

Static ORB Extractor for Visual SLAM

1st Chethan Chinder Chandrappa
Electrical and Computer Engineering
University of California Los Angeles
Los Angeles, USA
chinderc@g.ucla.edu

2nd Srinath Naik jmeera
Computer Science
University of California Los Angeles
Los Angeles, USA
srinath@cs.ucla.edu

Abstract—Despite significant progress in development of feature based visual SLAM systems, they exhibit poor performance when dynamic objects are present in the scene. This under performance can be attributed to using all the extracted feature points for location and pose estimation irrespective of whether they correspond to static/dynamic objects. In presence of dynamic objects, the calculated projection matrix to map from one frame to the next will be inaccurate due to excessive shift of features, thereby leading to poor estimation.

In this paper, we propose Reprojection Error Map for finding the dynamic mask in the scene using the ORB keypoint correspondences. We show the 3 methodologies in the design process and their drawbacks. We experimented all our approaches on dynamic object based scenes in TUM and KITTI datasets.

Index Terms—vSLAM, dynamic objects, ORB Keypoints, Reprojection Error Map.

I. INTRODUCTION

SLAM is the core of most of the autonomous robotic systems which move and interact with the environment. Accurate localization and mapping are essentials for safe and efficient functioning of these systems. Over the years, feature based visual SLAM systems[1] were developed, which can be used on plethora of visual SLAM tasks in real time due to their fast performance and easy adaptation. Recently, Deep Learning based visual odometry systems like DeepVO[2], UndeepVO[3], TartanVO[4] were built by training on huge data-sets. However, these are limited to few use cases as they are trained for some specific scenarios and are slow in general. This leaves the feature based SLAM methods remain state of the art till date.

Feature based SLAM methods extract key features in each frame, compute correspondences in subsequent frames and use projective transformation to estimate the relative camera pose and location. Although these methods achieve considerable performance in some of the applications, it still remains a challenge to create SLAM robust to dynamic environments. In dynamic environments, the inherent movements of the objects in the scene will cause an irregular shift in the key features corresponding to dynamic objects, resulting in an inaccurate projection matrix, leading to poor localization. Moreover, SLAM systems build a global map of the environment as time progresses, the map points corresponding to dynamic objects will accumulate, consuming space and increasing time for computations which are critical in real time systems.

Also, as the same environment in a particular frame might have different dynamic objects at various times, it is possible that loop closure detection will fail, resulting in improper optimization.

We have tested ORB-SLAM2 on few sequences of TUM dataset which contain dynamic objects. As it can be clearly seen from Fig. 1 and Fig. 2, which depict Absolute Trajectory Error and Relative Pose Error respectively, the estimated trajectory and poses significantly differ from the ground truth. Building a SLAM system robust to dynamic objects will improve tracking accuracy and create better or more representative map of the environment. Our contribution in this paper is two fold. First we experiment with different methodologies to generate a dynamic mask in a dynamic scene. Then we integrate the above system into ORB-SLAM2 front end to extract only static ORB features from each frame. Finally, we evaluate our approach and end-to-end system on dynamic object based scenes.

In the rest of this paper, the structure is as follows. Section II provides details regarding some recent work relevant to this area. Section III provides a detailed explanation on various approaches/design choices we have made as we progressed towards the objective of creating a robust dynamic object detector. Then, section IV provides details of our experimentation and results. Finally in section V we give a conclusion and possible directions for improvements and further research.

II. RELATED WORK

Creating SLAM robust to dynamic objects in the scene is an open research problem and the approaches for tackling it can be classified into two. Traditional methods in which pixel motion is estimated from video using geometric constraints such as reprojection, optical flow etc. to decide upon dynamic areas. [5] relies on EpicFlow[6] to get dense pixel matching between frames, then it uses a probabilistic inference to conclude the dynamic nature based on prior state and motion estimation in current frame. This method will have advantage on abrupt movements of camera which are of short time and does not contribute directly to the object motion. However, optical flow methods suffer from variations in brightness across the frames.

