# Comparison of modern open-source visual SLAM approaches

**Dinar Sharafutdinov · Mark Griguletskii · Pavel Kopanev · Mikhail Kurenkov · Gonzalo Ferrer · Aleksey Burkov · Aleksei Gonnochenko · Dzmitry Tsetserukou**

**Abstract** SLAM is one of the most fundamental areas of research in robotics and computer vision. State of the art solutions has advanced significantly in terms of accuracy and stability. Unfortunately, not all the approaches are available as open-source solutions and free to use. The results of some of them are difficult to reproduce, and there is a lack of comparison on common datasets. In our work, we make a comparative analysis of state-of-the-art open-source methods. We assess the algorithms based on accuracy, computational performance, robustness, and fault tolerance. Moreover, we present a comparison of datasets as well as an analysis of algorithms from a practical point of view. The findings of the work raise several crucial questions for SLAM researchers.

**Keywords** SLAM · VIO · Benchmarking · SLAM comparison

## 1 Introduction

SLAM or simultaneous localization and mapping is an important task in computer vision and robotics. Robots need to estimate their position, trajectory, and surrounding map during operation. The problem becomes more difficult for an unknown environment where there is no prior knowledge. SLAM algorithms only require differences between sequential positions of the robot and measurements between robot and landmarks. Thus, the goal of SLAM is to construct a consistent map of the environment, find a robot trajectory and localize it on the map.
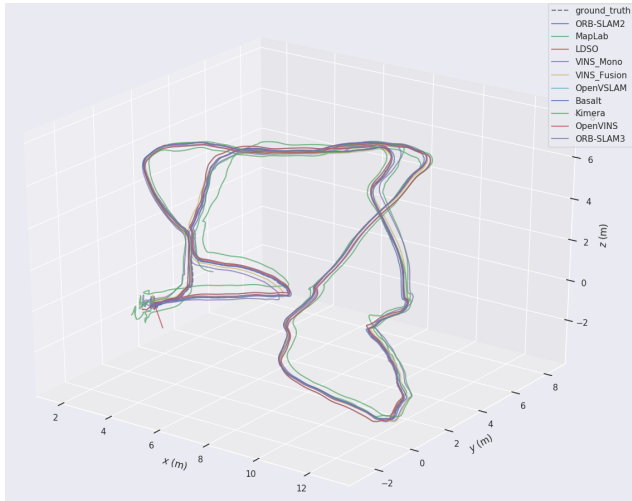
There are many open-source solutions which have been developed for the last 20 years. However, it is a problem of selecting an opens-source SLAM system among the options available. Different algorithms have different implementations and were tested on different datasets. However, good common evaluation of SLAM systems is essential for robotics engineers. They must known weakness and strengths to properly adjust SLAM solutions in their systems.

There are several approaches to solve the visual SLAM problem. Main methods are feature-based (Mur-Artal et al., 2015) and direct methods (Engel et al., 2014). Feature-based methods use key-points, they save it to local submaps or key-frames and then make bundle adjustment or graph optimization by, for example, Levenberg–Marquardt algorithm. Whereas direct-SLAM methods save all pixels of images on local maps and use photometric losses for measurement error estimations. Modern SLAM approaches are very different and usually highly modular. They consist of front-end and back-end. Moreover, some methods solve the SLAM problem with dynamic objects.
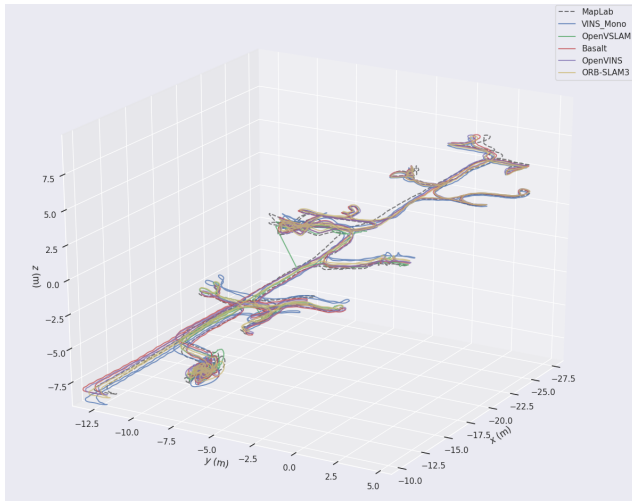
During last years, SLAM system for different sensors apart RGB cameras have been developed, for example, for RGB-D cameras and stereo cameras (Mur-Artal and Tardós, 2017). IMU measurements also have been integrated into SLAM systems (Qin et al., 2018). The combination of different information about the environment leads to a good solution of SLAM problem. Thus,

Dinar Sharafutdinov E-mail: dinar.sharafutdinov@skoltech.ru · Mark Griguletskii E-mail: mark.griguletskii@skoltech.ru · Pavel Kopanev E-mail: pavel.kopanev@skoltech.ru · Mikhail Kurenkov E-mail: mikhail.kurenkov@skoltech.ru · Gonzalo Ferrer E-mail: G.Ferrer@skoltech.ru · Dzmitry Tsetserukou E-mail: D.Tsetserukou@skoltech.ru
Skolkovo Institute of Science and Technology

Dinar Sharafutdinov · Mark Griguletskii · Pavel Kopanev · Aleksey Burkov E-mail: burkov.a.m@sberbank.ru · Aleksei Gonnochenko E-mail: gonnochenko.a.s@sberbank.ru
Sberbank Robotics Lab

(a) EuRoC MAV mh_5



(b) TUM VI corridor_1

Fig. 1: Examples of resulting trajectories

SLAM approaches must be implemented and compared on the datasets with different sensors.

Modern SLAM technologies achieve significant results but several challenges exist (Cesar et al., 2016). The main problems are long-term stability, scalability, sensor noises and miscalibrations, computation and memory requirements, dynamic objects, ambiguous scenes, accuracy requirements. Since one of the most popular areas of SLAM usage is mobile robotics with the outdoor environment, critical issues are dynamic objects and map size. Thus, modern SLAM solutions should take these challenges into account. Moreover, modern datasets and benchmarks need to be more focused on these challenging environments.

Nowadays, a lot of technical researches and papers suffer a common issue with the reproducibility of their methods and results. This is called the "reproducibility crisis". Frequently, it is very complicated to repeat the same results as described in paper or documentation. But this might be critical for researchers and engineers. Fast reproducibility allows to test and evaluate modern SLAM method in their environments and choose the appropriate setup for their task.

There are a lot of datasets and benchmarks for SLAM evaluation. Some of them are well described in a relevant work (Liu et al., 2021). A dataset should reflect SLAM challenges, for example, dynamic objects, ambiguous scenes, and huge spaces.

The purpose of the current paper is to give a general introduction to algorithmic state-of-the-art solutions for visual SLAM, to reproduce these algorithms on the same platform, make a comparison of the approaches on open datasets, and give practical advice for SLAM researchers. Figure 1 illustrates some trajectories for visualisation purpose and alignment consistency.

The contribution of the article is

1. Reproduction of the modern solutions on the same platform
2. Comparison of the modern solutions on the open datasets
3. Analysis of advantages and disadvantages of the methods
4. Practical hints for researchers for choosing SLAM methods coupled with published GitHub repository[1] with Dockers of the algorithms

## 2 Related work

For localization and mapping, many sensors are used as input data. For example, 3D LIDARs, 2D LIDARs, depth cameras, stereo cameras, global and rolling shutters cameras. LIDARs are quite precise but very expansive. Moreover, it needs to emit light signals to estimate distance. Therefore, LIDARs are active sensors. One of the most beneficial sensors for robots is the rolling shutter camera. It is cheap, dense and passive. They do not need to emit light and only receive surrounding information. Global shutter cameras might be more convenient, but the price is higher. Stereo and depth cameras such as D435 are also appropriate and very popular. All of them can be used in SLAM pipelines, but they have different features which are important while solving the SLAM problem. Apart from visual sensors, inertial measurements are also being used as an input: wheel odometry and IMU measurements are actively utilized in robotics.

---

[1] SLAM-Dockers        https://github.com/KopanevPavel/SLAM-Dockers

SLAM has been extensively studied during the last decades. Many reviews of SLAM solutions have been published. For example, Cadena et al. (Cesar et al., 2016) presented a good overview of modern SLAM technologies and the challenges which SLAM methods face. According to another review (Chen et al., 2020), today SLAM is going to spatial artificial intelligence age. This means that more and more deep learning techniques are used for solving SLAM problems.

## 2.1 Sparse visual SLAM

The history of feature-based SLAM or sparse visual SLAM began with the history of visual SLAM. One of the first successful approaches which solved the problem of simultaneous localization and mapping is a method presented by Davison et al. (Davison et al., 2007). In this work the MonoSLAM algorithm had been introduced. It was able to find the camera position from the input images. The authors managed to speed up the solution of the SfM (Structure from Motion) problem by using methods based on probabilistic filters. The algorithm can obtain the position of point landmarks and the trajectory of the camera in real-time.

There are several disadvantages of MonoSLAM method. The first problem is filtering approach which badly scales. The other disadvantage of MonoSLAM algorithm was a necessity to update the trajectory and the feature positions sequentially in the same thread. To solve this problem, the Klein et al. (Klein and Murray, 2007) proposed the PTAM (Parallel Tracking and Mapping) algorithm which achieves real-time performance on an algorithm that executes full batch updates. Additionally, authors of the work (Strasdat et al., 2010) showed that the solutions based on bundle adjustment (Klein and Murray, 2007) are superior to filter-based ones.

The key contribution in the field of solving the problem of visual SLAM was made by ORB-SLAM algorithm which was described in this article (Mur-Artal et al., 2015). The authors of the paper integrated all previous methods together in the same frame and combined the key elements of visual SLAM methods: Bags of Binary Words (Galvez-López and Tardos, 2012) visual terrain recognition algorithm, ORB (Rublee et al., 2011) algorithm for quick detection and description of key points and g2o (Kummerle et al., 2011) optimizer. Figure 2 shows general scheme of the method. There are also other backend optimizers which researchers highly used. For example, GTSAM (Dellaert and Kaess, 2012) and Ceres (Agarwal et al.). However, this article focuses more on different frontend parts of SLAM systems.

ORB-SLAM has three parallel processes. The first process is the construction of the local camera trajectory by matching the observed key points to the local map. The second process builds a local map and solves the local bundle adjustment problem. And the last process finds loop closures. Moreover, it can trigger a full optimization of the entire camera trajectory. The authors of ORB-SLAM also proposed a second version of their algorithm: ORB-SLAM2 (Mur-Artal and Tardós, 2017). This algorithm allows to use stereo and depth cameras.

ORB-SLAM method suffers from drifts of a scale. In order to reveal the geometric scale another source of information is often used, for example, an IMU sensor. The article (Qin and Shen, 2017a) introduces VINS algorithm which uses Kalman filter to merge sequential images and IMU data to calculate odometry. As a logical continuation of (Qin et al., 2019), the authors started to use global optimization to relocate the camera.

