

# Dynamic Object Detection in Range data using Spatiotemporal Normals

Raphael Falque<sup>†</sup>, Cedric Le Gentil<sup>†</sup>, and Fouad Sukkar

Robotics Institute at the University of Technology Sydney

raphael.guenot-falque@uts.edu.au, cedric.legentil@uts.edu.au, fouad.sukkar@uts.edu.au

<sup>†</sup> Both authors contributed equally to this paper

## Abstract

On the journey to enable robots to interact with the real world where humans, animals, and unpredictable elements are acting as independent agents; it is crucial for robots to have the capability to detect dynamic objects. In this paper, we argue that the detection of dynamic objects can be solved by computing the spatiotemporal normals of a point cloud. In our experiments, we demonstrate that this simple method can be used robustly for LiDAR and depth cameras with performances similar to the state of the art while offering a significantly simpler method.

## 1 Introduction

Automation is a cornerstone of modern societies. For example, one can think of the automatic production lines of most car manufacturers. However, to this day, it is still rare to encounter autonomous robots outside of controlled environments such as large-scale factories, warehouses, etc. A limiting factor is the ability of robots to adapt to uncontrolled situations or environments. Such ability is crucial to ensure safe operation especially when humans are present in the environment. Accordingly, as perception is the root of any autonomous system, the ability to detect dynamic objects in the surroundings of a robot is an essential step toward the democratisation of robots in society [Cadena *et al.*, 2016]. Additionally, in the context of localisation and mapping, it has been shown that the detection of dynamic elements leads to nearly 40% less odometry error [Pfreundschuh *et al.*, 2021]. As illustrated in Figure 1, this paper presents a method for dynamic object detection in data collected with range sensors such as lidars or color-depth (RGBD) cameras.

With the rise of the deep learning dominance in the field of computer vision [Krizhevsky *et al.*, 2012], the research community has naturally focused on the semantic

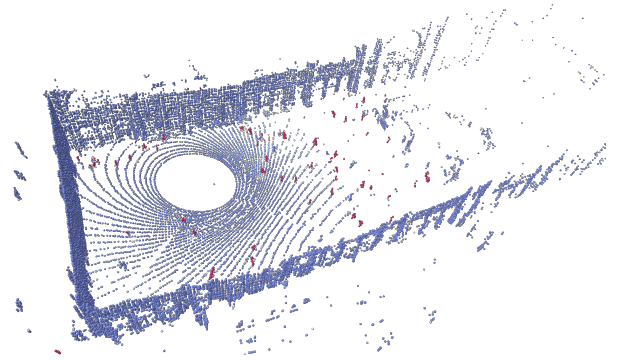


Figure 1: Illustration of the proposed method for detecting dynamic pedestrians in a large outdoor environment. The points that are detected as dynamic are shown in red and the static points are in blue.

segmentation of data as a proxy for classifying dynamic objects. The most straightforward approach is to feed the sensor output into a 2D-Convolutional Neural Network (CNN) and learn the dynamic components from large labelled datasets. [Chen *et al.*, 2019] proposed to use such method with the CNN Rangenet++ [Milioto *et al.*, 2019] and the KITTI Vision Benchmark dataset [Geiger *et al.*, 2012]. This work was later extended by stacking several consecutive scans in the network input and by benchmarking different CNN architectures [Chen *et al.*, 2021]. In [Dai *et al.*, 2018], lidar data is upsampled into image-like data to fine-tune a CNN object detector trained with standard images [Redmon and Farhadi, 2017]. As an alternative, in cases where the depth information is acquired with RGBD cameras, the color (RGB) information can be used for the semantic segmentation and then transferred onto the map Region-based Convolutional Neural Network (R-CNN) architecture [Henein *et al.*, 2020].

