
AI-Survey for Self-Flying Vehicles: Exploring the Challenges of Deep Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Everyone is talking about intuitive and automated transportation. An important
2 and very challenging part of this research field are autonomous unmanned aerial
3 vehicles (UAV) such as automated air taxis with a vertical take-off and landing
4 (VTOL) capability. On one hand autonomous VTOLs will redesign our personal
5 understanding of urban mobility, on the other hand automated UAVs will drastically
6 change any kind of delivery or transportation services and much more. However,
7 when studying computer vision and machine learning problems for UAVs or VTOLs
8 it becomes increasingly difficult to stay up-to-date. We provide a survey for the
9 topic of automated flights focusing on challenging Deep Learning problems with
10 a state-of-the-art overview. We give an outline of possible sensor set-ups and AI
11 based pipelines with leading results on established data sets. Finally we point out
12 currently missing investigations.

13 1 Introduction

14 Autonomous flying is a rapidly advancing application area with a lot of opportunities for Deep
15 Learning or Machine Learning based approaches. In common, two different pipelines can be
16 distinguished:

- 17 1. The mediated perception approach which semantically reasons the scene [12, 11, 24] and
18 determines the flight control decision based on it.
- 19 2. The end-to-end approach that learns the flying controls based on human behavior in and
20 end-to-end manner [16, 2, 28].

21 Fig. 1 gives an overview of both pipelines where exemplary possible applications are shown. (a)
22 SLAM is crucial for the local map and the vehicle pose within the environmental model. (b) Scene
23 Understanding is essential to interpret the environment, e.g. to detect static and dynamic objects
24 and their locations such as point wise classifications. (c) Sensor-Fusion is important to exploit the
25 strengths of the different sensor types like classification for cameras, reconstruction for Lidar or
26 dynamics for Radar. (d) End-2-End flying learns all decisions within a single network and can be
27 treated as alternative approach. Compared to other kinds of automated vehicles, Autonomous Flying
28 (AF) has specific challenges that characterize the use cases for Deep Learning:

- 29 • Scale Ambiguity: The 6DoF viewpoint ability for aerial vehicles impedes basic geometrical
30 tasks like visual depth estimation or visual reconstruction in comparison to 3DoF use cases
31 for ground robots and cars.
- 32 • Data Availability: Public data sets are rare compared to other computer vision tasks.
- 33 • Constraint Hardware: Applications have to run on a limited hardware with low energy
34 consumption.

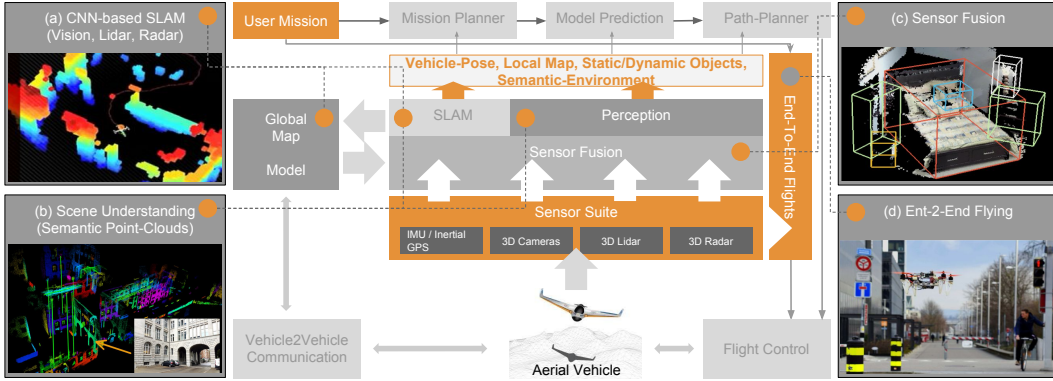


Figure 1: The principles of automated flying. The diagram outlines the state-of-the-art workflow. It all starts with the dedicated user mission. However, vehicle sensor data is essential to develop an environmental model for decision making and path planning. Several sensors like cameras, Lidars or Radars are crucial. In general two different paths are distinguished: 1. The mediated perception approach; 2. End-2-End Flying; Due to the complexity of the different tasks, leading approaches are mainly based on Machine Learning or Deep Learning, in particular Convolutions Neural Networks (CNN). (a-d) Illustrate example functions based on Deep Learning and their specific role within the pipeline [27, 4, 19, 7, 1].

