stereo Visual Odometry: Stereo visual odometry methods exploit the known baseline between the cameras of the stereo camera rig for estimating scale. Therefore, stereo methods are typically able to outperform monocular methods on the KITTI dataset (see Table 13.1 and Table 13.2).

Cvisic and Petrovic [143] decouple estimation of rotation and translation as translation is dependent on the scene depth while rotation is not. They estimate rotation using the five-point algorithm [491] and translation using the three-point method. Buczko and Willert [84] exploit the same idea and propose to use an initial rotation estimation to decouple rotational and translational optical flow. In contrast, Wang et al. [689] tackle the visual odometry problem with a direct method by combining static stereo with multi-view stereo as in [182, 183]. In contrast to [182, 183], they extend the energy function instead of relying on filtering approaches to update the geometry and provide an efficient bundle adjustment procedure for real-time optimization. One weakness of direct methods is that they often get stuck in local optima, especially in case of large motions. Zhu [788] address this problem with a dual Jacobian scheme for multi-scale pyramid optimization. This allows them to avoid local optima and obtain more accurate camera pose estimations that are closer to the optimal solution. In addition, they introduce a gradient-based feature representation, which improves robustness against illumination changes.

Lenac et al. [398] propose a filtering-based SLAM approach that leverages

	Method	Translation	Rotation	Runtime
1.	V-LOAM [769]	0.56 %	0.0013 [deg/m]	0.1 s / 2 cores
2.	LOAM [768]	0.59 %	0.0014 [deg/m]	0.1 s / 2 cores
3.	IMLS-SLAM [162]	0.69 %	0.0018 [deg/m]	1.25 s / 1 core
4.	MC2SLAM [486]	0.69 %	0.0016 [deg/m]	0.1 s / 4 cores
5.	LIMO2-GP [251]	0.84 %	0.0022 [deg/m]	0.2 s / 2 cores
6.	LIMO2 [251]	0.86 %	0.0022 [deg/m]	0.2 s / 2 cores
7.	CPFG-slam [320]	0.87 %	0.0025 [deg/m]	0.03 s / 4 cores
8.	LIMO [251]	0.93 %	0.0026 [deg/m]	0.2 s / 2 cores
9.	DEMO [767]	1.14 %	0.0049 [deg/m]	0.1 s / 2 cores
10.	STEAM-L WNOJ [645]	1.22 %	0.0058 [deg/m]	0.2 s / 1 core

Table 13.3: **KITTI Odometry LiDAR Leaderboard.** The numbers show relative translational errors and relative rotational errors, averaged over all subsequences of length 100 meters to 800 meters. Accessed on: April 2019.

a novel filtering solution on Lie groups. Combined with the visual odometry method proposed in [143], they are ranked second in stereo visual odometry. Cvišić et al. [142] improve the feature selection approach suggested in [143] with an age-based weighting factor suggested in [241] that gives higher weight to features that are horizontally closer to the image center. This allows them to better handle calibration errors and outperform all stereo-based methods (Table 13.2) while obtaining results competitive with LiDAR-based techniques.

Krešo and Šegvić [359] observed that camera calibration is critical for visual odometry and that the remaining calibration errors in pre-calibrated systems like KITTI have adversarial effects on the estimation results. They, therefore, propose to explicitly correct the calibration of the camera by exploiting ground truth motion which they use to recover a deformation field by optimizing the reprojection error of point feature correspondences in neighboring stereo frames.

LiDAR-based Odometry: Motivated by the impact of small calibration errors on the depth estimation of stereo-based methods [359], Gräter et al. [251] leverages depth information obtained from LiDAR for monocular visual odometry. Rejecting outliers based on a local plane assumption and fusing depth similar to [143, 85], they obtain competitive results (Table 13.3).

In contrast, Neuhaus et al. [486] directly address the SLAM problem by integrating LiDAR data with inertial measurements. The integration of IMU data allows them to cope with high-frequency motion, e.g., in off-road environments.

Inspired by RGB-D methods [488], Deschaud [162] uses an implicit surface representation [141] of the map for aligning new scans in a LiDAR SLAM

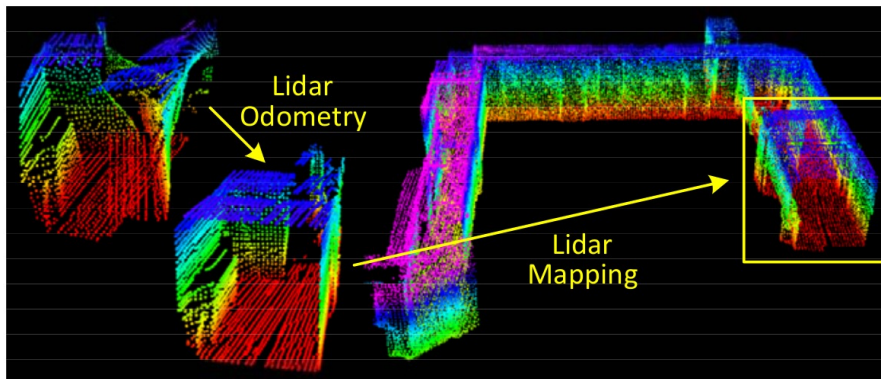


Figure 13.10: LOAM by [768] matches two consecutive LiDAR scans (LiDAR Odometry) and registers the new scan to a map (LiDAR Mapping). Figure courtesy of [768] © 2014 RSS.

approach. In combination with a specific sampling strategy based on LiDAR scans, they achieve results similar to [486].

The best performing methods on KITTI use 3D point clouds from LiDAR for ego-motion estimation (Table 13.3). Zhang and Singh [768] split the SLAM problem into LiDAR-based odometry at high frequency with low accuracy and LiDAR-mapping at low frequency with high accuracy, as illustrated in Figure 13.10. Their LiDAR-based odometry approach matches two consecutive LiDAR scans, whereas their LiDAR-based mapping approach matches and registers the new scan to a map, resulting in low drift and low computational complexity at the same time. Zhang and Singh [769] extend this work by combining visual odometry at high frequency with LiDAR-mapping at low frequency, which allows them to further improve upon their results (Table 13.3).

13.8 Discussion

While localization approaches are still missing an established unified benchmark for fair comparison and evaluation of methods, a new benchmark¹ based on multiple diverse datasets has recently been proposed by Sattler et al. [571]. Based on these results, it can be concluded that current techniques still fail to perform well in challenging real-world conditions, as identified in [96, 571]. One possible direction towards higher recall and more robustness is to incorporate deep CNN features encoding high-level information. For

¹<https://www.visuallocalization.net>

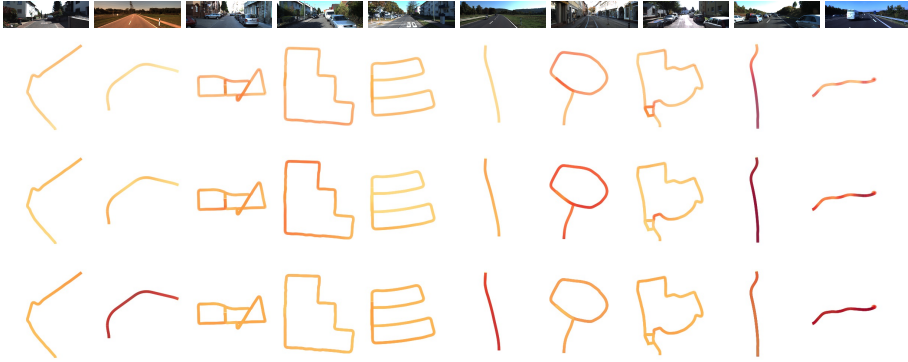


Figure 13.11: **KITTI Odometry.** From top-to-bottom: example image from the sequence, average translational error, average rotational error and speed. Averages are computed over 400 meter long trajectories and for the 15 best performing methods published on the KITTI website. Darker colors (i.e., red) indicate larger errors or higher speed.

instance, Schönberger et al. [590] and Radwan et al. [535] demonstrate that localization accuracy in challenging conditions can benefit from a semantic understanding of the environment.

In ego-motion estimation, monocular visual odometry methods can not yet compete with approaches using 3D information on the KITTI dataset. While LiDAR provides the richest source of information, stereo-based methods also achieve competitive results. In Figure 13.11, we visualize the average translational and rotational errors of the best performing visual odometry methods on the KITTI benchmark. The second row shows the translational error, and the third row shows the rotational error while the last row shows the speed. The highest translational and rotational errors are usually observed in case of strong turns. Furthermore, the error is correlated with speed and the amount of independently moving objects in the scene, which causes a decrease in the number of matched features in the background. While large errors can be observed for crowded highway scenes (second from right), only moderate errors occur when the highway is empty (right and second from left). Larger errors can also be observed in very narrow environments (fourth from right) where feature displacements are large. Overall, the most accurate motion estimation is achieved using 3D information. However, it is remarkable that state-of-the-art stereo-based methods achieve competitive results using cheap passive stereo sensors in comparison to more expensive LiDAR scanners.

Chapter 14

Scene Understanding

14.1 Problem Definition

One of the basic requirements of autonomous driving is to fully understand the surrounding area, such as a complex traffic scene. The complex task of outdoor scene understanding involves several sub-tasks such as depth estimation, scene categorization, object detection and tracking, event categorization, and more. Each of these tasks describes a particular aspect of a scene. It can be beneficial to model some of these aspects jointly in order to exploit the complementary nature of the different cues in the scene and to obtain a more holistic understanding. The goal of most scene understanding models is to obtain a rich but compact representation of the scene including its elements, e.g., layout, traffic participants, and their relation with each other.

In contrast to modeling these problems in 2D, 3D reasoning allows geometric scene understanding and results in a more informative representation of the scene in the form of 3D object models, layout elements, and occlusion relationships. In this section, we will focus on a subset of 3D scene understanding techniques that are particularly relevant to the autonomous driving task, excluding works on scene graph estimation or image tagging. One specific challenge in this context is the interpretation of urban and sub-urban traffic scenarios. Compared to highways and rural roads, urban scenarios comprise dynamic objects, a large degree of variability in the geometric layout of roads and crossroads, and an increased level of difficulty due to ambiguous visual features, occlusions, and challenging illumination conditions.

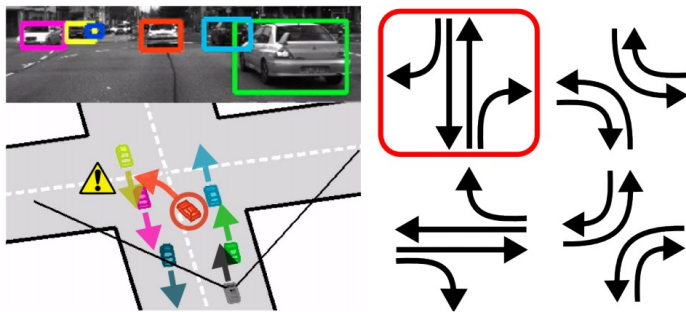


Figure 14.1: **Scene Understanding using Traffic Patterns.** Zhang et al. [766] propose to explicitly account for traffic patterns to improve scene layout and activity estimation results (right, correct situation marked in red). Figure courtesy of Zhang et al. [766] © 2013 IEEE.

14.2 Methods

While early work in computer vision [556, 269, 80, 493] already tackled the scene understanding problem from various perspectives, e.g., using a block world assumption [556] or via bottom-up top-down inference [493], most approaches relied on heuristics rather than learning and were not able to generalize to complex real-world scenes. In contrast, modern approaches try to learn complex relationships directly from data. In their pioneering work, Hoiem et al. [297] infer the overall 3D structure of an outdoor scene from a single image. The surface layout is represented as a set of coarse geometric classes with certain orientations such as support, vertical, and sky. These elements are inferred by learning an appearance-based model for each class. Oliveira et al. [495] propose a time-varying 3D representation using a set of planar polygons as primitives. Given 3D LiDAR point clouds, they find the support plane using RANSAC followed by a clustering of inliers to separate instances.

14.2.1 Road Topology and Traffic Participants

For autonomous driving, understanding the road topology and other traffic participants in the scene is of utmost importance. Ess et al. [193] use semantic segmentation as an intermediate representation to extract the road topology and to detect crosswalks and other traffic participants. In addition, their intermediate representation simultaneously encodes the spatial layout of the scene. Wojek and Schiele [709] detect vehicles and track them with a temporal filter based on a linear motion model. They also estimate the camera motion and propagate it to the next frame using a dynamic Conditional Random

Field model for joint labeling of object and scene classes. However, [193, 709] only infer a topological model of the scene and not a geometric model.

Wojek et al. [707] extend [709] to a probabilistic 3D scene model that encompasses multi-class object detection, object tracking, scene labeling, and reasoning about geometric relations. Geiger et al. [236] jointly reason about the 3D scene layout of intersections as well as the location and orientation of vehicles in the scene. They present a probabilistic generative model capturing the scene topology, geometry, and traffic activities by leveraging vehicle tracks, semantic labels, scene flow and occupancy grids.

Apart from 3D primitive-based representations, there exist other ways of representing a street scene. A more fine-grained model of the road is proposed by Topfer et al. [657]. The complex road scene is hierarchically decomposed into roads, lanes, and finally road-edges and lane-markings. This allows them to infer a more expressive model of the road compared to [236]. Seff and Xiao [595] define a list of road layout attributes such as the number of lanes, drivable directions, distance to intersections, etc. They first automatically collect a large-scale dataset for these attributes by leveraging existing street view image databases and online navigation maps (e.g., OpenStreetMap). Based on this dataset, they train a deep convolutional network to predict each attribute from a single street view image.

14.2.2 Physical and Temporal Relationships

While the detection of traffic participants is addressed in our review on object detection (Chapter 5) and object tracking (Chapter 6) approaches, scene understanding systems aim at integrating object detection and tracking with physical constraints and model the temporal behavior and relationship between traffic participants and the scene. Pellegrini et al. [505] model interactions between pedestrians (social behavior) and the scene (collisions) in a multi-target tracking formulation. Kuettel et al. [365] model spatio-temporal dependencies of moving agents in complex dynamic scenes by learning co-occurring activities and temporal rules between them. However, both approaches assume a static observer and a long observation period, i.e., the scene must be observed for a significant period of time before making a decision, therefore it is not applicable to autonomous systems.

In contrast, [711, 712, 766] consider a moving vehicle as observer and construct expressive 3D scene models by reasoning about occlusions and traffic patterns. Wojek et al. [711, 712] integrate multiple object part detectors [707] into the 3D scene model for explicit object-object occlusion reasoning (Figure 14.2). In addition, they enforce physically plausible trajectories by pruning geometrically infeasible detections. Zhang et al. [766] propose a more expressive generative model of 3D urban scenes similar to [236]. While the independent tracklets in [236] can lead to implausible inference results, they

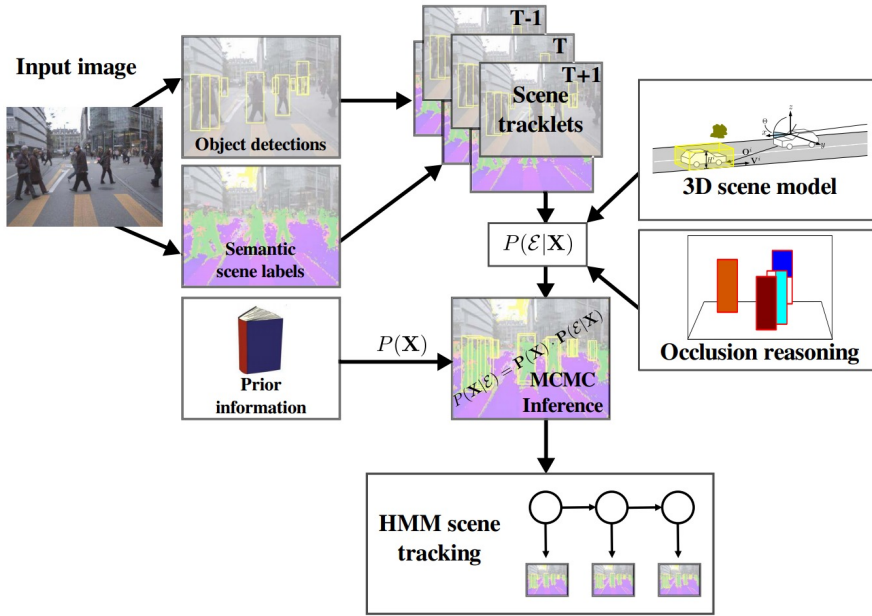


Figure 14.2: **Physical Relationships for Scene Understanding.** Overview of combined object detection and tracking system with explicit occlusion reasoning by Wojek et al. [712]. Figure courtesy of Wojek et al. [712] © 2013 IEEE.

reason about high-level semantics in the form of traffic patterns to avoid this problem (Figure 14.1) and force the solution to conform to traffic rules. This allows them to significantly improve scene estimation and vehicle-to-lane association results. Wang et al. [694] propose a top-view representation for complex road scenes that can be inferred from a single camera using a deep neural network.

14.3 Discussion

While early work on scene understanding struggled to infer expressive models of the real world, learning-based approaches led to models with increasing expressivity, ranging from simple 2D models to represent road topologies and objects [193, 709], to more complex 3D models [495, 236] which also incorporate physical [712, 694] and temporal [711, 712, 766] constraints. As motivated in [595], more expressive models can reduce the dependency on high definition maps. However, the level of expressiveness needed in autonomous driving re-

mains an open question and the accuracy achieved by state-of-the-art scene understanding models is still limited. In addition, a unified evaluation of scene understanding approaches is difficult due to the varying complexity of models and the different challenges they tackle.

Chapter 15

End-to-End Learning for Autonomous Driving

15.1 Problem Definition

Current state-of-the-art autonomous driving systems in industry are composed of numerous modules, e.g., detection (of traffic signs, lights, cars, pedestrians), segmentation (of lanes, facades), motion estimation, tracking of traffic participants, reconstruction etc. The results from these components are then typically combined in a planning module that feeds the control. However, this requires robust solutions to many open challenges in scene understanding in order to solve the problem of manipulating the car direction and speed. Furthermore, auxiliary loss functions are required to train each module (e.g., object detection, semantic segmentation) independently, hence ignoring the actual goals of the driving task which include travel time, safety, and comfort.

As an alternative, several methods consider autonomous driving as an end-to-end learning problem. In these approaches, the tasks of perception, planning, and control are combined, and a single model is trained end-to-end using a deep neural network. Most end-to-end autonomous driving systems map from sensory inputs, such as front-facing camera images, directly to driving actions such as steering angle.

15.2 Methods

End-to-end driving methods are typically trained from expert demonstrations to learn a driving policy that imitates the behavior of an expert or using reinforcement learning to explore the environment by trial and error (often

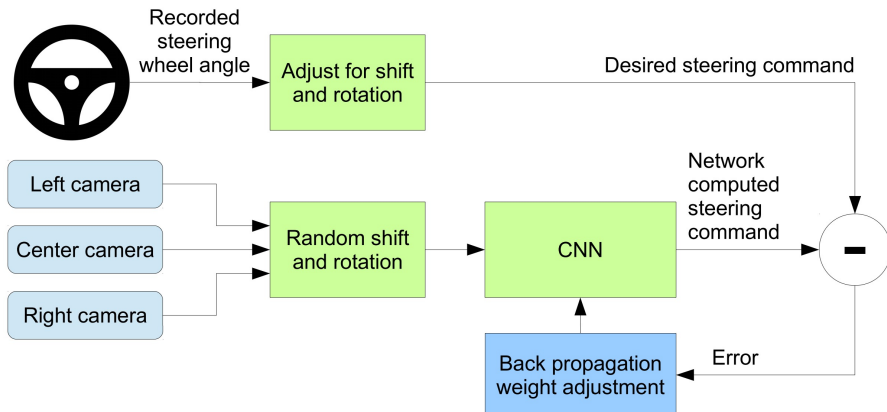


Figure 15.1: **End-to-end Learning for Lane Following.** A block diagram of an end-to-end model for lane following proposed by Bojarski et al. [60]. Conditioned on the image, a CNN estimates a steering command which is compared to the expert command for tuning the CNN weights in order to bring the CNN output closer to the desired output. Figure courtesy of Bojarski et al. [60]

in simulation). In the following sections, we first introduce the most relevant approaches proposed in the literature. We then discuss methods that combine ideas from behavior cloning and reinforcement learning. Finally, we discuss approaches that propose intermediate representations and demonstrate how driving models can be transferred from simulation to the real-world.

15.2.1 Behavior Cloning

Behavior cloning approaches learn to map sensor observations, such as RGB images, to desired driving behavior by learning to clone the behavior of an expert. Thus, these approaches fall into the category of supervised learning techniques. Most commonly, a deep neural network is employed to represent the mapping from observations to expert actions. In the 1980s, Pomerleau [526] propose ALVINN, the first demonstration of imitation learning for self-driving vehicles using a small fully connected neural network. 30 years later, Bojarski et al. [60] propose a deeper end-to-end deep convolutional neural network for lane following, illustrated in Figure 15.1, that maps images from the front-facing camera of a car to steering angles, given expert data. Xu et al. [735] propose an alternative approach and exploit large scale online datasets from uncalibrated sources to learn a driving model. Specifically, they formulate autonomous driving as a future ego-motion prediction problem. They

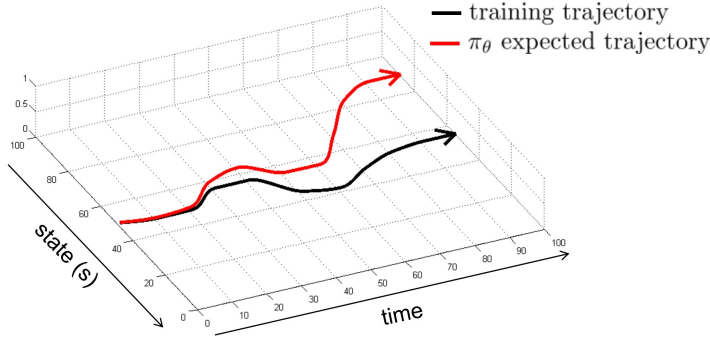


Figure 15.2: **Covariate Shift.** Covariate shift between expert demonstrations and trajectory generated by a behavior cloning policy. Errors by the behavior cloning policy π_θ compound when drifting away from expert demonstrations. Figure courtesy of Levine [402].

claim that predicting ego-motion instead of vehicle control allows their approach to generalize better to new platforms. Their deep learning architecture combines FCNs and LSTMs, and learns to predict the motion path given the current state of the agent.

Another problem with behavior cloning approach is that the training data is collected using an off-policy expert teacher, i.e., the training data is collected by rolling out the expert policy, which is different from the policy being learned. As collecting expert demonstrations for all possible situations is not practical, the training trajectories do not cover all possible states. At test time, the rollout of the behavior cloning policy thus causes it to move to a different distribution of states compared to the one it was trained on. Due to this covariate shift between the training and test time trajectories, the behavior cloning agent’s errors compound when drifting away from the expert demonstrations, as illustrated in Figure 15.2. In other words, the vehicle is likely to encounter new situations it has not been trained for and therefore acts wrongly.

In contrast, in on-policy rollout, training data is collected using the current policy being learned. Ross and Bagnell [559] propose DAgger to alleviate covariate shift by iteratively collecting corrective expert actions for the states visited by rolling out the currently learned driving policy. The driving policy parameters are then trained using the data collected on-policy. However, doing on-policy rollouts with an imperfect policy has the disadvantage of drifting and potentially reaching dangerous states, thus requiring a simulator for safe training. Laskey et al. [383] claim to provide a safer way of generating training data using expert policy with small amounts of noise injected

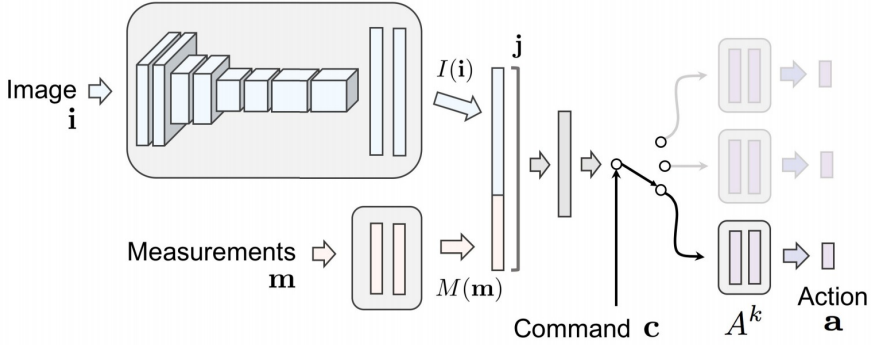


Figure 15.3: **Goal-conditional Behavior Cloning.** Architecture of goal-conditional end-to-end behavior cloning for autonomous driving proposed by Codevilla et al. [129]. The goal command acts as a switch that selects between specialized sub-policies that correspond to different commands such as lane following, turning left or turning right. Figure courtesy of Codevilla et al. [129] © 2018 IEEE.

to approximate the errors of on-policy rollout. They achieve this by iterating between learning a noise model that minimizes the covariate shift and generating data for training the behavior cloning agent.

