

vSlam-Zero

Partenza, Garrett
Khoury College of Computer Sciences
Northeastern University
Boston, USA
partenza.g@northeastern.edu

Abstract—SLAM-Zero is a toy implementation of visual-based monocular simultaneous localization and mapping, where an agent builds a three-dimensional map of its environment while moving through it. This project represents an effort to grasp the fundamental mathematical foundations behind epipolar geometry while gaining experience with open-cv in python.

Index Terms—simultaneous localization and mapping, visual odometry, computer vision

I. INTRODUCTION

SLAM stands for simultaneous localization and mapping, where an agent attempts to construct or update a map of an unexplored environment whilst simultaneously keeping track of the agent's location within it. The inspiration for this project came after watching George Hotz's eight-hour long Twitch stream where he developed a similar implementation. Major resources consulted during the development process notably include a visual slam textbook open source on GitHub, both original ORB-SLAM papers, and the Stanford epipolar geometry notes [1] [2] [3] [4]. The remainder of this article discusses related works, my methodology and the fundamental mathematics of visual-based monocular SLAM, experimental results, relevant files within the code repository, as well as a conclusion focusing on the projects weaknesses and future improvements.

II. RELATED WORK

The pinnacle paper within this domain of related works is the original Orb-Slam papers [2] [3], which was arguably the first viable feature based visual monocular SLAM system for versatile environments. Furthermore, the authors methodology achieved state-of-the-art camera localization precision all the while using the same features for tracking, mapping, relocalization and loop closing. The authors combination of a broad array of popular computer vision techniques served as a foundation for my initial research stages.

While the Orb-Slam papers provided a good high-level understanding of monocular vision based SLAM, a more granular understanding was necessary for implementing SLAM from scratch in Python. In this regard, arguably the most helpful related work in this field used during the development of this project was a Chinese computer vision textbook entitled Basic Knowledge on Visual SLAM: From Theory to Practice [1]. While the original publication uses C++ to provide examples, the authors go into mathematical depths for each step of the

SLAM process, specifically focusing on a visual monocular based approach.

III. METHODOLOGY

The algorithm in my approach follows that of a three step approach. First, given a pair of images taken sequentially of the environment, generate feature matches. Second, use the feature matches to calculate the pose transformation, notably the rotation and translation of the camera. Finally, use the pose transformation to triangulate the relevant key points. These three steps are described in greater detail in the following paragraphs.

A. Feature Matching

Feature matching is the process of detecting key pixels in two images and using descriptors to find corresponding matches. Features are pairs of key points and descriptors, where key points are the two dimensional image coordinates and descriptors are a feature vector which attempts to describe the key point in n-dimensional space. Key points are detected using good features to track, where key points are assumed to be the prominent corners in an image [5]. Next, orb features are computed using rotated BRIEF algorithm as in [6]. Matches between key points within both images are derived by brute force with the KNN algorithm, where the top two matches are retained. The ratio test is applied to the top two matches in order increase quality as described in [7].

B. Pose Transformation

Pose transformation is the process of using matches between two images to find the camera rotation and translation. To do this, the fundamental matrix is calculated using the eight-point algorithm, whereby its derivation can be reduced to solving a system of linear equations given eight or more matches. This process can be described in the following equalities, where F is the fundamental matrix and x is a set of corresponding key points.

$$x^T F x = 0 \quad (1)$$

$$F = \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \quad (2)$$

$$uu'f_{11} + vu'f_{12} + uf_{31} + vv'f_{22} + uv'f_{23} + vf_{32} + ff_{33} = 0 \quad (3)$$

Next, the intrinsic parameter matrix is used in combination with the fundamental matrix to calculate the essential matrix.

$$K^T \cdot F \cdot K \quad (4)$$

Finally, the rotation matrix R and translation vector t are calculated by decomposing the essential matrix and conducting a clarity check to ensure all points have positive depth. Together R and t create a tuple that performs a change of basis from the first camera's coordinate system to the second camera's coordinate system.