While these methods allow the detection of dynamic objects, they are incapable of distinguishing moving objects from static objects such as parked cars. For this reason, a part of the literature is studying the detection

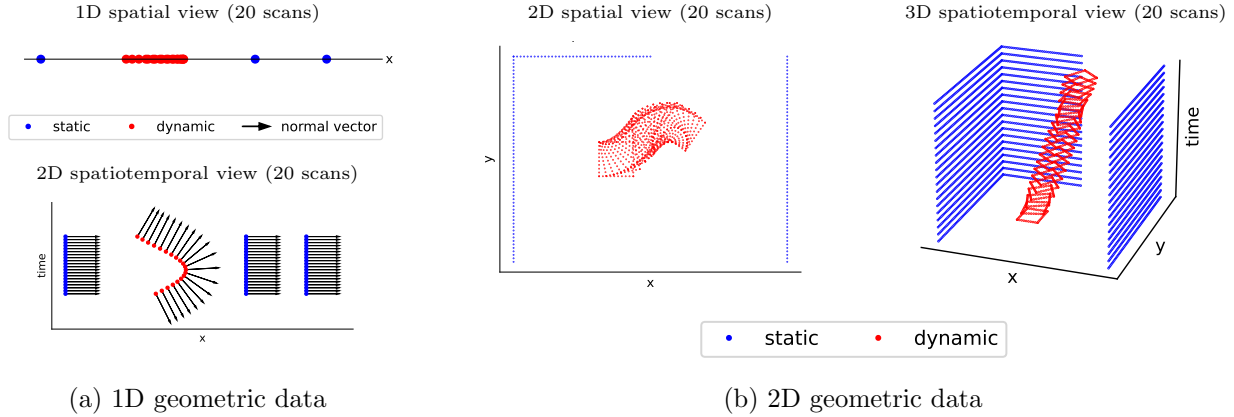


Figure 2: Spatiotemporal view of 1D (a) and 2D (b) geometric data in the presence of a dynamic object (20 range scans through time). It illustrates the link between the spatiotemporal normal vectors and the velocity of an object: for static objects the temporal component of the normals is null (the normals have been omitted in (b) for the sake of readability).

of dynamic objects without consideration of the semantic information.

More recently, several approaches are trying to promote the idea of detecting dynamic objects in the map space in contrast to the sensor space. These approaches generally perform better as the map aggregates long-term information and the classification is reduced to finding areas that were once *free* and are now *occupied* (i.e., all the newly occupied areas are parts of dynamic objects). A recent example of this strategy has been proposed by Schmid et al. where the Voxelblox framework [Oleynikova et al., 2017] has been extended for dynamic object detection [Schmid et al., 2023]. Voxelblox builds an efficient voxelization of the map with a hash map and Dynablox adds the mapping of the free-space confidence. Similarly, Mersch et al. proposed to use a sparse 4D (i.e., [X, Y, Z, time]) CNN to predict dynamic object and fuses the predictions of unoccupied space into the map with a Bayesian filter [Benedikt et al., 2023]. Given the great performance of these occupancy-based approaches, they can be used as an offline tool for labelling large datasets which can later be used for the training of real-time deep learning architectures [Pfreundschuh et al., 2021].

In this paper, we propose a simple method based on the concept of spatiotemporal normals vectors for dynamic object detection. The proposed method requires no training, has no requirement for any specific framework, and has a very limited number of parameters. More precisely, the contributions of this work are as follows:

- A simple yet effective method to detect dynamic objects in range data based on the analysis of the normal vectors in the spatiotemporal space.

- The open-source real-time implementation of our approach is available at [https://github.com/UTS-RI/dynamic\\_object\\_detection](https://github.com/UTS-RI/dynamic_object_detection).

Please note that the proposed method, similarly to [Schmid et al., 2023], requires the input 3D scans to be already registered in a common referential frame. The registration is outside the scope of this paper but can be performed simply using known kinematics if the sensor is mounted on a robot arm, or estimated with methods like [Le Gentil et al., 2021] or [Campos et al., 2021] in the context of localisation and mapping.