35 Due to those challenging circumstances our short survey will cover an overview of public aerial data
 36 sets for specific tasks with currently leading applications. We give an overview of possible sensor
 37 setups, specific work-flows for sensor fusion and point out there strengths and weaknesses. The main
 38 part gives an overview of possible Deep Learning based applications for AF referencing exemplary
 39 state-of-the-art developments.

Aerial Data Sets							
↓ <i>Name/Task</i> →	Semantics	Objects	Odometry	Vision	Lidar	Radar	Size
Stanford Drone [19]	✗	✓(2D)	✗	✓	✗	✗	~69GB ¹
DOTA [26]	✗	✓(2D)	✗	✓	✗	✗	2806F ²
ISPRS [15]	✓(2/3D)	✗	✗	✓	✓	✗	~20GB ³
40 VisDrone2018 [30]	✗	✓(2D)	✗	✓	✗	✗	3190F ⁴
Inria Aerial [17]	✓(2D)	✗	✗	✓	✗	✗	360F ⁵
Drone Mapper	✗	✗	✗	✓	✓	✗	- ⁶
Zurich Micro [18]	✗	✗	✓(6DoF)	✓	✗	✗	~28GB ⁷
EuRoC MAV [3]	✗	✗	✓(6DoF)	✓	✗	✗	~20GB ⁸
41 Kitty [12]	✓(2D)	✓(2/3D)	✓(3DoF)	✓	✓	✗	8110F ⁹

¹StanfordD: Several video sequences with instance tracking containing 7 classes in 8 different scenes

²DOTA: 2806 images (scale invariant) with 15 different object classes.

³ISPRS: Three different scenes (Toronto, Potsdam and Vaihingen) containing Lidar and RGB images (~ 40 image pairs per scene) with Semantic Pixel Classification (6 classes)

⁴VisDrone: 3190 frames in video and image footage with object boxes and tracking instances (12 classes).

⁵InriaA: Two pixel-wise classes (building, background) covering around 810 km² in 5 different regions.

⁶DMapper: Commercial data from <https://dronemapper.com> with HD-Lidar with accompanied RGB.

⁷ZurichM: A total of 5'237'298 2D keypoint observations and 1'382'274 3D points in Zurich.

⁸EuRoC: Around 10 indoor scenes with a static laser observer for odometry estimations.

⁹Kitti: Automotive Dataset with 8110 images with 2D and 3D Multiclass (8 classes) Object-boxes using Stereo Vision and Lidar such as 3Dof odometry.

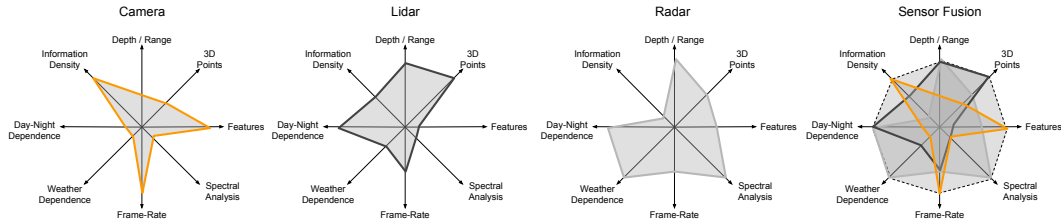


Figure 2: Sensor Fusion for Aerial Machine Learning. The figure shows four Star Plots analyzing the strengths and weaknesses of Camera, Lidar, Radar and Sensor Fusion. Individual strengths differ a lot. To benefit from all strengths Sensor Fusion is necessary. e.g. Object Classification can be trained easily using cameras [6] due to the good information density and the high value of visual features, whereas localization or reconstruction tasks benefit from Lidar sensing. Hence, 3D object detection mainly profits by fusion of cameras and Lidar, what can be proven by the Kitti leaderboard[12]. Radar has its advantages in the spectral analysis (2DFFT), i.e. it can directly measure the velocity of surrounding objects and many more tempo-spatial features. On the other hand Radar is resistant to weather or day/night conditions. Questionable is therefore rare usage of Radar data for ML in the domain of automated Flights.