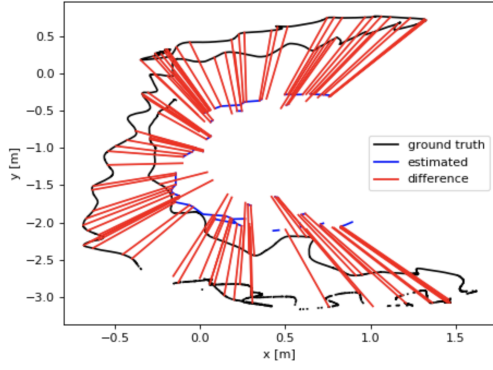


Fig. 1. Absolute Trajectory Error in a dynamic object scenario from TUM dataset

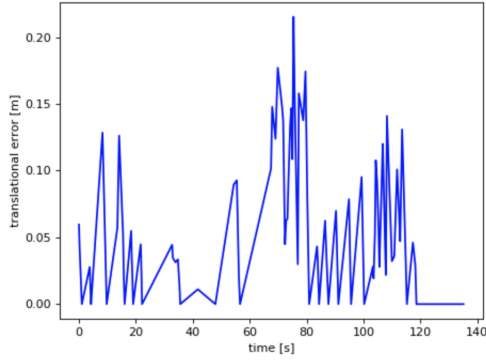


Fig. 2. Relative Pose Error in a dynamic object scenario from TUM dataset

Recently, we have Deep Learning based methods where pre-trained networks are used to segment the potential dynamic objects(eg: car, human) in the frame and then decide if these regions are actually moving. DS-SLAM[7] adopts SegNet[8] to get pixel wise semantic segmentation followed by moving consistency check between frames to determine dynamic points. Also, it creates a semantic map of the environment. DOT[9] uses Detectron2[10] network for semantic segmentation to generate dynamic object masks. Then, it performs consistency check between frames to decide upon which of these objects are actually moving. Finally, it also creates new masks based on the dynamic object motion estimation which is termed as 'Mask Propagation'. A major limitation of these methods are that they are not flexible and can only detect predetermined objects as dynamic. Moreover, they are slow in practice due to heavy computation.

Most of the above mentioned work integrates into front end of ORB-SLAM2 and report an improved performance on dynamic object scenarios of various data-sets.

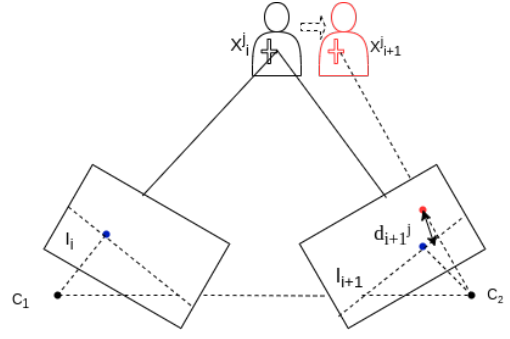


Fig. 3. Depiction of distance between epiline and corresponding keypoint

III. DESIGN PROCESS

A. Definitions

We denote I_i as a frame at an instant i . The temporal index is denoted by i . Each individual keypoint is referred by an index j . The keypoints used in the frame i are denoted by K_i^j . The fundamental matrix calculated between the frame F_i and F_{i+1} keypoints is denoted by $Fmat_i$. Similarly, the homography between the successive frames is denoted by H_i .

B. Method 1: Based on distance between epilines and keypoints

The feature keypoints and descriptors for the interest points are found using the ORB. The correspondences between the successive frames were found by applying the Flann based matcher on respective keypoints and descriptors. Each correspondence contains pair of Keypoints (K_i^j, K_{i+1}^j)

The 3D keypoint X_t^j is mapped on the F_i as K_i^j and on F_{i+1} as K_{i+1}^j . We can leverage the epipolar constraints on these keypoints for finding the dynamic keypoints. As shown in Fig. 3, distance d_{i+1}^j between the epiline l_{i+1}^j on frame F_{i+1} , calculated from corresponding keypoints K_i^j and Fundamental matrix $Fmat_i$, and corresponding keypoint K_{i+1}^j on frame F_{i+1} . If this distance d_j between the corresponding keypoints K_i^j and K_{i+1}^j in the consecutive frames is significant, we can consider the keypoint K_{i+1}^j on frame F_j belongs to dynamic object. We can see in the figure that a change in the position of the keypoint increases the distance d_{i+1}^j between the keypoint K_{i+1}^j and epiline l_{i+1}^j . If the distance is greater than the set threshold τ , it is considered as a keypoint belonging to a dynamic object and marked as an outlier or else an inlier. The mathematical formulation for the above defined epipolar constraint is defined below,

$$K_i^j F_i K_{i+1}^j > \tau$$

The Inlier set is then used for finding the camera poses $[R|T]$.