## 2 Method

### 2.1 Motivation and overview

In this paper, we consider a stream of range data through time (e.g., from a lidar or a depth camera). The goal is to classify the individual points of the incoming data as belonging to a dynamic object or not. A dynamic object is defined as an object with a non-null velocity in an earth-fixed reference frame  $\mathcal{F}_W$ . Accordingly, the proposed classification is agnostic to the nature of the objects present in the environment.

This work introduces a *dynamic score* estimated for each of the incoming points based on the computation of spatiotemporal normal vectors. In Figure 2, examples are provided to give an intuition on how the spatiotemporal normals relate to the points’ velocities. In the 1D scenario, the normal is directly linked to the velocity. In higher dimensions, the temporal component of the spatiotemporal normal corresponds to the projection of the actual velocity on the spatial normal to the surface.

An overview of the proposed pipeline is shown in Figure 3. The input of our method consists of point clouds

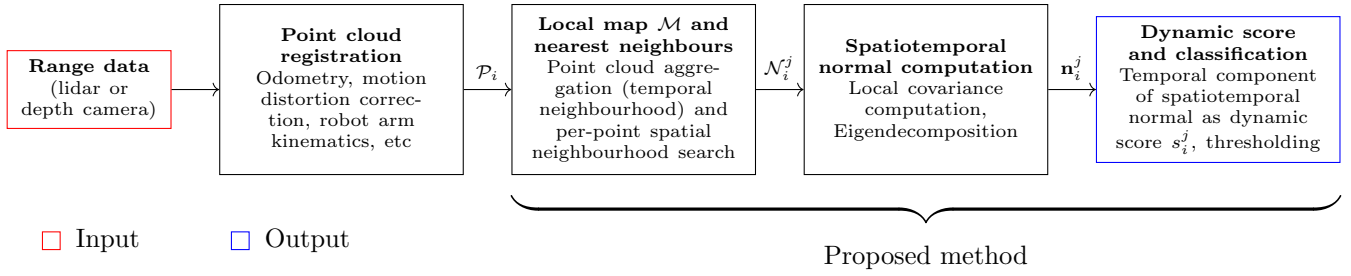


Figure 3: Overview of the proposed method for dynamic object detection in range data. The method relies on registered point clouds aggregated in local maps. The final *dynamic score* relies on the computation of the spatiotemporal normal vector for each of the points.

registered in a fixed reference frame. In the context of mapping, the registration can be performed based on odometry information (lidar, visual, etc.) while the kinematics of a robot arm can be leveraged for cobotics-manipulation scenarios. Similarly to standard normal vector estimation, the spatiotemporal normals are computed through a Principal Component Analysis (PCA) of the point’s neighbourhood. More specifically, the Eigendecomposition of the local covariance of the data is performed and the Eigenvector corresponding to the smallest Eigenvalue is used as the normal estimate. As discussed above, such a normal vector in the spatiotemporal space provides information regarding the point’s velocity. The dynamic score is then defined as the time component of the spatiotemporal normal. The following subsection presents a formal derivation of the proposed dynamic score.

## 2.2 Dynamic score via spatiotemporal

Let us consider an undistorted point cloud  $\mathcal{P}_i$  ( $i \in 1, \dots, N$ ) registered in a unique fixed frame  $\mathcal{F}_W$ . We denote  $\mathbf{p}_i^j$  the  $j^{th}$  point in  $\mathcal{P}_i$  and  $t_i^j$  the associated timestamp. The symbol  $\mathcal{M}$  refers to the aggregation of all the point clouds  $\mathcal{P}_i$  so that  $k \leq i \leq l$  (with  $k$  and  $l$  derived from the method’s parameters). To provide a dynamic score to a point  $\mathbf{p}_i^j$ , we first perform a neighbourhood search in  $\mathcal{M}$ . We define  $\mathcal{N}_i^j$  as the set of neighbour points. The local covariance  $\text{cov}_i^j$  is computed as