The development of ORB-SLAM2 can be considered as SOFT-SLAM (Cvišić et al., 2018), the authors of this algorithm made the following changes:

- Instead of the more computationally complex bundle adjustment for localization, SOFT visual odometry is used, which achieves an error of about 0.8% relative to the distance traveled.
- The streams responsible for mapping and odometry are separated from each other, thus the visual odometry stream is not blocked by the thread responsible for mapping, which leads to more stable processing time for incoming frames. In addition, unlike ORB-SLAM2, the algorithm is completely deterministic, i.e. always returns the same result for the same input.
- SOFT key points are used for loop closure (similar to ORB-SLAM2), which makes the system highly efficient, simple and reliable, while achieving sub pixel precision. Despite the fact that SOFT key points are not invariant with respect to rotation, the authors show on public datasets that in practice, loop closure occurs often enough that this drawback has little effect on the final result.

It should also be noted ORB-SLAM2-CNN (Tateno et al., 2017), which adds semantic information received from neural network to the ORB-SLAM2 framework.

Another cue point method worth mentioning is VINS-Mono (Qin and Shen, 2017b). This method is focused on processing monocular data augmented with data from IMU.

The article (Bustos et al., 2019) questions the visual SLAM approach based on global optimization of the entire trajectory and key point positions, that is, based

**TRACKING**

Figure: Tracking boxes — Stereo/RGB-D Frame → Pre-process Input → Pose Prediction (Motion Model) or Relocalization → Track Local Map → New KeyFrame Decision → KeyFrame

**PLACE RECOGNITION** — Visual Vocabulary, Recognition Database

**MAP** — MapPoints, KeyFrames, Covisibility Graph, Spanning Tree

**LOCAL MAPPING** — KeyFrame Insertion, Recent MapPoints Culling, New Points Creation, Local BA, Local KeyFrames Culling

Loop Correction — Optimize Essential Graph, Loop Fusion

Loop Detection — Compute SE3, Query Database

**FULL BA** — Update Map, Full BA
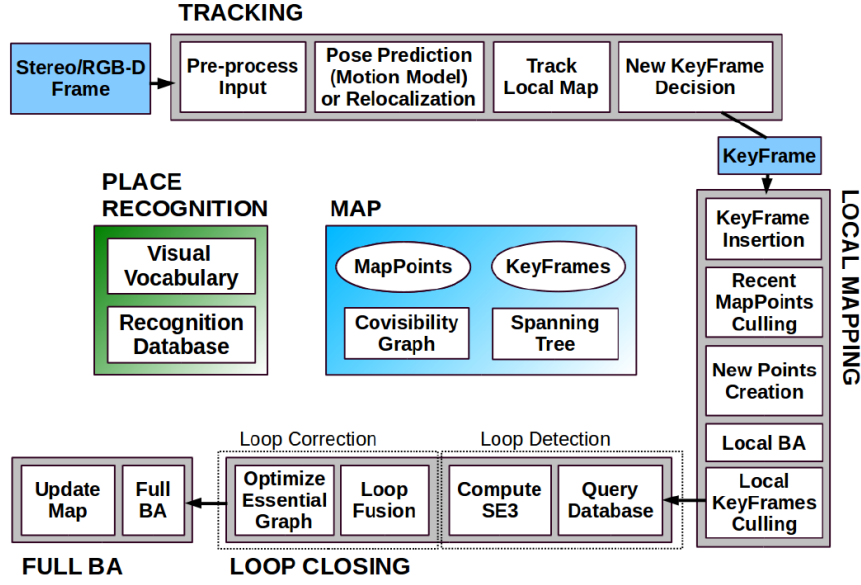
**LOOP CLOSING**

Fig. 2: The scheme of ORB-SLAM algorithm (Mur-Artal et al., 2015)

on bundle adjustment. As an alternative, the authors of the article propose the L-infinity SLAM method. The essence of this method is that the robot first separately finds the camera rotations by averaging the relative rotations, after which it solves the optimization problem to find the three-dimensional position of the cameras and key points with known camera orientations. One can also use the rotation averaging described in the publication (Bustos et al., 2019) for direct methods, and as additional restrictions for other ways to solve the visual SLAM problem.

2.2 Dense visual SLAM

Unlike methods for solving the problem of visual SLAM, which are based on key points, dense or direct methods do not use algorithmic features in their work. They use photometric error minimization for each pair of images to find the orientation and position of the camera in 3D space.

The first work that suggested using direct methods was DTAM (Newcombe et al., 2011b). For each image, the DTAM algorithm generated a dense 3D map for each pixel. The next algorithm was KineticFusion (Newcombe et al., 2011a), which uses a depth camera to build a dense 3D map. The authors use TSDF (truncated signed distancefunction) to describe each pixel, and ICP (iterative closest point) algorithm to map each depth image to a map.

The next milestone in the development of direct methods was the ElasticFusion (Whelan et al., 2016)

algorithm. This algorithm uses a direct representation of the surfaces of RGB-D cameras. For data fusion, this algorithm uses a non-rigid deformation model. Also in this algorithm there is no global optimization of the graph of camera positions, and most of the steps take place on a graphics accelerator (GPU).

In parallel with the development of KineticFusion and ElasticFusion, two other teams were developing semi-dense visual odometry (SVO) (Engel et al., 2013; Forster et al., 2014). This algorithm can run in real time on a robot processor (CPU). The SVO algorithm does not use all the pixels in the image, but only those that have a negligible gradient. The development of ideas from the publication (Engel et al., 2013) leads to the creation of a key work in this area - LSD-SLAM (Engel et al., 2014). LSD-SLAM uses $sim(3)$ - a metric to solve a problem with an uncertain scale, and also uses probabilistic inference to determine the error in constructing three-dimensional maps. The LSD-SLAM algorithm consists of three steps:

1. Get the image and determine the local offset relative to the current keyframe.
2. Update or create a new keyframe. If the current keyframe is updated, probabilistically merge the stereo depth data into the current keyframe. In case of creating a new keyframe, propagate the previous depth map to the new keyframe.
3. Updating the global map using position graph optimization. The edges in the graph are searched using the sim (3) metric, and the optimization is performed using the g2o library.

The LSD-SLAM development team continued the development of direct methods and proposed the DSO (direct sparse odometry) (Engel et al., 2017) method. Their method optimizes not only the photometric error, but also the geometric error simultaneously. Also, instead of assuming smoothness, the authors use probabilistic pixel sampling. In LDSO (Gao et al., 2018), they added graph optimization to get a complete solution to the visual SLAM problem. In fact, LDSO is an combination of the ORB-SLAM and LSD-SLAM approaches.

Separately, it should be noted that the authors of DSO in 2019 published a modification of DSO (Schubert et al., 2019), adapted for processing frames obtained using a camera with a floating shutter. In addition to the photometric bundle adjustment, IMU data are also used as additional constraints for the optimization problem. Due-to the fact that there is no selection of key points, this method (like DSO) can work not only with pixels representing edges and corners, but with any pixels that have a sufficient gradient.

Many improvements to direct methods have also been proposed recently. For example, CodeSLAM (Bloesch et al., 2018) suggests using autoencoders to train a more compact map representation for direct methods. Also, in the work KO-fusion (Houseago et al., 2019) it was proposed to use odometry and kinematics data obtained directly from the robot arm.

## 2.3 Dynamic visual SLAM

Usually, SLAM solutions suppose that scene is almost static or with a low level of dynamics. The simplest way of discarding outliers is to use RANSAC (Fischler and Bolles, 1987). In reality, almost in every scenario, there will be lots of dynamic objects. Both indoor and outdoor scenes are full of people, animals, cars, pedestrians, bicycles and other dynamic objects. Therefore, it leads to changes in feature maps. This especially important for loop closure detection. Additionally, in the case of a high level of dynamics, most slam solutions may show inconsistent results because of a lack of reliable features. There are several problems with dynamic objects. Firstly, the algorithms should detect the dynamic object and do not use it for the trajectory estimation as well as for the mapping process. A possible solution here is to delete these objects from the image or replace them with a background. Popular tools for that are multi-view geometry(Yu et al., 2018), semantic segmentation or detection neural networks(Bescos et al., 2018), scene(Bârsan et al., 2018) or optical flow and detection of background or foreground (Yang et al., 2019). Secondly, there is a need of building a consistent map with or without dynamic objects. Solving the problem

of lack of data due to the deletion of dynamic objects is an important aspect. Additional sensors, objects completion and reconstruction in combination with different map construction algorithms might help (Dou et al., 2016)(Rünz and Agapito, 2017). Dynamic SLAM algorithms often represent complex systems consisting of many parts and working not in real-time.

## 2.4 VIO

To improve localization accuracy, several modern SLAM techniques use IMU sensors (Maddern et al., 2011) and wheel encoders (Yang et al., 2019). This allows to take into account the position of the localized object in the intervals between frames received from the video camera. IMU preintegration(Forster et al., 2016) approach has shown its efficiency at short time intervals. Additionally, IMU is widely used for scale correction (Maddern et al., 2011) if only one camera sensor is present. Visual Odometry methods for position estimation are sensitive to environmental conditions such as lighting changes, surrounding texture, the presence of water or snow (Poddar et al., 2018). These changing conditions can lead to poor feature-tracking. In this case fusion with IMU helps to propagate position estimation relying not only on camera data.

There are two most popular IMU integration strategies for SLAM approaches utilizing commercial sensors: the loosely-coupled and the tightly-coupled(Falco et al., 2017). For loosely-coupled VIO system both inertial and visual information are seen as independent measurements. Tightly-coupled system considers the interaction of all measurements of sensors information before pose estimation. An example of VIO-based algorithm is VINS-Mono (Maddern et al., 2011). VINS-Mono is a tightly-coupled visual-inertial system. It utilizes IMU preintegration between camera frames. IMU has a higher update rate than the camera and preintegration avoids adding additional nodes in the pose graph keeping it light. For optimization, VINS-Mono uses Ceres Solver suitable for modeling and solving large, complicated optimization problems. IMU preintegration result is used as predicted pose value which is than compared with the pose calculated by utilizing camera (measured value). Also VINS-Mono has loop-closure capabilities.

## 3 Experimental setup

There were several main goals in the experimental stage. First of all, our goal is to test algorithms that are available to everyone. We took only those SLAM

solutions that have a link to the page with code and instructions. Secondly, we weren't targeting on doing experiments with all the possible solutions. Instead, it was decided to take one or a few algorithms from each group of SLAM methods. We wanted to see results of algorithms which work with monocular, stereo and RGB-D cameras. Another essential part is understanding how additional data such as IMU or wheel encoder will improve the results. The last point here is to try solutions with deep learning and with various front-end parts. Thirdly, we had the following main use case that suits for broad possible robotics tasks. Wheel robot with cameras and other sensors is moving in an outside or inside environment. It moves along dynamic objects and features-less areas. The robot also visits the same places several times (loop closures), and the traversal could be relatively long (30 minutes). These factors affect the datasets that we used for tests significantly.