42 2 Learning with Aerial Data

43 2.1 Public data sets

44 Different kinds of aerial data sets were established as it became important solving aerial computer
 45 vision tasks. To the best of our knowledge we summarized the most influential data sets in Tab. 1. At
 46 the moment, the main focus of research is aerial perception (e.g. multi-class object detection and
 47 tracking) and localization (e.g. odometry prediction) predominantly using camera inputs. All eight
 48 mentioned aerial data sets use cameras, only two use Lidar and none of them provide Radar ground
 49 truth. For comparison we mention the most comprehensive automotive data set Kitti [12]. Even Kitti
 50 does not provide public Radar data. We must conclude missing ground truth 3D boxes for aerial data
 51 and any kind of semantic Lidar annotations. Additionally, no one uses cameras with a large Field of
 52 View (FoV) or a stitched construction to cover 360 degrees of the vehicle.

53 2.2 3D Environmental Sensing

54 Lidar, Camera and Radar have different strengths and weaknesses that are important for solving Aerial
 55 Deep Learning Tasks. For a robust solution using Machine Learning Sensor Fusion is inevitable.
 56 Fig. 2 points out the advantages of Sensor Fusion. To our surprise, Radar is rarely used in perceptual
 57 fusion concepts, although it has standalone properties, like spectral analysis or weather resistance.
 58 We recommend a full fusion concept. Since, high quality data is inevitable for any kind of machine
 59 learning approach, we summarize the following Deep Learning challenges for our survey:

- 60 • Public Radar data (2D, 3D or Semantic ground truth) is missing.
- 61 • Additional ground truth for Lidar is (2D, 3D or Semantic) missing.
- 62 • Cameras are mainly used with a small FoV not covering 360 degrees.
- 63 • Highly redundant (minimum 3 sensors types) data sets are missing

64 3 Deep Learning based Autonomous Flying

65 Fig. 1 shows the basic principle of AF. We point out opportunities using DL in four different algorithm
 66 groups in the field of DL, whereas basic function (e.g. Semantic Segmentation) can be part of several
 67 groups (e.g. Semantic Maps):

68 **3.1 Localization, Mapping and Reconstruction**

69 **3.1.1 Visual Odometry**

70 Dense Tracking and Mapping (DTAM) [21] was the first published method estimating odometry
 71 with simultaneous mapping. Here, a key frame based minimization of the photo-metric error was
 72 introduced. The following cost function was used:

$$\mathbf{C}_r = \frac{1}{\|I(r)\|} \sum_{m \in I(r)} \|\mathbf{I}_r(\mathbf{u}) - \mathbf{I}_m(\mathbf{v})\|. \quad (1)$$

73 Currently, still traditional cost minimization is state-of-the-art. Recently, Direct Sparse Odometry
 74 (DSO) was published by Engel et al. [9] with leading results on Kitti [12]. The global cost takes
 75 geometric attributes (lens distortion, exposure time) is designed as:

$$\mathbf{C}_r = \frac{1}{\|I(r)\|} \sum_{m \in I(r)} \|\mathbf{I}_r(\mathbf{u}) - b_r - \frac{t_r e^{a_r}}{t_m e^{a_m}} \mathbf{I}_m(\mathbf{v}) - b_m\|. \quad (2)$$

76 Recently, Delmerico et al. [5] published a comprehensive UAV benchmark for traditional visual
 77 odometry estimation using the EuRoC [3] (6Dof, see section 2.2). The ablation study focuses on
 78 real-time capacity and accuracy. Most accurate method ODROID is based on key frame based
 79 optimization like DTAM (1).