Besides the drifting problem during test time, behavior cloning-based driving systems have other limitations. Sensor input alone is often not sufficient to uniquely infer control. Consider intersections, for example, where multiple possible actions are valid (left, right, straight). Without conditioning on the goal, all three options are acceptable. Thus, some of the behavior cloning agents, such as the one by Bojarski et al. [60], require human intervention for lane changes or turns. To alleviate this limitation, Codevilla et al. [129] propose a conditional imitation learning framework to learn a driving policy for steering and throttle control from a high-level navigational input in addition to the observations from the camera (Figure 15.3). The high-level navigational input represents the driver’s intention, such as the direction to take at the next intersection, which cannot be recovered from sensory input alone.

Codevilla et al. [130] identify other limitations of behavior cloning approaches related to generalization performance. They observe that in contrast to typical supervised learning tasks, the generalization performance for behavior cloning does not scale with training data. Moreover, they identify significant variance in performance when varying the model initialization or the order in which training examples are sampled from the dataset.

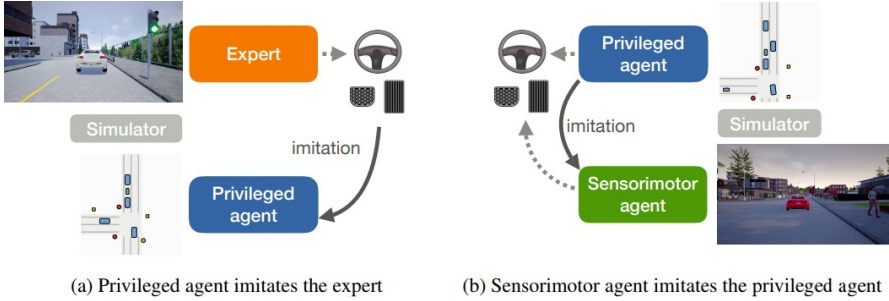


Figure 15.4: **Learning by Cheating.** Chen et al. [105] propose to first learn an agent with privileged information (a) which afterwards teaches an agent without access to privileged information (b) that learns to imitate the privileged agent. Figure courtesy of Chen et al. [105] © 2019 CoRL

Chen et al. [105] show that imitation learning can be simplified by decomposing it into two stages, as illustrated in Figure 15.4. They first train an agent that has access to privileged information. This privileged agent cheats by observing the ground-truth layout of the environment and the positions of all traffic participants. In the second stage, the privileged agent acts as a teacher that trains a purely vision-based sensorimotor agent. The resulting sensorimotor agent does not have access to any privileged information and does not cheat. They demonstrate that this approach substantially outperforms the state of the art on the CARLA benchmark and the recent NoCrash benchmark, attaining the best performance to date.

15.2.2 Reinforcement Learning

Approaches based on reinforcement learning (RL) learn to drive by training an agent that tries to maximize a user defined reward which the agent receives while interacting with the environment. In the autonomous driving application, the reward is defined by specifying the driving agent’s preferences and goals. Dosovitskiy et al. [175] propose a reinforcement learning method that trains a deep network based on a reward function provided by the CARLA simulator which combines speed, distance traveled towards the goal, collision damage, overlap with sidewalk and overlap with the opposite lane. For training the agent, they use the asynchronous advantage actor-critic (A3C) algorithm [354] which uses the value function learned by the critic to update the actor’s policy. Dosovitskiy et al. [175] observe that the RL agent performs significantly worse compared to a behavior cloning agent trained using conditional imitation learning [175] despite the fact that the RL agent was trained on a significantly larger set of visual observations. Recently, Kendall

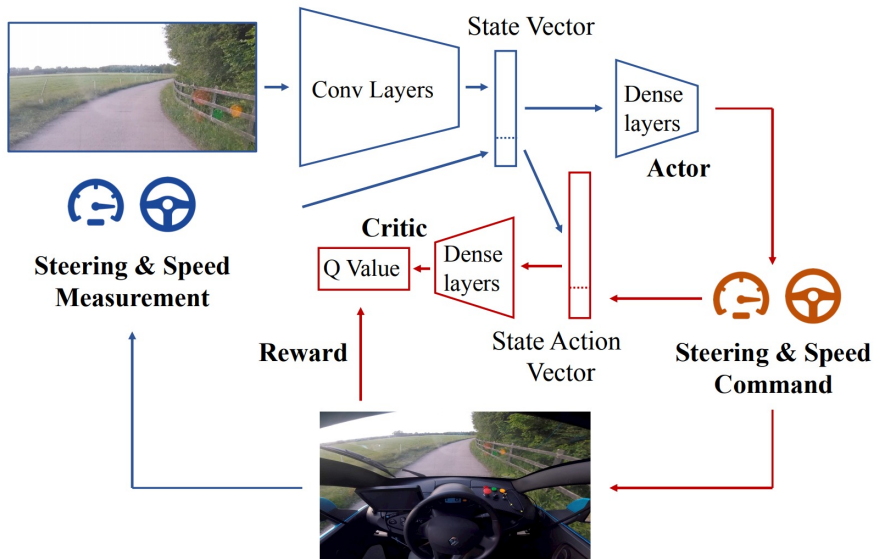


Figure 15.5: **Reinforcement Learning.** Block diagram of the autonomous system proposed by Kendall et al. [338]. The system is trained end-to-end using only the reward from the environment. The value function learned by the critic network is used to update the actor’s policy network parameters to increase the reward and improve the policy’s performance. Figure courtesy of Kendall et al. [338].

et al. [338] showed first promise in learning to drive in the real-world using a reinforcement learning agent (Figure 15.5). They use the deep deterministic policy gradients algorithm for training the RL agent and define the reward as the distance traveled by the vehicle without the safety driver taking control.

The aforementioned methods are trained using model-free reinforcement learning. The disadvantage of model-free methods is that they are often data inefficient and require a large number of interactions with the environment. In contrast, model-based reinforcement learning approaches learn a model of the environment dynamics from observational data and then exploit this model for training a driving policy. Model-based methods have been shown to significantly reduce the number of environment interactions required to learn an effective policy. However, model-based methods also typically require an interactive environment as a dynamics model trained on a fixed set of demonstrations may make incorrect predictions outside the training domain. The interactive training environment is however not practical in the real-world where such interactions are expensive and dangerous. To alleviate this

problem, Henaff et al. [288] propose to train a model-based policy which is encouraged to produce actions which the forward dynamics model is confident about. They achieve this by training the policy network to minimize an uncertainty cost which represents the mismatch between the states it induces and the states in the trained data.

15.2.3 Combined Methods

Behavior cloning methods are easy to train in a supervised fashion. However, they are poor at exploring the environment and therefore require extensive on-policy data augmentation using methods like DAgger [559]. RL approaches, in contrast, do not require per-frame supervision and are better at exploration. However, they are inefficient to train and require a simulator or non-practical trial and error runs in a real environment as well as careful design of the reward function. Therefore, several methods have been proposed to combine the strengths of both approaches.

Liang et al. [416] propose an approach to alleviate the low exploration efficiency of RL for large action space. They achieve this by constraining the policy search space by initializing the weights of the policy network of an RL algorithm by a network trained to clone the expert behavior. They observe significant improvements on the CARLA benchmark over agents trained using RL from scratch. Li et al. [407] propose an approach that learns to clone only the best behaviors of several sub-optimal teachers. They estimate the best teacher by estimating the value function of each sub-optimal teacher. The sub-optimal teachers are defined using several simple controllers over the planner output. Therefore, they do not require expert teachers for labeling data and allow for better exploration compared to learning from a single expert teacher. In addition, learning from multiple sub-optimal teachers leads to faster training compared to pure RL agents as exploration only happens from feasible states. The requirement to specify the reward function limits the practical use of Reinforcement Learning. An accurate specification of the reward requires tedious and computationally inefficient hyper-parameter tuning. Sharifzadeh et al. [607] propose to learn the unknown reward function of the driving behavior from expert demonstrations by applying Inverse Reinforcement Learning (IRL). In contrast to behavior cloning approaches that directly learn the observation-control mapping in a supervised fashion, Inverse Reinforcement Learning approaches claim to offer better generalization by learning a reward function that explains the expert behavior.

15.2.4 Intermediate Representations

Instead of directly learning a mapping from pixels to actions, Chen et al. [104] present an approach which first estimates a small number of human

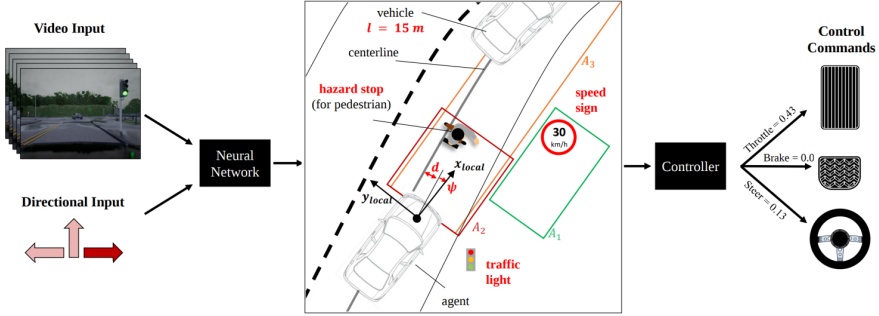


Figure 15.6: **Conditional Affordance Learning.** The input video and the high-level directional commands are fed into a neural network which predicts a set of affordances such as presence of red traffic lights or distance to the lane center. These affordances are used by a controller to compute the control output. Figure courtesy of Sauer et al. [573] © 2018 IEEE.

interpretable, pre-defined affordance measures such as the angle of the car relative to the road, the distance to the lane markings, and the distance to cars in the current and adjacent lane. These predicted affordances are then mapped to car actions using a rule-based controller to enable autonomous driving in the TORCS car racing simulation [725]. The advantage of mid-level representations is that the network predicting the mid-level representations can be trained and validated before deploying them. In addition, the mid-level representations are more interpretable compared to traditional behavior cloning approaches. Similarly, Sauer et al. [573] estimate several affordances from sensor inputs in order to drive a car, as illustrated in Figure 15.5. In contrast to Chen et al. [104], they consider the more challenging scenario of urban driving using the CARLA simulator [175]. In CARLA, the agent needs to obey traffic rules such as speed limits, red lights, avoid colliding with obstacles on the road and navigate at junctions with multiple possible driving directions. Sauer et al. [573] realize their driving agent by expanding the set of affordances to cover the most important aspects of urban environments. Similar to Chen et al. [104], they use a rule-based controller to map affordances to vehicle controls.

Recently, Zhou et al. [783] studied the significance of using intermediate representations pursued in computer vision research such as depth, segmentation, optical flow for improving several sensorimotor tasks such as urban driving. They observed that an agent that takes as input one or more of these intermediate representations along with the image learns significantly better sensorimotor control than an agent which uses just the raw image as input. They observed significant improvements even when the intermediate

representations were noisy predictions by a simple deep network. Bansal et al. [32] propose a perception module that translates raw sensor observations to a mid-level representation. Their representation includes a top-down rendering of the environment where 2D boxes of vehicles are drawn along with a rendering of the road information and traffic light states. They use this mid-level representation as input to a recurrent neural network (RNN) which outputs the control command. Similarly, Wang et al. [686] infer the depth and poses of the objects present in the scene from front-facing camera images and project the objects into an overhead view. They train a behavior cloning agent over the concatenation of front-facing and overhead images and observe improved performance over an agent trained only on front-facing images. In the same spirit, Müller et al. [478] train a driving policy in CARLA with mid-level representations as input. Specifically, they used binary segmentation estimated from a scene segmentation network as input to the driving policy network and observed improvements over an agent trained on raw camera images.

Similar to the aforementioned methods, Mehta et al. [452] also propose to use intermediate visual affordances such as “distance to intersection”, and action primitives such as “slow down” as input to the driving policy network. However, in contrast to the aforementioned works, they predict visual affordances and action primitives as an auxiliary task to the driving control task. They claim that predicting representations which are crucial for the driving decision allow the policy network to learn superior internal representations leading to more efficient training and better generalization. Kendall et al. [338] studied the importance of using an intermediate representation for state representation instead of raw pixels for learning a reinforcement learning-based driving policy. They observed significant improvements in data efficiency in training the driving policy using a compressed representation of the raw image, obtained using a Variational Autoencoder (VAE)

15.2.5 Transferring from Simulation to the Real World

One major limitation of reinforcement learning is the necessity of a simulation environment for trial and error. Thus, during training only synthetic data is considered and the models usually do not generalize to real data. To address this problem, Pan et al. [498] propose to transfer a reinforcement learning agent trained in a virtual environment to the real-world. More specifically, they learn an image translation network to translate non-realistic simulated images to realistic images. Their translation network is composed of two conditional GANs, the first for segmenting virtual images from the simulator, and the second for translating the segmented images to their realistic counterparts. In the same spirit, Bewley et al. [51] propose to train an image-to-image translation network for transferring a driving policy from simulation to real-world without any real-world control labels (Figure 15.5). In contrast to

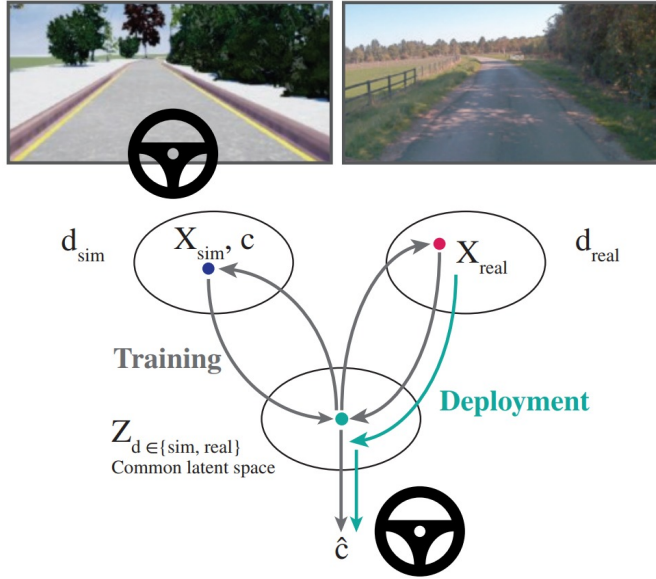


Figure 15.7: **From Simulation to the Real World.** Bewley et al. [51] proposed a model for end-to-end driving by learning to translate between simulated and real-world images, jointly learning a control policy from the common latent space Z using expert labels in simulation. Their method does not require real-world control labels and is able to learn a policy which can be transferred with improved generalization to real-world driving. Figure courtesy of Bewley et al. [51]

Pan et al. [498], which uses an explicit semantic segmentation as intermediate representation, they use an implicit latent structure as intermediate representation. They propose two autoencoder-like networks for translating between domains where a common latent space is learned through direct and cyclic losses. Their control network is trained using behavior cloning by passing the latent code as input to the control network.

15.3 Datasets

As behavior cloning approaches can be trained on offline expert demonstrations, several public datasets have been introduced in the last few years to train and evaluate such methods. The comma.ai dataset [567] provides 7.25 hours of driving data with a camera in the windshield, capturing images of the road at 20Hz. The dataset also provides observations from several other

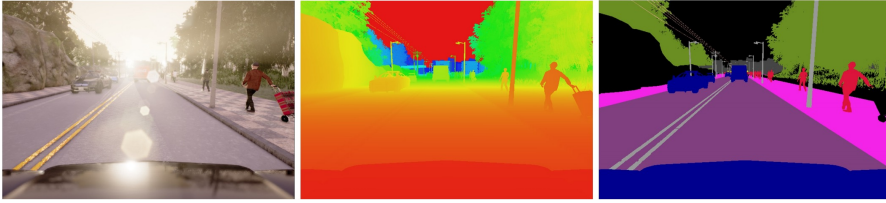


Figure 15.8: **Sensor Modalities in CARLA.** From left to right: RGB image, ground-truth depth and ground-truth semantic segmentation. Additional sensor models can be plugged in via the provided API. Figure courtesy of Dosovitskiy et al. [175] © 2017 IEEE.

sensors such as car speed, steering angle, GPS, gyroscope, and IMU. However, the dataset only contains training examples for highway scenarios, and is therefore not suitable for learning a driving policy which operates in more challenging situations such as in cities. The Berkeley DeepDrive Video dataset [735] comprises 10,000 hours of driving in cities, on highways, in towns, and in rural areas. The dataset has been recorded using forward-facing dash cameras along with observations from sensors such as GPS, IMU, gyroscope, and magnetometer. As discussed in the previous section, online rollouts of the driving policy is an important requirement for training and evaluation of most end-to-end learning methods. However, deploying a partially trained model in a real environment to collect training data is both dangerous and impractical. Therefore, realistic driving simulators are a key requirement for training and evaluating these models.

As one of the first open-source simulators, the TORCS racing car simulator [725] has been used for learning and evaluation of road lane following by Chen et al. [104]. However, the TORCS environments is simplistic, lacking complexities such as traffic participants, junctions, etc.

In contrast, CARLA [175] provides a more realistic, complex and flexible open-source simulator for autonomous driving that enables training and validation in urban driving conditions. It provides high quality images along with ground-truth depth and semantic segmentation as pseudo-sensors, as illustrated in Figure 15.8. In order to replicate the complex nature of urban driving, the environments in CARLA exhibit realistic urban street layouts with traffic rules, intersections, buildings, pedestrians, street signs and other traffic participants. The simulator also provides different weather and lighting conditions in order to evaluate the generalization ability of the driving agent. CARLA also provides a benchmark based on four increasingly difficult driving tasks and is actively expanded in terms of the environments, assets and agents it provides.

However, existing real-world datasets and synthetic simulators often fail

to capture the long tail of the distribution which covers important but rare situations. These rare events can only be effectively captured with a large fleet of vehicles that log these situations in real-world driving. Tesla’s Autopilot system [647] is a dormant logging-only mode that can be queried for multiple instances of rare failure situations so that the model can be trained to avoid such failures. In addition, Shadow Mode allows Tesla to validate the Autopilot system running in the background in real situations. However, data from Tesla vehicles are proprietary and hence not released to other companies or public research institutions.

15.4 Metrics

There are no standard metrics and benchmarks for autonomous driving and thus most methods usually evaluate on their own set of metrics and datasets. The most popular benchmark CARLA [175] evaluates on two metrics. First, the percentage of successfully completed episodes under the four different conditions provided by CARLA. And second, the average distance (in kilometers) driven between two infractions. Infractions include driving on the opposite lane, driving on the sidewalk, colliding with other vehicles, colliding with pedestrians, and hitting static objects. Codevilla et al. [128] use CARLA to analyze the correlation between offline and online metrics for evaluation of autonomous driving agents. They observe that offline metrics such as the squared or absolute error of the steering angle are poorly correlated with online metrics such as the success rate of reaching the goal. Their work highlights the tension between imitation learning and reinforcement learning. While reinforcement learning allows to train for the desired goal, training an imitation learning agent is significantly easier and does not require potentially unsafe exploration.

15.5 Discussion

A common characteristic of most end-to-end driving methods is the need to collect online training data. While behavior cloning methods have shown promising results by learning purely from expert demonstrations, minimizing the covariate shift between the expert trajectories and the agent’s policy is still an open problem. Similarly, reinforcement learning approaches require millions of trial and error runs and thus can only be safely applied in simulation environments. Moreover, as shown by [128], offline evaluation metrics are poorly correlated with online driving performance. Therefore, safe training and validation of end-to-end learning models require further development of realistic simulators such as CARLA [175] on which these methods can be

trained before transferring the resulting policies to the real-world. Flexibility is another desirable characteristic when developing simulators: the resulting simulations should allow for highly diverse and complex scenarios, yet also model the long tail of the data distribution to capture rare events. While realistic simulation environments are important, there will likely remain a domain gap between simulated and real data. Therefore, another critical direction of future research is the design of end-to-end learning methods which can be robustly transferred from simulated environments to the real-world. Furthermore, the lack of interpretability of end-to-end driving networks prevents deeper insights into the modes of operation (in particular legal relevant failure cases) and thus requires further investigation.

Chapter 16

Conclusion

This book provides a comprehensive survey on problems, datasets, and methods in computer vision for autonomous vehicles. Towards this goal, we considered the historically most relevant literature as well as the state of the art on several relevant topics, including recognition, reconstruction, motion estimation, tracking, scene understanding, and end-to-end learning. We discussed open problems and current research challenges in each of these areas and also provided a novel in-depth analysis of the KITTI benchmark.

While self-driving vehicles have a long history, it remains difficult to make predictions when self-driving vehicles will hit the consumer market. Traditionally, the problems involved in achieving or surpassing human-level performance on this task have been underestimated. Difficulties include the high accuracy that needs to be attained, the robustness required for safe self-driving as well as adverse weather conditions (snow, rain, night). Furthermore, most self-driving systems rely on accurate HD maps for localization and detection of static infrastructure, which are hard to create and to maintain up-to-date. In addition, some of the most challenging scenarios are less structured (parking areas, complex roundabouts) and thus need to be mastered without HD maps. Pedestrians pose another challenge to self-driving vehicles as their behavior is often erratic, and communication with them can be key for making a driving decision. Other challenges include complex planning tasks such as merging into traffic and negotiating with other vehicles. Further, several ethical and legal questions need to be addressed before self-driving vehicles can be deployed in large numbers on public roads.

From a technical perspective, modular pipelines offer the advantage of parallelization, interpretability, and ease of introducing prior knowledge. However, human-engineered modules often rely on heuristics or intuitions, which may be inaccurate or wrong. Learning driving policies from data is an attractive alternative, however bridging the gap to modular and interpretable

systems as well as attaining human-level performance remain unsolved problems to date. A particularly challenging problem is generalization to unseen environments and to handle rare events for which little data is available.

We are at an exciting time where self-driving technology receives considerable attention and progress is fast. At the same time, it is of prime importance that we stay objective and cautious with the claims that we make in order not to gamble people's trust in this new technology or put people's lives at stake. Writing a survey on this rapidly evolving field was a major tour de force. We are well aware that some of the approaches surveyed in this work might be outdated in the near future. However, some of the works presented in this survey will stand the test of time and will be remembered as landmarks in the development of autonomous vehicles. We hope that this survey, in combination with our online navigation tool¹, will become useful references, encourage new research, and ease the entry for beginners starting in this exciting field.

16.1 Acknowledgement

We thank Raghudeep Gadde, Varun Jampani, Yiyi Liao, Despoina Paschali-dou, Jörg Stückler, Torsten Sattler, Siyu Tang, and Osman Ulusoy for sharing their expert knowledge and giving us valuable feedback on early versions of the draft. We also highly appreciate the help of Davide Scaramuzza, Daniel Maturana, and many others in the community for their feedback and suggestions for related work. Finally, we would like to thank all researchers who gave us permission to use the figures from their papers and greatly appreciate the support of Benjamin Coors and Jonas Wulff who provided additional illustrations of their work.

¹http://www.cvlibs.net/projects/autonomous_vision_survey

Bibliography

- [1] David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler. “Efficient Interactive Annotation of Segmentation Datasets with Polygon-RNN++”. In: (2018).
- [2] Daimler AG. *Bosch and Daimler. Metropolis in California to become a pilot city for automated driving*. <https://www.daimler.com/innovation/case/autonomous/pilot-city-for-automated-driving.html>. Online: accessed 18-October-2019. 2019.
- [3] Sameer Agarwal, Noah Snavely, Ian Simon, Steven M. Seitz, and Richard Szeliski. “Building Rome in a Day”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2009.
- [4] Hamed Habibi Aghdam, Elnaz Jahani Heravi, and Domenec Puig. “A practical approach for detection and classification of traffic signs using Convolutional Neural Networks”. In: *Robotics and Autonomous Systems (RAS)* 84 (2016), pp. 97–112.
- [5] Jose M. Alvarez, Theo Gevers, and Antonio M. Lopez. “3D Scene Priors for Road Detection”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2010.
- [6] José Manuel Álvarez, Theo Gevers, Yann LeCun, and Antonio M. López. “Road Scene Segmentation from a Single Image”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [7] José Manuel Álvarez and Antonio M. López. “Road Detection Based on Illuminant Invariance”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 12.1 (2011), pp. 184–193.
- [8] Mohamed Aly. “Real time detection of lane markers in urban streets”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2008.
- [9] P. Anandan. “A computational framework and an algorithm for the measurement of visual motion”. In: *International Journal of Computer Vision (IJCV)* 2.3 (1989), pp. 283–310.