## 3.1 Choice of algorithms

### 3.1.1 Sparse and dense SLAM

ORB-SLAM2 (Mur-Artal and Tardós, 2017) was used as a baseline sparse SLAM solution. It is quite popular among researchers, and it proved efficiency in many scenarios. The second version is an expansion of the first hence we did not use ORB-SLAM1 (Mur-Artal et al., 2015). Peculiarities of the algorithm were covered in the literature review section.

Another algorithm here is an OpenVSLAM (Sumikura et al., 2019). Generally, it is still ORB-SLAM but with additional features from ProSLAM, and UcoSLAM. The main value of this solution is practical extensions such as support of the various type of camera models (perspective, fisheye, equirectangular), map saving and creation features, and finally convenient documentation and UI.

The most modern approach that we considered is recently published ORB-SLAM3 (Campos et al., 2020). This sparse SLAM system is a natural development of the previous version of ORB-SLAM. Now it is a visual-inertial approach with an updated place recognition module (high-recall place recognition with geometrical and local consistency check) as well as with the support of pinhole and fisheye camera models. Moreover, it allows localization in a multi-map setup.

Also we have chosen one of the dense SLAM methods for comparison. LDSO (Gao et al., 2018) is graph-based dense SLAM system with loop closure detection. It is based on Direct Sparse Odometry (DSO) (Engel et al., 2018). LDSO is robust and suitable for feature-poor environments utilizing any image pixel with sufficient intensity gradient. It uses feature-based bag-of-words for loop closure detection and was validated on several open-source datasets.

### 3.1.2 Visual-Inertial Odometry and SLAM

Typical Visual-Inertial Odometry system implies usage of a camera and IMU sensors. OpenVINS (Geneva et al., 2020) is Multi-State Constraint Kalman Filter (MSCKF) based VIO estimator. It uses IMU in propagation step of the filter and camera data in the update step. In current implementation OpenVINS also has secondary loosely coupled loop closure thread based on VINS-Fusion (Qin et al., 2019). Additionally OpenVINS has interface wrapper for exporting visual-inertial runs into the ViMap structure taken by maplab (Schneider et al., 2018).

VINS-Mono (Qin et al., 2018) and VINS-Fusion are graph-based VIO approaches. VINS-Fusion is an extension of VINS-Mono and supports multiple visual-inertial sensor types (mono camera + IMU, stereo cameras + IMU, even stereo cameras only). Also it supports global sensors like GPS and Barometer and has global graph optimization module. Additionally VINS-Fusion supports loop closure.

Basalt (Usenko et al., 2019) is other graph-based VIO approach. It has a lot of similarities with previously described VINS-Fusion (KLT feature tracking and Gauss-Newton non-linear optimization). But regrading map optimization this approach utilizes ORB features.
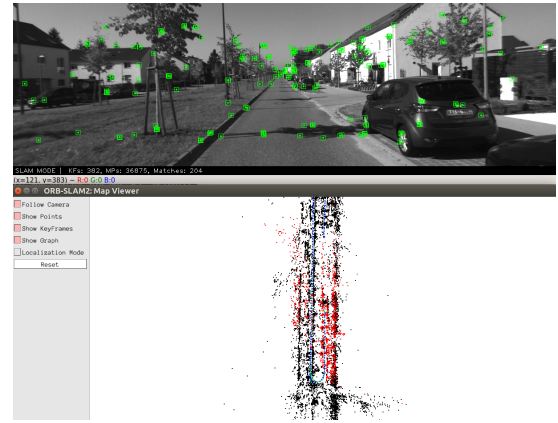
Kimera (Antoni et al., 2020) is an open source C ++ library that implements one way of solving SLAM problem in real time. The modular structure includes: a VIO module with GTSAM-based (Dellaert and Kaess, 2012) VIO approach, using IMU-preintegration (Forster et al., 2016) and structureless vision factors (Carlone et al., 2014) for fast and accurate state assessment, a robust position graph optimizer (RPGO) for global trajectory estimation which adds a robustness layer that avoids SLAM failures due to perceptual aliasing, and relieves the user from time-consuming parameter tuning, a lightweight 3D mesh module (Kimera-Mesher) for fast 3D mesh reconstruction (Rosinol et al., 2019) and obstacle avoidance and a dense 3D semantic reconstruction module (Kimera-Semantics) that builds a more-accurate global 3D mesh using a volumetric approach (Oleynikova et al., 2017), and semantically annotates the 3D mesh using 2D pixel-wise semantic segmentation based on deep learning methods.
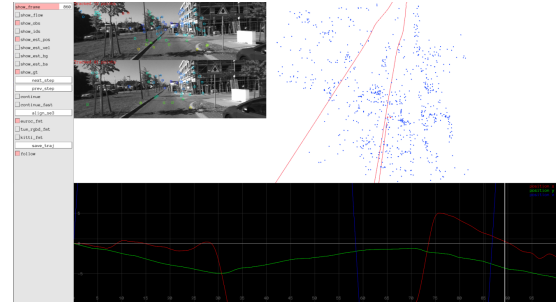
*3.1.3 Dynamic SLAM*

In this section, we have started with three algorithms: DynaSLAM (Bescos et al., 2018), DynSLAM (Bârsan et al., 2018) and DRE-SLAM (Yang et al., 2019). All of them have theoretical potential and relatively popular implementations. DynaSLAM is a successor of ORB-SLAM2 with added front-end part that allows working in a dynamic environment. The solution consists of the segmentation of dynamic objects using CNN in monocular and stereo cases, and a combination of deep neural methods and multi-view geometry in the RGBD case. The important feature is that with DynaSLAM, it is possible to detect a priory dynamic objects, objects that may be static at the moment but the dynamic in essence. Pixel-wise semantic segmentation of potentially movable objects is possible with the help of Mask R-CNN. Then combining segmentation with multi-view geometry allows finding even static objects moved by dynamic ones (e.g. book in the hands of a person). After the deletion of dynamic objects, there are empty regions on images. To solve that, the authors use background inpainting by using information from previous views.

DynSLAM is a more complex solution that solves even more tasks for large scale dynamic environments. Taking stereo images as the input, they compute depth maps (using ELAS or DispNet) and sparse scene flow (libviso2). Multi-task Network Cascades allow finding dynamic objects. Using obtained information, they do 3D objects tracking and reconstruction. The final result is a static map without dynamic objects with the help of InfiniTAM. The last one, DRE-SLAM, is a bit different from the previous ones. Apart from image data, it also uses two wheel encoders to handle the lack of features in dynamic scenes. It extracts ORB features from RGB-D images, uses YOLOv3 for dynamic object detection and then applies multiview constraint-based pixels culling. Loop closure detection implemented through a bag-of-words approach. The final map is OctoMap which is constructed by fusing sub-OctoMaps.
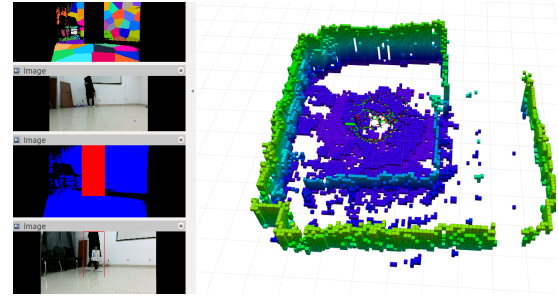
Even though DynaSLAM and DynSLAM have developed and published solutions, we were unable to run them. GitHub pages full of issues and implementations are not supported by authors, which makes solving varying problems very difficult. On the other hand, DRE-SLAM ran smoothly on ROS but it has examples only for data collected by authors. Comparison with other algorithms on popular datasets was difficult because of the lack of wheel encoder data in many of them. Still, we compared it with several other algorithms on the DRE dataset.



(a) ORB-SLAM2



(b) Basalt



(c) DRE-SLAM

Fig. 3: GUI examples

## 3.2 A practical overview of open source solutions

During the process of work with algorithms, we identified a number of aspects that must be taken into account when dealing with SLAM solutions.

First, all algorithms were run on Linux-based systems (most often Ubuntu). Launching on Windows seems to be quite time-consuming and, in general, all examples and instructions from the authors are made for Ubuntu or macOS. Moreover, it was found that a number of algorithms require older versions of Ubuntu (16.04 or 18.04). In some cases, this can be critical for a successful launch.

Secondly, algorithms have different interfaces. The most unified way is to use the Robot Operating System

(ROS). In this case, everything depends on the support of the developers of the algorithm. The presence of a separate community, pipeline for data usage, and many examples make ROS a convenient option for launching. Other options are graphical interfaces used by developers (e.g. Pangolin viewer). The quality of the viewer implementation and interface depends on the developer and may differ from one algorithm to another.

Thirdly, the points mentioned above make it especially important to have support from the authors of the SLAM systems. It can be in form of answers to issues on GitHub, availability of documentation and examples of launching. In some cases, the authors have their own website with documentation, or even a channel on Slack to discuss the algorithm and emerging problems (Sumikura et al., 2019).

Fourth, the most convenient way to run SLAM algorithms is to use a ready-made Docker image. Just a few authors provide such a solution. To simplify the launch of algorithms and to the benefit of the community, we make the repository with docker images publicly available (Kopanev, 2021) and also upload them to Docker Hub.
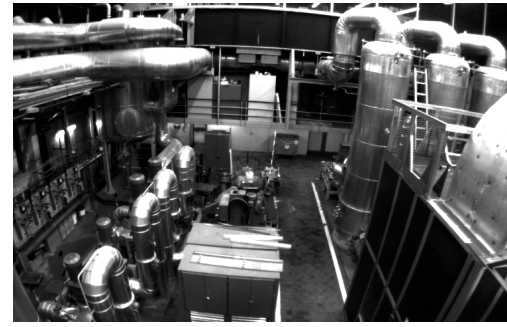
A full comparison of practical aspects could be found in Table 1. We qualitatively measured such aspects of original implementations of algorithms as availability of documentation, support from the authors (possibility to get an answer to the question in several days), examples on popular datasets, the convenience of the interface and change of parameters of algorithms as well as camera parameters, and presence of Docker and ROS wrapper.

## 3.3 Datasets description

Datasets are actively used for SLAM algorithms validation and comparison. They should include ground truth information and have a redundant amount of sensors in order to test different SLAM approaches based on different sensors stacks.

In our work we reviewed visual and visual-inertial SLAM approaches, and therefore one of the main dataset choice factors was camera images presence. Datasets were examined for the existence of loop closures and additional sensors (IMU, GPS, or wheel encoders). Here we do a short description of all the datasets that we considered. Our final choice for the experiments are five datasets: EuRoC MAV, TUM VI, KITI, Open Loris, and DRE dataset. Camera data is shown in Fig. 4.