C. Method 2: Based on Motion Vector Clustering

In the initial frame F_i , A good set of keypoints K_i^j is found using Harris corner detectors. These keypoints were

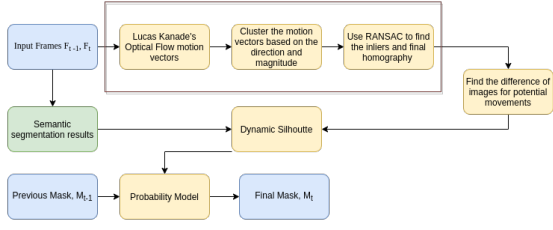


Fig. 4. Increase image size :: pipeline for finding the dynamic mask using Motion vector clustering

tracked in the successive frames using the Iterative Lucas-Kanade method. In optical flow, an assumption of brightness consistency across the temporal axis is made. The following equation represents the optical flow equation applied on successive frames.

$$\nabla I \cdot \vec{V} = -I_t$$

Here, \vec{V} represents the motion flow in second frame, ∇I represents the spatial intensity gradient, and I_t is Intensity gradient along the temporal axis. we have two unknowns (V_x, V_y) and one equation. we need more equations to solve this problem. To solve this ambiguity, Lucas Kanades algorithm assumes the local neighbourhood of a pixel have the same optical flow. After having equations for the local neighborhood of the pixels, least squares can be used to solve the flow vectors and where the keypoints are in the successive frames. Likewise we can track all the keypoints in successive frames, and track them.

After tracking these keypoints in the successive frames, a new set of keypoint correspondences (K_i^j, K_{i+1}^j) were found which are present in both frames.

As described in the fig 2, Motion vectors \vec{v}_i^j are calculated using these pair of keypoint correspondences K_i^j, K_{i+1}^j . Each motion vector consist of direction and magnitude. If camera is moving, we can see that most of the motion vectors have the similar direction and magnitude of the motion vector depends on the distance of the object from the camera for a rigid body transformation. After getting the set of motion vectors $\{\vec{v}_i^j\}$ between correspondences, we will find the clusters among the motion vectors which are similar. We used K-means clustering to find the similar of motion vectors of cluster. The number of clusters can basically set to number of objects which can be moved. We will obtain the number of objects in the scene after panoptic segmentation of the scenes. We clustered the motion vectors based on the number of objects. We will randomly pick the clusters and use of the keypoints related to the motion vectors in the cluster to find the homography H_i^j . After repeating the above steps for some iterations, from the set of homographies H_i^j which are obtained from random clusters, we will find the best homography, inliers and outliers using ransac method to find the best homography which minimizes reprojection error between keypoints.

After finding the best homography, We will apply this homography on the first frame F_i to get projected frame F'_i .

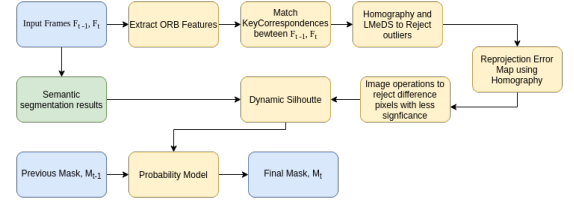


Fig. 5. Increase size::Pipeline for finding the dynamic Mask using Reprojection Error Map

Take the difference between the F'_i and F_{i+1} . After finding the difference image D_i , we applied the opening operation on it to reject the weak edges and get the strong differences candidates. We can then use the semantic segmentation results and match them with the difference candidates to find the potential dynamic regions S_i . Flood fill algorithm used to get dynamic by taking the difference candidates as the candidates. We can maintain the probability model to get the robust dynamic object foreground M_i over the frames. We initialize the probability map with 0.5 with all the pixels for first frame.

$$P_{i+1} = (1 - \alpha)P_i + \alpha S_i$$

$$M_{i+1}(x, y) = P_{i+1}(x, y) > \tau_p$$

Above equations represent the probability model for finding the final dynamic mask. We give the weight of α to new potential dynamic region and $1 - \alpha$ to the previous probability map. Using the equation, we calculate the probability map at time step $i + 1$. Thresholding the probability map by τ_p will give us the dynamic mask M_i . We can then use this inverted mask to get the static ORB keypoints.

D. Method 3: Based on Reprojection Error Map

The ORB keypoints and descriptors are extracted and matched to get the correspondences between the two frames as explained in method 1. After getting the correspondences, we will use the least median squared errors for finding the best homography for the set of correspondences between the frames. The least median squared errors are used because of their robustness to outliers. In the least median squared error problem, the minimum value for the median squared errors of homography over the sampled correspondences is found.

$$H_i = \min \text{med}_{1,2..i \leq p} H_i(K_t, K_{t+1})$$

We can use this homography matrix to find the projection F'_i of F_i on $i+1$. Then, we can find the reprojection error map by subtracting the two images. The following steps of finding the difference of images and curating the potential dynamic candidates using the probability model as described in method 2. Pipeline for the third method is as shown in the fig 3.