- [10] Henrik Andreasson and Achim J. Lilienthal. “6D scan registration using depth-interpolated local image features”. In: *Robotics and Autonomous Systems (RAS)*. 2010.
- [11] M. Andriluka, S. Roth, and B. Schiele. “People-Tracking-by-Detection and People-Detection-by-Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [12] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. “Monocular 3D Pose Estimation and Tracking by Detection”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2010.
- [13] Anton Andriyenko and Konrad Schindler. “Multi-target tracking by continuous energy minimization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.
- [14] Anton Andriyenko, Konrad Schindler, and Stefan Roth. “Discrete-continuous optimization for multi-target tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [15] Dragomir Anguelov, Carole Dulong, Daniel Filip, Christian Frueh, Stéphane Lafon, Richard Lyon, Abhijit S. Ogale, Luc Vincent, and Josh Weaver. “Google Street View: Capturing the World at Street Level”. In: *IEEE Computer* 43.6 (2010), pp. 32–38.
- [16] R. Arandjelović, P. Gronat, A. Torii, T. Pajdla, and J. Sivic. “NetVLAD: CNN architecture for weakly supervised place recognition”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [17] Pablo Andrés Arbeláez, Jordi Pont-Tuset, Jonathan T. Barron, Ferran Marqués, and Jitendra Malik. “Multiscale Combinatorial Grouping”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [18] Anurag Arnab and Philip H. S. Torr. “Pixelwise Instance Segmentation with a Dynamically Instantiated Network”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [19] Maryam Babaei, Ali Athar, and Gerhard Rigoll. “Multiple People Tracking Using Hierarchical Deep Tracklet Re-Identification”. In: *arXiv.org* (2018).
- [20] H. Badino, U. Franke, and R. Mester. “Free Space Computation Using Stochastic Occupancy Grids and Dynamic Programming”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV) Workshops*. 2007.
- [21] Hernan Badino, Uwe Franke, and David Pfeiffer. “The Stixel World - A Compact Medium Level Representation of the 3D-World”. In: *Proc. of the DAGM Symposium on Pattern Recognition (DAGM)*. 2009.

- [22] Hernan Badino, Daniel Huber, and Takeo Kanade. “Real-Time Topometric Localization”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. May 2012.
- [23] Vijay Badrinarayanan, Ignas Budvytis, and Roberto Cipolla. “Mixture of trees probabilistic graphical model for video segmentation”. In: *International Journal of Computer Vision (IJCV)* 110.1 (2014), pp. 14–29.
- [24] Vijay Badrinarayanan, Fabio Galasso, and Roberto Cipolla. “Label Propagation in Video Sequences”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2010.
- [25] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. “SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 39.12 (2017), pp. 2481–2495.
- [26] Min Bai, Wenjie Luo, Kaustav Kundu, and Raquel Urtasun. “Exploiting Semantic Information and Deep Matching for Optical Flow”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [27] Min Bai and Raquel Urtasun. “Deep Watershed Transform for Instance Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 2858–2866.
- [28] Christian Bailer, Bertram Taetz, and Didier Stricker. “Flow Fields: Dense Correspondence Fields for Highly Accurate Large Displacement Optical Flow Estimation”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [29] Simon Baker, Daniel Scharstein, J. Lewis, Stefan Roth, Michael Black, and Richard Szeliski. “A Database and Evaluation Methodology for Optical Flow”. In: *International Journal of Computer Vision (IJCV)* 92 (2011), pp. 1–31.
- [30] Vassileios Balntas. *SILDa: A Multi-Task Dataset for Evaluating Visual Localization*. <https://medium.com/scape-technologies/silda-a-multi-task-dataset-for-evaluating-visual-localization-7fc6c2c56c74>. Online: accessed 17-June-2019. 2019.
- [31] Vassileios Balntas, Shuda Li, and Victor Prisacariu. “RelocNet: Continuous Metric Learning Relocalisation Using Neural Nets”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 782–799.
- [32] Mayank Bansal, Alex Krizhevsky, and Abhijit Ogaler. “ChauffeurNet: Learning to Drive by Imitating the Best and Synthesizing the Worst”. In: *arXiv.org* (2018).

- [33] Mayank Bansal, Harpreet S. Sawhney, Hui Cheng, and Kostas Daniilidis. “Geo-localization of street views with aerial image databases”. In: *Proc. of the International Conf. on Multimedia (ICM)*. 2011.
- [34] S.Y. Bao, M. Chandraker, Yuanqing Lin, and S. Savarese. “Dense Object Reconstruction with Semantic Priors”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [35] Nick Barnes, Alexander Zelinsky, and Luke S Fletcher. “Real-time speed sign detection using the radial symmetry detector”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 9.2 (2008), pp. 322–332.
- [36] Aseem Behl, Omid Hosseini Jafari, Siva Karthik Mustikovela, Hassan Abu Alhaija, Carsten Rother, and Andreas Geiger. “Bounding Boxes, Segmentations and Object Coordinates: How Important is Recognition for 3D Scene Flow Estimation in Autonomous Driving Scenarios?” In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [37] Aseem Behl, Despoina Paschalidou, Simon Donne, and Andreas Geiger. “PointFlowNet: Learning Representations for Rigid Motion Estimation from Point Clouds”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [38] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Juergen Gall. “A Dataset for Semantic Segmentation of Point Cloud Sequences”. In: *arXiv.org* (2019).
- [39] Jens Behley, Volker Steinhage, and Armin B. Cremers. “Laser-based Segment Classification Using a Mixture of Bag-of-Words”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2013.
- [40] Jens Behley, Volker Steinhage, and Armin B. Cremers. “Performance of histogram descriptors for the classification of 3D laser range data in urban environments”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2012.
- [41] Jorge Beltrán, Carlos Guindel, Francisco Miguel Moreno, Daniel Cruzado, Fernando García, and Arturo de la Escalera. “BirdNet: A 3D Object Detection Framework from LiDAR Information”. In: *Proc. IEEE Conf. on Intelligent Transportation Systems (ITSC)*. 2018, pp. 3517–3523.
- [42] Rodrigo Benenson, Markus Mathias, Radu Timofte, and Luc J. Van Gool. “Pedestrian detection at 100 frames per second”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.

- [43] Rodrigo Benenson, Mohamed Omran, Jan Hendrik Hosang, and Bernt Schiele. “Ten Years of Pedestrian Detection, What Have We Learned?”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [44] Jérôme Berclaz, Francois Fleuret, and Pascal Fua. “Multiple Object Tracking using Flow Linear Programming”. In: *Performance Evaluation of Tracking and Surveillance* (2009).
- [45] Jerome Berclaz, Francois Fleuret, Engin Turetken, and Pascal Fua. “Multiple Object Tracking Using K-Shortest Paths Optimization”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 33.9 (2011), pp. 1806–1819.
- [46] Gedas Bertasius, Lorenzo Torresani, and Jianbo Shi. “Object Detection in Video with Spatiotemporal Sampling Networks”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [47] Massimo Bertozzi, Luca Bombini, Alberto Broggi, Michele Buzzoni, Elena Cardarelli, Stefano Cattani, Pietro Cerri, Alessandro Coati, Stefano Debattisti, Andrea Falzoni, Rean Isabella Fedriga, Mirko Felisa, Luca Gatti, Alessandro Giacomazzo, Paolo Grisleri, Maria Chiara Laghi, Luca Mazzei, Paolo Medici, Matteo Panciroli, Pier Paolo Porta, Paolo Zani, and Pietro Versari. “VIAC: An out of ordinary experiment”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2011, pp. 175–180.
- [48] Massimo Bertozzi, Alberto Broggi, and Alessandra Fascioli. “Vision-based intelligent vehicles: State of the art and perspectives”. In: *Robotics and Autonomous Systems (RAS)* 32.1 (2000), pp. 1–16.
- [49] P.J. Besl and H.D. McKay. “A method for registration of 3D shapes”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 14 (1992), pp. 239–256.
- [50] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft. “Simple online and realtime tracking”. In: *Proc. IEEE International Conf. on Image Processing (ICIP)*. 2016.
- [51] Alex Bewley, Jessica Rigley, Yuxuan Liu, Jeffrey Hawke, Richard Shen, Vinh Dieu Lam, and Alex Kendall. “Learning to Drive from Simulation without Real World Labels”. In: *arXiv.org abs/1812.03823* (2018).
- [52] Rahul Bhotika, David J. Fleet, and Kiriakos N. Kutulakos. “A Probabilistic Theory of Occupancy and Emptiness”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2002.
- [53] Jonathan Binas, Daniel Neil, Shih-Chii Liu, and Tobi Delbrück. “DDD17: End-To-End DAVIS Driving Dataset”. In: *Proc. of the International Conf. on Machine learning (ICML) Workshops* (2017).

- [54] Michael J. Black and P. Anandan. “A framework for the robust estimation of optical flow”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 1993.
- [55] Maros Blaha, Christoph Vogel, Audrey Richard, Jan D. Wegner, Thomas Pock, and Konrad Schindler. “Large-Scale Semantic 3D Reconstruction: An Adaptive Multi-Resolution Model for Multi-Class Volumetric Labeling”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [56] José-Luis Blanco-Claraco, Francisco-Ángel Moreno-Dueñas, and Javier González Jiménez. “The Málaga urban dataset: High-rate stereo and LiDAR in a realistic urban scenario”. In: *International Journal of Robotics Research (IJRR)* 33.2 (2014), pp. 207–214.
- [57] Erik Bochinski, Volker Eiselein, and Thomas Sikora. “High-Speed tracking-by-detection without using image information”. In: *Proc. of International Conf. on Advanced Video and Signal Based Surveillance (AVSS)*. 2017, pp. 1–6.
- [58] András Bódis-Szomorú, Hayko Riemenschneider, and Luc Van Gool. “Efficient Volumetric Fusion of Airborne and Street-Side Data for Urban Reconstruction”. In: *Proc. of the International Conf. on Pattern Recognition (ICPR)*. 2016.
- [59] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter V. Gehler, Javier Romero, and Michael J. Black. “Keep It SMPL: Automatic Estimation of 3D Human Pose and Shape from a Single Image”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [60] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. “End to End Learning for Self-Driving Cars”. In: *arXiv.org* 1604.07316 (2016).
- [61] Amol Borkar, Monson Hayes, and Mark T. Smith. “A Novel Lane Detection System With Efficient Ground Truth Generation”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 13.1 (2012), pp. 365–374.
- [62] Jean-Yves Bouguet. *Camera Calibration Toolbox for Matlab*. 2010. URL: http://www.vision.caltech.edu/bouguetj/calib_doc.
- [63] Zeyd Boukhers, Kimiaki Shirahama, and Marcin Grzegorzek. “Less restrictive camera odometry estimation from monocular camera”. In: *Multimedia Tools Appl.* 77.13 (2018), pp. 16199–16222.

- [64] Yuri Boykov, Olga Veksler, and Ramin Zabih. “Fast Approximate Energy Minimization via Graph Cuts”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 23 (1999), p. 2001.
- [65] Bert De Brabandere, Davy Neven, and Luc Van Gool. “Semantic Instance Segmentation with a Discriminative Loss Function”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops* (2017).
- [66] Samarth Brahmabhatt, Jinwei Gu, Kihwan Kim, James Hays, and Jan Kautz. “Geometry-Aware Learning of Maps for Camera Localization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 2616–2625.
- [67] D. Braid, A. Broggi, and G. Schmiedel. “The TerraMax Autonomous Vehicle”. In: *Journal of Field Robotics (JFR)* (2006).
- [68] Markus Braun, Sebastian Krebs, Fabian Flohr, and Dariu M. Gavrilă. “The EuroCity Persons Dataset: A Novel Benchmark for Object Detection”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* (2019).
- [69] Markus Braun, Qing Rao, Yikang Wang, and Fabian Flohr. “Pose-RCNN: Joint object detection and pose estimation using 3D object proposals”. In: *Proc. IEEE Conf. on Intelligent Transportation Systems (ITSC)*. IEEE. 2016, pp. 1546–1551.
- [70] Technische Universität Braunschweig. *Project Stadtpilot*. <https://www.tu-braunschweig.de/stadtpilot>. Online: accessed 18-October-2019. 2010.
- [71] Kristian Bredies, Karl Kunisch, and Thomas Pock. “Total Generalized Variation”. In: *Journal of Imaging Sciences (SIAM)* 3.3 (2010), pp. 492–526.
- [72] Michael D. Breitenstein, Fabian Reichlin, Bastian Leibe, Esther Koller-Meier, and Luc J. Van Gool. “Online Multiperson Tracking-by-Detection from a Single, Uncalibrated Camera”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 33.9 (2011), pp. 1820–1833.
- [73] Michael D. Breitenstein, Fabian Reichlin, Bastian Leibe, Esther Koller-Meier, and Luc Van Gool. “Robust Tracking-by-Detection using a Detector Confidence Particle Filter”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2009.
- [74] William Brendel, Mohamed R. Amer, and Sinisa Todorovic. “Multi-Object Tracking as Maximum Weight Independent Set”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.

- [75] A. Broggi, M. Bertozzi, A. Fascioli, and G. Conte. *Automatic Vehicle Guidance: the Experience of the ARGO Vehicle*. Singapore: World Scientific, 1999.
- [76] A. Broggi, M. Bertozzi, A. Fascioli, and M. Sechi. “Shape-based Pedestrian Detection”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2000.
- [77] Alberto Broggi, Pietro Cerri, Stefano Debattisti, Maria Chiara Laghi, Paolo Medici, Daniele Molinari, Matteo Panciroli, and Antonio Prioretti. “PROUD - Public Road Urban Driverless-Car Test”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 16.6 (2015), pp. 3508–3519.
- [78] Alberto Broggi, Pietro Cerri, Paolo Medici, Pier Paolo Porta, and Guido Ghisio. “Real time road signs recognition”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. IEEE. 2007, pp. 981–986.
- [79] Alberto Broggi, Paolo Medici, Elena Cardarelli, Pietro Cerri, Alessandro Giacomazzo, and Nicola Finardi. “Development of the control system for the Visslab Intercontinental Autonomous Challenge”. In: *Proc. IEEE Conf. on Intelligent Transportation Systems (ITSC)*. 2010, pp. 635–640.
- [80] Rodney A. Brooks. “Model-Based Three Dimensional Interpretations of Two Dimensional Images”. In: *Proc. of the International Joint Conf. on Artificial Intelligence (IJCAI)*. 1981, pp. 619–624.
- [81] T. Brox and J. Malik. “Large Displacement Optical Flow: Descriptor Matching in Variational Motion Estimation”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 33 (Mar. 2011), pp. 500–513.
- [82] Marcus A. Brubaker, Andreas Geiger, and Raquel Urtasun. “Map-Based Probabilistic Visual Self-Localization”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 38.4 (2016), pp. 652–665.
- [83] A. Bruhn and J. Weickert. “A Confidence Measure for Variational Optic flow Methods”. In: *Geometric Properties for Incomplete Data*. Ed. by Reinhard Klette, Ryszard Kozera, Lyle Noakes, and Joachim Weickert. Dordrecht: Springer Netherlands, 2006, pp. 283–298. ISBN: 978-1-4020-3858-7. DOI: 10.1007/1-4020-3858-8_15.
- [84] Martin Buczko and Volker Willert. “Flow-Decoupled Normalized Reprojection Error for Visual Odometry”. In: *Proc. IEEE Conf. on Intelligent Transportation Systems (ITSC)*. 2016.

- [85] Martin Buczko and Volker Willert. “How to distinguish inliers from outliers in visual odometry for high-speed automotive applications”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2016.
- [86] Martin Buczko and Volker Willert. “Monocular Outlier Detection for Visual Odometry”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2017, pp. 739–745.
- [87] Martin Buczko, Volker Willert, Julian Schwehr, and Jürgen Adamy. “Self-Validation for Automotive Visual Odometry”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2018, pp. 1–6.
- [88] Ignas Budvytis, Vijay Badrinarayanan, and Roberto Cipolla. “Label propagation in complex video sequences using semi-supervised learning”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2010.
- [89] M. Buehler, K. Iagnemma, and S. Singh. *The 2005 darpa grand challenge: The great robot race*. Vol. 36. Springer, 2007.
- [90] Martin Buehler, Karl Iagnemma, and Sanjiv Singh. “The DARPA Urban Challenge”. In: *DARPA Challenge*. Advanced Robotics 56 (2009).
- [91] Roland Bunschoten, Ben J. A. Kröse, and Nikos A. Vlassis. “Robust scene reconstruction from an omnidirectional vision system”. In: *IEEE Trans. on Robotics and Automation (TRA)* 19.2 (2003), pp. 351–357.
- [92] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. “A naturalistic open source movie for optical flow evaluation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [93] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. “nuScenes: A multimodal dataset for autonomous driving”. In: *arXiv.org* (2019).
- [94] Zhaowei Cai, Quanfu Fan, Rogério Schmidt Feris, and Nuno Vasconcelos. “A Unified Multi-scale Deep Convolutional Neural Network for Fast Object Detection”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [95] Federico Camposeco, Andrea Cohen, Marc Pollefeys, and Torsten Sattler. “Hybrid scene Compression for Visual Localization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [96] Nicholas Carlevaris-Bianco, Arash K. Ushani, and Ryan M. Eustice. “University of Michigan North Campus long-term vision and lidar dataset”. In: *International Journal of Robotics Research (IJRR)* 35.9 (2016), pp. 1023–1035.

- [97] João Carreira, Rui Caseiro, Jorge P. Batista, and Cristian Sminchisescu. “Semantic Segmentation with Second-Order Pooling”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012, pp. 430–443.
- [98] João Carreira and Cristian Sminchisescu. “CPMC: Automatic Object Segmentation Using Constrained Parametric Min-Cuts”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 34.7 (2012), pp. 1312–1328.
- [99] Jan Cech, Jordi Sanchez-Riera, and Radu P. Horaud. “Scene Flow Estimation by Growing Correspondence Seeds”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.
- [100] GM Heritage Center. *Self-Driving Cars, in 1956?* https://www.gmheritagecenter.com/featured/Autonomous_Vehicles.html. Online: accessed 18-October-2019. 2017.
- [101] Dan Cernea. *OpenMVS: Open Multiple View Stereovision*. <http://cdcseacave.github.io/openMVS>. Online: accessed 23-April-2019. 2015.
- [102] Florian Chabot, Mohamed Chaouch, Jaonary Rabarisoa, Céline Teulière, and Thierry Chateau. “Deep MANTA: A Coarse-to-fine Many-Task Network for joint 2D and 3D vehicle analysis from monocular image”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 2040–2049.
- [103] Jia-Ren Chang and Yong-Sheng Chen. “Pyramid Stereo Matching Network”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 5410–5418.
- [104] Chenyi Chen, Ari Seff, Alain L. Kornhauser, and Jianxiong Xiao. “DeepDriving: Learning Affordance for Direct Perception in Autonomous Driving”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015, pp. 2722–2730.
- [105] Dian Chen, Brady Zhou, Vladlen Koltun, and Philipp Krähenbühl. “Learning by Cheating”. In: *Proc. Conf. on Robot Learning (CoRL)*. 2019.
- [106] Liang-Chieh Chen, Maxwell D. Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jonathon Shlens. “Searching for Efficient Multi-Scale Architectures for Dense Image Prediction”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2018, pp. 8713–8724.

- [107] Liang-Chieh Chen, Sanja Fidler, Alan L. Yuille, and Raquel Urtasun. “Beat the MTurkers: Automatic Image Labeling from Weak 3D Supervision”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [108] Liang-Chieh Chen, Alexander Hermans, George Papandreou, Florian Schroff, Peng Wang, and Hartwig Adam. “MaskLab: Instance Segmentation by Refining Object Detection With Semantic and Direction Features”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 4013–4022.
- [109] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. “DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 40.4 (2018), pp. 834–848.
- [110] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. “Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs”. In: *Proc. of the International Conf. on Learning Representations (ICLR)*. 2015.
- [111] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. “Rethinking Atrous Convolution for Semantic Image Segmentation”. In: *arXiv.org abs/1706.05587* (2017).
- [112] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. “Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 833–851.
- [113] Qifeng Chen and Vladlen Koltun. “Full Flow: Optical Flow Estimation By Global Optimization over Regular Grids”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [114] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew G Berneshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. “3D Object Proposals for Accurate Object Class Detection”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2015.
- [115] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Huimin Ma, Sanja Fidler, and Raquel Urtasun. “3D Object Proposals Using Stereo Imagery for Accurate Object Class Detection”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 40.5 (2018), pp. 1259–1272.

- [116] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. “Multi-View 3D Object Detection Network for Autonomous Driving”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [117] Ian Cherabier, Christian Häne, Martin R. Oswald, and Marc Pollefeys. “Multi-Label Semantic 3D Reconstruction Using Voxel Blocks”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2016.
- [118] Ian Cherabier, Johannes Schönberger, Martin Oswald, Marc Pollefeys, and Andreas Geiger. “Learning Priors for Semantic 3D Reconstruction”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [119] W. Choi, C. Pantofaru, and S. Savarese. “A General Framework for Tracking Multiple People from a Moving Camera”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 35.7 (2013), pp. 1577–1591.
- [120] Wongun Choi. “Near-Online Multi-target Tracking with Aggregated Local Flow Descriptor”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [121] Wongun Choi and Silvio Savarese. “A Unified Framework for Multi-Target Tracking and Collective Activity Recognition”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [122] Francois Chollet. “Xception: Deep Learning with Depthwise Separable Convolutions”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 1800–1807.
- [123] Peng Chu, Heng Fan, Chiu C. Tan, and Haibin Ling. “Online Multi-Object Tracking With Instance-Aware Tracker and Dynamic Model Refreshment”. In: *Proc. of the IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2019, pp. 161–170.
- [124] Peng Chu and Haibin Ling. “FAMNet: Joint Learning of Feature, Affinity and Multi-dimensional Assignment for Online Multiple Object Tracking”. In: *arXiv.org abs/1904.04989* (2019).
- [125] Qi Chu, Wanli Ouyang, Hongsheng Li, Xiaogang Wang, Bin Liu, and Nenghai Yu. “Online Multi-object Tracking Using CNN-Based Single Object Tracker with Spatial-Temporal Attention Mechanism”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 4846–4855.
- [126] Dan C. Ciresan, Ueli Meier, Jonathan Masci, and Jürgen Schmidhuber. “Multi-column deep neural network for traffic sign classification”. In: *Neural Networks* 32 (2012), pp. 333–338.

- [127] Dan C Cirezan, Ueli Meier, Jonathan Masci, and Jürgen Schmidhuber. “A committee of neural networks for traffic sign classification”. In: *International Joint Conference on Neural Networks (IJCNN)*. 2011, pp. 1918–1921.
- [128] Felipe Codevilla, Antonio M. Lopez, Vladlen Koltun, and Alexey Dosovitskiy. “On Offline Evaluation of Vision-based Driving Models”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [129] Felipe Codevilla, Matthias Miiller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. “End-to-End Driving Via Conditional Imitation Learning”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2018, pp. 1–9.
- [130] Felipe Codevilla, Eder Santana, Antonio M. López, and Adrien Gaidon. “Exploring the Limitations of Behavior Cloning for Autonomous Driving”. In: *arXiv.org abs/1904.08980* (2019).
- [131] Robert T. Collins. “A Space-Sweep Approach to True Multi-Image Matching”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 1996, pp. 358–363.
- [132] Timothy F Cootes, Christopher J Taylor, David H Cooper, and Jim Graham. “Active shape models-their training and application”. In: *Computer Vision and Image Understanding (CVIU)* 61.1 (1995), pp. 38–59.
- [133] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. “The Cityscapes Dataset for Semantic Urban Scene Understanding”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [134] Marius Cordts, Lukas Schneider, Markus Enzweiler, Uwe Franke, and Stefan Roth. “Object-Level Priors for Stixel Generation”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2014.
- [135] N. Cornelis, B. Leibe, K. Cornelis, and L. J. Van Gool. “3D Urban Scene Modeling Integrating Recognition and Reconstruction”. In: *International Journal of Computer Vision (IJCV)* 78.2-3 (July 2008), pp. 121–141.
- [136] Mitsubishi Motors Corporation. *Mitsubishi Motors Develops ‘New Driver Support System’*. <https://www.mitsubishi-motors.com/en/corporate/pressrelease/corporate/detail429.html>. Online: accessed 17-May-2019. 1998.