A commonly used option for evaluating visual and visual-inertial SLAM algorithms is the EuRoC MAV



(a) EuRoC MAV



(b) TUM VI



(c) KITTI



(d) OpenLoris-Scene



(e) DRE dataset

Fig. 4: Examples of images from the datasets

dataset (Burri et al., 2016), but its image resolution and bit depth is not quite state-of-the-art anymore. Also this dataset is mostly suitable for Micro Air Vehicles and does not include odometry readings suitable for some SLAM approaches for mobile wheeled robots. SLAM article authors usually test and compare their algorithms on EuRoC dataset primarily but in many cases SLAM algorithms show good performance on this dataset and bad on other datasets. That will be shown later.

TUM VI (Schubert et al., 2018) is other popular dataset for evaluating visual and visual-inertial

Table 1: Practical comparison of algorithms

| Framework | Documentation | Support | Usage examples on popular datasets | Convenience of usage | Docker | ROS |
|---|---|---|---|---|---|---|
| ORB-SLAM2 | github | - | + | + | - | + |
| MapLab | github | - | + | + | - | + |
| LDSO | github | + | + | + | - | - |
| VINS-Mono | github | - | + | + | + | + |
| VINS-Fusion | github | - | + | + | + | + |
| OpenVSLAM | web-page, github | + | + | + | + | + |
| Basalt | github | + | + | + | + | - |
| Kimera | github | + | EuRoC only | + | + | + |
| OpenVINS | web-page, github | + | + | + | - | + |
| ORB-SLAM3 | github | + | + | + | - | + |
| DRE | github | - | - | + | - | + |

SLAM algorithms. It includes more varied scenes including indoor and outdoor environments and longer sequences than EuRoC dataset. Dataset provides time-synchronized camera images with 1024x1024 resolution and IMU measurements. However, this dataset provides accurate pose ground truth from a motion capture system only at the start and end of the sequences. There is no ground truth data for the outdoor part of the sequence.

OpenLoris-Scene dataset (Shi et al., 2020) is a dataset which aims to help evaluate the maturity of SLAM and scene understanding algorithms for real-world deployment, by providing visual, inertial and odometry data recorded with real robots in real scenes, and ground-truth robot trajectories acquired by motion capture system or high-resolution LiDARs. The important peculiarity of this dataset is quite hard real-life conditions with recordings of the same scenes in different lighting levels and with dynamic content.

TUM RGB-D SLAM Dataset (Sturm et al., 2012b) is a dataset which includes RGB-D data and ground-truth data with the goal to establish a novel benchmark for the evaluation of visual odometry and visual SLAM systems. Dataset contains the color and depth images of a Microsoft Kinect sensor and accelerometer data. The ground-truth trajectory was obtained from a high-accuracy motion-capture system.

KITTI (Geiger et al., 2012) vision dataset was collected using an autonomous driving platform and scenes are captured by driving around the mid-size city of Karlsruhe, in rural areas and on highways. Dataset has a benchmark suite which includes odometry, object detection and tracking, road (Fritsch et al., 2013), stereo and flow benchmarks. Odometry dataset was chosen for visual SLAM evaluation. This benchmark consists of 22 stereo sequences, saved in loss-less *png* format and has a leaderboard. For this benchmark it is possible to provide results using monocular or stereo visual odometry,

laser-based SLAM or algorithms that combine visual and LIDAR information. The data was recorded using an eight core i7 computer.

KAIST Urban Data Set (Jeong et al., 2019) is other dataset collected using car platform. The vehicle was equipped with two 2D and two 3D LiDARs to collect data on the surrounding environment. Additionally sensor suite included a stereo camera installed facing the front of the vehicle, GPS, VRS GPS, IMU, Fiber Optic Gyro (FOG) and altimeter. All sensor information was provided in a raw file format with timestamps. Three PCs were used to collect the data. The system clocks of the three PCs were periodically synchronized using the Chrony library (Curnow). Each PC used an i7 processor, 512 GB solid state drive (SSD), and 32 GB DDR4 memory.

Oxford RobotCar Dataset (Maddern et al., 2017) contains over 100 repetitions of a consistent route through Oxford, UK, captured over a period of over a year. The dataset captures many different combinations of weather, traffic and pedestrians, along with longer term changes such as construction and roadworks.

RPNG OpenVINS Dataset was developed by OpenVINS (Geneva et al., 2020) authors. It includes ArUco datasets built using Synchronized Visual-Inertial Sensor System (Nikolic et al., 2014) in indoor environment with ArUco markers. Also this dataset includes two long sequences (2.3 and 7.4 km) build using ironsides (Zhang et al., 2017) visual-inertial sensor.

DRE dataset (Yang et al., 2019) is a dataset collected using Redbot robot. It was made by DRE SLAM (Yang et al., 2019) authors for their algorithm evaluation. One of the key features of this dataset is presence of odometry readings from wheel encoders. Also dataset includes scenes with low and high dynamic objects.

Segway DRIVE benchmark (Huai et al., 2019) includes datasets collected by Segway delivery robots deployed in real office buildings and shopping malls. Each

robot was equipped with a global-shutter fisheye camera, a consumer-grade IMU synced to the camera on chip, two low-cost wheel encoders, and a removable high-precision lidar.

The Zurich Urban Micro Aerial Vehicle Dataset (Majdik et al., 2017) was collected using camera equipped Micro Aerial Vehicle (MAV) flying within urban streets at low altitudes (i.e., 5-15 meters above the ground). The 2 km dataset consists of time synchronized aerial high-resolution images, GPS and IMU sensor data, ground-level street view images, and ground truth data.

The Malaga Stereo and Laser Urban Dataset (Blanco-Claraco et al., 2014) was gathered entirely in urban scenarios with a car equipped with several sensors, including one stereo camera (Bumblebee2) and five laser scanners. One distinctive feature of this dataset is high-resolution stereo images.

The University of Michigan North Campus Long-Term Vision and LIDAR Dataset (Carlevaris-Bianco et al., 2016) consists of omnidirectional imagery, 3D lidar, planar lidar, GPS, and proprioceptive sensors for odometry collected using a Segway robot.

The following Table 2 shows sensors used in the described earlier datasets. Datasets main and additional features are shown in the Tables 3 and 4. As it could be seen from the tables, there is no perfect dataset that matches different demands, such as indoor/outdoor data, changing conditions, dynamic content, several types of sensors, etc. As well as there is no unified way of data organisation. Several standards of ground truth formats exist (TUM, KITTY, EuRoC), but only this. It means that there is always a need to use several datasets and tune data formats for the algorithms which are itself demands special conditions for data.

### 3.4 Metrics

In general, the result of the SLAM system is a trajectory and a map of the surrounding area. Despite the fact that it is theoretically possible to compare the resulting maps, in practice, obtaining reference maps is not a trivial task. In addition, different systems use different map formats, making this task even more difficult. For this reason, most often comparisons are made using the resulting trajectories.

However, it should be noted that a good trajectory does not necessarily mean a good map. For example, a map might be sparse with low number of features or landmarks whereas a trajectory is precise. Small errors on the map, which have little effect on localization accuracy, can interfere with the normal functioning of the robot, for example, by blocking movement in a doorway or corridor.

Readings from additional sensors can be used to obtain reference trajectories. For outdoor datasets, high-precision GNSS systems are used in conjunction with the IMU sensor (Geiger et al., 2013), while indoor optical tracking systems (Sturm et al., 2012a) can be used.

Relative Positional Error (RPE) (Kummerle et al., 2011) measures the trajectory drift over a fixed time interval $\Delta$ in a reference frame, thus it does not accumulate previous errors. Additionally, many researchers tend to separate translational and rotational errors because it gives more clear understanding of SLAM algorithm quality and it is not plausible to mix angles with poses (degrees with meters). For pose predictions of SLAM algorithm $\mathbf{P_1}, ..., \mathbf{P_n} \in \mathrm{SE}(3)$, and the ground-true poses $\mathbf{Q_1}, ..., \mathbf{Q_n} \in \mathrm{SE}(3)$ the relative error for each point in time interval is defined as in (Kasar, 2018):

$$\mathbf{E}_i = \left( \mathbf{Q}_i^{-1} \mathbf{Q}_{i+\Delta} \right)^{-1} \left( \mathbf{P}_i^{-1} \mathbf{P}_{i+\Delta} \right)^{-1} \tag{1}$$

Thus, for a sequence of $n$ measurements and their corresponding positions, we get $m = n - \Delta$ relative errors. Using the error data, the root mean square error is calculated:

$$\mathrm{RMS}\left(\mathbf{E}_{1:n}, \Delta\right) = \left( \frac{1}{m} \sum_{i=1}^{m} \|\mathbf{E}_i\|^2 \right)^{1/2} \tag{2}$$

In practise, it is useful to know Relative Rotational and Relative Translation errors for time interval $\Delta$. It allows to estimate quality of SLAM algorithm for particular situation without an influence of previous movements. Thus, rotational and translational parts can be calculated separately:

$$\mathrm{RMS}\left(\mathbf{E}_{1:n}^{trans}, \Delta\right) = \left( \frac{1}{m} \sum_{i=1}^{m} \|\mathbf{E}_i^{trans}\|^2 \right)^{1/2} \tag{3}$$

$$\mathrm{RMS}\left(\mathbf{E}_{1:n}^{rot}, \Delta\right) = \left( \frac{1}{m} \sum_{i=1}^{m} \|\mathbf{E}_i^{rot}\|^2 \right)^{1/2} \tag{4}$$

where $\mathbf{E}_i^{trans}$ refers to the translation components of the Relative Pose Error and $\mathbf{E}_i^{rot}$ refers to the rotational part.

In some cases, instead of the mean square, the median or means of other orders can be used. For example, the arithmetic mean will be less sensitive to anomalies than the root mean square, but more sensitive than the median.