C. Triangulation

Triangulation in visual-based monocular slam involves using R , t , to find the three dimensional world point of a key point matched between two images. To do this, the projection matrix is created by horizontally stacking R and t , before using least squares to cast the image point into three dimensional space as in [8]. Finally, in order to improve the mapping quality, the reprojection error of each world point is calculated by reprojecting the point back onto the image frame. Points that create large reprojection errors are masked.

IV. EXPERIMENTS AND RESULTS

Empirically evaluating the results of my system pose a challenge. Being that my system is concerned with the mapping of environments from driving videos of dash cameras mounted onto cars, we do not know the exact world points of features within the vehicles frame of view. As a result, we can not calculate a metric difference between my systems predictions and ground truth. While it would be possible to use a labeled SLAM dataset unrelated to driving footage that does have world point ground truths, I was not able to familiarize, download, and integrate one of these datasets within the short time span of this project. As a result, evaluation of my system was done visually, whereby features of the last frame and current frame are plotted onto the original image, and a three-dimensional scatter plot is displayed side-by-side and paused for review at each frame. An example of one of these plot pairs are show below in Figure 1 and Figure 2. In Figure 1, the green dots represent the current frame feature and the blue dots represent the previous frame feature. Lines are drawn between matching features (note these lines are small because the feature matching process of my system is very good, and thus there are not many outliers). In Figure 2, the single green dot represents the current camera location and the blue dots represent the 3d world coordinates of the current frame with the origin at the camera. It is clearly visible that the drawn world points within the three-dimensional plot correspond nicely to the detected features in the image frame.

V. RELEVANT FILES

Below outlines the relevant files in my code repository along with their functionalities.

- 1) slam.py
 - a) reads video and executes slam



Fig. 1. Current image frame with plotted feature matches.

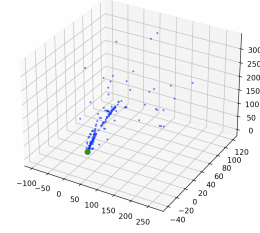


Fig. 2. Current image frame converted to three dimensions using vSlam-Zero.

- b) plots map
- 2) extract.py
 - a) calculate key points and descriptors
 - b) retain previous frame key points and descriptors
- 3) pose.py
 - a) calculate R and t
- 4) triangle.py
 - a) triangulate world points

VI. SUMMARY

Considering the complexities of epipolar geometry and the limited time I had to implemented the project, I am overall satisfied with my algorithms performance. One of the weaknesses of my implementation is that it does not include a back end optimization such as in the original ORB-SLAM paper. As a result, if world points are continuously plotted without resetting the axis after each frame, pose transformation error grows linearly and results in a erroneous mapping. Given more time, this would be my focus for improvement.

REFERENCES

- [1] Xiang Gao, Tao Zhang, Yi Liu, and Qinrui Yan. 14 lectures on visual slam: From theory to practice. *Publishing House of Electronics Industry*, 2017.
- [2] Raul Mur-Artal, J. M. M. Montiel, and Juan D. Tardos. Orb-slam: A versatile and accurate monocular slam system. *IEEE Transactions on Robotics*, 31(5):1147–1163, Oct 2015.
- [3] Raul Mur-Artal and Juan D. Tardos. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, Oct 2017.
- [4] Kenji Hata and Silvio Savarese. Computer science 231 a course notes 3: Epipolar geometry. 2011.
- [5] Jianbo Shi and Carlo Tomasi. Good features to track. pages 593–600, 1994.

- [6] Michael Calonder, Vincent Lepetit, Christoph Strecha, and Pascal Fua. Brief: Binary robust independent elementary features. In Kostas Daniilidis, Petros Maragos, and Nikos Paragios, editors, *Computer Vision – ECCV 2010*, pages 778–792, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.
- [7] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91–110, November 2004.
- [8] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004.