- [137] Arthur Daniel Costea, Andra Petrovai, and Sergiu Nedevschi. “Fusion Scheme for Semantic and Instance-level Segmentation”. In: *Proc. IEEE Conf. on Intelligent Transportation Systems (ITSC)*. 2018, pp. 3469–3475.
- [138] Arthur Daniel Costea, Robert Varga, and Sergiu Nedevschi. “Fast Boosting Based Detection Using Scale Invariant Multimodal Multiresolution Filtered Features”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 993–1002.
- [139] David J. Crandall, Andrew Owens, Noah Snavely, and Dan Huttenlocher. “Discrete-continuous optimization for large-scale structure from motion”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011, pp. 3001–3008.
- [140] Mark Cummins and Paul Newman. “FAB-MAP: Probabilistic Localization and Mapping in the Space of Appearance”. In: *International Journal of Robotics Research (IJRR)* 27.6 (2008), pp. 647–665.
- [141] Brian Curless and Marc Levoy. “A Volumetric Method for Building Complex Models from Range Images”. In: *ACM Trans. on Graphics*. 1996.
- [142] Igor Cvišić, Josip Cescic, Ivan Markovic, and Ivan Petrovic. “Soft-slam: Computationally efficient stereo visual slam for autonomous uavs”. In: *Journal of Field Robotics (JFR)* (2017).
- [143] Igor Cvisic and Ivan Petrovic. “Stereo odometry based on careful feature selection and tracking”. In: *Proc. European Conf. on Mobile Robotics (ECMR)*. 2015.
- [144] Angela Dai and Matthias Nießner. “3DMV: Joint 3D-Multi-view Prediction for 3D Semantic Scene Segmentation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 458–474.
- [145] Jifeng Dai, Kaiming He, Yi Li, Shaoqing Ren, and Jian Sun. “Instance-Sensitive Fully Convolutional Networks”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016, pp. 534–549.
- [146] Jifeng Dai, Kaiming He, and Jian Sun. “Convolutional feature masking for joint object and stuff segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 3992–4000.
- [147] Jifeng Dai, Kaiming He, and Jian Sun. “Instance-Aware Semantic Segmentation via Multi-Task Network Cascades”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [148] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. “R-FCN: Object Detection via Region-based Fully Convolutional Networks”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2016.

- [149] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. “Deformable Convolutional Networks”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)* 1703.06211 (2017).
- [150] N. Dalal and B. Triggs. “Histograms of oriented gradients for human detection”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2005.
- [151] A. Dame, V.A. Prisacariu, C.Y. Ren, and I. Reid. “Dense Reconstruction Using 3D Object Shape Priors”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [152] DARPA. *The DARPA Grand Challenge: Ten Years Later*. <https://www.darpa.mil/news-events/2014-03-13>. Online: accessed 18-June-2019. 2014.
- [153] Afshin Dehghan, Shayan Modiri Assari, and Mubarak Shah. “GMMCP Tracker: Globally Optimal Generalized Maximum Multi Clique Problem for Multiple Object Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [154] Afshin Dehghan, Yicong Tian, Philip H. S. Torr, and Mubarak Shah. “Target Identity-Aware Network Flow for Online Multiple Target Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [155] Joerg Deigmoeller and Julian Eggert. “Stereo Visual Odometry Without Temporal Filtering”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2016.
- [156] Amaël Delaunoy and Marc Pollefeys. “Photometric Bundle Adjustment for Dense Multi-view 3D Modeling”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [157] Amaël Delaunoy and Emmanuel Prados. “Gradient Flows for Optimizing Triangular Mesh-based Surfaces: Applications to 3D Reconstruction Problems Dealing with Visibility”. In: *International Journal of Computer Vision (IJCV)* 95.2 (2011), pp. 100–123.
- [158] Frank Dellaert, Dieter Fox, Wolfram Burgard, and Sebastian Thrun. “Monte Carlo Localization for Mobile Robots”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 1999.
- [159] Frank Dellaert and Michael Kaess. “Square Root SAM: Simultaneous Localization and Mapping via Square Root Information Smoothing”. In: *International Journal of Robotics Research (IJRR)* (2006).
- [160] Jia Deng, Wei Dong, Richard Socher, Li-jia Li, Kai Li, and Li Fei-fei. “Imagenet: A large-scale hierarchical image database”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.

- [161] Chaitanya Desai, Deva Ramanan, and Charless C. Fowlkes. “Discriminative Models for Multi-Class Object Layout”. In: *International Journal of Computer Vision (IJCV)* 95.1 (2011), pp. 1–12.
- [162] Jean-Emmanuel Deschaud. “IMLS-SLAM: Scan-to-Model Matching Based on 3D Data”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2018, pp. 2480–2485.
- [163] Ayush Dewan, Tim Caselitz, Gian Diego Tipaldi, and Wolfram Burgard. “Rigid scene flow for 3D LiDAR scans”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2016.
- [164] E. D. Dickmanns, R. Behringer, D. Dickmanns, T. Hildebrandt, M. Maurer, F. Thomanek, and J. Schiehlen. “The seeing passenger car ‘VaMoRs-P’”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 1994.
- [165] E. D. Dickmanns and B. D. Mysliwetz. “Recursive 3-D road and relative ego-state recognition”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 14.2 (Feb. 1992), pp. 199–213.
- [166] Ernst D. Dickmanns. *Dynamic Machine Vision*. <http://dyna-vision.de/>. Online: accessed 18-June-2019. 1995.
- [167] Ernst D. Dickmanns and Volker Graefe. “Dynamic monocular machine vision”. In: *Machine Vision and Applications (MVA)* 1.4 (1988), pp. 223–240.
- [168] Ernst D. Dickmanns, Birger D. Mysliwetz, and Thomas Christians. “An integrated spatio-temporal approach to automatic visual guidance of autonomous vehicles”. In: *IEEE Trans. on Systems, Man and Cybernetics (TSMC)* 20.6 (1990), pp. 1273–1284.
- [169] P. Dollar, C. Wojek, B. Schiele, and P. Perona. “Pedestrian Detection: An Evaluation of the State of the Art”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)*. Vol. 99. 2011.
- [170] Piotr Dollár, Ron Appel, Serge J. Belongie, and Pietro Perona. “Fast Feature Pyramids for Object Detection”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 36.8 (2014), pp. 1532–1545.
- [171] Piotr Dollár, Zhuowen Tu, Pietro Perona, and Serge Belongie. “Integral channel features”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. BMVC Press, 2009.
- [172] Piotr Dollar, Christian Wojek, Bernt Schiele, and Pietro Perona. “Pedestrian Detection: A Benchmark”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.

- [173] Piotr Dollár, Christian Wojek, Bernt Schiele, and Pietro Perona. “Pedestrian Detection: An Evaluation of the State of the Art”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 34.4 (2012), pp. 743–761.
- [174] A. Dosovitskiy, P. Fischer, E. Ilg, P. Haeusser, C. Hazirbas, V. Golkov, P. v.d. Smagt, D. Cremers, and T. Brox. “FlowNet: Learning Optical Flow with Convolutional Networks”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [175] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. “CARLA: An Open Urban Driving Simulator”. In: *Proc. Conf. on Robot Learning (CoRL)*. 2017.
- [176] Amnon Drory, Carsten Haubold, Shai Avidan, and Fred A. Hamprecht. “Semi-Global Matching: A Principled Derivation in Terms of Message Passing”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2014.
- [177] Xinxin Du, Marcelo H. Ang, Sertac Karaman, and Daniela Rus. “A General Pipeline for 3D Detection of Vehicles”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2018, pp. 3194–3200.
- [178] Liuyun Duan and Florent Lafarge. “Towards Large-Scale City Reconstruction from Satellites”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [179] Renaud Dubé, Daniel Dugas, Elena Stumm, Juan Nieto, Roland Siegwart, and Cesar Cadena. “SegMatch: Segment based place recognition in 3D point clouds”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. IEEE. 2017, pp. 5266–5272.
- [180] J. Engel, J. Sturm, and D. Cremers. “Semi-Dense Visual Odometry for a Monocular Camera”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2013.
- [181] Jakob Engel, Vladlen Koltun, and Daniel Cremers. “Direct Sparse Odometry”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 40 (2018), pp. 611–625.
- [182] Jakob Engel, Thomas Schöps, and Daniel Cremers. “LSD-SLAM: Large-Scale Direct Monocular SLAM”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [183] Jakob Engel, Jörg Stückler, and Daniel Cremers. “Large-scale direct SLAM with stereo cameras”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2015.

- [184] Martin Engelcke, Dushyant Rao, Dominic Zeng Wang, Chi Hay Tong, and Ingmar Posner. “Vote3Deep: Fast Object Detection in 3D Point Clouds Using Efficient Convolutional Neural Networks”. In: *arXiv.org* 609.06666 (2016).
- [185] Martin Engelcke, Dushyant Rao, Dominic Zeng Wang, Chi Hay Tong, and Ingmar Posner. “Vote3Deep: Fast object detection in 3D point clouds using efficient convolutional neural networks”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2017, pp. 1355–1361.
- [186] M. Enzweiler and D. M. Gavrilu. “Monocular Pedestrian Detection: Survey and Experiments”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 31 (2009), pp. 2179–2195.
- [187] M. Enzweiler and D.M. Gavrilu. “A mixed generative-discriminative framework for pedestrian classification”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [188] Markus Enzweiler and Dariu M. Gavrilu. “A Multilevel Mixture-of-Experts Framework for Pedestrian Classification”. In: *IEEE Trans. on Image Processing (TIP)* 20.10 (2011), pp. 2967–2979.
- [189] Friedrich Erbs, Beate Schwarz, and Uwe Franke. “From stixels to objects - A conditional random field based approach”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2013.
- [190] Friedrich Erbs, Beate Schwarz, and Uwe Franke. “Stixmentation - Probabilistic Stixel based Traffic Scene Labeling”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2012.
- [191] A. Ess, B. Leibe, K. Schindler, and L. Van Gool. “Robust multi-person tracking from a mobile platform”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 31 (2009), pp. 1831–1846.
- [192] A. Ess, B. Leibe, K. Schindler, and L. Van Gool. “A Mobile Vision System for Robust Multi-Person Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [193] A. Ess, T. Mueller, H. Grabner, and L. van Gool. “Segmentation-Based Urban Traffic Scene Understanding”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2009.
- [194] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. “The Pascal Visual Object Classes (VOC) Challenge”. In: *International Journal of Computer Vision (IJCV)* 88.2 (2010), pp. 303–338.

- [195] Nolang Fanani, Matthias Ochs, Henry Bradler, and Rudolf Mester. “Keypoint trajectory estimation using propagation based tracking”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2016, pp. 933–939.
- [196] Nolang Fanani, Alina Sturck, Marc Barnada, and Rudolf Mester. “Multimodal scale estimation for monocular visual odometry”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2017, pp. 1714–1721.
- [197] Nolang Fanani, Alina Sturck, Matthias Ochs, Henry Bradler, and Rudolf Mester. “Predictive monocular odometry (PMO): What is possible without RANSAC and multiframe bundle adjustment?”. In: *Image and Vision Computing (IVC)* 68 (2017), pp. 3–13.
- [198] Gunnar Farneback. “Two-Frame Motion Estimation Based on Polynomial Expansion”. In: *Scandinavian Conference on Image Analysis (SCIA)*. 2003.
- [199] Olivier D. Faugeras and Renaud Keriven. “Variational principles, surface evolution, PDEs, level set methods, and the stereo problem”. In: *IEEE Trans. on Image Processing (TIP)* 7.3 (1998), pp. 336–344.
- [200] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. “Detect to Track and Track to Detect”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [201] Pedro F. Felzenszwalb, David A. McAllester, and Deva Ramanan. “A discriminatively trained, multiscale, deformable part model”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [202] Pedro Felzenszwalb and Daniel Huttenlocher. “Efficient Belief Propagation for Early Vision”. In: *International Journal of Computer Vision (IJCV)* 70.1 (Oct. 2006), pp. 41–54.
- [203] J. Ferryman and A. Shahrokni. “PETS2009: Dataset and challenge”. In: *Performance Evaluation of Tracking and Surveillance*. 2009, pp. 1–6.
- [204] G. Floros and B. Leibe. “Joint 2D-3D temporally consistent semantic segmentation of street scenes”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [205] Jan-Michael Frahm, Pierre Fite-Georgel, David Gallup, Tim Johnson, Rahul Raguram, Changchang Wu, Yi-Hung Jen, Enrique Dunn, Brian Clipp, Svetlana Lazebnik, and Marc Pollefeys. “Building Rome on a Cloudless Day”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.

- [206] Jan-Michael Frahm, Marc Pollefeys, Svetlana Lazebnik, David Gallup, Brian Clipp, Rahul Raguram, Changchang Wu, Christopher Zach, and Tim Johnson. “Fast robust large-scale mapping from video and internet photo collections”. In: *ISPRS Journal of Photogrammetry and Remote Sensing (JPRS)* 65 (2010), pp. 538–550.
- [207] U. Franke, S. Mehring, A. Suissa, and S. Hahn. “The Daimler-Benz steering assistant: a spin-off from autonomous driving”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 1994.
- [208] Uwe Franke, Dariu Gavrilă, Steffen Görzig, Frank Lindner, Frank Paetzold, and Christian Wöhler. “Autonomous Driving Goes Downtown”. In: *Intelligent Systems (IS)* 13.6 (1998), pp. 40–48.
- [209] Uwe Franke, Clemens Rabe, Hernán Badino, and Stefan Gehrig. “6D-Vision: Fusion of Stereo and Motion for Robust Environment Perception”. In: *Proc. of the DAGM Symposium on Pattern Recognition (DAGM)*. 2005.
- [210] Friedrich Fraundorfer and Davide Scaramuzza. “Visual Odometry: Part II - Matching, Robustness, and Applications.” In: *Robotics and Automation Magazine (RAM)* (2011).
- [211] Jannik Fritsch, Tobias Kuehnl, and Andreas Geiger. “A New Performance Measure and Evaluation Benchmark for Road Detection Algorithms”. In: *Proc. IEEE Conf. on Intelligent Transportation Systems (ITSC)*. 2013.
- [212] Davi Frossard and Raquel Urtasun. “End-to-end Learning of Multi-sensor 3D Tracking by Detection”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2018.
- [213] Duncan P. Frost, Olaf Kähler, and David W. Murray. “Object-aware bundle adjustment for correcting monocular scale drift”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2016.
- [214] Christian Früh, Siddharth Jain, and Avidesh Zakhori. “Data Processing Algorithms for Generating Textured 3D Building Facade Meshes from Laser Scans and Camera Images”. In: *International Journal of Computer Vision (IJCV)* 61.2 (2005), pp. 159–184.
- [215] Simon Fuhrmann, Fabian Langguth, and Michael Goesele. “MVE - A Multi-View Reconstruction Environment”. In: *Eurographics Workshop on Graphics and Cultural Heritage (GCH)*. 2014, pp. 11–18.

- [216] Paul Timothy Furgale, Ulrich Schwesinger, Martin Ruflı, Wojciech Derendarz, Hugo Grimmett, Peter Mühlfellner, Stefan Wonneberger, Julian Timpner, Stephan Rottmann, Bo Li, Bastian Schmidt, Thien-Nghia Nguyen, Elena Cardarelli, Stefano Cattani, Stefan Bruning, Sven Horstmann, Martin Stellmacher, Holger Mielenz, Kevin Köser, Markus Beermann, Christian Hane, Lionel Heng, Gim Hee Lee, Friedrich Fraundorfer, Rene Iser, Rudolph Triebel, Ingmar Posner, Paul Newman, Lars C. Wolf, Marc Pollefeys, Stefan Brosig, Jan Effertz, Cédric Pradalier, and Roland Siegwart. “Toward automated driving in cities using close-to-market sensors: An overview of the V-Charge Project”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2013.
- [217] Yasutaka Furukawa and Jean Ponce. “Accurate, Dense, and Robust Multi-View Stereopsis”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 32.8 (2010), pp. 1362–1376.
- [218] Raghudeep Gadde, Varun Jampani, Martin Kiefel, Daniel Kappler, and Peter V. Gehler. “Superpixel Convolutional Networks Using Bilateral Inceptions”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [219] Raghudeep Gadde, Varun Jampani, Renaud Marlet, and Peter V. Gehler. “Efficient 2D and 3D Facade Segmentation Using Auto-Context”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 40.5 (2018), pp. 1273–1280.
- [220] David Gadot and Lior Wolf. “PatchBatch: a Batch Augmented Loss for Optical Flow”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [221] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. “Virtual Worlds as Proxy for Multi-Object Tracking Analysis”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [222] Guillermo Gallego, Henri Rebecq, and Davide Scaramuzza. “A Unifying Contrast Maximization Framework for Event Cameras, With Applications to Motion, Depth, and Optical Flow Estimation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [223] Carolina Galleguillos, Andrew Rabinovich, and Serge J. Belongie. “Object categorization using co-occurrence, location and appearance”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [224] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. “Massively Parallel Multiview Stereopsis by Surface Normal Diffusion”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.

- [225] D. Gallup, J. M. Frahm, P. Mordohai, and M. Pollefeys. “Variable baseline/resolution stereo”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [226] David Gallup, Jan-Michael Frahm, Philippos Mordohai, Qingxiong Yang, and Marc Pollefeys. “Real-time plane-sweeping stereo with multiple sweeping directions”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2007.
- [227] David Gallup, Jan-Michael Frahm, and Marc Pollefeys. “Piecewise planar and non-planar stereo for urban scene reconstruction”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2010.
- [228] Dorian Gálvez-López and Juan D. Tardós. “Bags of Binary Words for Fast Place Recognition in Image Sequences”. In: *IEEE Trans. on Robotics* 28.5 (2012), pp. 1188–1197.
- [229] Alberto Garcia-Garcia, Francisco Gomez-Donoso, José García Rodríguez, Sergio Orts-Escolano, Miguel Cazorla, and Jorge Azorín López. “PointNet: A 3D Convolutional Neural Network for real-time object class recognition”. In: *International Joint Conference on Neural Networks (IJCNN)*. 2016.
- [230] Álvaro Arcos García, Juan Antonio Álvarez-García, and Luis Miguel Soria-Morillo. “Evaluation of deep neural networks for traffic sign detection systems”. In: *Neurocomputing* 316 (2018), pp. 332–344.
- [231] D. M. Gavrila and S. Munder. “Multi-cue pedestrian detection and tracking from a moving vehicle”. In: *International Journal of Computer Vision (IJCV)* 73 (2007), pp. 41–59.
- [232] Daniel Gehrig, Henri Rebecq, Guillermo Gallego, and Davide Scaramuzza. “Asynchronous, Photometric Feature Tracking Using Events and Frames”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [233] Stefan K. Gehrig, Felix Eberli, and Thomas Meyer. “A Real-Time Low-Power Stereo Vision Engine Using Semi-Global Matching.” In: *Proc. of the International Conf. on Computer Vision Systems (ICVS)*. 2009.
- [234] Andreas Geiger. “Monocular road mosaicing for urban environments”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2009.
- [235] Andreas Geiger, Martin Lauer, Frank Moosmann, Benjamin Ranft, Holger Rapp, Christoph Stiller, and Julius Ziegler. “Team AnnieWAY’s entry to the Grand Cooperative Driving Challenge 2011”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 13.3 (Sept. 2012), pp. 1008–1017.

- [236] Andreas Geiger, Martin Lauer, Christian Wojek, Christoph Stiller, and Raquel Urtasun. “3D Traffic Scene Understanding from Movable Platforms”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 36.5 (2014), pp. 1012–1025.
- [237] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. “Vision meets Robotics: The KITTI Dataset”. In: *International Journal of Robotics Research (IJRR)* 32.11 (2013), pp. 1231–1237.
- [238] Andreas Geiger, Philip Lenz, and Raquel Urtasun. “Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [239] Andreas Geiger, Frank Moosmann, Omer Car, and Bernhard Schuster. “Automatic Calibration of Range and Camera Sensors using a single Shot”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2012.
- [240] Andreas Geiger, Martin Roser, and Raquel Urtasun. “Efficient Large-Scale Stereo Matching”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2010.
- [241] Andreas Geiger, Julius Ziegler, and Christoph Stiller. “StereoScan: Dense 3D Reconstruction in Real-time”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2011.
- [242] David Geronimo, Antonio M. Lopez, Angel D. Sappa, and Thorsten Graf. “Survey on Pedestrian Detection for Advanced Driver Assistance Systems”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 32.7 (2010), pp. 1239–1258.
- [243] Christopher Geyer and Kostas Daniilidis. “A unifying theory for central panoramic systems and practical implications”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2000.
- [244] Golnaz Ghiasi and Charless C. Fowlkes. “Laplacian Pyramid Reconstruction and Refinement for Semantic Segmentation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [245] J. Giebel, D.M. Gavrila, and C. Schnörr. “A Bayesian Framework for Multi-cue 3D Object Tracking”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2004.
- [246] Ross B. Girshick. “Fast R-CNN”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [247] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. “Rich feature hierarchies for accurate object detection and semantic segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.

- [248] Georgia Gkioxari, Bharath Hariharan, Ross Girshick, and Jitendra Malik. “Using k-Poselets for Detecting People and Localizing Their Keypoints”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [249] José-Joel González-Barbosa and Simon Lacroix. “Fast Dense Panoramic Stereovision”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2005, pp. 1210–1215.
- [250] A. González, D. Vázquez, A. M. Lóopez, and J. Amores. “On-Board Object Detection: Multicue, Multimodal, and Multiview Random Forest of Local Experts”. In: *IEEE Trans. on Cybernetics* (2016).
- [251] Johannes Gräter, Alexander Wilczynski, and Martin Lauer. “LIMO: Lidar-Monocular Visual Odometry”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2018, pp. 7872–7879.
- [252] Hugo Grimmett, Mathias Bürki, Lina María Paz, Pedro Pinies, Paul Timothy Furgale, Ingmar Posner, and Paul Newman. “Integrating metric and semantic maps for vision-only automated parking”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2015.
- [253] P. Grisleri and I. Fedriga. “The BRAiVE platform”. In: *Proc. of the IFAC Symposium on Intelligent Autonomous Vehicles (IFAC)*. 2010.
- [254] Carlos Guindel, David Martín, and José María Armingol. “Fast Joint Object Detection and Viewpoint Estimation for Traffic Scene Understanding”. In: *Proc. IEEE Intelligent Transportation Systems Magazine (ITSM)* 10.4 (2018), pp. 74–86.
- [255] Jose Guivant and Eduardo Nebot. *Victoria Park Dataset*. http://www-personal.acfr.usyd.edu.au/nebot/victoria_park.htm. Online: accessed 8-April-2019. 2006.
- [256] Fatma Güney and Andreas Geiger. “Deep Discrete Flow”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2016.
- [257] Fatma Güney and Andreas Geiger. “Displets: Resolving Stereo Ambiguities using Object Knowledge”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [258] Bertan Günyel, Rodrigo Benenson, Radu Timofte, and Luc J. Van Gool. “Stixels Motion Estimation without Optical Flow Computation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.

- [259] Norbert Haala and Karl-Heinrich Anders. “Acquisition of 3D urban models by analysis of aerial images, digital surface models, and existing 2D building information”. In: *Proc. of the SPIE Conf. Integrating Photogrammetric Techniques with Scene Analysis and Machine Vision III*. Vol. 3072. International Society for Optics and Photonics. 1997, pp. 212–222.
- [260] Timo Hackel, Nikolay Savinov, Lubor Ladicky, Jan Dirk Wegner, Konrad Schindler, and Marc Pollefeys. “Semantic3D.net: A new Large-scale Point Cloud Classification Benchmark”. In: *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences (APRS)* (2017).
- [261] Timo Hackel, Jan D. Wegner, and Konrad Schindler. “Fast semantic segmentation of 3d point clouds with strongly varying density”. In: *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences (APRS)* III-3 (2016), pp. 177–184.
- [262] Christian Haene, Nikolay Savinov, and Marc Pollefeys. “Class Specific 3D Object Shape Priors Using Surface Normals”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [263] Christian Haene, Christopher Zach, Andrea Cohen, Roland Angst, and Marc Pollefeys. “Joint 3D Scene Reconstruction and Class Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [264] Christian Haene, Christopher Zach, Bernhard Zeisl, and Marc Pollefeys. “A Patch Prior for Dense 3D Reconstruction in Man-Made Environments”. In: *Proc. of the International Conf. on 3D Digital Imaging, Modeling, Data Processing, Visualization and Transmission (THREE-DIMPVT)*. 2012.
- [265] Lars Hammarstrand, Fredrik Kahl, Will Maddern, Tomas Pajdla, Marc Pollefeys, Torsten Sattler, Josef Sivic, Erik Stenborg, Carl Toft, and Akihiko Torii. *Workshop on Long-Term Visual Localization under Changing Conditions*. <https://sites.google.com/view/ltml2019/>. Online: accessed 17-June-2019. 2019.
- [266] Xufeng Han, Thomas Leung, Yangqing Jia, Rahul Sukthankar, and Alexander C. Berg. “MatchNet: Unifying Feature and Metric Learning for Patch-Based Matching”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. June 2015.
- [267] Christian Häne, Lionel Heng, Gim Hee Lee, Alexey Sizov, and Marc Pollefeys. “Real-Time Direct Dense Matching on Fisheye Images Using Plane-Sweeping Stereo”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2014.