Table 2: Sensor setup for dataset collection

| Dataset name | Camera | GPS | IMU | Lidar | Wheel odometry | Additional sensors |
|---|---|---|---|---|---|---|
| EuRoC MAV | Aptina MT9V034 global shutter, WVGA monochrome (2×20 Hz) | x | ADIS16448 (200 Hz) | x | x | x |
| TUM VI | 1 stereo gray (20 Hz) | x | BMI160 (200 Hz) | x | x | x |
| KITTI | 2x Point Grey FL2-14S3C-C global shutter (15 Hz) | OXTS RT 3003 (250 Hz) | OXTS RT 3003 (250 Hz) | Velodyne HDL-64E (5–20 Hz) | x | x |
| OpenLoris-Scene | RealSense D435i (30 Hz) RealSense T265 (30 Hz) | x | RealSense (60-400 Hz) | Hokuyo UTM30LX (40 Hz) | Odometer (20 Hz) | x |
| DRE dataset | Kinect 2.0 (20 Hz) | x | x | x | Wheel encoders (100 Hz) | x |
| KAIST | 2x FLIR FL3-U3-20E4C-C global shutter (2x10 Hz) | U-Blox EVK-7P (10 Hz) | Xsens MTi-300 (200 Hz) | Velodyne VLP-16 (10 Hz), SICK LMS-511 (100 Hz) | 2x RLS LM13 encoders (100 Hz) | VRS–GPS, Three-axis FOG, Altimeter |
| Malaga Urban | 1 stereo RGB (20 Hz) | DELUO (1 Hz) | xSens MTi (100 Hz) | 3x Hokuyo UTM-30LX (25 Hz), 2x SICK LMS-200 (75 Hz) | x | x |
| Oxford RobotCar | BBX3-13S2C-38 global shutter (16 Hz), 3x GS2-FW-14S5C-C global shutter (11.1 Hz) | NovAtel SPAN-CPT (50 Hz) | NovAtel SPAN-CPT (50 Hz) | 2 x SICK LMS-151 (50 Hz) 1 x SICK LD-MRS (12.5 Hz) | x | x |
| RPNG OpenVINS | Aptina MT9V034 (60 Hz) | x | ADIS16448 (200 Hz) | x | x | x |
| Segway DRIVE | RealSence ZR300 (30 Hz) | x | BMI055 (250 Hz) | Hokuyo UTM30LX (40 Hz) | x | x |
| TUM RGB-D | Kinect (30 Hz) | x | x | x | x | x |
| UMich NCLT | 6 RGB (omni) (5 Hz) | Garmin 18x (5 Hz) | Microstrain 3DM-GX3-45 (100 Hz) | Velodyne HDL-32E (10 Hz) Hokuyo UTM-30LX (40 Hz) Hokuyo URG-04LX (10 Hz) | Wheel encoders | FOG, RTK GPS |
| Zurich Urban MAV | RGB rolling shutter (30 Hz) | GPS | IMU (10 Hz) | x | x | x |

Table 3: Dataset main features

| Dataset name | Year | ROS bag | Ground truth | Indoor / Outdoor |
|---|---|---|---|---|
| EuRoC MAV | 2016 | + | Motion capture system (acc. ≈ 1 mm) | Indoor |
| TUM VI | 2018 | + | Motion capture system (acc. ≈ 1 mm) | Indoor / Outdoor |
| KITTI | 2013 | + (parser is provided) | OXTS RT 3003 inertial navigation system (INS) (acc. < 10 cm) | Outdoor |
| OpenLoris-Scene | 2019 | + | OptiTrack motion capture system / LIDAR SLAM (acc. < 10 cm) | Indoor |
| DRE dataset | 2019 | + | Downview camera and markers (acc. unknown) | Indoor |
| KAIST | 2019 | + (rosbag player) | LIDAR SLAM (acc. unknown) | Outdoor |
| Malaga Urban | 2014 | - | GPS (low acc.) | Outdoor |
| Oxford RobotCar | 2017 | - | SPAN GNSS Inertial Navigation System (acc. < 10 cm) | Outdoor |
| RPNG OpenVINS | 2019 | + | GPS in outdoor scenes (low acc.) | Indoor / Outdoor |
| Segway DRIVE | 2019 | + | LIDAR SLAM (acc. < 10 cm) | Indoor |
| TUM RGB-D | 2012 | + | Motion capture system (acc. ≈ 1 mm) | Indoor |
| UMich NCLT | 2015 | + (parser is provided) | Fused GPS/IMU/LIDAR (acc. ≈ 10 cm) | Outdoor |
| Zurich Urban MAV | 2017 | + (parser is provided) | Pix4D SLAM (acc. unknown) | Outdoor |

For SLAM the global consistency of the resulting trajectory is also important, which can be estimated by calculating the difference between the corresponding estimated and reference poses (Absolute error) in global frame. Since the frames of reference for the estimated and reference trajectories may differ from each other, the first step is to align them with each other. This can be done using the analytical solution provided in (Horn, 1987). The solution finds a transformation $\mathbf{S} \in$ SE(3) that minimizes the root-mean-square difference between $\mathbf{P}_{1:n}$ and $\mathbf{Q}_{1:n}$. Using this transformation, the difference between the poses can be calculated as:

$$\mathbf{F}_i = \mathbf{Q}_i^{-1}\mathbf{S}\mathbf{P}_i \qquad (5)$$

$$\mathrm{RMS}\left(\mathbf{F}_{1:n}\right) = \left(\frac{1}{m}\sum_{i=1}^{m}\|\mathbf{F}_i\|^2\right)^{1/2} \qquad (6)$$

$$\mathrm{RMS}\left(\mathbf{F}_{1:n}\right) = \left(\frac{1}{m}\sum_{i=1}^{m}\|\mathbf{F}_i\|^2\right)^{1/2} \qquad (7)$$

In practice, both metrics show strong correlation with each other. But we are interested in comparison of different trajectories with one reference Ground-True trajectory for the whole length. Therefore, the more intuitive Absolute Error is used most often as well as in our results section.

For the calculation of trajectories and metrics, comparison EVO tool (Grupp, 2017) was used. Exhaustive documentation, visualisation possibilities, and easy to use interface make it very convenient to compare a large number of results. Another possible way to track the performance is to use the KITTI Vision Benchmark Suite (Geiger et al., 2012) web page. It allows seeing the translation and rotation error of an algorithm in a web-table with more than a hundred different solutions, but it is limited to one dataset.

Table 4: Dataset additional features

| Dataset name | Variable weather conditions | Dynamic objects | Memory consumption | Path length |
|---|---|---|---|---|
| EuRoC MAV | - | - | 1-2.5 Gb | 30-130 m |
| TUM VI | - | + | 3-60 Gb | <20 km |
| KITTI | + | + | >20 Gb | <40 km |
| OpenLoris-Scene | - | + | 6-33 Gb | ? |
| DRE dataset | - | Static (ST), low dynamic (LD), high dynamic (HD) | 4-8 Gb | ? |
| KAIST | - | + | 1-30 Gb | 1-30 km |
| Malaga Urban | - | + | 1-33 Gb | <36.8 km |
| Oxford RobotCar | + | + | 10-500 Gb | ? |
| RPNG OpenVINS | - | - | 1-2.6 Gb | 27-105 m and 2.3, 7.4 km |
| Segway DRIVE | - | + | 1-20 Gb | 50 km |
| TUM RGB-D | - | Separate category: dynamic objects | 0.2-2.5 Gb | 0.4 km |
| UMich NCLT | + | + | 80-110 Gb | 1-7.5 km |
| Zurich Urban MAV | - | + | 28 Gb | 2 km |

## 4 Results

The following tables represent the Absolute Pose Error metric (maximum and root mean squared values) for the algorithms we are considering. We present metrics both for position error in meters and for rotation error in degrees. We mark the lowest error result among all algorithms by bold green numbers and by italic purple numbers the second result. The first three tables represent results on popular and very common datasets to the robotics community. Not all algorithms present in all tables. Absence means the inability to run on certain types of data, problems with parameters or algorithms themselves. Experiments on Open LORIS data and DRE dataset have been carried out for checking results in harsh conditions (sharp turns, textureless data, dynamic content). If a SLAM system could not initialise or lose trajectory we mark it as "X". The algorithms in the tables placed in the order of their release date. The examples of obtained trajectories of algorithms on EuRoC MAV mh5 and TUM VI corridor 1 sequence illustrated in Fig. 1.

Table 5 illustrates Basalt to be a winner on EuRoC V1_01_easy sequence and OpenVSLAM on MH_05_difficult. Worth reminding, "easy" EuRoC sequences differs from "difficult" one by jiggles and light conditions because of quadcopter data recording. OpenVSLAM showed better results on worse conditions. In general, most of the algorithms show low errors and comparable results. Also, we can notice a much bigger error of LDSO, the only Dense SLAM method that we look at.

Table 6 represents errors on TUM VI data sequence. Unfortunately, not all methods are capable to work with fish-eye data. The leader here is ORB-SLAM3 with errors, an order of magnitude less than others. It is worth noting, that this approach is 2 years newer than MapLab and VINS Mono - the second leaders, whereas all other algorithms in the table have similar release time. VINS Mono has two times higher error than MapLab on corridor_1 sequence whereas it significantly wins on room1 and magistrale_1 sequences. The length of magistrale_1 sequence is much bigger (around 1 km) than others which makes it especially difficult for the algorithms. Even though OpenVSLAM shows good results on other sequences, here it is used to get lost, which led to huge rotational errors. The important peculiarity of this dataset is a handheld camera and lots of motion which makes IMU data especially important here.

Table 7 represents errors on KITTI dataset. Here we see only the algorithms that could work with stereo data. ORB-SLAM2 and OpenVSLAM show pretty equal results with low errors (RMSE < 1 meter) on the majority of sequences. VINS Fusion and Basalt have several times worse results. Meanwhile, LDSO suffers much on KITTI sequences. We can notice that all methods have difficulties on sequence 02, which has an irregular shape and many smooth curves.

We summarised all the results on three main datasets in Fig. 5 and Fig. 6. Here we used one sequence from each dataset and normalized error values. In terms of translational error, ORB-SLAM3 shows its superiority on three datasets. MapLab and LDSO on the con-

Table 5: Results of SLAM algorithms tested on EUROC dataset

| Framework | V1_01 | | | | MH_05 | | | |
|---|---|---|---|---|---|---|---|---|
| | position(m) | | rotation(deg) | | position(m) | | rotation(deg) | |
| | max | rms | max | rms | max | rms | max | rms |
| ORB-SLAM2 | 0.16 | 0.08 | 13.97 | *4.70* | 0.15 | 0.05 | 20.93 | 6.60 |
| MapLab | 0.30 | 0.15 | 10.03 | 6.13 | 1.13 | 0.59 | 2.99 | 1.27 |
| LDSO | 0.16 | 0.08 | 10.34 | **4.68** | 1.49 | 0.08 | 28.35 | 20.45 |
| VINS-Mono | 0.21 | 0.09 | 6.77 | 6.09 | 0.48 | 0.30 | 2.58 | *0.79* |
| VINS-Fusion | 0.24 | 0.07 | 10.19 | 5.60 | 0.40 | 0.19 | 4.81 | 1.87 |
| OpenVSLAM | 0.15 | 0.08 | 14.13 | 4.73 | *0.14* | **0.04** | 20.85 | 6.54 |
| Basalt | **0.05** | **0.03** | **5.99** | 5.36 | 0.17 | 0.09 | **1.42** | **0.67** |
| Kimera | 0.08 | 0.04 | 6.34 | 5.59 | 0.58 | 0.34 | 3.82 | 2.34 |
| OpenVINS | 0.11 | 0.05 | *6.05* | 5.45 | 0.48 | 0.16 | 3.36 | 1.30 |
| ORB-SLAM3 | *0.07* | **0.03** | 6.39 | 5.94 | **0.11** | **0.04** | *1.68* | 0.93 |