80 **3.1.2 Unsupervised Odometry and Depth Estimation**

81 To our surprise, Deep Learning is currently not dominating odometry challenges. However, promising
 82 results are recently published. GeoNet by Yin et al. [29] minimizes an additive cost function that is
 83 completely consisting of geometric unsupervised terms, i.e. a joint estimation of monocular depth,
 84 optical flow and egomotion. The overall cost is used to train a combination of CNNs. The full
 85 pipeline can be divided into a Rigid-Structure-Decoder such as a Non-Rigid-Motion Localizer. The
 86 loss is composed by:

$$\mathcal{L} = \sum \sum [\mathcal{L}_{rw} + \mathcal{L}_{ds} + \mathcal{L}_{fw} + \mathcal{L}_{fs} + \mathcal{L}_{gc}] \quad (3)$$

87 \mathcal{L}_{rw} (warping loss) and \mathcal{L}_{ds} (depth smoothness) define the rigid decoder. \mathcal{L}_{fw} , \mathcal{L}_{fs} and \mathcal{L}_{gc} describe
 88 the non-rigid motion localizer. The method outperforms significantly ORB-Slam on single Kitti
 89 Traces for trajectory accuracy (RMSE) and demonstrates the power of unsupervised Deep Learning.

90 **3.1.3 Competitive Learning of Odometry and Depth**

91 Recently, generative adversarial networks (GAN) outperformed lots of generative computer vision
 92 tasks. Milz et al. [20] used a cGAN doing Image-to-Image translation, i.e. Pix2Pix by Isola et al.
 93 [14], performing aerial depth estimation using Lidar ground truth. The overall loss minimizes the
 94 following term:

$$\mathcal{L} = \mathbb{E}_{x,y} \{\log(D(x,y))\} + \mathbb{E}_{x,z} \{\log(1 - D(x,G(x,z)))\} + \lambda \cdot \mathbb{E}_{x,y,z} \{\|y - G(x,z)\|_1\} \quad (4)$$

95 The method is composed by a generative G and a descriptive network D (see Fig. 3) In order to
 96 create more and more accurate data, the loss of G is reduced, whereas a training step of D results in
 97 an increase of the partial loss $(1 - D)$ ideally. Hence, a competitive loss is the result. The advantage
 98 of the approach is, that the overall loss design is learned by the network itself.

99 Ranjan et al. goes a step further and combines a competitive such as a collaborative loss to an overall
 100 cost, which is composed by camera motion, monocular depth, optical flow such as motion estimation

⁹See reference [21] for detailed explanation of (1) and [9] for a detailed explanation of (2)

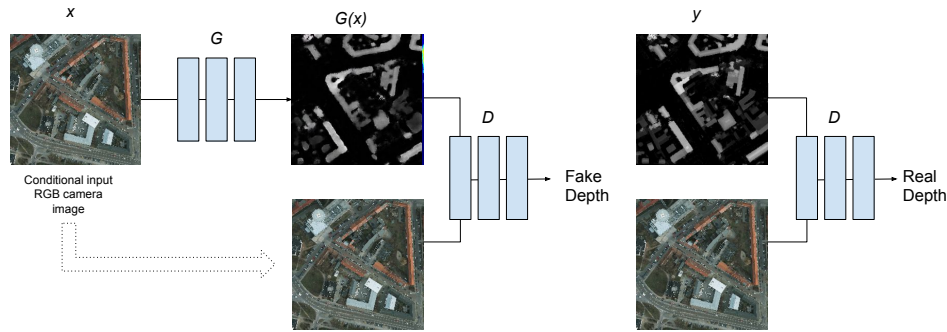


Figure 3: Competitive monocular Depth estimation using conditional GANs. The figure shows Milz et al. [25] implementation of the cGAN playing the minimax game. A generator G is used to create a fake image $G(x)$ (Depth reconstruction) based on the conditional input camera image x . The discriminator D tries to distinguish between a real Depth map $D(y)$ and fake image $D(G(y))$. The method shows promising results on the ISPRS data set[15].