- [268] Christian Häne, Torsten Sattler, and Marc Pollefeys. “Obstacle detection for self-driving cars using only monocular cameras and wheel odometry”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2015.
- [269] Allen Hanson. *Computer vision systems*. Elsevier, 1978.
- [270] Bharath Hariharan, Pablo Andrés Arbeláez, Ross B. Girshick, and Jitendra Malik. “Hypercolumns for object segmentation and fine-grained localization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 447–456.
- [271] Bharath Hariharan, Pablo Andrés Arbeláez, Ross B. Girshick, and Jitendra Malik. “Simultaneous Detection and Segmentation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [272] Bharath Hariharan, Pablo Arbelaez, Lubomir D. Bourdev, Subhransu Maji, and Jitendra Malik. “Semantic contours from inverse detectors”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2011, pp. 991–998.
- [273] R. I. Hartley. “In defence of the 8-point algorithm”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 1995.
- [274] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Second. Cambridge University Press, 2004.
- [275] Wilfried Hartmann, Silvano Galliani, Michal Havlena, Luc Van Gool, and Konrad Schindler. “Learned Multi-Patch Similarity”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [276] Nick Hawes. *Driving the revolution*. <https://www.birmingham.ac.uk/news/thebirminghambrief/items/2016/11/driving-the-revolution.aspx>. Online: accessed 18-October-2019. 2016.
- [277] Zeeshan Hayder, Xuming He, and Mathieu Salzmann. “Boundary-Aware Instance Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 587–595.
- [278] James Hays and Alexei A. Efros. “im2gps: estimating geographic information from a single image”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [279] Caner Hazirbas, Lingni Ma, Csaba Domokos, and Daniel Cremers. “FuseNet: Incorporating Depth into Semantic Segmentation via Fusion-Based CNN Architecture”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2016.
- [280] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. “Mask R-CNN”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 2980–2988.

- [281] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep Residual Learning for Image Recognition”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [282] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Identity Mappings in Deep Residual Networks”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016, pp. 630–645.
- [283] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [284] Xuming He, Richard S. Zemel, and Miguel A. Carreira-Perpinan. “Multiscale Conditional Random Fields for Image Labeling”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2004.
- [285] Xuming He, Richard S. Zemel, and Debajyoti Ray. “Learning and Incorporating Top-Down Cues in Image Segmentation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2006.
- [286] David J. Heeger. “Optical flow using spatiotemporal filters”. In: *International Journal of Computer Vision (IJCV)* 1.4 (1988), pp. 279–302.
- [287] Janne Heikkila and Olli Silven. “A Four-step Camera Calibration Procedure with Implicit Image Correction”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 1997.
- [288] Mikael Henaff, Alfredo Canziani, and Yann LeCun. “Model-Predictive Policy Learning with Uncertainty Regularization for Driving in Dense Traffic”. In: *arXiv.org abs/1901.02705* (2019).
- [289] Lionel Heng, Paul Timothy Furgale, and Marc Pollefeys. “Leveraging Image-based Localization for Infrastructure-based Calibration of a Multi-camera Rig”. In: *Journal of Field Robotics (JFR)* 32.5 (2015), pp. 775–802.
- [290] Lionel Heng, Bo Li, and Marc Pollefeys. “CamOdoCal: Automatic intrinsic and extrinsic calibration of a rig with multiple generic cameras and odometry”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2013.
- [291] Roberto Henschel, Laura Leal-Taixé, Daniel Cremers, and Bodo Rosenhahn. “Fusion of Head and Full-Body Detectors for Multi-Object Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*. 2018, pp. 1428–1437.

- [292] Roberto Henschel, Zhou Y, and Bodo Rosenhahn. “Fusion of Head and Full-Body Detectors for Multi-Object Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*. 2019.
- [293] H. Hirschmüller and D. Scharstein. “Evaluation of Cost Functions for Stereo Matching”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. June 2007, pp. 1–8.
- [294] Heiko Hirschmüller. “Stereo Processing by Semiglobal Matching and Mutual Information”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 30.2 (2008), pp. 328–341.
- [295] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Computation* 9.8 (1997), pp. 1735–1780.
- [296] D. Hoiem, A. Efros, and M. Hebert. “Putting Objects in Perspective”. In: *International Journal of Computer Vision (IJCV)* 80 (2008), pp. 3–15.
- [297] Derek Hoiem, Alexei A. Efros, and Martial Hebert. “Recovering Surface Layout from an Image”. In: *International Journal of Computer Vision (IJCV)* 75.1 (Oct. 2007), pp. 151–172.
- [298] Berthold K. P. Horn and Brian G. Schunck. “Determining Optical Flow”. In: *Artificial Intelligence (AI)* 17.1-3 (1981), pp. 185–203.
- [299] Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipsing, and Christian Igel. “Detection of traffic signs in real-world images: The German traffic sign detection benchmark”. In: *International Joint Conference on Neural Networks (IJCNN)*. 2013, pp. 1–8.
- [300] Hanzhang Hu, Daniel Munoz, J. Andrew Bagnell, and Martial Hebert. “Efficient 3-D scene analysis from streaming data”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2013, pp. 2297–2304.
- [301] Hou-Ning Hu, Yen-Chen Lin, Ming-Yu Liu, Hsien-Tzu Cheng, Yung-Ju Chang, and Min Sun. “Deep 360 Pilot: Learning a Deep Agent for Piloting through 360° Sports Video”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [302] Xiaowei Hu, Xuemiao Xu, Yongjie Xiao, Hao Chen, Shengfeng He, Jing Qin, and Pheng-Ann Heng. “SINet: A Scale-insensitive Convolutional Neural Network for Fast Vehicle Detection”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* (2018).
- [303] C. Huang, B. Wu, and R. Nevatia. “Robust object tracking by hierarchical association of detection responses”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2008.

- [304] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. “Densely Connected Convolutional Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 2261–2269.
- [305] Jing Huang and Suya You. “Point Cloud Labeling using 3D Convolutional Neural Network”. In: *Proc. of the International Conf. on Pattern Recognition (ICPR)*. 2016.
- [306] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. “DeepMVS: Learning Multi-View Stereopsis”. In: *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. 2018, pp. 2821–2830.
- [307] Xinyu Huang, Xinjing Cheng, Qichuan Geng, Binbin Cao, Dingfu Zhou, Peng Wang, Yuanqing Lin, and Ruigang Yang. “The ApolloScape Dataset for Autonomous Driving”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [308] Frédéric Huguet and Frédéric Devernay. “A Variational Method for Scene Flow Estimation from Stereo Sequences”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2007.
- [309] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. “LiteFlowNet: A Lightweight Convolutional Neural Network for Optical Flow Estimation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [310] Junhwa Hur and Stefan Roth. “MirrorFlow: Exploiting Symmetries in Joint Optical Flow and Occlusion Estimation”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [311] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. “FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2017).
- [312] Sergey Ioffe and Christian Szegedy. “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”. In: *Proc. of the International Conf. on Machine learning (ICML)*. 2015.
- [313] Arnold Irschara, Christopher Zach, Jan-Michael Frahm, and Horst Bischof. “From structure-from-motion point clouds to fast location recognition”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009, pp. 2599–2606.

- [314] Varun Jampani, Martin Kiefel, and Peter V. Gehler. “Learning Sparse High Dimensional Filters: Image Filtering, Dense CRFs and Bilateral Neural Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [315] Joel Janai, Fatma Güney, Anurag Ranjan, Michael Black, and Andreas Geiger. “Unsupervised Learning of Multi-Frame Optical Flow with Occlusions”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [316] Joel Janai, Fatma Güney, Jonas Wulff, Michael Black, and Andreas Geiger. “Slow Flow: Exploiting High-Speed Cameras for Accurate and Diverse Optical Flow Reference Data”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [317] Michal Jancosek and Tomáš Pajdla. “Multi-view reconstruction preserving weakly-supported surfaces”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.
- [318] Simon Jégou, Michal Drozdal, David Vázquez, Adriana Romero, and Yoshua Bengio. “The One Hundred Layers Tiramisu: Fully Convolutional DenseNets for Semantic Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*. 2017, pp. 1175–1183.
- [319] Rasmus Ramsbøl Jensen, Anders Lindbjerg Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. “Large Scale Multi-view Stereopsis Evaluation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [320] Kaijin Ji, Huiyan Chen, Huijun Di, Jianwei Gong, Guangming Xiong, Jianyong Qi, and Tao Yi. “CPFG-SLAM: a Robust Simultaneous Localization and Mapping based on LIDAR in Off-Road Environment”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2018, pp. 650–655.
- [321] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. “SurfaceNet: An End-to-end 3D Neural Network for Multiview Stereopsis”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [322] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross B. Girshick, Sergio Guadarrama, and Trevor Darrell. “Caffe: Convolutional Architecture for Fast Feature Embedding”. In: *Proc. of the International Conf. on Multimedia (ICM)*. 2014.
- [323] H. Jiang, S. Fels, and J. J. Little. “A linear programming approach for multiple object tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2007.

- [324] Junqi Jin, Kun Fu, and Changshui Zhang. “Traffic sign recognition with hinge loss trained convolutional neural networks”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 15.5 (2014), pp. 1991–2000.
- [325] S. J. Julier and J. K. Uhlmann. “A counter example to the theory of simultaneous localization and map building”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2001.
- [326] Sang-Il Jung and Ki-Sang Hong. “Deep network aided by guiding network for pedestrian detection”. In: *Pattern Recognition Letters* 90 (2017), pp. 43–49.
- [327] Michael Kaess, Hordur Johannsson, Richard Roberts, Viorela Ila, John J. Leonard, and Frank Dellaert. “iSAM2: Incremental Smoothing and Mapping Using the Bayes Tree”. In: *International Journal of Robotics Research (IJRR)* 31 (2 2012), pp. 217–236.
- [328] Michael Kaess, Ananth Ranganathan, and Frank Dellaert. “iSAM: Incremental Smoothing and Mapping”. In: *IEEE Trans. on Robotics* 24.6 (2008), pp. 1365–1378.
- [329] R. S. Kaminsky, Noah Snavely, Steven M. Seitz, and Richard Szeliski. “Alignment of 3D point clouds to overhead images”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*. 2009.
- [330] Kai Kang, Hongsheng Li, Tong Xiao, Wanli Ouyang, Junjie Yan, Xihui Liu, and Xiaogang Wang. “Object Detection in Videos with Tubelet Proposal Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [331] Kai Kang, Wanli Ouyang, Hongsheng Li, and Xiaogang Wang. “Object Detection from Video Tubelets with Convolutional Neural Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [332] Yue Kang, Hang Yin, and Christian Berger. “Test Your Self-Driving Algorithm: An Overview of Publicly Available Driving Datasets and Virtual Testing Environments”. In: *Proc. IEEE Transactions on Intelligent Vehicles (T-IV)* 4.2 (2019), pp. 171–185.
- [333] Abhishek Kar, Christian Häne, and Jitendra Malik. “Learning a Multi-View Stereo Machine”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2017.
- [334] A. Kassir and T. Peynot. “Reliable automatic camera-laser calibration”. In: *Proc. IEEE Australasian Conf. on Robotics and Automation (ACRA)*. 2010.

- [335] Christoph Gustav Keller, Markus Enzweiler, Marcus Rohrbach, David Fernández Llorca, Christoph Schnörr, and Darius M. Gavrila. “The Benefits of Dense Stereo for Pedestrian Detection”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 12.4 (2011), pp. 1096–1106.
- [336] Alex Kendall, Yarin Gal, and Roberto Cipolla. “Multi-Task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 7482–7491.
- [337] Alex Kendall, Matthew Grimes, and Roberto Cipolla. “PoseNet: A Convolutional Network for Real-Time 6-DOF Camera Relocalization”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [338] Alex Kendall, Jeffrey Hawke, David Janz, Przemyslaw Mazur, Daniele Reda, John Mark Allen, Vinh Dieu Lam, Alex Bewley, and Amar Shah. “Learning to Drive in a Day”. In: *arXiv.org abs/1807.00412* (2018).
- [339] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, and Peter Henry. “End-to-End Learning of Geometry and Context for Deep Stereo Regression”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [340] Christian Kerl, Jürgen Sturm, and Daniel Cremers. “Robust odometry estimation for RGB-D cameras”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2013, pp. 3748–3754.
- [341] Margret Keuper, Siyu Tang, Bjoern Andres, Thomas Brox, and Bernt Schiele. “Motion Segmentation & Multiple Object Tracking by Correlation Co-Clustering”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* (2018).
- [342] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg. “Multiple Hypothesis Tracking Revisited”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [343] Chanh Kim, Fuxin Li, and James M. Rehg. “Multi-Object Tracking with Neural Gating Using Bilinear LSTM”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [344] Alexander Kirillov, Ross B. Girshick, Kaiming He, and Piotr Dollár. “Panoptic Feature Pyramid Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2019).
- [345] Alexander Kirillov, Kaiming He, Ross B. Girshick, Carsten Rother, and Piotr Dollár. “Panoptic Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019, pp. 9404–9413.

- [346] Alexander Kirillov, Evgeny Levinkov, Bjoern Andres, Bogdan Savchynskyy, and Carsten Rother. “InstanceCut: From Edges to Instances with MultiCut”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 7322–7331.
- [347] Bernd Kitt, Andreas Geiger, and Henning Lategahn. “Visual Odometry based on Stereo Image Sequences with RANSAC-based Outlier Rejection Scheme”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2010.
- [348] Georg Klein and David W. Murray. “Parallel Tracking and Mapping for Small AR Workspaces”. In: *Proc. of the International Symposium on Mixed and Augmented Reality (ISMAR)*. 2007, pp. 225–234.
- [349] Reinhard Klette. *Vision-based Driver Assistance Systems*. Tech. rep. CITR, Auckland, New Zealand, 2015.
- [350] Bryan Matthew Klingner, David Martin, and James Roseborough. “Street View Motion-from-Structure-from-Motion”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2013, pp. 953–960.
- [351] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. “Tanks and Temples: Benchmarking Large-Scale Scene Reconstruction”. In: *ACM Trans. on Graphics* 36.4 (2017).
- [352] Pushmeet Kohli, Lubor Ladicky, and Philip H. S. Torr. “Robust Higher Order Potentials for Enforcing Label Consistency”. In: *International Journal of Computer Vision (IJCV)* 82.3 (2009), pp. 302–324.
- [353] Vladimir Kolmogorov. “Convergent Tree-Reweighted Message Passing for Energy Minimization”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 28.10 (2006), pp. 1568–1583.
- [354] Vijay R. Konda and John N. Tsitsiklis. “Actor-Critic Algorithms”. In: *Advances in Neural Information Processing Systems (NIPS)*. 1999, pp. 1008–1014.
- [355] Claudia Kondermann, Daniel Kondermann, Bernd Jähne, and Christoph S. Garbe. “An Adaptive Confidence Measure for Optical Flows Based on Linear Subspace Projections”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2007, pp. 132–141.
- [356] Claudia Kondermann, Rudolf Mester, and Christoph S. Garbe. “A Statistical Confidence Measure for Optical Flows”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2008.

- [357] Daniel Kondermann, Rahul Nair, Katrin Honauer, Karsten Krispin, Jonas Andrusis, Alexander Brock, Burkhard Gussfeld, Mohsen Rahimimoghaddam, Sabine Hofmann, Claus Brenner, and Bernd Jahne. “The HCI Benchmark Suite: Stereo and Flow Ground Truth With Uncertainties for Urban Autonomous Driving”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*. 2016.
- [358] Philipp Krähenbühl and Vladlen Koltun. “Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2011.
- [359] Ivan Krešo and Siniša Šegvić. “Improving the Egomotion Estimation by Correcting the Calibration Bias”. In: *Proc. of the Conf. on Computer Vision Theory and Applications (VISAPP)*. 2015.
- [360] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2012.
- [361] Till Kroeger, Radu Timofte, Dengxin Dai, and Luc Van Gool. “Fast Optical Flow Using Dense Inverse Search”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016, pp. 471–488.
- [362] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L. Waslander. “Joint 3D Proposal Generation and Object Detection from View Aggregation”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2018, pp. 1–8.
- [363] Jason Ku, Alex D. Pon, and Steven L. Waslander. “Monocular 3D Object Detection Leveraging Accurate Proposals and Shape Reconstruction”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2019).
- [364] T. Kuehnl, F. Kummert, and J. Fritsch. “Spatial Ray Features for Real-Time Ego-Lane Extraction”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2012.
- [365] Daniel Kuettel, Michael D. Breitenstein, Luc Van Gool, and Vittorio Ferrari. “What’s going on?: Discovering Spatio-Temporal Dependencies in Dynamic Scenes”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2010.
- [366] Andreas Kuhn, Heiko Hirschmüller, Daniel Scharstein, and Helmut Mayer. “A TV Prior for High-Quality Scalable Multi-View Stereo Reconstruction”. In: *International Journal of Computer Vision (IJCV)* 124.1 (2017), pp. 2–17.
- [367] Sanjiv Kumar and Martial Hebert. “A Hierarchical Field Framework for Unified Context-Based Classification”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2005.

- [368] Rainer Kümmerle, Giorgio Grisetti, Hauke Strasdat, Kurt Konolige, and Wolfram Burgard. “G²o: A general framework for graph optimization”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2011, pp. 3607–3613.
- [369] Abhijit Kundu, Yin Li, Frank Dellaert, Fuxin Li, and James M. Rehg. “Joint Semantic Segmentation and 3D Reconstruction from Monocular Video”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [370] Abhijit Kundu, Yin Li, and James M. Rehg. “3D-RCNN: Instance-Level 3D Object Reconstruction via Render-and-Compare”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 3559–3568.
- [371] Abhijit Kundu, Vibhav Vineet, and Vladlen Koltun. “Feature Space Optimization for Semantic Video Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [372] Georg Kuschke and Daniel Cremers. “Fast and Accurate Large-scale Stereo Reconstruction using Variational Methods”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV) Workshops*. 2013.
- [373] Kiriakos N. Kutulakos and Steven M. Seitz. “A Theory of Shape by Space Carving”. In: *International Journal of Computer Vision (IJCV)* 38.3 (2000), pp. 199–218.
- [374] Jan Kybic and Claudia Nieuwenhuis. “Bootstrap optical flow confidence and uncertainty measure”. In: *Computer Vision and Image Understanding (CVIU)* 115.10 (2011), pp. 1449–1462.
- [375] Patrick Labatut, Jean-Philippe Pons, and Renaud Keriven. “Efficient Multi-View Reconstruction of Large-Scale Scenes using Interest Points, Delaunay Triangulation and Graph Cuts”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2007, pp. 1–8.
- [376] Ankit Laddha, Mehmet Kemal Kocamaz, Luis E. Navarro-Serment, and Martial Hebert. “Map-Supervised Road Detection”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2016.
- [377] L. Ladicky, C. Russell, P. Kohli, and P. Torr. “Associative hierarchical CRFs for object class image segmentation”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2009.
- [378] Lubor Ladicky, Chris Russell, Pushmeet Kohli, and Philip H.S. Torr. “Graph cut based inference with co-occurrence statistics”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.

- [379] Lubor Ladicky, Christopher Russell, Pushmeet Kohli, and Philip H. S. Torr. “Associative Hierarchical Random Fields”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 36.6 (2014), pp. 1056–1077.
- [380] Florent Lafarge, Xavier Descombes, Josiane Zerubia, and Marc Pierrot Deseilligny. “Structural Approach for Building Reconstruction from a Single DSM”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 32.1 (2010), pp. 135–147.
- [381] Florent Lafarge, Renaud Keriven, Mathieu Bredif, and Hoang-Hiep Vu. “A Hybrid Multiview Stereo Algorithm for Modeling Urban Scenes.” In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 35.1 (2013), pp. 5–17.
- [382] Florent Lafarge and Clément Mallet. “Creating Large-Scale City Models from 3D-Point Clouds: A Robust Approach with Hybrid Representation”. In: *International Journal of Computer Vision (IJCV)* 99.1 (2012), pp. 69–85.
- [383] Michael Laskey, Jonathan Lee, Roy Fox, Anca D. Dragan, and Ken Goldberg. “DART: Noise Injection for Robust Imitation Learning”. In: *Proc. Conf. on Robot Learning (CoRL)*. 2017, pp. 143–156.
- [384] Henning Lategahn, Andreas Geiger, and Bernd Kitt. “Visual SLAM for Autonomous Ground Vehicles”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2011.
- [385] M. Lauer. “Grand Cooperative Driving Challenge 2011”. In: *Proc. IEEE Intelligent Transportation Systems Magazine (ITSM)* 3.3 (2011), pp. 38–40.
- [386] Hei Law and Jia Deng. “CornerNet: Detecting Objects as Paired Key-points”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 765–781.
- [387] Nam Le, Alexander Heili, and Jean-Marc Odobez. “Long-Term Time-Sensitive Costs for CRF-Based Tracking by Detection”. In: *Proc. of the European Conf. on Computer Vision (ECCV) Workshops*. 2016.
- [388] Laura Leal-Taixé, Cristian Canton-Ferrer, and Konrad Schindler. “Learning by Tracking: Siamese CNN for Robust Target Association”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops*. 2016.
- [389] Laura Leal-Taixé, Anton Milan, Ian D. Reid, Stefan Roth, and Konrad Schindler. “MOTChallenge 2015: Towards a Benchmark for Multi-Target Tracking”. In: *arXiv.org* (2015).

- [390] Laura Leal-Taixé, Anton Milan, Ian D. Reid, Stefan Roth, and Konrad Schindler. “MOTChallenge 2015: Towards a Benchmark for Multi-Target Tracking”. In: *arXiv.org* 1504.01942 (2015).
- [391] Gim Hee Lee, Friedrich Fraundorfer, and Marc Pollefeys. “Motion Estimation for Self-Driving Cars with a Generalized Camera”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [392] Gim Hee Lee, Friedrich Fraundorfer, and Marc Pollefeys. “Structureless pose-graph loop-closure with a multi-camera system on a self-driving car”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2013.
- [393] Gim Hee Lee, Marc Pollefeys, and Friedrich Fraundorfer. “Relative Pose Estimation for a Multi-camera System with Known Vertical Direction”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [394] Seokju Lee, Jun-Sik Kim, Jae Shin Yoon, Seunghak Shin, Oleksandr Bailo, Namil Kim, Tae-Hee Lee, Hyun Seok Hong, Seung-Hoon Han, and In So Kweon. “VPGNet: Vanishing Point Guided Network for Lane and Road Marking Detection and Recognition”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 1965–1973.
- [395] B. Leibe, N. Cornelis, K. Cornelis, and L. Van Gool. “Dynamic 3D Scene Analysis from a Moving Vehicle”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2007.
- [396] B. Leibe, A. Leonardis, and B. Schiele. “Robust Object Detection with Interleaved Categorization and Segmentation”. In: *International Journal of Computer Vision (IJCV)* 77.1-3 (2008), pp. 259–289.
- [397] B. Leibe, K. Schindler, N. Cornelis, and L. Van Gool. “Coupled Detection and Tracking from Static Cameras and Moving Vehicles”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 30.10 (2008), pp. 1683–1698.
- [398] Kruno Lenac, Josip Cesic, Ivan Markovic, and Ivan Petrovic. “Exactly sparse delayed state filter on Lie groups for long-term pose graph SLAM”. In: *International Journal of Robotics Research (IJRR)* 37.6 (2018), pp. 585–610.
- [399] Philip Lenz, Andreas Geiger, and Raquel Urtasun. “FollowMe: Efficient Online Min-Cost Flow Tracking with Bounded Memory and Computation”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.