$$\text{cov}_i^j = \frac{1}{\|\mathcal{N}_i^j\|} \sum_{\mathbf{p}_u^v \in \mathcal{N}_i^j} \left( \begin{bmatrix} \mathbf{p}_u^v \\ t_u^v \end{bmatrix} - \mathbf{m}_i^j \right) \left( \begin{bmatrix} \mathbf{p}_u^v \\ t_u^v \end{bmatrix} - \mathbf{m}_i^j \right)^\top \quad (1)$$

with

$$\mathbf{m}_i^j = \frac{1}{\|\mathcal{N}_i^j\|} \sum_{\mathbf{p}_u^v \in \mathcal{N}_i^j} \begin{bmatrix} \mathbf{p}_u^v \\ t_u^v \end{bmatrix}. \quad (2)$$

The dynamic score  $s_i^j$  is defined as the absolute value of the temporal component of the Eigenvector associated with the smallest Eigenvalue of  $\text{cov}_i^j$ . A low dynamic

score corresponds to low velocity in the 3D space. Thus, the point classification as static or dynamic is performed by thresholding  $s_i^j$ .

## 2.3 Implementation

The proposed method has been implemented to run in real-time in C++ and written as a Robot Operating System (ROS) node that *subscribes* to a point cloud  $\mathcal{P}_{in}$  and *publishes* two point clouds, one dynamic  $\mathcal{P}_{out}^{dyn}$  and one static  $\mathcal{P}_{out}^{sta}$ .

The implementation first performs a voxel grid down-sampling, with a voxel size  $d_v$ , using a hash map [Nießner *et al.*, 2013] and aggregates the point clouds in the map  $\mathcal{M}$ . The map is defined as a sliding window over the last  $2N + 1$  point clouds, where similarly to a First In, First Out (FIFO) pile, each time a new point cloud is aggregated into  $\mathcal{M}$ , the oldest point cloud is then removed. A kd-tree [Blanco and Rai, 2014] is built with  $\mathcal{M}$  and for each of the points in the  $(N + 1)^{th}$  last point cloud, a query of the neighbours within the radius  $d_r$  is then performed to compute the spatiotemporal normal using PCA. The dynamic points are then detected by using a threshold  $thr$  over the temporal component of the computed normal.

Eventually, our method relies on only four parameters: the voxel size  $d_v$ , the number of clouds  $N$ , the radius for the neighbour search  $d_r$ , and the dynamic score threshold  $thr$ . All these parameters have a physical meaning and are easy to tune depending on the application requirements. The parameters on the voxel size,  $d_v$ , and the number of aggregated point clouds,  $N$ , have a direct relationship with the computational time and can be used to balance sensors with high resolution and high framerates. The radius used for finding the points’ neighbours,  $d_r$ , has a marginal impact on the computational time and should be chosen with respect to the scale of the scanned environment. Finally, the threshold,  $thr$ , should be chosen depending on how important it is to remove dynamic points.

### 3 Experiments

The experiment section is separated into two parts, first, we test the proposed method on an established dataset which has been used consistently in the recent state of the art [Pfreundschuh *et al.*, 2021; Schmid *et al.*, 2023], secondly, we show a practical study case on how the dynamic coefficient can be used in human-robot collaboration scenario. Different parameters are used for each scenario as the experiments are performed with sensors that have radically different ranges (LiDAR and RGBD camera).

#### 3.1 The DOALS dataset

We applied our method on the Urban Dynamic Objects LiDAR (DOALS) dataset [Pfreundschuh *et al.*, 2021]. The dataset contains two scenarios: the first one is made of a simulated small town with moving elements (e.g., cars, planes, pedestrians, animals, etc.) that provide proper ground truth for dynamic elements; the second scenario is from real-world scenes recorded around the Zurich metropolitan area. An Ouster OS1 64 LiDAR, producing clouds with 131.072 points at 10 Hz, was used for the data collection and manual annotation was provided for some of the timestamps. In this experiment, we use the following parameters:  $N = 10$ ,  $d_r = 0.3$ , and  $thr = 0.25$  and the voxel size  $d_v$  is set proportionally to the scale of the point cloud (i.e., the diagonal of its bounding box) as follows:

$$d_v = \frac{\text{scale}(\mathcal{P}_i)}{600}. \quad (3)$$

A qualitative evaluation of our method is shown in Figure 4 where we compare the dynamic coefficient to the prediction of Dynablox and the ground truth annotations.