Table 6: Results of SLAM algorithms tested on TUM VI dataset

| Framework | Corridor 1 | | | | Magistrale 1 | | | | Room 1 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | position(m) | | rotation(deg) | | position(m) | | rotation(deg) | | position(m) | | rotation(deg) | |
| | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms |
| MapLab | *0.78* | 0.27 | 4.18 | 2.85 | 2.72 | 0.88 | 13.92 | 6.37 | 0.27 | 0.14 | 6.05 | 2.10 |
| VINS-Mono | 1.74 | 0.44 | **1.40** | *0.47* | **0.35** | **0.24** | **3.64** | **1.06** | *0.12* | *0.05* | **3.79** | 1.74 |
| OpenVSLAM | 0.91 | *0.23* | 135.99 | 43.61 | 1.26 | 0.51 | 154.69 | 41.44 | 0.13 | 0.07 | 102.54 | 32.72 |
| Basalt | 1.34 | 0.33 | 4.25 | 2.32 | 1.91 | 0.68 | *5.28* | *1.47* | 0.23 | 0.09 | 5.92 | *1.25* |
| OpenVINS | 1.22 | 0.36 | 5.62 | 3.43 | 2.68 | 0.88 | 30.25 | 20.02 | 0.20 | 0.07 | 9.69 | 5.25 |
| ORB-SLAM3 | **0.02** | **0.01** | *1.48* | **0.45** | *0.57* | *0.33* | 6.66 | 3.69 | **0.02** | **0.01** | *5.74* | **0.42** |

Table 7: Results of SLAM algorithms tested on KITTI dataset

| Framework | 00 | | | | 02 | | | | 05 | | | | 06 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | position | | rotation | | position | | rotation | | position | | rotation | | position | | rotation | |
| | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms |
| ORB-SLAM2 | 2.82 | *0.89* | **6.67** | *0.74* | *15.76* | 6.60 | 4.37 | 1.47 | *0.94* | *0.39* | *2.16* | *0.38* | *1.14* | *0.62* | 0.76 | **0.39** |
| LDSO | 24.18 | 10.91 | 9.66 | 2.40 | 91.37 | 22.30 | 5.56 | 2.41 | 14.23 | 4.47 | 2.37 | 0.87 | 25.92 | 13.69 | 0.84 | 0.59 |
| VINS-Fusion | 13.80 | 5.20 | 7.93 | 3.22 | 47.90 | 18.17 | 9.31 | 4.78 | 16.78 | 4.79 | 7.49 | 3.02 | 7.14 | 2.50 | 6.93 | 4.81 |
| OpenVSLAM | **2.41** | **0.88** | 6.94 | 0.75 | **12.88** | **5.53** | **4.18** | **1.20** | 1.19 | 0.46 | **2.15** | *0.38* | 1.31 | 0.79 | 1.20 | 0.69 |
| Basalt | 5.09 | 3.38 | 7.25 | 0.88 | 19.49 | 9.11 | 4.88 | 1.72 | 4.96 | 2.48 | 2.30 | 0.86 | 4.00 | 2.11 | 3.17 | 2.78 |
| ORB-SLAM3 | *2.57* | 0.90 | *6.85* | **0.73** | 16.27 | *6.35* | *4.27* | *1.43* | **0.76** | **0.36** | 2.18 | **0.37** | **0.67** | **0.39** | 0.76 | *0.41* |

trary have the worst results. Other approaches demonstrate even performance. ORB-SLAM2 and OpenVS-LAM show low errors on KITTI, and Basalt on EUROC. Rotational errors on EUROC data is on a similar level of magnitude for all the algorithms. This could be explained by the nature of data from a micro aerial vehicle. ORB-SLAM3, Basalt, Vins-Mono and MapLab have the lowest rotational errors on the TUM VI sequence. Whereas OpenVSLAM the only visual method on TUM VI shows the highest error. Most of the approaches show comparable errors on KITTI except LDSO and VINS-FUSION.

Experiments on the Open LORIS dataset aim to compare successful algorithms from previous runs that use data from different types of sensors. We have monocular RGB-D ORB-SLAM2, stereo VINS-Mono with IMU data, and OpenVSLAM with images from a stereo fish-eye camera but without IMU. Table 8 illustrates that, in general, all three algorithms have similar performance on several sequences. But if we will look at details, we can notice that on the most crowded and dynamic sequences (cafe2 and home1), ORB-SLAM2 is losing trajectory at some moment. Moreover, monocular VINS mono overall has higher errors than stereo fish-eye OpenVSLAM. LDSO could not initialise on these data. Also, the wheel encoder data is in an inappropriate data format for the DRE-SLAM.

In Table 9 we compared results of ORB-SLAM2 (mono and rgb-d versions) with DRE-SLAM on DRE dataset with dynamic content. We can clearly see that DRE SLAM (Yang et al., 2019) gives the best results both on low dynamic and high dynamic sequences in comparison with ORB-SLAM2.
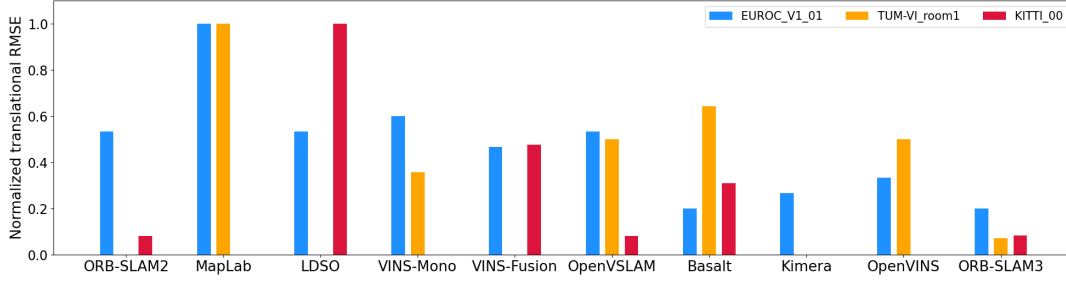
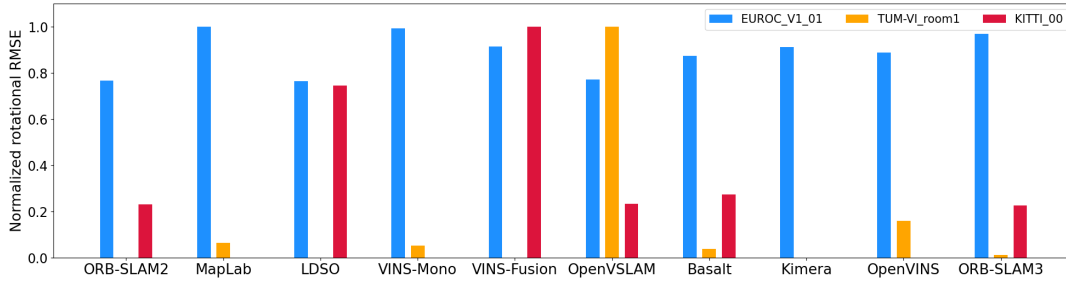Fig. 5: Comparison of translational errors on EuRoC, TUM VI and KITTI datasets



Fig. 6: Comparison of rotational errors on EuRoC, TUM VI and KITTI datasets

Table 8: Results of SLAM algorithms tested on Open LORIS dataset

| Framework | office5 | | | | office7 | | | | cafe1 | | | | cafe2 | | | | home1 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | position | | rotation | | position | | rotation | | position | | rotation | | position | | rotation | | position | | rotation | |
| | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms | max | rms |
| **ORB-SLAM2** | 0.30 | **0.21** | 7.45 | 2.97 | **0.09** | **0.04** | 2.88 | **1.45** | 0.36 | 0.21 | 5.04 | 2.72 | X | X | X | X | X | X | X | X |
| **VINS-Mono** | 0.32 | 0.22 | 5.83 | 3.38 | 0.40 | 0.12 | 3.75 | 2.42 | 1.31 | 0.42 | 5.40 | 3.47 | 0.75 | **0.15** | **4.27** | **2.79** | 1.09 | 0.55 | 16.01 | 8.86 |
| **OpenVSLAM** | **0.28** | 0.23 | **5.12** | 1.26 | 0.10 | 0.05 | **2.86** | 1.55 | **0.14** | **0.10** | **3.35** | **0.77** | **0.42** | 0.22 | 6.37 | 4.35 | **0.50** | **0.36** | **9.36** | **2.30** |

Table 10 indicates CPU load and RAM memory usage. All algorithms were tested on the system with Intel Core i7 8565U CPU and 16 GB of RAM. VINS-Mono demands fewer resources than others on the TUM VI corridor (indoor) environment but rather mediocre results on the outdoor sequence. Basalt shows a small need for memory both on indoor and outdoor sequences. OpenVINS consumes the smallest amount of CPU resources on TUM VI outdoors. All other algorithms on the outdoors sequence show significant demand for memory. The reason for that is the length of the sequence, leading to a huge amount of visual features. Kimera shows itself to be quite sufficient in memory consumption and CPU load on Euroc MH_05 sequence. Results of ORB-SLAM2 and LDSO are comparable and quite demanding both in memory and CPU load. DRE-SLAM expectedly demands lots of resources even for relatively short sequence. The possible reason is the use of several modules for a bunch of tasks.

## 5 Conclusions and recommendations

Currently, there is no perfect open-source out-of-the-box Visual SLAM solution for different environments and conditions. It is always needed to tune algorithm parameters for the particular task, preprocess data, solve dependency and performance issues. Some approaches (ORB-SLAM2/3, OpenVSLAM) are close to the perfect ones from a practical point of view. However, we faced problems with using more data sources or additional challenges such as dynamic objects detection.