- [400] Stefan Leutenegger, Paul Timothy Furgale, Vincent Rabaud, Margarita Chli, Kurt Konolige, and Roland Siegwart. “Keyframe-Based Visual-Inertial SLAM using Nonlinear Optimization”. In: *Proc. Robotics: Science and Systems (RSS)*. 2013.
- [401] Dan Levi, Noa Garnett, and Ethan Fetaya. “StixelNet: A Deep Convolutional Network for Obstacle Detection and Road Segmentation”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2015.
- [402] Sergey Levine. *CS 294: Deep Reinforcement Learning*. http://rail.eecs.berkeley.edu/deeprlcourse-fa17/f17docs/lecture_2_behavior_cloning.pdf. Online: accessed 18-October-2019. 2017.
- [403] Evgeny Levinkov, Jonas Uhrig, Siyu Tang, Mohamed Omran, Eldar Insafutdinov, Alexander Kirillov, Carsten Rother, Thomas Brox, Bernt Schiele, and Bjoern Andres. “Joint Graph Decomposition & Node Labeling: Problem, Algorithms, Applications”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 1904–1912.
- [404] J. Levinson, M. Montemerlo, and S. Thrun. “Map-Based Precision Vehicle Localization in Urban Environments”. In: *Proc. Robotics: Science and Systems (RSS)*. 2007.
- [405] J. Levinson and S. Thrun. “Robust vehicle localization in urban environments using probabilistic maps”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2010.
- [406] Bo Li, Tianlei Zhang, and Tian Xia. “Vehicle Detection from 3D Lidar Using Fully Convolutional Network”. In: *Proc. Robotics: Science and Systems (RSS)*. 2016.
- [407] Guohao Li, Matthias Mueller, Vincent Casser, Neil Smith, Dominik L. Michels, and Bernard Ghanem. “Teaching UAVs to Race With Observational Imitation Learning”. In: *arXiv.org abs/1803.01129* (2018).
- [408] Jun Li, Xue Mei, Danil V. Prokhorov, and Dacheng Tao. “Deep Neural Network for Structural Prediction and Lane Detection in Traffic Scene”. In: *IEEE Trans. on Neural Networks and Learning Systems* 28.3 (2017), pp. 690–703.
- [409] Qizhu Li, Anurag Arnab, and Philip H. S. Torr. “Weakly- and Semi-supervised Panoptic Segmentation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 106–124.
- [410] Xiangtai Li, Houlong Zhao, Lei Han, Yunhai Tong, and Kuiyuan Yang. “GFF: Gated Fully Fusion for Semantic Segmentation”. In: *arXiv.org* (2019).

- [411] Y. Li, N. Snavely, and D. P. Huttenlocher. “Location Recognition using Prioritized Feature Matching”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.
- [412] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. “Fully Convolutional Instance-Aware Semantic Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 4438–4446.
- [413] Yujia Li and Richard S. Zemel. “Mean-Field Networks”. In: *Proc. of the International Conf. on Machine learning (ICML) Workshops* (2014).
- [414] Yunpeng Li, David J. Crandall, and Daniel P. Huttenlocher. “Landmark classification in large-scale image collections.” In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2009.
- [415] Yunpeng Li, Noah Snavely, Dan Huttenlocher, and Pascal Fua. “World-wide Pose Estimation using 3D Point Clouds”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [416] Xiaodan Liang, Tairui Wang, Luona Yang, and Eric Xing. “CIRL: Controllable Imitative Reinforcement Learning for Vision-based Self-driving”. In: *arXiv.org abs/1807.03776* (2018).
- [417] Tsung-Yi Lin, Serge J. Belongie, and James Hays. “Cross-View Image Geolocalization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013, pp. 891–898.
- [418] Tsung-Yi Lin, Yin Cui, Serge J. Belongie, and James Hays. “Learning deep representations for ground-to-aerial geolocalization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [419] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. “Focal Loss for Dense Object Detection”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 2999–3007.
- [420] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. “Microsoft COCO: Common Objects in Context”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [421] Liu Liu, Hongdong Li, and Yuchao Dai. “Efficient Global 2D-3D Matching for Camera Localization in a Large-Scale 3D Map”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 2391–2400.
- [422] Shu Liu, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. “SGN: Sequential Grouping Networks for Instance Segmentation”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 3516–3524.

- [423] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. “Path Aggregation Network for Instance Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 8759–8768.
- [424] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. “SSD: Single Shot MultiBox Detector”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [425] Jonathan Long, Evan Shelhamer, and Trevor Darrell. “Fully Convolutional Networks for Semantic Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [426] H. Longuet-Higgins. “A Computer Algorithm for Reconstructing a Scene from Two Projections”. In: *Nature* 293 (1981), pp. 133–135.
- [427] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. “SMPL: A Skinned Multi-Person Linear Model”. In: *ACM Trans. on Graphics* (2015).
- [428] William E. Lorensen and Harvey E. Cline. “Marching Cubes: A High Resolution 3D Surface Construction Algorithm”. In: *ACM Trans. on Graphics*. 1987.
- [429] David G. Lowe. “Distinctive Image Features from Scale-Invariant Keypoints”. In: *International Journal of Computer Vision (IJCV)* 60.2 (2004), pp. 91–110.
- [430] Stephanie M. Lowry, Niko Sünderhauf, Paul Newman, John J. Leonard, David D. Cox, Peter I. Corke, and Michael J. Milford. “Visual Place Recognition: A Survey”. In: *IEEE Trans. on Robotics* 32.1 (2016), pp. 1–19.
- [431] Chuanhua Lu, Hideaki Uchiyama, Diego Thomas, Atsushi Shimada, and Rin-ichiro Taniguchi. “Sparse Cost Volume for Efficient Stereo Matching”. In: *Remote Sensing (RS)* 10.11 (2018), p. 1844.
- [432] W. Luo, A. Schwing, and R. Urtasun. “Efficient Deep Learning for Stereo Matching”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [433] Zhaoyang Lv, Chris Beall, Pablo Alcantarilla, Fuxin Li, Zsolt Kira, and Frank Dellaert. “A Continuous Optimization Approach for Efficient and Accurate Scene Flow”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [434] Simon Lynen, Torsten Sattler, Michael Bosse, Joel A. Hesch, Marc Pollefeys, and Roland Siegwart. “Get Out of My Lab: Large-scale, Real-Time Visual-Inertial Localization”. In: *Proc. Robotics: Science and Systems (RSS)*. 2015.

- [435] Liqian Ma, Siyu Tang, Michael J. Black, and Luc Van Gool. “Customized Multi-Person Tracker”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2018.
- [436] Wei-Chiu Ma, Shenlong Wang, Rui Hu, Yuwen Xiong, and Raquel Urtasun. “Deep Rigid Instance Scene Flow”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2019).
- [437] Oisín Mac Aodha, Ahmad Humayun, Marc Pollefeys, and Gabriel J. Brostow. “Learning a Confidence Measure for Optical Flow”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 35.5 (2013), pp. 1107–1120.
- [438] Will Maddern, Geoff Pascoe, Chris Linegar, and Paul Newman. “1 Year, 1000km: The Oxford RobotCar Dataset”. In: *International Journal of Robotics Research (IJRR)* (2016).
- [439] Saturnino Maldonado-Bascón, Sergio Lafuente-Arroyo, Pedro Gil-Jimenez, Hilario Gómez-Moreno, and Francisco López-Ferrerías. “Road-sign detection and recognition based on support vector machines”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 8.2 (2007), pp. 264–278.
- [440] Fabian Manhardt, Wadim Kehl, and Adrien Gaidon. “ROI-10D: Monocular Lifting of 2D Detection to 6D Pose and Metric Shape”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2019).
- [441] Vikash Mansinghka, Tejas Kulkarni, Yura Perov, and Josh Tenenbaum. “Approximate Bayesian Image Interpretation using Generative Probabilistic Graphics Programs”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2013.
- [442] Ana I. Maqueda, Antonio Loquercio, Guillermo Gallego, Narciso N. García, and Davide Scaramuzza. “Event-Based Vision Meets Deep Learning on Steering Prediction for Self-Driving Cars”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [443] Andelo Martinović, Jan Knopp, Hayko Riemenschneider, and Luc Van Gool. “3D All The Way: Semantic Segmentation of Urban Scenes from Start to End in 3D”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [444] Markus Mathias, Andelo Martinovic, and Luc Van Gool. “ATLAS: A Three-Layered Approach to Facade Parsing”. In: *International Journal of Computer Vision (IJCV)* 118.1 (2016), pp. 22–48.

- [445] Markus Mathias, Radu Timofte, Rodrigo Benenson, and Luc Van Gool. “Traffic sign recognition - How far are we from the solution?”. In: *International Joint Conference on Neural Networks (IJCNN)*. IEEE. 2013, pp. 1–8.
- [446] Larry H. Matthies, Richard Szeliski, and Takeo Kanade. “Incremental estimation of dense depth maps from image sequences”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 1988, pp. 366–374.
- [447] Gellert Mattyus, Shenlong Wang, Sanja Fidler, and Raquel Urtasun. “Enhancing Road Maps by Parsing Aerial Images Around the World”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [448] Gellert Mattyus, Shenlong Wang, Sanja Fidler, and Raquel Urtasun. “HD Maps: Fine-Grained Road Segmentation by Parsing Ground and Aerial Images”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [449] Daniel Maturana and Sebastian Scherer. “VoxNet: A 3D Convolutional Neural Network for real-time object recognition”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2015.
- [450] N. Mayer, E. Ilg, P. Haeusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. “A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [451] John McCormac, Ankur Handa, Andrew J. Davison, and Stefan Leutenegger. “SemanticFusion: Dense 3D semantic mapping with convolutional neural networks”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2017, pp. 4628–4635.
- [452] Ashish Mehta, Adithya Subramanian, and Anbumani Subramanian. “Learning End-to-end Autonomous Driving using Guided Auxiliary Supervision”. In: *arXiv.org abs/1808.10393* (2018).
- [453] C. Mei and P. Rives. “Single View Point Omnidirectional Camera Calibration from Planar Grids”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2007.
- [454] Simon Meister, Junhwa Hur, and Stefan Roth. “UnFlow: Unsupervised Learning of Optical Flow with a Bidirectional Census Loss”. In: *Proc. of the Conf. on Artificial Intelligence (AAAI)*. 2018.
- [455] Moritz Menze and Andreas Geiger. “Object Scene Flow for Autonomous Vehicles”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.

- [456] Moritz Menze, Christian Heipke, and Andreas Geiger. “Discrete Optimization for Optical Flow”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2015.
- [457] Moritz Menze, Christian Heipke, and Andreas Geiger. “Joint 3D Estimation of Vehicles and Scene Flow”. In: *Proc. of the ISPRS Workshop on Image Sequence Analysis (ISA)*. 2015.
- [458] Moritz Menze, Christian Heipke, and Andreas Geiger. “Object Scene Flow”. In: *ISPRS Journal of Photogrammetry and Remote Sensing (JPRS)* 140 (2018), pp. 60–76.
- [459] Paul Merrell, Amir Akbarzadeh, Liang Wang, Philippos Mordohai, Jan-Michael Frahm, Ruigang Yang, David Nistér, and Marc Pollefeys. “Real-Time Visibility-Based Fusion of Depth Maps”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2007, pp. 1–8.
- [460] Cade Metz. *A Toaster on Wheels to Deliver Groceries? Self-Driving Tech Tests Practical Uses*. <https://www.nytimes.com/2018/12/18/technology/driverless-mini-car-deliver-groceries.html>. Online: accessed 18-October-2019. 2018.
- [461] Branislav Micusik and Jana Kosecka. “Piecewise planar city 3D modeling from street view panoramic sequences.” In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.
- [462] A. Milan, S. Roth, and K. Schindler. “Continuous Energy Minimization for Multitarget Tracking”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 36.1 (2014), pp. 58–72.
- [463] Anton Milan, Laura Leal-Taixé, Ian D. Reid, Stefan Roth, and Konrad Schindler. “MOT16: A Benchmark for Multi-Object Tracking”. In: *arXiv.org* 1603.00831 (2016).
- [464] Anton Milan, Seyed Hamid Rezaatofghi, Anthony R. Dick, Ian D. Reid, and Konrad Schindler. “Online Multi-Target Tracking Using Recurrent Neural Networks”. In: *Proc. of the Conf. on Artificial Intelligence (AAAI)*. 2017, pp. 4225–4232.
- [465] Anton Milan, Konrad Schindler, and Stefan Roth. “Detection- and Trajectory-Level Exclusion in Multiple Object Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [466] Mohammad Hossein Mirabdollah and Bärbel Mertsching. “Fast Techniques for Monocular Visual Odometry”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2015.
- [467] Mohammad Hossein Mirabdollah and Bärbel Mertsching. “On the Second Order Statistics of Essential Matrix Elements”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2014.

- [468] Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, Denis Teplyashin, Karl Moritz Hermann, Mateusz Malinowski, Matthew Koichi Grimes, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, and Raia Hadsell. “The StreetLearn Environment and Dataset”. In: *arXiv.org* abs/1903.01292 (2019).
- [469] Dennis Mitzel and Bastian Leibe. “Taking Mobile Multi-object Tracking to the Next Level: People, Unknown Objects, and Carried Items”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [470] Rahul Mohan. “Deep Deconvolutional Networks for Scene Parsing”. In: *arXiv.org* 1411.4101 (2014).
- [471] Michael Montemerlo, Sebastian Thrun, Daphne Koller, and Ben Wegbreit. “FastSLAM: A Factored Solution to the Simultaneous Localization and Mapping Problem”. In: *Artificial Intelligence (AI)*. 2002.
- [472] J. Montoya, J. D. Wegner, L. Ladicky, and K. Schindler. “Semantic segmentation of aerial images in urban areas with class-specific higher-order cliques”. In: *In ISPRS Conf. Photogrammetric Image Analysis (PIA)*. 2015.
- [473] Department of Motor Vehicles CA. *Report of Traffic Collision Involving an Autonomous Vehicle*. https://www.dmv.ca.gov/portal/dmv/detail/vr/autonomous/autonomousveh_ol316. Online: accessed 17-May-2019. 2019.
- [474] Pierre Moulon, Pascal Monasse, Renaud Marlet, et al. *OpenMVG. An Open Multiple View Geometry library*. <https://github.com/openMVG/openMVG>. Online: accessed 23-April-2019. 2012.
- [475] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Košecká. “3D Bounding Box Estimation Using Deep Learning and Geometry”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE. 2017, pp. 5632–5640.
- [476] Elias Mueggler, Guillermo Gallego, and Davide Scaramuzza. “Continuous-Time Trajectory Estimation for Event-based Vision Sensors”. In: *Proc. Robotics: Science and Systems (RSS)*. 2015.
- [477] Elias Mueggler, Henri Rebecq, Guillermo Gallego, Tobi Delbrück, and Davide Scaramuzza. “The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and SLAM”. In: *International Journal of Robotics Research (IJRR)*. 2017.
- [478] Matthias Müller, Alexey Dosovitskiy, Bernard Ghanem, and Vladlen Koltun. “Driving Policy Transfer via Modularity and Abstraction”. In: *arXiv.org* abs/1804.09364 (2018).

- [479] Daniel Munoz, J. Andrew Bagnell, and Martial Hebert. “Stacked Hierarchical Labeling”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.
- [480] Daniel Munoz, J. Andrew Bagnell, Nicolas Vandapel, and Martial Hebert. “Contextual Classification with Functional Max-Margin Markov Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.
- [481] Raul Mur-Artal, J. M. M. Montiel, and Juan D. Tardós. “ORB-SLAM: A Versatile and Accurate Monocular SLAM System”. In: *IEEE Trans. on Robotics* 31.5 (2015), pp. 1147–1163.
- [482] Przemyslaw Musialski, Peter Wonka, Daniel G. Aliaga, Michael Wimmer, Luc J. Van Gool, and Werner Purgathofer. “A Survey of Urban Reconstruction”. In: *Computer Graphics Forum* 32.6 (2013), pp. 146–177.
- [483] Heesoo Myeong, Ju Yong Chang, and Kyoung Mu Lee. “Learning object relationships via graph-based context model”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [484] Michal Neoral and Jan Šochman. “Object Scene Flow with Temporal Consistency”. In: *Proc. of the Computer Vision Winter Workshop (CVWW)*. 2017.
- [485] Michal Neoral, Jan Sochman, and Jiri Matas. “Continual Occlusions and Optical Flow Estimation”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2018.
- [486] Frank Neuhaus, Tilman Koß, Robert Kohnen, and Dietrich Paulus. “MC2SLAM: Real-Time Inertial Lidar Odometry Using Two-Scan Motion Compensation”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2018, pp. 60–72.
- [487] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulò, and Peter Kotschieder. “The Mapillary Vistas Dataset for Semantic Understanding of Street Scenes”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [488] Richard A. Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J. Davison, Pushmeet Kohli, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. “KinectFusion: Real-time Dense Surface Mapping and Tracking”. In: *Proc. of the International Symposium on Mixed and Augmented Reality (ISMAR)*. 2011.
- [489] Richard A. Newcombe, Steven Lovegrove, and Andrew J. Davison. “DTAM: Dense tracking and mapping in real-time”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2011.

- [490] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. “Real-time 3D Reconstruction at Scale using Voxel Hashing”. In: *ACM Trans. on Graphics*. 2013.
- [491] David Nistér. “An Efficient Solution to the Five-Point Relative Pose Problem”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 26.6 (2004), pp. 756–777.
- [492] Sang Min Oh, Sarah Tariq, Bruce N. Walker, and Frank Dellaert. “Map-based priors for localization”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2004.
- [493] Yuichi Ohta. *Knowledge-based interpretation of outdoor natural color scenes*. Vol. 4. Morgan Kaufmann, 1985.
- [494] Gabriel Oliveira, Wolfram Burgard, and Thomas Brox. “Efficient Deep Methods for Monocular Road Segmentation”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2016.
- [495] Viviane M. de Oliveira, Vítor Santos, Angel Domingo Sappa, Paulo Dias, and A. Paulo Moreira. “Incremental scenario representations for autonomous driving using geometric polygonal primitives”. In: *Robotics and Autonomous Systems (RAS)* 83 (2016), pp. 312–325.
- [496] Benedikt Ortelt, Christian Herrmann, Dieter Willersinn, and Jürgen Beyerer. “Foveal Vision for Instance Segmentation of Road Images”. In: *Proc. of the International Joint Conf. on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISIGRAPP)*. 2018, pp. 371–378.
- [497] Aljosa Osep, Wolfgang Mehner, Markus Mathias, and Bastian Leibe. “Combined image- and world-space tracking in traffic scenes”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2017.
- [498] Xinlei Pan, Yurong You, Ziyang Wang, and Cewu Lu. “Virtual to Real Reinforcement Learning for Autonomous Driving”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2017.
- [499] Gaurav Pandey, James R. McBride, and Ryan M. Eustice. “Ford campus vision and lidar data set”. In: *International Journal of Robotics Research (IJRR)* 30 (2011), pp. 1543–1552.
- [500] Jiahao Pang, Wenxiu Sun, Jimmy S. J. Ren, Chengxi Yang, and Qiong Yan. “Cascade Residual Learning: A Two-Stage Convolutional Neural Network for Stereo Matching”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 878–886.
- [501] Constantine Papageorgiou and Tomaso A. Poggio. “A Trainable System for Object Detection”. In: *International Journal of Computer Vision (IJCV)* 38.1 (2000), pp. 15–33.

- [502] Despoina Paschalidou, Ali Osman Ulusoy, Carolin Schmitt, Luc van Gool, and Andreas Geiger. “RayNet: Learning Volumetric 3D Reconstruction with Ray Potentials”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [503] Mark A. Paskin. “Thin Junction Tree Filters for Simultaneous Localization and Mapping”. In: *Proc. of the International Joint Conf. on Artificial Intelligence (IJCAI)*. 2003, pp. 1157–1166.
- [504] Rohan Paul and Paul Newman. “FAB-MAP 3D: Topological mapping with spatial and visual appearance”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2010.
- [505] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc J. Van Gool. “You’ll never walk alone: Modeling social behavior for multi-target tracking”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2009.
- [506] Bojan Pepik, Michael Stark, Peter V. Gehler, and Bernt Schiele. “Multi-View and 3D Deformable Part Models”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 37.11 (2015), pp. 2232–2245.
- [507] F. Pereira, J. Luft, G. Ilha, A. Sofiatti, and A. Susin. “Backward Motion for Estimation Enhancement in Sparse Visual Odometry”. In: *Workshop of Computer Vision (WVC)*. 2017, pp. 61–66.
- [508] A. G. A. Perera, C. Srinivas, A. Hoogs, G. Brooksby, and W. Hu. “Multi-object tracking through simultaneous long occlusions and split-merge conditions”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2006.
- [509] Mikael Persson, Tommaso Piccini, Michael Felsberg, and Rudolf Mester. “Robust stereo visual odometry from monocular techniques”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2015.
- [510] D. Pfeiffer and U. Franke. “Efficient representation of traffic scenes by means of dynamic stixels”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2010.
- [511] David Pfeiffer and Uwe Franke. “Towards a Global Optimal Multi-Layer Stixel Representation of Dense 3D Data”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2011.
- [512] Cuong Cao Pham and Jae Wook Jeon. “Robust object proposals re-ranking for object detection in autonomous driving using convolutional neural networks”. In: *Signal Processing: Image Communication (SPIC)* 53 (2017), pp. 110–122.

- [513] Sudeep Pillai, Rares Ambrus, and Adrien Gaidon. “SuperDepth: Self-Supervised, Super-Resolved Monocular Depth Estimation”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2019, pp. 9250–9256.
- [514] Peter Pinggera, Uwe Franke, and Rudolf Mester. “High-performance long range obstacle detection using stereo vision”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2015.
- [515] Peter Pinggera, Sebastian Ramos, Stefan Gehrig, Uwe Franke, Carsten Rother, and Rudolf Mester. “Lost and Found: detecting small road hazards for self-driving vehicles”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2016.
- [516] Pedro H. O. Pinheiro, Ronan Collobert, and Piotr Dollár. “Learning to Segment Object Candidates”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2015, pp. 1990–1998.
- [517] Pedro Oliveira Pinheiro, Tsung-Yi Lin, Ronan Collobert, and Piotr Dollár. “Learning to Refine Object Segments”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016, pp. 75–91.
- [518] Hamed Pirsiavash, Deva Ramanan, and Charless C. Fowlkes. “Globally-Optimal Greedy Algorithms for Tracking a Variable Number of Objects”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.
- [519] Leonid Pishchulin, Eldar Insafutdinov, Siyu Tang, Bjoern Andres, Mykhaylo Andriluka, Peter V. Gehler, and Bernt Schiele. “DeepCut: Joint Subset Partition and Labeling for Multi Person Pose Estimation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [520] Leonid Pishchulin, Arjun Jain, Mykhaylo Andriluka, Thorsten Thormaehlen, and Bernt Schiele. “Articulated People Detection and Pose Estimation: Reshaping the Future”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [521] Leonard Plotkin. “PyDriver: Entwicklung eines Frameworks für räumliche Detektion und Klassifikation von Objekten in Fahrzeugumgebung”. MA thesis. Karlsruhe Institute of Technology, 2015.
- [522] Tobias Pohlen, Alexander Hermans, Markus Mathias, and Bastian Leibe. “Full-Resolution Residual Networks for Semantic Segmentation in Street Scenes”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 3309–3318.
- [523] Thomas Pollard and Joseph L. Mundy. “Change Detection in a 3-d World”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2007.

- [524] M. Pollefeys. “Detailed Real-Time Urban 3D Reconstruction from Video”. In: *International Journal of Computer Vision (IJCV)* 78.2-3 (July 2008), pp. 143–167.
- [525] D. Pomerleau and T. Jochem. “Rapidly adapting machine vision for automated vehicle steering”. In: *EXPERT* (1996).
- [526] Dean Pomerleau. “ALVINN: An Autonomous Land Vehicle in a Neural Network”. In: *Advances in Neural Information Processing Systems (NIPS)*. 1988, pp. 305–313.
- [527] Dean Pomerleau and Todd Jochem. *Look, Ma, No Hands*. <https://www.cmu.edu/news/stories/archives/2015/july/look-ma-no-hands.html>. Online: accessed 18-June-2019. 2015.
- [528] Jordi Pont-Tuset, Pablo Arbelaez, Jonathan T. Barron, Ferran Marqués, and Jitendra Malik. “Multiscale Combinatorial Grouping for Image Segmentation and Object Proposal Generation”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 39.1 (2017), pp. 128–140.
- [529] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. “Frustum pointnets for 3d object detection from rgb-d data”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2018).
- [530] Zhen Qin and Christian R. Shelton. “Improving Multi-Target Tracking via Social Grouping”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [531] Weichao Qiu, Fangwei Zhong, Yi Zhang, Siyuan Qiao, Zihao Xiao, Tae Soo Kim, Yizhou Wang, and Alan Yuille. “UnrealCV: Virtual Worlds for Computer Vision”. In: *ACM Multimedia Open Source Software Competition* (2017).
- [532] L. H. Quam. “Hierarchical warp stereo”. In: *Image Understanding Workshop (IUW)*. 1984.
- [533] Clemens Rabe, Thomas Mueller, Andreas Wedel, and Uwe Franke. “Dense, Robust, and Accurate Motion Field Estimation from Stereo Image Sequences in Real-Time”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.
- [534] A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora, and S. Belongie. “Objects in Context”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2007.
- [535] Noha Radwan, Abhinav Valada, and Wolfram Burgard. “VLocNet++: Deep Multitask Learning for Semantic Visual Localization and Odometry”. In: *IEEE Robotics and Automation Letters (RA-L)* 3.4 (2018), pp. 4407–4414.