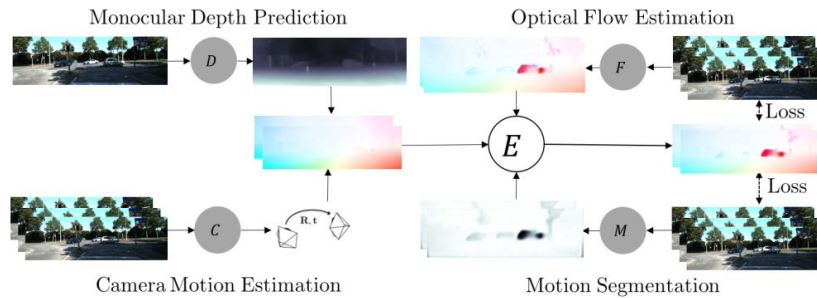


Figure 4: Collaborative and competitive odometry estimation and reconstruction. The figure is taken from [23] outlining the basic idea of the overall loss with promising results on Kitti.

101 3.1.4 Point-Cloud based SLAM using CNN based Semantic Points

102 SegMap by Dube et al. [8] uses Lidar based Point-Clouds to perform overall SLAM. The clue is an
 103 feature based global optimization function that is performed on semantic point clouds. The semantic
 104 point cloud classification is performed by a CNN. The model reduces drastically the number of tracked
 105 features and improves accuracy. The approach yields competitive results on Kitti (see. Fig.)

106 3.2 Perception and Scene Understanding

107 3.2.1 Visual Object Detection

108 The DOTA leader board [26] is good signpost for modeling visual object detectors. The currently
 109 leading approach is a mask R-CNN by He et al. [13]. The mask R-CNN performs instance object
 110 segmentation on DOTA with an overall mAP of 0.762.

Statistic	KITTI	Powerplant	Foundry
Duration (s)	114	850	1086
Number of robots	5	3	2
Number of segmented local cloud	557	758	672
Average number of segments per cloud	42.9	37.0	45.4
Bandwidth for transmitting local clouds (kB/s)	4814.7	1269.2	738.1
Bandwidth for transmitting segments (kB/s)	2626.6	219.4	172.2
Bandwidth for transmitting descriptors (kB/s)	60.4	9.5	8.1
Final map size with the SegMap descriptor (kB)	386.2	181.3	121.2
Number of successful localizations	113	27	85

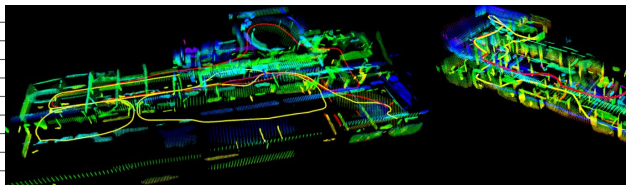


Figure 5: SegMap by Deube et al. performs localization and mapping based on Semantic Point Clouds sensed by Lidars. Results in the left table are promising ([7, 12]). The right area outlines qualitative odometry and mapping predictions by Dube et al.

Classes	IoU Aerial GAN
Impervious surfaces	79.4%
Building	87.1%
Low vegetation	67.3%
Tree	70.3%
Car	24.1%
Clutter/background	30.7%
Mean IoU	59.8%

Figure 6: IoU for the aerial GANeration approach by Milz et al. [20] in the domain of image to semantic segmentation translation (ISPRS dataset[15]). The right part shows qualitative results.

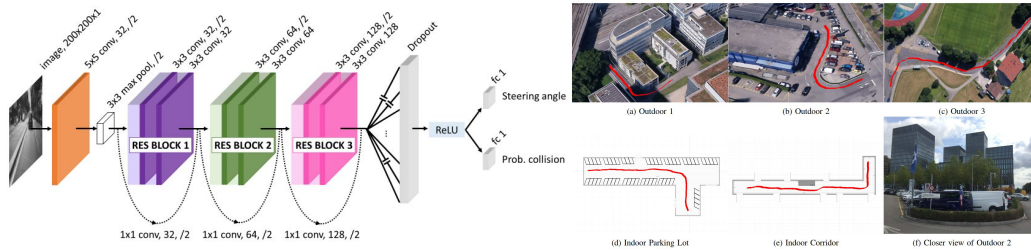


Figure 7: Dronet by Loquercio et al. (Parts of the Figure are taken from [16]). The left part shows the Resnet [10] architecture which is directly trained by the movement of observed agents in the urban area (cars, bicycles). The right part (a-e) shows qualitative movement results in different scenes.