Furthermore, to provide a quantitative evaluation, we compare the Intersection over Union (IoU) for LC Free space [Modayil and Kuipers, 2008], Dynablox, and the proposed method in Table 1 on all real scenes from the DOALS dataset<sup>1</sup>. The evaluation is performed on the original point cloud resolution for Dynablox and LC Free Space, while we upsample our dynamic prediction to the original point cloud resolution. The upsampling from the downsampled cloud to the full-resolution cloud is performed by a local search in the vicinity of the dynamic points. More specifically, for all dynamic points in the downsampled cloud, we search for their neighbours within a radius of 0.5m in the full-resolution cloud and define them as dynamic.

The results from Table 1 show that the proposed method has performances relatively similar to the state

Table 1: Quantitative evaluation of the Intersection over Union (IoU) [%] between LC Free space [Modayil and Kuipers, 2008], Dynablox [Schmid *et al.*, 2023], and the proposed method. A limit is applied to the range at 20m. SV, HG, and ND stand for the Shopville, Niederdorf, and Hauptgebaeude scenes respectively. We did not manage to reproduce exactly the same performances as the original paper for Dynablox which were slightly higher.

	Station	SV	HG	ND	all
LC Free space	0.49	0.32	0.25	0.18	0.31
Dynablox	0.81	0.84	0.82	0.81	0.82
ours	0.80	0.81	0.85	0.76	0.81

of the art while alleviating the requirement for complex frameworks.

#### 3.2 Study case: human robot interaction

Unlike traditional robotic arms, collaborative robots (cobots) are designed to be operating with humans in their surroundings. There are typically four types of cobot Human-Robot Interactions (HRI): physically separated, coexistence, cooperation, and collaboration [Guertler *et al.*, 2023]. For the latter three cases, the cobot must react and adapt to the human to safely carry out tasks. In the case of coexistence and cooperation, a common control strategy is to slow down the cobot based on its proximity to any humans [Marvel, 2013; Villani *et al.*, 2018]. For collaboration, our method could be used in conjunction with dynamic motion planning algorithms [Likhachev *et al.*, 2005; Ferguson and Stentz, 2006; Alwala and Mukadam, 2021] for detecting changes in the environment and enabling the cobot to react in real-time to these changes.

In this study case, we consider a shared workspace between a robot arm and a human such as the one shown in Figure 5(a). We demonstrate how an embedded perception system can allow a robot arm to be *aware* of its surroundings by detecting moving elements in the workspace. Our experimental setup consists of a UR5 robot arm equipped with an RGBD Intel Realsense d435 camera on its end effector. The arm moves according to randomly selected 6-DoF waypoints above a table containing multiple objects. While the robot is moving, a human operator displaces some of the objects in the workspace.

The RGBD camera point clouds are projected into the robot’s base referential frame according to the arm’s kinematics. An example of raw data and estimated dynamic score are shown in Figure 5(b) and (c). The parameters for this experiment are set as  $N = 4$ ,  $d_r = 0.3$ ,  $thr = 0.01$  and  $d_v = \frac{\text{scale}(\mathcal{P}_i)}{100}$ . Additionally, the real-time

<sup>1</sup>We did not manage to run Dynablox on the virtual dataset.