- [536] René Ranftl, Kristian Bredies, and Thomas Pock. “Non-local Total Generalized Variation for Optical Flow Estimation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [537] Rene Ranftl, Thomas Pock, and Horst Bischof. “Minimizing TGV-based Variational Models with Non-Convex Data terms”. In: *Proc. of the International Conf. on Scale Space and Variational Methods in Computer Vision (SSVM)*. 2013.
- [538] Anurag Ranjan and Michael Black. “Optical Flow Estimation using a Spatial Pyramid Network”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [539] Henri Rebecq, Timo Horstschaefer, Guillermo Gallego, and Davide Scaramuzza. “EVO: A Geometric Approach to Event-based 6-DOF Parallel Tracking and Mapping in Real-time”. In: *IEEE Robotics and Automation Letters (RA-L)*. 2016.
- [540] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. “You Only Look Once: Unified, Real-Time Object Detection”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 779–788.
- [541] Joseph Redmon and Ali Farhadi. “YOLO9000: Better, Faster, Stronger”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [542] Vladimir Reilly, Haroon Idrees, and Mubarak Shah. “Detection and Tracking of Large Number of Targets in Wide Area Surveillance”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.
- [543] Jimmy Ren, Xiaohao Chen, Jianbo Liu, Wenxiu Sun, Jiahao Pang, Qiong Yan, Yu-Wing Tai, and Li Xu. “Accurate single stage detector using recurrent rolling convolution”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [544] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2015.
- [545] Shaoqing Ren, Kaiming He, Ross B. Girshick, Xiangyu Zhang, and Jian Sun. “Object Detection Networks on Convolutional Feature Maps”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 39.7 (2017), pp. 1476–1481.
- [546] Zhile Ren, Deqing Sun, Jan Kautz, and Erik B. Sudderth. “Cascaded Scene Flow Prediction Using Semantic Segmentation”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2017, pp. 225–233.

- [547] Maria I Restrepo, Ali O Ulusoy, and Joseph L Mundy. “Evaluation of feature-based 3-D registration of probabilistic volumetric scenes”. In: *ISPRS Journal of Photogrammetry and Remote Sensing (JPRS)* (2014).
- [548] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. “EpicFlow: Edge-Preserving Interpolation of Correspondences for Optical Flow”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [549] Christian Richardt, Hyeonwoo Kim, Levi Valgaerts, and Christian Theobalt. “Dense Wide-Baseline Scene Flow From Two Handheld Video Cameras”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2016.
- [550] Stephan R. Richter, Zeeshan Hayder, and Vladlen Koltun. “Playing for Benchmarks”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [551] Stephan R. Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. “Playing for Data: Ground Truth from Computer Games”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [552] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. “OctNetFusion: Learning Depth Fusion from Data”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2017.
- [553] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. “OctNet: Learning Deep 3D Representations at High Resolutions”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [554] Hayko Riemenschneider, András Bódis-Szomorú, Julien Weissenberg, and Luc Van Gool. “Learning Where to Classify in Multi-view Semantic Segmentation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [555] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. “Performance Measures and a Data Set for Multi-Target, Multi-Camera Tracking”. In: *Proc. of the European Conf. on Computer Vision (ECCV) Workshops*. 2016.
- [556] Lawrence G. Roberts. “Machine perception of three-dimensional solids”. PhD thesis. Massachusetts Institute of Technology, 1963.
- [557] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. “U-Net: Convolutional Networks for Biomedical Image Segmentation”. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. 2015.

- [558] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio Lopez. “The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [559] Stéphane Ross and Drew Bagnell. “Efficient Reductions for Imitation Learning”. In: *Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [560] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. “GrabCut: interactive foreground extraction using iterated graph cuts”. In: *ACM Trans. on Graphics*. Vol. 23. 3. 2004, pp. 309–314.
- [561] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski. “ORB: an efficient alternative to SIFT or SURF”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2011, pp. 2564–2571.
- [562] Leonid I Rudin, Stanley Osher, and Emad Fatemi. “Nonlinear total variation based noise removal algorithms”. In: *Physica D: Nonlinear Phenomena* 60.1-4 (1992), pp. 259–268.
- [563] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. “Learning representations by back-propagating errors”. In: *Nature* 323 (1986), pp. 533–536.
- [564] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. “ImageNet Large Scale Visual Recognition Challenge”. In: *International Journal of Computer Vision (IJCV)* 115.3 (2015), pp. 211–252.
- [565] Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. “Tracking the Untrackable: Learning to Track Multiple Cues with Long-Term Dependencies”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [566] Soham Saha, Girish Varma, and C. V. Jawahar. “Improved Visual Relocalization by Discovering Anchor Points”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2018, p. 164.
- [567] Eder Santana and George Hotz. “Learning a Driving Simulator”. In: *arXiv.org abs/1608.01230* (2016).
- [568] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. “From Coarse to Fine: Robust Hierarchical Localization at Large Scale”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.

- [569] T. Sattler, B. Leibe, and L. Kobbelt. “Efficient Effective Prioritized Matching for Large-Scale Image-Based Localization”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* PP.99 (2016), pp. 1–1.
- [570] Torsten Sattler, Michal Havlena, Filip Radenovic, Konrad Schindler, and Marc Pollefeys. “Hyperpoints and Fine Vocabularies for Large-Scale Location Recognition”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [571] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, Fredrik Kahl, and Tomás Pajdla. “Benchmarking 6DOF Outdoor Visual Localization in Changing Conditions”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018, pp. 8601–8610.
- [572] Torsten Sattler, Qunjie Zhou, Marc Pollefeys, and Laura Leal-Taixé. “Understanding the Limitations of CNN-based Absolute Camera Pose Regression”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [573] Axel Sauer, Nikolay Savinov, and Andreas Geiger. “Conditional Affordance Learning for Driving in Urban Environments”. In: *Proc. Conf. on Robot Learning (CoRL)*. 2018.
- [574] Manolis Savva, Angel X. Chang, and Pat Hanrahan. “Semantically-Enriched 3D Models for Common-sense Knowledge”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) Workshops* (2015).
- [575] Rohan Saxena, René Schuster, Oliver Wasenmüller, and Didier Stricker. “PWOC-3D: Deep Occlusion-Aware End-to-End Scene Flow Estimation”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)* (2019).
- [576] D. Scaramuzza and R. Siegwart. “Appearance-Guided Monocular Omnidirectional Visual Odometry for Outdoor Ground Vehicles”. In: *IEEE Trans. on Robotics* 24.5 (2008), pp. 1015–1026.
- [577] Davide Scaramuzza and Friedrich Fraundorfer. “Visual Odometry [Tutorial]”. In: *Robotics and Automation Magazine (RAM)* 18.4 (2011), pp. 80–92.
- [578] Davide Scaramuzza, Friedrich Fraundorfer, and Roland Siegwart. “Real-time monocular visual odometry for on-road vehicles with 1-point RANSAC”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2009.

- [579] Davide Scaramuzza and Agostino Martinelli. “A Toolbox for Easily Calibrating Omnidirectional Cameras”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2006.
- [580] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nesic, Xi Wang, and Porter Westling. “High-Resolution Stereo Datasets with Subpixel-Accurate Ground Truth”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2014.
- [581] Daniel Scharstein and Richard Szeliski. “A taxonomy and evaluation of dense two-frame stereo correspondence algorithms”. In: *International Journal of Computer Vision (IJCV)* 47 (2002), pp. 7–42.
- [582] Daniel Scharstein and Richard Szeliski. “High-Accuracy Stereo Depth Maps Using Structured Light”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2003.
- [583] Konstantin Schauwecker, Reinhard Klette, and Andreas Zell. “A new feature detector and stereo matching method for accurate high-performance sparse stereo matching”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2012, pp. 5171–5176.
- [584] Samuel Scheidegger, Joachim Benjaminsson, Emil Rosenberg, Amrit Krishnan, and Karl Granström. “Mono-Camera 3D Multi-Object Tracking Using Deep Learning Detections and PMBM Filtering”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2018, pp. 433–440.
- [585] Lukas Schneider, Marius Cordts, Timo Rehfeld, David Pfeiffer, Markus Enzweiler, Uwe Franke, Marc Pollefeys, and Stefan Roth. “Semantic Stixels: Depth is not enough”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2016.
- [586] Miriam Schönbein and Andreas Geiger. “Omnidirectional 3D Reconstruction in Augmented Manhattan Worlds”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2014.
- [587] Miriam Schönbein, Tobias Strauss, and Andreas Geiger. “Calibrating and Centering Quasi-Central Catadioptric Cameras”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2014.
- [588] Johannes Lutz Schönberger and Jan-Michael Frahm. “Structure-from-Motion Revisited”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [589] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. “Pixelwise View Selection for Unstructured Multi-View Stereo”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.

- [590] Johannes Schönberger, Marc Pollefeys, Andreas Geiger, and Torsten Sattler. “Semantic Visual Localization”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [591] Thomas Schöps, Johannes Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. “A Multi-View Stereo Benchmark with High-Resolution Images and Multi-Camera Videos”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [592] Markus Schreiber, Carsten Knöppel, and Uwe Franke. “LaneLoc: Lane marking based localization using highly accurate maps”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2013.
- [593] Samuel Schulter, Paul Vernaza, Wongun Choi, and Manmohan Chandraker. “Deep Network Flow for Multi-Object Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [594] Sean Scott. *Meet Scout*. <https://blog.aboutamazon.com/transportation/meet-scout>. Online: accessed 18-October-2019. 2019.
- [595] Ari Seff and Jianxiong Xiao. “Learning from Maps: Visual Common Sense for Autonomous Driving”. In: *arXiv.org* 1611.08583 (2016).
- [596] Steven M. Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. “A Comparison and Evaluation of Multi-View Stereo Reconstruction Algorithms”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2006.
- [597] Akihito Seki and Marc Pollefeys. “Patch Based Confidence Prediction for Dense Disparity Map”. In: *Proc. of the British Machine Vision Conf. (BMVC)*. 2016.
- [598] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. “Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 618–626.
- [599] Sunando Sengupta, Eric Greveson, Ali Shahrokni, and Philip HS Torr. “Urban 3D Semantic Modelling Using Stereo Vision”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2013.
- [600] Sunando Sengupta, Paul Sturgess, Lubor Ladicky, and Philip H. S. Torr. “Automatic dense visual semantic mapping from street-level imagery.” In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2012.

- [601] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. “OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks”. In: *Proc. of the International Conf. on Learning Representations (ICLR)*. 2014.
- [602] Pierre Sermanet, Koray Kavukcuoglu, Soumith Chintala, and Yann LeCun. “Pedestrian Detection with Unsupervised Multi-stage Feature Learning”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [603] Pierre Sermanet and Yann LeCun. “Traffic sign recognition with multi-scale Convolutional Networks.” In: *International Joint Conference on Neural Networks (IJCNN)*. 2011, pp. 2809–2813.
- [604] Laura Sevilla-Lara, Deqing Sun, Varun Jampani, and Michael J. Black. “Optical Flow with Semantic Segmentation and Localized Layers”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [605] Khurram Shafique, Mun Wai Lee, and Niels Haering. “A Rank Constrained Continuous Formulation of Multi-Frame Multi-Target Tracking Problem”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [606] Qi Shan, Changchang Wu, Brian Curless, Yasutaka Furukawa, Carlos Hernández, and Steven M. Seitz. “Accurate Geo-Registration by Ground-to-Aerial Image Matching”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2014.
- [607] S. Sharifzadeh, I. Chiotellis, R. Triebel, and D. Cremers. “Learning to Drive using Inverse Reinforcement Learning and Deep Q-Networks”. In: *Advances in Neural Information Processing Systems (NIPS) Workshops*. 2016.
- [608] Sarthak Sharma, Junaid Ahmed Ansari, J. Krishna Murthy, and K. Madhava Krishna. “Beyond Pixels: Leveraging Geometry and Shape Cues for Online Multi-Object Tracking”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2018.
- [609] A. Shashua, Y. Gdalyahu, and G. Hayun. “Pedestrian detection for driving assistance systems: Single-frame classification and system level performance”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2004.
- [610] Han Shen, Lichao Huang, Chang Huang, and Wei Xu. “Tracklet Association Tracker: An End-to-End Learning-based Association Approach for Multi-Object Tracking”. In: *arXiv.org* (2018).

- [611] H. Sheng, Y. Zhang, J. Chen, Z. Xiong, and J. Zhang. “Heterogeneous Association Graph Fusion for Target Association in Multiple Object Tracking”. In: *IEEE Transactions on Circuits and Systems for Video Technology* (2019).
- [612] J. Shotton, J. Winn, C. Rother, and A. Criminisi. “TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context”. In: *International Journal of Computer Vision (IJCV)* 81 (2009), pp. 2–23.
- [613] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. “Scene Coordinate Regression Forests for Camera Relocalization in RGB-D Images”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [614] Guang Shu. “Part-based Multiple-Person Tracking with Partial Occlusion Handling”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [615] Eero P. Simoncelli, Edward H. Adelson, and David J. Heeger. “Probability distributions of optical flow”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 1991.
- [616] K. Simonyan and A. Zisserman. “Very Deep Convolutional Networks for Large-Scale Image Recognition”. In: *Proc. of the International Conf. on Learning Representations (ICLR)*. 2015.
- [617] M. Smith, I. Baldwin, W. Churchill, R. Paul, and P. Newman. “The new college vision and laser data set”. In: *International Journal of Robotics Research (IJRR)* 28 (2009), pp. 595–599.
- [618] R. Smith, Matthew Self, and Peter Cheeseman. “Estimating uncertain spatial relationships in robotics”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 1987.
- [619] Noah Snavely, Steven M. Seitz, and Richard Szeliski. “Modeling the World from Internet Photo Collections”. In: *International Journal of Computer Vision (IJCV)* 80.2 (2008), pp. 189–210.
- [620] Noah Snavely, Steven M. Seitz, and Richard Szeliski. “Photo Tourism: Exploring Photo Collections in 3D”. In: *ACM Trans. on Graphics*. 2006, pp. 835–846.
- [621] Shiyu Song and Manmohan Chandraker. “Robust Scale Estimation in Real-Time Monocular SFM for Autonomous Driving”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- [622] Xiao Song, Xu Zhao, Hanwen Hu, and Liangji Fang. “EdgeStereo: A Context Integrated Residual Pyramid Network for Stereo Matching”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2018.

- [623] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. “The German Traffic Sign Recognition Benchmark: A multi-class classification competition”. In: *International Joint Conference on Neural Networks (IJCNN)*. 2011, pp. 1453–1460.
- [624] Rainer Stiefelhagen, Keni Bernardin, Rachel Bowers, John Garofolo, Djamel Mostefa, and Padmanabhan Soundararajan. “The CLEAR 2006 Evaluation”. In: *CLEAR*. 2007.
- [625] Hauke Strasdat, Andrew J. Davison, J. M. M. Montiel, and Kurt Konolige. “Double window optimisation for constant time visual SLAM”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2011, pp. 2352–2359.
- [626] Hauke Strasdat, J. M. M. Montiel, and Andrew J. Davison. “Scale Drift-Aware Large Scale Monocular SLAM”. In: *Proc. Robotics: Science and Systems (RSS)*. 2010.
- [627] Christoph Strecha, Wolfgang von Hansen, Luc J. Van Gool, Pascal Fua, and Ulrich Thoennessen. “On benchmarking camera calibration and multi-view stereo for high resolution imagery”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [628] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. “A benchmark for the evaluation of RGB-D SLAM systems”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2012, pp. 573–580.
- [629] Frédéric Suard, Alain Rakotomamonjy, Abdelaziz Bensrhair, and Alberto Broggi. “Pedestrian detection using infrared images and histograms of oriented gradients”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. IEEE. 2006, pp. 206–212.
- [630] T. Suleymanov, L. M. Paz, P. Piniés, G. Hester, and P. Newman. “The path less taken: A fast variational approach for scene segmentation used for closed loop control”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2016.
- [631] Deqing Sun, Stefan Roth, and Michael J. Black. “A Quantitative Analysis of Current Practices in Optical Flow Estimation and the Principles Behind Them”. In: *International Journal of Computer Vision (IJCV)* 106.2 (2014), pp. 115–137.
- [632] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. “Models Matter, So Does Training: An Empirical Study of CNNs for Optical Flow Estimation”. In: *arXiv.org* (2018).

- [633] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. “PWC-Net: CNNs for Optical Flow Using Pyramid, Warping, and Cost Volume”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [634] Min Sun and Silvio Savarese. “Articulated part-based model for joint object detection and pose estimation.” In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2011.
- [635] Linus Svärm, Olof Enqvist, Fredrik Kahl, and Magnus Oskarsson. “City-Scale Localization for Cameras with Known Vertical Direction”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 39.7 (2017), pp. 1455–1461.
- [636] Linus Svärm, Olof Enqvist, Magnus Oskarsson, and Fredrik Kahl. “Accurate Localization and Pose Estimation for Large 3D Models”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014, pp. 532–539.
- [637] Tomáš Svoboda and Tomáš Pajdla. “Epipolar Geometry for Central Catadioptric Cameras”. In: *International Journal of Computer Vision (IJCV)* 49.1 (2002), pp. 23–37.
- [638] Chris Sweeney. *Theia Multiview Geometry Library: Tutorial & Reference*. <http://theia-sfm.org>. Online: accessed 23-April-2019. 2016.
- [639] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. “Going deeper with convolutions”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [640] Richard Szeliski. *Computer Vision - Algorithms and Applications*. Texts in Computer Science. Springer, 2011.
- [641] Peng Tang, Chunyu Wang, Xinggang Wang, Wenyu Liu, Wenjun Zeng, and Jingdong Wang. “Object Detection in Videos by High Quality Object Linking”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* (2019).
- [642] Siyu Tang, Bjoern Andres, Mykhaylo Andriluka, and Bernt Schiele. “Multi-person Tracking by Multicut and Deep Matching”. In: *Proc. of the European Conf. on Computer Vision (ECCV) Workshops*. 2016.
- [643] Siyu Tang, Bjoern Andres, Mykhaylo Andriluka, and Bernt Schiele. “Subgraph Decomposition for Multi-Target Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [644] Siyu Tang, Mykhaylo Andriluka, Bjoern Andres, and Bernt Schiele. “Multiple People Tracking by Lifted Multicut and Person Re-identification”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 3701–3710.

- [645] Tim Yuqing Tang, David Juny Yoon, and Timothy D. Barfoot. “A White-Noise-on-Jerk Motion Prior for Continuous-Time Trajectory Estimation on SE(3)”. In: *IEEE Robotics and Automation Letters (RA-L)* 4.2 (2019), pp. 594–601.
- [646] Tatsunori Tani, Sudeep N. Sinha, and Yoichi Sato. “Fast Multi-frame Stereo Scene Flow with Motion Segmentation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 6891–6900.
- [647] Tesla. *Introducing Software Version 9.0*. <https://www.tesla.com/blog/introducing-software-version-9?redirect=no>. Online: accessed 8-June-2019. 2018.
- [648] Tesla. *Tesla Autopilot*. <https://www.tesla.com/autopilot>. Online: accessed 18-October-2019. 2014.
- [649] Charles Thorpe, Martial H. Hebert, Takeo Kanade, and Steven A. Shafer. “Vision and Navigation for the Carnegie-Mellon Navlab”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 10.3 (May 1988), pp. 362–372.
- [650] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics*. The MIT Press, 2005.
- [651] Sebastian Thrun, Michael Montemerlo, Hendrik Dahlkamp, David Stavens, Andrei Aron, James Diebel, Philip Fong, John Gale, Morgan Halpenny, Gabriel Hoffmann, Kenny Lau, Celia M. Oakley, Mark Palatucci, Vaughan R. Pratt, Pascal Stang, Sven Strohband, Cedric Dupont, Lars-Erik Jendrossek, Christian Koelen, Charles Markey, Carlo Rummel, Joe van Niekerk, Eric Jensen, Philippe Alessandrini, Gary R. Bradski, Bob Davies, Scott Ettinger, Adrian Kaehler, Ara V. Nefian, and Pamela Mahoney. “Stanley: The robot that won the DARPA Grand Challenge”. In: *Journal of Field Robotics (JFR)* 23.9 (2006), pp. 661–692.
- [652] Wei Tian, Martin Lauer, and Long Chen. “Online Multi-Object Tracking Using Joint Domain Information in Traffic Scenarios”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* (2019).
- [653] Yicong Tian, Afshin Dehghan, and Mubarak Shah. “On Detection, Data Association and Segmentation for Multi-target Tracking”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* PP (2018), pp. 1–1.
- [654] LLC TIME USA. *Science: Radio Auto*. <http://content.time.com/time/magazine/article/0,9171,720720,00.html>. Online: accessed 18-October-2019. 1925.

- [655] Radu Timofte and Luc Van Gool. “Sparse Flow: Sparse Matching for Small to Large Displacement Optical Flow”. In: *Proc. of the IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2015.
- [656] Federico Tombari, Samuele Salti, and Luigi di Stefano. “Unique Signatures of Histograms for Local Surface Description”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.
- [657] Daniel Topfer, Jens Spehr, Jan Effertz, and Christoph Stiller. “Efficient Road Scene Understanding for Intelligent Vehicles Using Compositional Hierarchical Models”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 16.1 (2015), pp. 441–451.
- [658] A. Torii, R. Arandjelović, J. Sivic, M. Okutomi, and T. Pajdla. “24/7 Place Recognition by View Synthesis”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [659] Andy Tsai, Anthony Yezzi Jr, William Wells, Clare Tempany, Dewey Tucker, Ayres Fan, W Eric Grimson, and Alan Willsky. “A shape-based approach to the segmentation of medical imagery using level sets”. In: *Medical Imaging* 22.2 (2003), pp. 137–154.
- [660] Sik-Ho Tsang. *Review: DenseNet - Dense Convolutional Network (Image Classification)*. <https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>. Online: accessed 8-June-2019. 2018.
- [661] Zhuowen Tu and Xiang Bai. “Auto-Context and Its Application to High-Level Vision Tasks and 3D Brain Image Segmentation”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 32.10 (2010), pp. 1744–1757.
- [662] Stepan Tulyakov, Anton Ivanov, and Francois Fleuret. “Practical Deep Stereo (PDS): Toward applications-friendly deep stereo matching”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2018, pp. 5875–5885.
- [663] Uber. *Advanced Technologies Group*. <https://www.uber.com/de/de/atg/>. Online: accessed 18-October-2019. 2015.
- [664] Jonas Uhrig, Marius Cordts, Uwe Franke, and Thomas Brox. “Pixel-Level Encoding and Depth Layering for Instance-Level Semantic Labeling”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2016, pp. 14–25.
- [665] Jasper RR Uijlings, Koen EA van de Sande, Theo Gevers, and Arnold WM Smeulders. “Selective search for object recognition”. In: *International Journal of Computer Vision (IJCV)* 104.2 (2013), pp. 154–171.

- [666] Ali Osman Ulusoy, Michael Black, and Andreas Geiger. “Semantic Multi-view Stereo: Jointly Estimating Objects and Voxels”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [667] Ali Osman Ulusoy, Andreas Geiger, and Michael J. Black. “Towards Probabilistic Volumetric Reconstruction using Ray Potentials”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2015.
- [668] S. Uras, F. Girosi, A. Verri, and V. Torre. “A computational approach to motion perception”. In: *Biological Cybernetics* 60.2 (1988), pp. 79–87.
- [669] Vladyslav C. Usenko, Jakob Engel, Jörg Stückler, and Daniel Cremers. “Direct visual-inertial odometry with stereo cameras”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2016, pp. 1885–1892.
- [670] Abhinav Valada, Rohit Mohan, and Wolfram Burgard. “Self-Supervised Model Adaptation for Multimodal Semantic Segmentation”. In: *arXiv.org* (2018).
- [671] Abhinav Valada, Johan Vertens, Ankit Dhall, and Wolfram Burgard. “AdapNet: Adaptive semantic segmentation in adverse environmental conditions”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2017, pp. 4644–4651.
- [672] Julien PC Valentin, Sunando Sengupta, Jonathan Warrell, Ali Shahrokni, and Philip HS Torr. “Mesh based semantic modelling for indoor and outdoor scenes”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [673] S. Vedula, S. Baker, P. Rander, R. Collins, and T. Kanade. “Three-dimensional scene flow”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 1999.
- [674] Jakob J. Verbeek and Bill Triggs. “Scene Segmentation with CRFs Learned from Partially Labeled Images”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2007, pp. 1553–1560.
- [675] Yannick Verdie and Florent Lafarge. “Detecting parametric objects in large scenes by Monte Carlo sampling”. In: *International Journal of Computer Vision (IJCV)* 106.1 (2014), pp. 57–75.
- [676] Antoni Rosinol Vidal, Henri Rebecq, Timo Horstschaefer, and Davide Scaramuzza. “Ultimate SLAM? Combining Events, Images, and IMU for Robust Visual SLAM in HDR and High-Speed Scenarios”. In: *IEEE Robotics and Automation Letters (RA-L)* (2018).