Visual-inertial approaches benefit from using IMU data and show competitive results on the considered datasets. VINS-Mono, Basalt and OpenVINS show quite competitive results. The drawback of VIO is the need for an additional data source, but from the analysis of datasets, we can conclude that it is common to have IMU data. At the same time, it can be seen that visual-only ORB-SLAM2, which was published

Table 9: Results of SLAM algorithms tested on DRE dataset

| Framework | Low Dynamic (ST2) | | | | High Dynamic (HD2) | | | |
|---|---|---|---|---|---|---|---|---|
| | position | | rotation | | position | | rotation | |
| | max | rms | max | rms | max | rms | max | rms |
| ORB-SLAM2 mono | 0.27 | 0.09 | 3.95 | 1.62 | 1.12 | 0.31 | 39.77 | 16.49 |
| ORB-SLAM2 rgb-d | 0.16 | 0.09 | 5.14 | 2.05 | 1.23 | 0.70 | 45.55 | 19.30 |
| DRE-SLAM | 0.05 | 0.01 | 2.51 | 0.85 | 0.26 | 0.06 | 19.85 | 4.94 |

Table 10: CPU Load and Memory Usage

| Framework | Dataset | Memory RAM (Mb) | | | CPU Load (%) | | |
|---|---|---|---|---|---|---|---|
| | | Average | Max | Std | Average | Max | Std |
| Maplab | TUM VI corridor | 526.5 | 721.7 | 36.1 | 85.2 | 523.2 | 34.9 |
| VINS-Mono | | 87.0 | 92.1 | 1.1 | 68.5 | 118.8 | 36.2 |
| Basalt | | 107.8 | 109.2 | 1.3 | 418.3 | 560.1 | 56.9 |
| OpenVINS | | 1686.2 | 3329.1 | 952.5 | 85.2 | 523 | 34.9 |
| **Maplab | TUM VI outdoors | 1587.5 | 4587.4 | 682.6 | 160.6 | 792.8 | 93.2 |
| VINS-Mono | | 2851.0 | 4524.2 | 1285.9 | 182.5 | 306.7 | 94.2 |
| OpenVSLAM | | 2639.6 | 4363.4 | 1071.3 | 230.6 | 310.5 | 38.8 |
| Basalt | | 136.5 | 161.5 | 8.0 | 537.7 | 676.9 | 57.9 |
| OpenVINS | | 6490.8 | 12819.5 | 3694.5 | 100.5 | 147.8 | 14.3 |
| ORB-SLAM2 | Euroc MH_05 | 705.9 | 869.9 | 102.7 | 187.2 | 249.8 | 37.0 |
| LDSO | | 946.3 | 1082.1 | 102.9 | 160.1 | 263.9 | 50.6 |
| Kimera | | 250.5 | 273.6 | 46.4 | 80.5 | 120.8 | 17.6 |
| DRE-SLAM | DRE dataset hd2 | 5010.7 | 6288.1 | 1777.0 | 200.0 | 303.5 | 77.3 |

four years ago, still can fight for 1-2 place. The same applies to the similar but refined OpenVSLAM, which got one of the best results on many datasets except TUM VI. Kimera represents a nice modular concept with great opportunities but low usability. Thus, it was hard to launch it on other datasets except for EuRoC. The only dense approach (LDSO) that we tested showed far worse results compared to sparse algorithms. The most modern ORB-SLAM3 shows its superiority in different setups and seems to be a nice option to continue experiments. It tends to be a leader for the majority of datasets. Unfortunately, we could not run many dynamic SLAM approaches because of various issues. DRE-SLAM showed perfect results on dynamic data with the help of wheel encoders. But the problem is that wheel encoders data is difficult to find in datasets and it is an additional constraint to the setup during real-life data recording. The overall conclusion of results is that for different data various approaches could show good results and there is no striking leader.

Another important point is that we found many SLAM solutions which have promising results in the paper but completely unusable GitHub repositories with a lack of support and examples on the benchmark datasets. This refers to the reproducibility crisis which we stated in the introduction part of the paper.

Support and relevance of the implementation are important, as well as detailed instructions by the authors. Quickly updating packages cause dependency and repeatability issues even for a three-year-old solution.

The same is true for datasets. The absence of a complete dataset that could allow comparing different types of algorithms makes it difficult to understand which approach is the best one. Moreover, there is a need to unify data organisation both in datasets and in SLAM algorithms. Different formats of the ground truth data add a mess to work. For example, the KITTI format without timestamps makes it unusable for comparison of trajectories in other formats because of timestamps absence. Nevertheless, there is a significant choice of datasets suitable for comparing distinctive types of algorithms or specific environmental conditions.

We can conclude that a user-friendly open-source slam solution is rather rare. One of them was OpenVS-LAM. But during the paper writing, it has been deleted from the GitHub[2] due to the concerns on similarities of the source code to the ORB-SLAM2. Indeed, it can be noted that the resulting trajectories and errors are

---

[2] https://github.com/xdspacelab/openvslam/wiki/Termination-of-the-release

very similar in these two algorithms. Even though the National Institute of Advanced Industrial Science and Technology did not find the copyright infringement incidences, they decided to terminate the release of Open-VSLAM to avoid any risk of possible copyright issues.

Finally, this would be very convenient if every research work or SLAM approach had a Docker container to run and check the results described in the proposed method as it might solve reproducibility issues. To prove the practical consistency of an algorithm, it would be helpful if the latter had a ROS wrapper to check the proposed method on a real robot. Also, this would be meaningful if robotics datasets had as many sensors as possible in order to check more approaches for a better understanding of benefits and disadvantages.

# References

S. Agarwal, K. Mierle, and Others. Ceres solver. http://ceres-solver.org.

R. Antoni, A. Marcus, C. Yun, and C. Luca. Kimera: an open-source library for real-time metric-semantic localization and mapping. In *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2020. URL https://github.com/MIT-SPARK/Kimera.

I. A. Bârsan, P. Liu, M. Pollefeys, and A. Geiger. Robust dense mapping for large-scale dynamic environments. In *International Conference on Robotics and Automation (ICRA)*, 2018.

B. Bescos, J. M. Facil, J. Civera, and J. Neira. Dynaslam: Tracking, mapping, and inpainting in dynamic scenes. *IEEE Robotics and Automation Letters*, 3(4):4076–4083, Oct 2018. ISSN 2377-3774. doi: 10.1109/lra.2018.2860039. URL http://dx.doi.org/10.1109/LRA.2018.2860039.

J.-L. Blanco-Claraco, F.-A. Moreno-Dueñas, and J. González-Jiménez. The málaga urban dataset: High-rate stereo and lidar in a realistic urban scenario. *Int. J. Rob. Res.*, 33(2):207–214, Feb. 2014. ISSN 0278-3649. doi: 10.1177/0278364913507326. URL https://doi.org/10.1177/0278364913507326.

M. Bloesch, J. Czarnowski, R. Clark, S. Leutenegger, and A. J. Davison. CodeSLAM - Learning a Compact, Optimisable Representation for Dense Visual SLAM. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2560–2568. IEEE, jun 2018. ISBN 978-1-5386-6420-9. doi: 10.1109/CVPR.2018.00271. URL http://arxiv.org/abs/1804.00874https://ieeexplore.ieee.org/document/8578369/.

M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, and R. Siegwart. The euroc micro aerial vehicle datasets. *The International Journal of Robotics Research*, 2016. doi: 10.1177/0278364915620033. URL http://ijr.sagepub.com/content/early/2016/01/21/0278364915620033.abstract.

A. P. Bustos, T.-J. Chin, A. Eriksson, and I. Reid. Visual SLAM: Why Bundle Adjust? In *2019 International Conference on Robotics and Automation (ICRA)*, volume 2019-May, pages 2385–2391. IEEE, may 2019. ISBN 978-1-5386-6027-0. doi: 10.1109/ICRA.2019.8793749. URL https://ieeexplore.ieee.org/document/8793749/.

C. Campos, R. Elvira, J. J. G/´omez, J. M. M. Montiel, and J. D. Tardós. ORB-SLAM3: An accurate open-source library for visual, visual-inertial and multi-map SLAM. *arXiv preprint arXiv:2007.11898*, 2020.

N. Carlevaris-Bianco, A. K. Ushani, and R. M. Eustice. University of michigan north campus long-term vision and lidar dataset. *Int. J. Rob. Res.*, 35(9):1023–1035, Aug. 2016. ISSN 0278-3649. doi: 10.1177/0278364915614638. URL https://doi.org/10.1177/0278364915614638.

L. Carlone, Z. Kira, C. Beall, V. Indelman, and F. Dellaert. Eliminating conditionally independent sets in factor graphs: A unifying perspective based on smart factors. *IEEE International Conference on Robotics and Automation (ICRA)*, 2014. URL https://ieeexplore.ieee.org/document/6907483.

C. Cesar, C. Luca, C. Henry, L. Yasir, S. Davide, N. Jose, R. Ian, and L. J. J. Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Transactions on Robotics*, 32(6):1309–1332, jun 2016. ISSN 15523098. doi: 10.1109/TRO.2016.2624754. URL http://arxiv.org/abs/1606.05830http://dx.doi.org/10.1109/TRO.2016.2624754.

C. Chen, B. Wang, C. X. Lu, N. Trigoni, and A. Markham. A Survey on Deep Learning for Localization and Mapping: Towards the Age of Spatial Machine Intelligence. *arXiv*, jun 2020. URL http://arxiv.org/abs/2006.12567.

R. Curnow. Chrony. URL https://chrony.tuxfamily.org.

I. Cvišić, J. Ćesić, I. Marković, and I. Petrović. SOFT-SLAM: Computationally efficient stereo visual simultaneous localization and mapping for autonomous unmanned aerial vehicles. *Journal of Field Robotics*, 35(4):578–595, jun 2018. ISSN 15564959. doi: 10.1002/rob.21762. URL http://doi.wiley.com/10.1002/rob.21762.

A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse. MonoSLAM: Real-Time Single Camera SLAM. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(6):1052–1067, jun 2007. ISSN 0162-8828. doi: 10.1109/TPAMI.2007.1049. URL `http://ieeexplore.ieee.org/document/4160954/`.

F. Dellaert and M. Kaess. Georgia tech smoothing and mapping (gtsam). *IEEE International Conference on Robotics and Automation (ICRA)*, 2012. URL `https://smartech.gatech.edu/handle/1853/45226`.

M. Dou, S. Khamis, Y. Degtyarev, P. Davidson, S. R. Fanello, A. Kowdle, S. O. Escolano, C. Rhemann, D. Kim, J. Taylor, et al. Fusion4d: Real-time performance capture of challenging scenes. *ACM Transactions on Graphics (TOG)*, 35(4):1–13, 2016.

J. Engel, J. Sturm, and D. Cremers. Semi-dense visual odometry for a monocular camera. In *Proceedings of the IEEE international conference on computer vision*, pages 1449–1456, 2013.

J. Engel, T. Schöps, and D. Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014.

J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):611–625, 2017.

J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, mar 2018.

G. Falco, M. Pini, and G. Marucco. Loose and tight gnss/ins integrations: Comparison of performance assessed in real urban scenarios. *Sensors*, 17(2), 2017. ISSN 1424-8220. doi: 10.3390/s17020255. URL `https://www.mdpi.com/1424-8220/17/2/255`.

M. A. Fischler and R. C. Bolles. Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography. In *Readings in Computer Vision*, pages 726–740. Elsevier, jan 1987. doi: 10.1016/b978-0-08-051581-6.50070-2.

C. Forster, M. Pizzoli, and D. Scaramuzza. Svo: Fast semi-direct monocular visual odometry. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 15–22. IEEE, 2014.

C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza. On-manifold preintegration theory for fast and accurate visual-inertial navigation. *IEEE Trans. Robotics: 33(1):1-21*, 2016. URL `https://arxiv.org/abs/1512.02363`.

J. Fritsch, T. Kuehnl, and A. Geiger. A new performance measure and evaluation benchmark for road detection algorithms. In *International Conference on Intelligent Transportation Systems (ITSC)*, 2013.

D. Galvez-López and J. D. Tardos. Bags of Binary Words for Fast Place Recognition in Image Sequences. *IEEE Transactions on Robotics*, 28(5):1188–1197, oct 2012. ISSN 1552-3098. doi: 10.1109/TRO.2012.2197158. URL `http://ieeexplore.ieee.org/document/6202705/`.