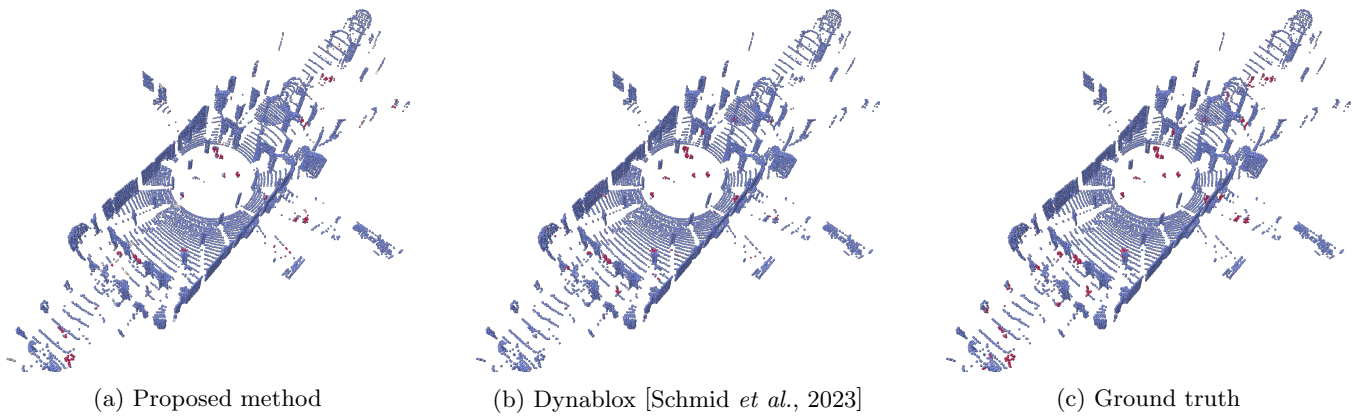


Figure 4: Comparison between the proposed method (a), Dynablox (b) and the ground truth labels (c) on the *Hauptgebaeude* scene from the Urban Dynamic Objects LiDAR (DOALS) Dataset [Pfreundschuh *et al.*, 2021].

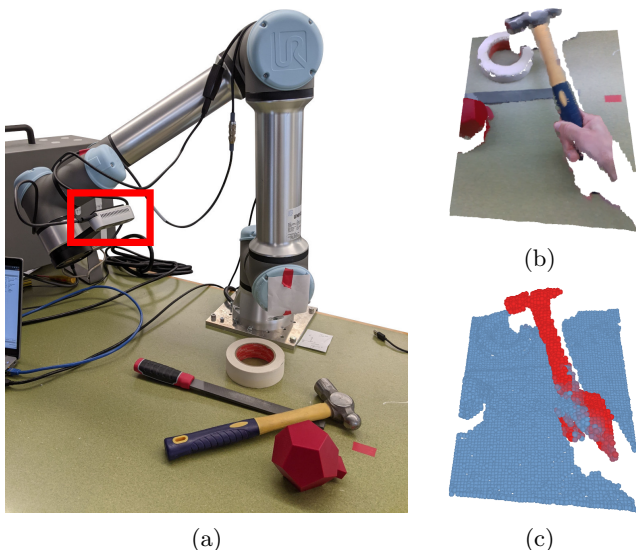


Figure 5: HRI experiment with a shared workspace between a human operator and a robot arm. The real-time detection of the dynamic motion can help the robot to avoid harming the human operator. The shared workspace is shown in (a), with the scanned objects, the robot arm, and the mounted RGBD Realsense camera. The capture RGBD data are displayed in (b), and the dynamic classification of the points is shown in (c) with the dynamic points in red and the static points in blue.

detection of the moving elements is showcased in the attached video where both the registered point clouds and their associated dynamic scores are displayed.

## 4 Limitations

One of the major limitations of the proposed approach is the difficulty of detecting dynamic points close to the ground as part of the ground is picked in the neighbour search and changes the direction of the spatiotemporal normal. Such a problem could be avoided by removing

the ground which is a common approach in the literature [Petrovskaya and Thrun, 2009; Postica *et al.*, 2016; Arora *et al.*, 2021; Arora *et al.*, 2023].