- [677] Sudheendra Vijayanarasimhan and Kristen Grauman. “Active Frame Selection for Label Propagation in Videos”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [678] Vibhav Vineet, Ondrej Miksik, Morten Lidegaard, Matthias Niessner, Stuart Golodetz, Victor A. Prisacariu, Olaf Kahler, David W. Murray, Shahram Izadi, Patrick Perez, and Philip H. S. Torr. “Incremental Dense Semantic Stereo Fusion for Large-Scale Semantic Scene Reconstruction”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2015.
- [679] P. A. Viola, M. J. Jones, and D. Snow. “Detecting pedestrians using patterns of motion and appearance”. In: *International Journal of Computer Vision (IJCV)* 63(2) (2005), pp. 153–161.
- [680] Paul A. Viola and Michael J. Jones. “Robust Real-Time Face Detection”. In: *International Journal of Computer Vision (IJCV)* 57.2 (2004), pp. 137–154.
- [681] Christoph Vogel, Stefan Roth, and Konrad Schindler. “An Evaluation of Data Costs for Optical Flow”. In: *Proc. of the German Conference on Pattern Recognition (GCPR)*. 2013.
- [682] Christoph Vogel, Konrad Schindler, and Stefan Roth. “3D scene flow estimation with a piecewise rigid scene model”. In: *International Journal of Computer Vision (IJCV)* 115.1 (2015), pp. 1–28.
- [683] George Vosselman, Sander Dijkman, et al. “3D building model reconstruction from point clouds and ground plans”. In: *Proc. of the ISPRS Workshop Land Surface Mapping and Characterization Using Laser Altimetry* 34.3/W4 (2001), pp. 37–44.
- [684] Florian Walch, Caner Hazirbas, Laura Leal-Taixé, Torsten Sattler, Sebastian Hilsenbeck, and Daniel Cremers. “Image-Based Localization Using LSTMs for Structured Feature Correlation”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 627–637.
- [685] S. Walk, N. Majer, K. Schindler, and B. Schiele. “New features and insights for pedestrian detection”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2010.
- [686] Dequan Wang, Coline Devin, Qi-Zhi Cai, Philipp Krähenbühl, and Trevor Darrell. “Monocular Plan View Networks for Autonomous Driving”. In: *arXiv.org abs/1905.06937* (2019).
- [687] Dominic Zeng Wang and Ingmar Posner. “Voting for Voting in On-line Point Cloud Object Detection”. In: *Proc. Robotics: Science and Systems (RSS)*. 2015.

- [688] Gaoang Wang, Yizhou Wang, Haotian Zhang, Renshu Gu, and Jenq-Neng Hwang. “Exploit the Connectivity: Multi-Object Tracking with TrackletNet”. In: *arXiv.org abs/1811.07258* (2018).
- [689] Rui Wang, Martin Schwörer, and Daniel Cremers. “Stereo DSO: Large-Scale Direct Sparse Visual Odometry with Stereo Cameras”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017, pp. 3923–3931.
- [690] Shenlong Wang, Simon Suo, Wei-Chiu Ma, Andrei Pokrovsky, and Raquel Urtasun. “Deep Parametric Continuous Convolutional Neural Networks”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. June 2018.
- [691] Shiyao Wang, Yucong Zhou, Junjie Yan, and Zhidong Deng. “Fully Motion-Aware Network for Video Object Detection”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018.
- [692] Xiaoyu Wang, Ming Yang, Shenghuo Zhu, and Yuanqing Lin. “Regionlets for Generic Object Detection”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 37.10 (2015), pp. 2071–2084.
- [693] Yang Wang, Yi Yang, Zhenheng Yang, Liang Zhao, Peng Wang, and Wei Xu. “Occlusion Aware Unsupervised Learning of Optical Flow”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [694] Ziyang Wang, Buyu Liu, Samuel Schuster, and Manmohan Chandraker. “A Parametric Top-View Representation of Complex Road Scenes”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [695] Waymo. *Be an early rider*. <https://waymo.com/apply>. Online: accessed 18-October-2019. 2019.
- [696] A. Wedel, T. Brox, T. Vaudrey, C. Rabe, U. Franke, and D. Cremers. “Stereoscopic scene flow computation for 3D motion understanding”. In: *International Journal of Computer Vision (IJCV)* 95.1 (2011), pp. 29–51.
- [697] A. Wedel, C. Rabe, H. Badino, H. Loose, U. Franke, and D. Cremers. “B-Spline Modeling of Road Surfaces with an Application to Free Space Estimation”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* 10.4 (2009), pp. 572–583.
- [698] Andreas Wedel, Clemens Rabe, Tobi Vaudrey, Thomas Brox, Uwe Franke, and Daniel Cremers. “Efficient Dense Scene Flow from Sparse or Dense Stereo Data”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2008.

- [699] Jan D. Wegner, Steven Branson, David Hall, Konrad Schindler, and Pietro Perona. “Cataloging Public Objects Using Aerial and Street-Level Images - Urban Trees”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. June 2016.
- [700] Jan Dirk Wegner, Javier A. Montoya-Zegarra, and Konrad Schindler. “A Higher-Order CRF Model for Road Network Extraction”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [701] Jan Dirk Wegner, Javier Alexander Montoya-Zegarra, and Konrad Schindler. “Road networks as collections of minimum cost paths”. In: *ISPRS Journal of Photogrammetry and Remote Sensing (JPRS)* 108.Complete (2015), pp. 128–137.
- [702] D. Wei, C. Liu, and W.T. Freeman. “A Data-driven Regularization Model for Stereo and Flow”. In: *Proc. of the International Conf. on 3D Vision (3DV)*. 2014.
- [703] Philippe Weinzaepfel, Jérôme Revaud, Zaïd Harchaoui, and Cordelia Schmid. “DeepFlow: Large Displacement Optical Flow with Deep Matching”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2013.
- [704] Thomas Whelan, Stefan Leutenegger, Renato F. Salas-Moreno, Ben Glocker, and Andrew J. Davison. “ElasticFusion: Dense SLAM Without A Pose Graph”. In: *Proc. Robotics: Science and Systems (RSS)*. 2015.
- [705] H. Winner, S. Hakuli, F. Lotz, C. Singer, Andreas Geiger, et al. *Handbook of Driver Assistance Systems*. Springer Vieweg, 2015.
- [706] Christian Wöhler and Joachim K. Anlauf. “A time delay neural network algorithm for estimating image-pattern shape and motion”. In: *Image and Vision Computing (IVC)* (1999).
- [707] C. Wojek, S. Roth, K. Schindler, and B. Schiele. “Monocular 3D Scene Modeling and Inference: Understanding Multi-Object Traffic Scenes”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2010.
- [708] C. Wojek, S. Walk, and B. Schiele. “Multi-Cue Onboard Pedestrian Detection”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.
- [709] Christian Wojek and Bernt Schiele. “A Dynamic Conditional Random Field Model for Joint Labeling of Object and Scene Classes”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2008.
- [710] Christian Wojek and Bernt Schiele. “A Performance Evaluation of Single and Multi-feature People Detection”. In: *Proc. of the DAGM Symposium on Pattern Recognition (DAGM)*. 2008.

- [711] Christian Wojek, Stefan Walk, Stefan Roth, and Bernt Schiele. “Monocular 3D Scene Understanding with Explicit Occlusion Reasoning”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.
- [712] Christian Wojek, Stefan Walk, Stefan Roth, Konrad Schindler, and Bernt Schiele. “Monocular Visual Scene Understanding: Understanding Multi-Object Traffic Scenes”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 35.4 (2013), pp. 882–897.
- [713] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. “Simple online and realtime tracking with a deep association metric”. In: *Proc. IEEE International Conf. on Image Processing (ICIP)*. 2017, pp. 3645–3649.
- [714] Mark Wolff, Robert T. Collins, and Yanxi Liu. “Regularity-Driven Facade Matching Between Aerial and Street Views”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. June 2016.
- [715] Oliver Woodford, Philip Torr, Ian Reid, and Andrew Fitzgibbon. “Global Stereo Reconstruction under Second-Order Smoothness Priors”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 31 (2009), pp. 2115–2128.
- [716] Scott Workman, Richard Souvenir, and Nathan Jacobs. “Wide-Area Image Geolocalization with Aerial Reference Imagery”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [717] B. Wu and R. Nevatia. “Detection and Tracking of Multiple, Partially Occluded Humans by Bayesian Combination of Edgelet Part Detectors”. In: *International Journal of Computer Vision (IJCV)* 75.2 (2007), pp. 247–266.
- [718] Changchang Wu. *VisualSFM: A visual structure from motion system*.
- [719] Tao Wu and Ananth Ranganathan. “A practical system for road marking detection and recognition”. In: *Proc. IEEE Intelligent Vehicles Symposium (IV)*. 2012, pp. 25–30.
- [720] Tianfu Wu, Bo Li, and Song-Chun Zhu. “Learning And-Or Model to Represent Context and Occlusion for Car Detection and Viewpoint Estimation”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 38.9 (2016), pp. 1829–1843.
- [721] Zheng Wu, Thomas H. Kunz, and Margrit Betke. “Efficient track linking methods for track graphs using network-flow and set-cover techniques”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011, pp. 1185–1192.

- [722] Zheng Wu, Ashwin Thangali, Stan Sclaroff, and Margrit Betke. “Coupling Detection and Data Association for Multiple Object Tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [723] Zifeng Wu, Chunhua Shen, and Anton van den Hengel. “Wider or Deeper: Revisiting the ResNet Model for Visual Recognition”. In: *Pattern Recognition* 90 (2019), pp. 119–133.
- [724] Jonas Wulff and Michael J. Black. “Efficient Sparse-to-Dense Optical Flow Estimation using a Learned Basis and Layers”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [725] Bernhard Wymann, Eric Espié, Christophe Guionneau, Christos Dimitrakakis, Rémi Coulom, and Andrew Sumner. “Torcs, the open racing car simulator”. In: *arXiv.org* (2015).
- [726] Yu Xiang, Alexandre Alahi, and Silvio Savarese. “Learning to Track: Online Multi-object Tracking by Decision Making”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [727] Yu Xiang, Wongun Choi, Yuanqing Lin, and Silvio Savarese. “Data-driven 3d voxel patterns for object category recognition”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [728] Yu Xiang, Wongun Choi, Yuanqing Lin, and Silvio Savarese. “Subcategory-Aware Convolutional Neural Networks for Object Proposals and Detection”. In: *Proc. of the IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2017, pp. 924–933.
- [729] Jianxiong Xiao, Tian Fang, Peng Zhao, Maxime Lhuillier, and Long Quan. “Image-based street-side city modeling”. In: *ACM Trans. on Graphics* 28.5 (2009), 114:1–114:12.
- [730] Jianxiong Xiao and Long Quan. “Multiple view semantic segmentation for street view images.” In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2009.
- [731] Jun Xie, Martin Kiefel, Ming-Ting Sun, and Andreas Geiger. “Semantic Instance Annotation of Street Scenes by 3D to 2D Label Transfer”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [732] Junliang Xing, Haizhou Ai, and Shihong Lao. “Multi-object tracking through occlusions by local tracklets filtering and global tracklets association with detection responses”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.

- [733] Xuehan Xiong, Daniel Munoz, J. Andrew Bagnell, and Martial Hebert. “3-D scene analysis via sequenced predictions over points and regions”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2011, pp. 2609–2616.
- [734] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. “UPSNet: A Unified Panoptic Segmentation Network”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [735] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. “End-to-End Learning of Driving Models from Large-Scale Video Datasets”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 3530–3538.
- [736] Qingshan Xu and Wenbing Tao. “Multi-Scale Geometric Consistency Guided Multi-View Stereo”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [737] K. Yamaguchi, D. McAllester, and R. Urtasun. “Robust Monocular Epipolar Flow Estimation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2013.
- [738] Koichiro Yamaguchi, Tamir Hazan, David McAllester, and Raquel Urtasun. “Continuous Markov Random Fields for Robust Stereo Estimation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [739] Koichiro Yamaguchi, David McAllester, and Raquel Urtasun. “Efficient joint segmentation, occlusion labeling, stereo and flow estimation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [740] Yan Yan, Yuxing Mao, and Bo Li. “SECOND: Sparsely Embedded Convolutional Detection”. In: *Sensors* 18.10 (2018), p. 3337.
- [741] Bin Yang, Junjie Yan, Zhen Lei, and Stan Z Li. “Craft objects from images”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 6043–6051.
- [742] Bo Yang, Chang Huang, and R. Nevatia. “Learning affinities and dependencies for multi-target tracking using a CRF model”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2011.
- [743] Bo Yang and Ram Nevatia. “An online learned CRF model for multi-target tracking”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.

- [744] Fan Yang, Wongun Choi, and Yuanqing Lin. “Exploit All the Layers: Fast and Accurate CNN Object Detector with Scale Dependent Pooling and Cascaded Rejection Classifiers”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [745] Guorun Yang, Hengshuang Zhao, Jianping Shi, Zhidong Deng, and Jiaya Jia. “SegStereo: Exploiting Semantic Information for Disparity Estimation”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 660–676.
- [746] Nan Yang, Rui Wang, Jörg Stückler, and Daniel Cremers. “Deep Virtual Stereo Odometry: Leveraging Deep Depth Prediction for Monocular Direct Sparse Odometry”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 835–852.
- [747] Tingting Yang, Xiang Long, Arun Kumar Sangaiah, Zhigao Zheng, and Chao Tong. “Deep detection network for real-life traffic sign in vehicular networks”. In: *Computer Networks* 136 (2018), pp. 95–104.
- [748] Yi Yang, Simon Baker, Anitha Kannan, and Deva Ramanan. “Recognizing proxemics in personal photos”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012, pp. 3522–3529.
- [749] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. “MVSNet: Depth Inference for Unstructured Multi-view Stereo”. In: *Proc. of the European Conf. on Computer Vision (ECCV)* (2018).
- [750] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. “Recurrent MVSNet for High-resolution Multi-view Stereo Depth Inference”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2019).
- [751] Ju Hong Yoon, Chang-Ryeol Lee, Ming-Hsuan Yang, and Kuk-Jin Yoon. “Online Multi-object Tracking via Structural Constraint Event Aggregation”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [752] Ju Hong Yoon, Ming-Hsuan Yang, Jongwoo Lim, and Kuk-Jin Yoon. “Bayesian Multi-object Tracking Using Motion Context from Multiple Objects”. In: *Proc. of the IEEE Winter Conference on Applications of Computer Vision (WACV)*. 2015.
- [753] Fengwei Yu, Wenbo Li, Quanquan Li, Yu Liu, Xiaohua Shi, and Junjie Yan. “POI: Multiple Object Tracking with High Performance Detection and Appearance Feature”. In: *Proc. of the European Conf. on Computer Vision (ECCV) Workshops*. 2016.
- [754] Fisher Yu and Vladlen Koltun. “Multi-Scale Context Aggregation by Dilated Convolutions”. In: *Proc. of the International Conf. on Learning Representations (ICLR)*. 2016.

- [755] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. “BDD100K: A Diverse Driving Video Database with Scalable Annotation Tooling”. In: *arXiv.org* (2018).
- [756] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. *Berkeley DeepDrive*. <https://bdd-data.berkeley.edu/>. Online: accessed 05-June-2019. 2019.
- [757] Fisher Yu, Jianxiong Xiao, and Thomas A. Funkhouser. “Semantic alignment of LiDAR data at city scale”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 1722–1731.
- [758] Jason J. Yu, Adam W. Harley, and Konstantinos G. Derpanis. “Back to Basics: Unsupervised Learning of Optical Flow via Brightness Constancy and Motion Smoothness”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2016.
- [759] C. Zach, T. Pock, and H. Bischof. “A Duality Based Approach for Realtime TV-L1 Optical Flow”. In: *Proc. of the DAGM Symposium on Pattern Recognition (DAGM)*. 2007, pp. 214–223.
- [760] Christopher Zach, Thomas Pock, and Horst Bischof. “A duality based approach for realtime TV-L1 optical flow”. In: *Pattern Recognition Letters*. Springer Berlin Heidelberg, 2007, pp. 214–223.
- [761] Christopher Zach, Thomas Pock, and Horst Bischof. “A Globally Optimal Algorithm for Robust TV-L1 Range Image Integration.” In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2007.
- [762] Amir Roshan Zamir, Afshin Dehghan, and Mubarak Shah. “GMCP-Tracker: Global Multi-Object Tracking Using Generalized Minimum Clique Graphs”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2012.
- [763] Jure Žbontar and Yann LeCun. “Stereo Matching by Training a Convolutional Neural Network to Compare Image Patches”. In: *Journal of Machine Learning Research (JMLR)* 17.65 (2016), pp. 1–32.
- [764] Matthew D. Zeiler, Graham W. Taylor, and Rob Fergus. “Adaptive deconvolutional networks for mid and high level feature learning”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2011.
- [765] Bernhard Zeisl, Torsten Sattler, and Marc Pollefeys. “Camera Pose Voting for Large-Scale Image-Based Localization”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [766] Hongyi Zhang, Andreas Geiger, and Raquel Urtasun. “Understanding High-Level Semantics by Modeling Traffic Patterns”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2013.

- [767] Ji Zhang, Michael Kaess, and Sanjiv Singh. “Real-time Depth Enhanced Monocular Odometry”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2014.
- [768] Ji Zhang and Sanjiv Singh. “LOAM: Lidar Odometry and Mapping in Real-time”. In: *Proc. Robotics: Science and Systems (RSS)*. 2014.
- [769] Ji Zhang and Sanjiv Singh. “Visual-lidar odometry and mapping: low-drift, robust, and fast”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2015.
- [770] L. Zhang, Y. Li, and R. Nevatia. “Global Data Association for Multi-Object Tracking Using Network Flows”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2008.
- [771] Ning Zhang, Jeff Donahue, Ross B. Girshick, and Trevor Darrell. “Part-Based R-CNNs for Fine-Grained Category Detection”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2014.
- [772] Qilong Zhang and R. Pless. “Extrinsic calibration of a camera and laser range finder”. In: *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*. 2004.
- [773] Richard Zhang, Stefan A. Candra, Kai Vetter, and Avideh Zakhori. “Sensor fusion for semantic segmentation of urban scenes”. In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*. 2015, pp. 1850–1857.
- [774] Shanshan Zhang, Rodrigo Benenson, Mohamed Omran, Jan Hendrik Hosang, and Bernt Schiele. “How Far are We from Solving Pedestrian Detection?” In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [775] Shiquan Zhang, Xu Zhao, Liangji Fang, Haiping Fei, and Haitao Song. “Led: Localization-Quality Estimation Embedded Detector”. In: *Proc. IEEE International Conf. on Image Processing (ICIP)*. 2018, pp. 584–588.
- [776] Yimeng Zhang and Tsuhan Chen. “Efficient inference for fully-connected CRFs with stationarity”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [777] Ziyu Zhang, Sanja Fidler, and Raquel Urtasun. “Instance-Level Segmentation for Autonomous Driving with Deep Densely Connected MRFs”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [778] Ziyu Zhang, Alexander G. Schwing, Sanja Fidler, and Raquel Urtasun. “Monocular Object Instance Segmentation and Depth Ordering with CNNs”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.

- [779] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. “Pyramid Scene Parsing Network”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 6230–6239.
- [780] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip H. S. Torr. “Conditional Random Fields as Recurrent Neural Networks”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [781] Yantao Zheng, Ming Zhao, Yang Song, Hartwig Adam, Ulrich Bud-demeier, Alessandro Bissacco, Fernando Brucher, Tat-Seng Chua, and Hartmut Neven. “Tour the world: Building a web-scale landmark recog-nition engine.” In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2009.
- [782] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. “Learning Deep Features for Scene Recognition using Places Database”. In: *Advances in Neural Information Processing Sys-tems (NIPS)*. 2014.
- [783] Brady Zhou, Philipp Krähenbühl, and Vladlen Koltun. “Does com-puter vision matter for action?” In: *arXiv.org abs/1905.12887* (2019).
- [784] Chen Zhou, Fatma Güney, Yizhou Wang, and Andreas Geiger. “Ex-ploiting Object Similarity in 3D Reconstruction”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2015.
- [785] Yin Zhou and Oncel Tuzel. “VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection”. In: *Proc. IEEE Conf. on Com-puter Vision and Pattern Recognition (CVPR)*. 2018.
- [786] H. Zhu, K. V. Yuen, L. Mihaylova, and H. Leung. “Overview of En-vironment Perception for Intelligent Vehicles”. In: *IEEE Trans. on Intelligent Transportation Systems (TITS)* PP.99 (2017), pp. 1–18.
- [787] Ji Zhu, Hua Yang, Nian Liu, Minyoung Kim, Wenjun Zhang, and Ming-Hsuan Yang. “Online Multi-Object Tracking with Dual Match-ing Attention Networks”. In: *Proc. of the European Conf. on Computer Vision (ECCV)*. 2018, pp. 379–396.
- [788] Jianke Zhu. “Image Gradient-based Joint Direct Visual Odometry for Stereo Camera”. In: *Proc. of the International Joint Conf. on Artificial Intelligence (IJCAI)*. 2017, pp. 4558–4564.
- [789] Xizhou Zhu, Jifeng Dai, Lu Yuan, and Yichen Wei. “Towards High Performance Video Object Detection”. In: *Proc. IEEE Conf. on Com-puter Vision and Pattern Recognition (CVPR)*. 2018.

- [790] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. “Flow-Guided Feature Aggregation for Video Object Detection”. In: *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*. 2017.
- [791] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. “Deep Feature Flow for Video Recognition”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2017.
- [792] Yousong Zhu, Jinqiao Wang, Chaoyang Zhao, Haiyun Guo, and Hanqing Lu. “Scale-adaptive Deconvolutional Regression Network for Pedestrian Detection”. In: *Proc. of the Asian Conf. on Computer Vision (ACCV)*. 2016.
- [793] Zhe Zhu, Dun Liang, Songhai Zhang, Xiaolei Huang, Baoli Li, and Shimin Hu. “Traffic-sign detection and classification in the wild”. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 2110–2118.
- [794] Yueqing Zhuang, Li Tao, Fan Yang, Cong Ma, Ziwei Zhang, Huizhu Jia, and Xiaodong Xie. “RelationNet: Learning Deep-Aligned Representation for Semantic Image Segmentation”. In: *Proc. of the International Conf. on Pattern Recognition (ICPR)*. 2018, pp. 1506–1511.
- [795] Yueqing Zhuang, Fan Yang, Li Tao, Cong Ma, Ziwei Zhang, Yuan Li, Huizhu Jia, Xiaodong Xie, and Wen Gao. “Dense Relation Network: Learning Consistent and Context-Aware Representation for Semantic Image Segmentation”. In: *Proc. IEEE International Conf. on Image Processing (ICIP)*. 2018, pp. 3698–3702.
- [796] M.Z. Zia, M. Stark, B. Schiele, and K. Schindler. “Detailed 3D Representations for Object Recognition and Modeling”. In: *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)* 35.11 (Nov. 2013), pp. 2608–2623.
- [797] M.Zeeshan Zia, Michael Stark, and Konrad Schindler. “Towards Scene Understanding with Detailed 3D Object Representations”. In: *International Journal of Computer Vision (IJCV)* 112.2 (2015), pp. 188–203.
- [798] Julius Ziegler, Philipp Bender, Markus Schreiber, and Henning Lategahn. “Making Bertha Drive - An Autonomous Journey on a Historic Route”. In: *Proc. IEEE Intelligent Transportation Systems Magazine (ITSM)* 6.2 (2014), pp. 8–20.
- [799] Henning Zimmer, Andrés Bruhn, and Joachim Weickert. “Optic flow in harmony”. In: *International Journal of Computer Vision (IJCV)* 93.3 (2011), pp. 368–388.