X. Gao, R. Wang, N. Demmel, and D. Cremers. LDSO: Direct Sparse Odometry with Loop Closure. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2198–2204. IEEE, oct 2018. ISBN 978-1-5386-8094-0. doi: 10.1109/IROS.2018.8593376. URL `https://ieeexplore.ieee.org/document/8593376/`.

A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.

P. Geneva, K. Eckenhoff, W. Lee, Y. Yang, and G. Huang. Openvins: A research platform for visual-inertial estimation. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4666–4672, 2020. doi: 10.1109/ICRA40945.2020.9196524.

M. Grupp. evo: Python package for the evaluation of odometry and slam. `https://github.com/MichaelGrupp/evo`, 2017.

B. K. Horn. Closed-form solution of absolute orientation using unit quaternions. *Journal of the Optical Society of America A*, 4(4):629–642, 1987. doi: 10.1364/JOSAA.4.000629. URL `http://people.csail.mit.edu/bkph/papers/Absolute-OCR.pdf`.

C. Houseago, M. Bloesch, and S. Leutenegger. Ko-fusion: dense visual slam with tightly-coupled kinematic and odometric tracking. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 4054–4060. IEEE, 2019.

J. Huai, Y. Qin, F. Pang, and Z. Chen. Segway drive benchmark: Place recognition and slam data collected by a fleet of delivery robots, 2019.

J. Jeong, Y. Cho, Y.-S. Shin, H. Roh, and A. Kim. Complex urban dataset with multi-level sensors from highly diverse urban environments. *The International Journal of Robotics Research*, page 0278364919843996, 2019.

A. Kasar. Benchmarking and comparing popular visual slam algorithms. *arXiv*, 2018.

G. Klein and D. Murray. Parallel Tracking and Mapping for Small AR Workspaces. In *2007 6th IEEE and ACM International Symposium on Mixed*

*and Augmented Reality*, pages 1–10. IEEE, nov 2007. ISBN 978-1-4244-1749-0. doi: 10.1109/ISMAR. 2007.4538852. URL http://ieeexplore.ieee.org/document/4538852/.

P. Kopanev. Slam-dockers. https://github.com/KopanevPavel/SLAM-Dockers, 2021.

R. Kummerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard. G2o: A general framework for graph optimization. In *2011 IEEE International Conference on Robotics and Automation*, pages 3607–3613. IEEE, may 2011. ISBN 978-1-61284-386-5. doi: 10.1109/ICRA.2011.5979949. URL http://ieeexplore.ieee.org/document/5979949/.

Y. Liu, Y. Fu, F. Chen, B. Goossens, W. Tao, and H. Zhao. Datasets and Evaluation for Simultaneous Localization and Mapping Related Problems: A Comprehensive Survey. *arXiv*, feb 2021. URL http://arxiv.org/abs/2102.04036.

W. Maddern, G. Pascoe, M. Gadd, D. Barnes, B. Yeomans, and P. Newman. Robust real-time visual odometry with a single camera and an imu. *Proceedings of The British Machine Vision Conference (BMVC), Dundee, UK*, 2011. URL https://www.researchgate.net/publication/256666913_Robust_Real-Time_Visual_Odometry_with_a_Single_Camera_and_an_IMU.

W. Maddern, G. Pascoe, C. Linegar, and P. Newman. 1 Year, 1000km: The Oxford RobotCar Dataset. *The International Journal of Robotics Research (IJRR)*, 36(1):3–15, 2017. doi: 10.1177/0278364916679498. URL http://dx.doi.org/10.1177/0278364916679498.

A. L. Majdik, C. Till, and D. Scaramuzza. The zurich urban micro aerial vehicle dataset. *Int. J. Rob. Res.*, 36(3):269–273, Mar. 2017. ISSN 0278-3649. doi: 10.1177/0278364917702237. URL https://doi.org/10.1177/0278364917702237.

R. Mur-Artal and J. D. Tardós. ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, 2017. doi: 10.1109/TRO.2017.2705103.

R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós. ORB-SLAM: a versatile and accurate monocular SLAM system. *IEEE Transactions on Robotics*, 31(5):1147–1163, 2015. doi: 10.1109/TRO.2015.2463671.

R. A. Newcombe, A. J. Davison, S. Izadi, P. Kohli, O. Hilliges, J. Shotton, D. Molyneaux, S. Hodges, D. Kim, and A. Fitzgibbon. KinectFusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, pages 127–136. IEEE, oct 2011a. ISBN 978-1-4577-2185-4. doi: 10.1109/

ISMAR.2011.6092378. URL http://ieeexplore.ieee.org/document/6162880/.

R. A. Newcombe, S. J. Lovegrove, and A. J. Davison. DTAM: Dense tracking and mapping in real-time. In *2011 International Conference on Computer Vision*, pages 2320–2327. IEEE, nov 2011b. ISBN 978-1-4577-1102-2. doi: 10.1109/ICCV.2011.6126513. URL http://ieeexplore.ieee.org/document/6126513/.

J. Nikolic, J. Rehder, M. Burri, P. Gohl, S. Leutenegger, P. T. Furgale, and R. Siegwart. A synchronized visual-inertial sensor system with fpga preprocessing for accurate real-time slam. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 431–437, 2014. doi: 10.1109/ICRA.2014.6906892.

H. Oleynikova, Z. Taylor, M. Fehr, R. Siegwart, and J. Nieto. Voxblox: Incremental 3d euclidean signed distance fields foron-board mav planning. *arXiv*, 2017. doi: 10.1109/IROS.2017.8202315.

S. Poddar, R. Kottath, and V. Karar. Evolution of visual odometry techniques, 2018.

T. Qin and S. Shen. Robust initialization of monocular visual-inertial estimation on aerial robots. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, volume 2017-Septe, pages 4225–4232. IEEE, sep 2017a. ISBN 978-1-5386-2682-5. doi: 10.1109/IROS.2017.8206284. URL https://ieeexplore.ieee.org/document/8206284/.

T. Qin and S. Shen. Robust initialization of monocular visual-inertial estimation on aerial robots. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, volume 2017-Septe, pages 4225–4232. IEEE, sep 2017b. ISBN 978-1-5386-2682-5. doi: 10.1109/IROS.2017.8206284. URL https://ieeexplore.ieee.org/document/8206284/.

T. Qin, P. Li, and S. Shen. Vins-mono: A robust and versatile monocular visual-inertial state estimator. *IEEE Transactions on Robotics*, 34(4):1004–1020, 2018. doi: 10.1109/TRO.2018.2853729.

T. Qin, S. Cao, J. Pan, and S. Shen. A general optimization-based framework for global pose estimation with multiple sensors, 2019.

A. Rosinol, T. Sattler, M. Pollefeys, and L. Carlone. Incremental visual-inertial 3d mesh generation with structural regularities. *IEEE Int. Conf. on Robotics and Automation (ICRA)*, 2019. URL https://arxiv.org/pdf/1903.01067.

E. Rublee, V. Rabaud, K. Konolige, and G. Bradski. ORB: An efficient alternative to SIFT or SURF. In *2011 International Conference on Com-*

*puter Vision*, pages 2564–2571. IEEE, nov 2011. ISBN 978-1-4577-1102-2. doi: 10.1109/ICCV.2011. 6126544. URL http://ieeexplore.ieee.org/ document/6126544/.

M. Rünz and L. Agapito. Co-fusion: Real-time segmentation, tracking and fusion of multiple objects. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4471–4478. IEEE, 2017.

T. Schneider, M. Dymczyk, M. Fehr, K. Egger, S. Lynen, I. Gilitschenski, and R. Siegwart. Maplab: An open framework for research in visual-inertial mapping and localization. *IEEE Robotics and Automation Letters*, PP:1–1, 01 2018. doi: 10.1109/LRA. 2018.2800113.

D. Schubert, T. Goll, N. Demmel, V. Usenko, J. Stückler, and D. Cremers. The tum vi benchmark for evaluating visual-inertial odometry. *arXiv*, 2018. doi: 10.1109/IROS.2018.8593419.

D. Schubert, N. Demmel, L. von Stumberg, V. Usenko, and D. Cremers. Rolling-shutter modelling for direct visual-inertial odometry. *arXiv preprint arXiv:1911.01015*, 2019.

X. Shi, D. Li, P. Zhao, Q. Tian, Y. Tian, Q. Long, C. Zhu, J. Song, F. Qiao, L. Song, Y. Guo, Z. Wang, Y. Zhang, B. Qin, W. Yang, F. Wang, R. H. M. Chan, and Q. She. Are we ready for service robots? the OpenLORIS-Scene datasets for lifelong SLAM. In *2020 International Conference on Robotics and Automation (ICRA)*, pages 3139–3145, 2020.

H. Strasdat, J. M. Montiel, and A. J. Davison. Real-time monocular SLAM: Why filter? In *Proceedings - IEEE International Conference on Robotics and Automation*, pages 2657–2664, 2010. ISBN 9781424450381. doi: 10.1109/ROBOT.2010.5509636.

J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A benchmark for the evaluation of rgb-d slam systems. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 573–580. IEEE, 2012a.

J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A benchmark for the evaluation of rgb-d slam systems. In *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012b.

S. Sumikura, M. Shibuya, and K. Sakurada. OpenVSLAM: A Versatile Visual SLAM Framework. In *Proceedings of the 27th ACM International Conference on Multimedia*, MM '19, pages 2292–2295, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6889-6. doi: 10.1145/3343031.3350539. URL http://doi.acm.org/10.1145/3343031.3350539.

K. Tateno, F. Tombari, I. Laina, and N. Navab. CNN-SLAM: Real-time dense monocular SLAM with learned depth prediction. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017-Janua:6565–6574, apr 2017. doi: 10.1109/CVPR.2017.695. URL http://ieeexplore.ieee.org/document/8100178/http://arxiv.org/abs/1704.03489.

V. Usenko, N. Demmel, D. Schubert, J. Stückler, and D. Cremers. Visual-inertial mapping with non-linear factor recovery. *arXiv*, 2019. doi: 10.1109/LRA.2019. 2961227.

T. Whelan, R. F. Salas-Moreno, B. Glocker, A. J. Davison, and S. Leutenegger. ElasticFusion: Real-time dense SLAM and light source estimation. *The International Journal of Robotics Research*, 35(14): 1697–1716, dec 2016. ISSN 0278-3649. doi: 10. 1177/0278364916669237. URL http://journals. sagepub.com/doi/10.1177/0278364916669237.

D. Yang, S. Bi, W. Wang, C. Yuan, X. Qi, and Y. Cai. Dre-slam: Dynamic rgb-d encoder slam for a differential-drive robot. *Remote Sensing*, 11(4):380, 2019.

C. Yu, Z. Liu, X.-J. Liu, F. Xie, Y. Yang, Q. Wei, and Q. Fei. Ds-slam: A semantic visual slam towards dynamic environments. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1168–1174. IEEE, 2018.

Z. Zhang, S. Liu, G. Tsai, H. Hu, C.-C. Chu, and F. Zheng. Pirvs: An advanced visual-inertial slam system with flexible sensor fusion and hardware co-design, 2017.