111 3.2.2 Semantic Segmentation

112 Aerial Semantic Segmentation was recently performed by Milz et al. using the ISPRS data set.
 113 Similar to section 3.1.3, the approach uses a cGAN to model the task as Image-to-Image translation
 114 problem. The results on the ISPRS are state-of-the-art. In Section 3.1.4 we have already referenced
 115 to semantic point cloud classification, which could be implicitly used for SLAM. As shown by Qi et
 116 al. [22] the overall idea is to approximate a symmetric f function on the point-set $x_{1..n}$ by applying
 117 local function h to get transformed elements of the data (5). This approximation is directly used in
 118 the overall loss to get a geometrical assessment and therefore a semantic segmentation of the points.

$$f(\{x_1, \dots, x_n\}) \approx g(h(x_1), \dots, h(x_n)) \quad (5)$$

119 3.3 Prediction, Planning and End-to-End flying

120 Prediction and planning for Aerial Vehicles are currently rarely solved using Deep Learning. Loquer-
 121 cio et al.[16] proposed an End-to-End approach imitating the movement of cars and bicycles using
 122 UAVs in Urban areas. The concept uses the ground truth motion of real cars/bicycles to train a CNN
 123 directly. The models architecture and qualitative results are shown in Fig.7

124 4 Conclusion

125 We have shown a compressed survey for AI based Autonomous Flights using Deep Learning for
 126 solving modular Tasks. We note, that DL has arrived in many parts like SLAM, perception, prediction
 127 or End-2-End flying. However, currently the main challenge is a comprehensive sensor redundant data
 128 set with three-dimensional ground truth (e.g. point semantics). To benefit from the strength of several
 129 sensor types. To our surprise, the main research focuses on cameras. Consequently, complex and
 130 comprehensive visual models are developed to perform tasks like reconstruction or depth estimation,
 131 e.g. competitive learning (section 3), which could be taken directly from Lidar or Radar. Hence, we
 132 highly recommend the usage of Lidar, Cameras and Radar.