In the current implementation, one of the main limitations for safety applications consists of the delay between the dynamic predictions and the data stream. More specifically, given a sensor that generates depth measurements with a frequency  $f$ , the delay is  $\frac{N}{f}$  seconds. This limitation can be mitigated by increasing the framerate from the sensor or by limiting the number of clouds used for the dynamic score prediction.

## 5 Conclusion

In this paper, we propose a simple yet effective method for detecting dynamic elements. Our method is on par with the state of the art in terms of performance while still keeping an extremely simple formulation which is a guarantee for robustness. Therefore thanks to its simplicity we do believe that the proposed approach can be used as a stepping stone for more complex frameworks. A typical application could be to use the dynamic score in the feature space of learning algorithms for end-to-end frameworks.

In future work, we are planning to explore the coupling of the dynamic perception system with the control part of a robot arm, to provide a safer and *aware* interaction for HRI. Furthermore, we want to integrate the dynamic object detection into a full Simultaneous Localisation And Mapping (SLAM) framework for the generation of clean maps.

## 6 Acknowledgment

Cedric Le Gentil is supported by the Australian Research Council Discovery Project under Grant DP210101336.



## References

- [Alwala and Mukadam, 2021] Kalyan Vasudev Alwala and Mustafa Mukadam. Joint sampling and trajectory optimization over graphs for online motion planning. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4700–4707. IEEE, 2021.
- [Arora et al., 2021] Mehul Arora, Louis Wiesmann, Xieyuanli Chen, and Cyrill Stachniss. Mapping the static parts of dynamic scenes from 3d lidar point clouds exploiting ground segmentation. In *2021 European Conference on Mobile Robots (ECMR)*, pages 1–6. IEEE, 2021.
- [Arora et al., 2023] Mehul Arora, Louis Wiesmann, Xieyuanli Chen, and Cyrill Stachniss. Static map generation from 3d lidar point clouds exploiting ground segmentation. *Robotics and Autonomous Systems*, 159:104287, 2023.
- [Benedikt et al., 2023] Mersch Benedikt, Guadagnino Tiziano, Chen Xieyuanli, Vizzo Ignacio, Behley Jens, and Stachniss Cyrill. Building Volumetric Beliefs for Dynamic Environments Exploiting Map-Based Moving Object Segmentation. *IEEE Robotics and Automation Letters (RA-L)*, 8(8):5180–5187, 2023.
- [Blanco and Rai, 2014] Jose Luis Blanco and Pranjali Kumar Rai. nanoflann: a C++ header-only fork of FLANN, a library for nearest neighbor (NN) with kd-trees. <https://github.com/jlblancoc/nanoflann>, 2014.
- [Cadena et al., 2016] Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif, Davide Scaramuzza, Jose Neira, Ian Reid, and John J. Leonard. Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Transactions on Robotics*, 32(6):1309–1332, 2016.
- [Campos et al., 2021] Carlos Campos, Richard Elvira, Juan J. Gómez, José M. M. Montiel, and Juan D. Tardós. ORB-SLAM3: An accurate open-source library for visual, visual-inertial and multi-map SLAM. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021.
- [Chen et al., 2019] Xieyuanli Chen, Andres Milioto, Emanuele Palazzolo, Philippe Giguere, Jens Behley, and Cyrill Stachniss. Suma++: Efficient lidar-based semantic slam. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4530–4537. IEEE, 2019.
- [Chen et al., 2021] Xieyuanli Chen, Shijie Li, Benedikt Mersch, Louis Wiesmann, Jürgen Gall, Jens Behley, and Cyrill Stachniss. Moving object segmentation in 3d lidar data: A learning-based approach exploiting sequential data. *IEEE Robotics and Automation Letters*, 6(4):6529–6536, 2021.
- [Dai et al., 2018] Benny Dai, Cedric Le Gentil, and Teresa Vidal-Calleja. Connecting the dots for real-time LiDAR-based object detection with YOLO. *Australasian Conference on Robotics and Automation, ACRA*, 2018-Decem, 2018.
- [Ferguson and Stentz, 2006] Dave Ferguson and Anthony Stentz. Using interpolation to improve path planning: The field d\* algorithm. *Journal of Field Robotics*, 23(2):79–101, 2006.
- [Geiger et al., 2012] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012.
- [Guertler et al., 2023] Matthias Guertler, Laura Tomidei, Nathalie Sick, Marc Carmichael, Gavin Paul, Annika Wambsganss, Victor Hernandez Moreno, and Sazzad Hussain. When is a robot a cobot? moving beyond manufacturing and arm-based cobot manipulators. In *2023 INTERNATIONAL CONFERENCE ON ENGINEERING DESIGN (ICED)*, volume 3, pages 3889–3898. Cambridge University Press, 2023.
- [Henein et al., 2020] Mina Henein, Jun Zhang, Robert Mahony, and Viorela Ila. Dynamic slam: The need for speed. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2123–2129. IEEE, 2020.
- [Krizhevsky et al., 2012] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [Le Gentil et al., 2021] Cedric Le Gentil, Teresa Vidal-Calleja, and Shoudong Huang. IN2LAAMA: INertial Lidar Localisation Autocalibration And MApping. *IEEE Transactions on Robotics*, 2021.
- [Likhachev et al., 2005] Maxim Likhachev, David I Ferguson, Geoffrey J Gordon, Anthony Stentz, and Sebastian Thrun. Anytime dynamic a\*: An anytime, replanning algorithm. In *ICAPS*, volume 5, pages 262–271, 2005.
- [Marvel, 2013] Jeremy A Marvel. Performance metrics of speed and separation monitoring in shared workspaces. *IEEE Transactions on automation Science and Engineering*, 10(2):405–414, 2013.
- [Milioto et al., 2019] Andres Milioto, Ignacio Vizzo, Jens Behley, and Cyrill Stachniss. Rangenet++: Fast and accurate lidar semantic segmentation. In