References

- 133
- 134 [1] Aseem Behl, Omid Hosseini Jafari, Siva Karthik Mustikovela, Hassan Abu Alhaija, Carsten Rother, and
135 Andreas Geiger. Bounding boxes, segmentations and object coordinates: How important is recognition
136 for 3d scene flow estimation in autonomous driving scenarios? In *International Conference on Computer
137 Vision (ICCV)*, 2017.
- 138 [2] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon Goyal,
139 Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving
140 cars. *arXiv preprint arXiv:1604.07316*, 2016.
- 141 [3] Michael Burri, Janosch Nikolic, Pascal Gohl, Thomas Schneider, Joern Rehder, Sammy Omari, Markus
142 Achtelik, and Roland Siegwart. The euroc micro aerial vehicle datasets. 35, 01 2016.
- 143 [4] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for
144 autonomous driving. *CoRR*, abs/1611.07759, 2016.
- 145 [5] Jeffrey A. Delmerico and Davide Scaramuzza. A benchmark comparison of monocular visual-inertial
146 odometry algorithms for flying robots. *2018 IEEE International Conference on Robotics and Automation
147 (ICRA)*, pages 2502–2509, 2018.
- 148 [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image
149 Database. In *CVPR09*, 2009.
- 150 [7] Renaud Dubé, Andrei Cramariuc, Daniel Dugas, Juan Nieto, Roland Siegwart, and Cesar Cadena. SegMap:
151 3d segment mapping using data-driven descriptors. In *Robotics: Science and Systems (RSS)*, 2018.
- 152 [8] Renaud Dubé, Daniel Dugas, Elena Stumm, Juan Nieto, Roland Siegwart, and Cesar Cadena. Segmatch:
153 Segment based place recognition in 3d point clouds. In *IEEE International Conference on Robotics and
154 Automation (ICRA)*, pages 5266–5272. IEEE, 2017.
- 155 [9] J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. In *arXiv:1607.02565*, July 2016.
- 156 [10] Clément Farabet, NYU EDU, Camille Couprie, Laurent Najman, and Yann LeCun. Scene parsing with
157 multiscale feature learning, purity trees, and optimal covers.
- 158 [11] Daniel Gehrig, Henri Rebecq, Guillermo Gallego, and Davide Scaramuzza. Asynchronous, photometric
159 feature tracking using events and frames. *CoRR*, abs/1807.09713, 2018.
- 160 [12] Andreas Geiger, Martin Lauer, Christian Wojek, Christoph Stiller, and Raquel Urtasun. 3d traffic scene
161 understanding from movable platforms. *IEEE transactions on pattern analysis and machine intelligence*,
162 36(5):1012–1025, 2014.
- 163 [13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. *CoRR*, abs/1703.06870,
164 2017.
- 165 [14] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional
166 adversarial networks. *arxiv*, 2016.
- 167 [15] K. Khoshelham, L. Díaz Vilariño, M. Peter, Z. Kang, and D. Acharya. The isprs benchmark on indoor
168 modelling. *ISPRS - International Archives of the Photogrammetry, Remote Sensing and Spatial Information
169 Sciences*, XLII-2/W7:367–372, 2017.
- 170 [16] Antonio Loquercio, Ana I. Maqueda, Carlos R. del-Blanco, and Davide Scaramuzza. Dronet: Learning to
171 fly by driving. *IEEE Robotics and Automation Letters*, 3(2):1088–1095, 2018.
- 172 [17] Emmanuel Maggiori, Yuliya Tarabalka, Guillaume Charpiat, and Pierre Alliez. Can semantic labeling
173 methods generalize to any city? the inria aerial image labeling benchmark. In *IEEE International
174 Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE, 2017.
- 175 [18] András L Majdik, Charles Till, and Davide Scaramuzza. The zurich urban micro aerial vehicle dataset. *Int.
176 J. Rob. Res.*, 36(3):269–273, March 2017.
- 177 [19] Huynh Manh and Gita Alaghband. Scene-1stm: A model for human trajectory prediction. *CoRR*,
178 abs/1808.04018, 2018.
- 179 [20] Stefan Milz. Aerial ganeration: Towards realistic data augmentation using conditional gan. In *2nd
180 International Workshop on Computer Vision for UAVs*, Sept. 2018.

- 181 [21] Richard A. Newcombe, Steven J. Lovegrove, and Andrew J. Davison. Dtam: Dense tracking and mapping
182 in real-time. In *Proceedings of the 2011 International Conference on Computer Vision, ICCV '11*, pages
183 2320–2327, Washington, DC, USA, 2011. IEEE Computer Society.
- 184 [22] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point
185 sets for 3d classification and segmentation. *CoRR*, abs/1612.00593, 2016.
- 186 [23] Anurag Ranjan, Varun Jampani, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J. Black. Adversarial
187 collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation.
188 *CoRR*, abs/1805.09806, 2018.
- 189 [24] Davide Scaramuzza and Roland Siegwart. Appearance-guided monocular omnidirectional visual odometry
190 for outdoor ground vehicles. *IEEE transactions on robotics*, 24(5):1015–1026, 2008.
- 191 [25] Min Wang, Baoyuan Liu, and Hassan Foroosh. Design of efficient convolutional layers using single
192 intra-channel convolution, topological subdivisoning and spatial “bottleneck” structure.
- 193 [26] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge J. Belongie, Jiebo Luo, Mihai Datcu, Marcello
194 Pelillo, and Liangpei Zhang. DOTA: A large-scale dataset for object detection in aerial images. *CoRR*,
195 abs/1711.10398, 2017.
- 196 [27] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box
197 estimation. *CoRR*, abs/1711.10871, 2017.
- 198 [28] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. End-to-end learning of driving models from
199 large-scale video datasets. *arXiv preprint arXiv:1612.01079*, 2016.
- 200 [29] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera
201 pose, 2018.
- 202 [30] Pengfei Zhu, Longyin Wen, Xiao Bian, Ling Haibin, and Qinghua Hu. Vision meets drones: A challenge.
203 *arXiv preprint arXiv:1804.07437*, 2018.