- 2019 *IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 4213–4220. IEEE, 2019.
- [Modayil and Kuipers, 2008] Joseph Modayil and Benjamin Kuipers. The initial development of object knowledge by a learning robot. *Robotics and autonomous systems*, 56(11):879–890, 2008.
- [Nießner *et al.*, 2013] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 32(6):1–11, 2013.
- [Oleynikova *et al.*, 2017] Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan Nieto. Voxblox: Incremental 3d euclidean signed distance fields for on-board mav planning. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- [Petrovskaya and Thrun, 2009] Anna Petrovskaya and Sebastian Thrun. Model based vehicle detection and tracking for autonomous urban driving. *Autonomous Robots*, 26(2-3):123–139, 2009.
- [Pfreundschuh *et al.*, 2021] Patrick Pfreundschuh, Hubertus FC Hendriks, Victor Reijgwart, Renaud Dubé, Roland Siegwart, and Andrei Cramariuc. Dynamic object aware lidar slam based on automatic generation of training data. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11641–11647. IEEE, 2021.
- [Postica *et al.*, 2016] Gheorghii Postica, Andrea Romanoni, and Matteo Matteucci. Robust moving objects detection in lidar data exploiting visual cues. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1093–1098. IEEE, 2016.
- [Redmon and Farhadi, 2017] Joseph Redmon and Ali Farhadi. YOLO9000: Better, faster, stronger. *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017-Janua:6517–6525, 2017.
- [Schmid *et al.*, 2023] Lukas Schmid, Olov Andersson, Aurelio Sulser, Patrick Pfreundschuh, and Roland Siegwart. Dynablox: Real-time detection of diverse dynamic objects in complex environments. *IEEE Robotics and Automation Letters*, 8(10):6259–6266, 2023.
- [Villani *et al.*, 2018] Valeria Villani, Fabio Pini, Francesco Leali, and Cristian Secchi. Survey on human–robot collaboration in industrial settings: Safety, intuitive interfaces and applications. *Mechatronics*, 55:248–266, 2018.