IV. EXPERIMENTS AND RESULTS

A. Distance between epilines and keypoints

In method 1, we have set the threshold τ as 3 pixels considering the central pixel as keypoint in frame F_{i+1} and farthest pixel is a corner pixel in the 5x5 neighborhood. We can



Fig. 6. Motion vector clustering for optical flow vectors on the successive images



Fig. 7. Potential dynamic mask after matching with segmentation results

see the following results, where blue colored points are marked as static keypoints and red colored points are marked as Dynamic keypoints. We observed threshold τ varies depending on the depth of keypoint from the camera. It is evident that, if the object is far from the distance between the keypoints and epilines reduces. We need to set the threshold based on the depth of the keypoint. Therefore, we concluded that we need the depth of the keypoint for efficiently setting the threshold for each of the keypoint. However, our depth data from the given dataset also had some noise. Also, sudden changes in the camera movement were resulting in bad correspondences because of blur in the motion as shown.

B. Based Motion vector Clustering

For the optical flow based method we set to extract 100 corners, which are then tracked in the successive frames. if the tracked corners falls below 50 percent, again new 100 corners are extracted from the frame and tracked. These tracked keypoints are then used for finding the homography using the minimum projection error. The number of classes are limited to minimum of 5 and number of dynamic classes. we found out from experiments that setting 5 classes could work best for clustering the motion vectors which can compensate for both depth problems as well as dynamic objects. Few of the results are shown in the following figures 6 to 8.

For the probability model, we have set the alpha to be 0.75 (add reference). We can observe that some of the static regions are present even after applying probability model fig 8. It is because of the motion vectors also depend on the depth of the key correspondences.

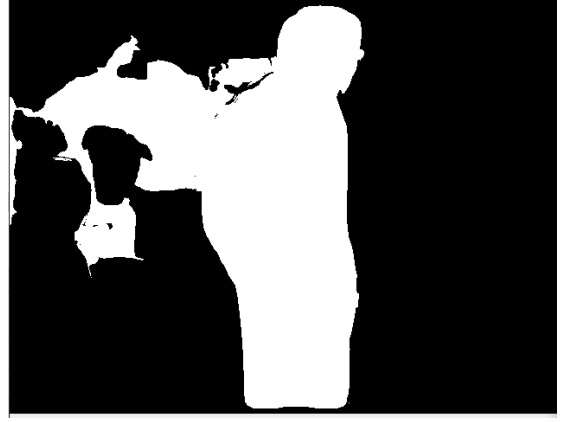


Fig. 8. Noisy dynamic region after probability update using motion vector clustering method



Fig. 9. Removal False positive dynamic regions using probability model for Reprojection error map method

C. Based Reprojection error map

We extracted 1000 ORB keypoints. We used Flann Based Matcher with minimum matches to be 50. In Lowe's ratio (add reference), we used min ratio 0.75 for filtering the bad correspondences. Similar to method, we initialized probability map with 0.5 for all pixel locations. We kept α equal to 0.75 for the probability model. The results of these experiments are shown from fig 9 and 10.

V. CONCLUSION & FUTURE WORK

Overall, we have proposed an out of the box dynamic object detection framework which is integrated into front end of SLAM to filter out keypoints corresponding to dynamic objects. This helps to create SLAM systems robust to dynamic environments. Our framework is independent and can be used on other tasks where dynamic objects are needed from videos, such as video tracking. We also realize potential research fron-



Fig. 10. Results of static ORB keypoints after using dynamic mask for Reprojection error map method.

tiers where further improvement can be achieved by tracking the dynamic objects and using their motion models to better estimate the camera pose and location.

ACKNOWLEDGMENT

Thanks to all the fellow classmates of ECE209AS for their feedback and suggestions on the project.

REFERENCES

- [1] Raul Mur-Artal and Juan D Tardós. “Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras”. In: *IEEE transactions on robotics* 33.5 (2017), pp. 1255–1262.
- [2] Sen Wang et al. “Deepvo: Towards end-to-end visual odometry with deep recurrent convolutional neural networks”. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2017, pp. 2043–2050.
- [3] Ruihao Li et al. “Undeepvo: Monocular visual odometry through unsupervised deep learning”. In: *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE. 2018, pp. 7286–7291.
- [4] Wenshan Wang, Yaoyu Hu, and Sebastian Scherer. “TartanVO: A Generalizable Learning-based VO”. In: *arXiv preprint arXiv:2011.00359* (2020).
- [5] Yuxiang Sun, Ming Liu, and Max Q-H Meng. “Motion removal for reliable RGB-D SLAM in dynamic environments”. In: *Robotics and Autonomous Systems* 108 (2018), pp. 115–128.
- [6] Jerome Revaud et al. “Epicflow: Edge-preserving interpolation of correspondences for optical flow”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1164–1172.
- [7] Chao Yu et al. “DS-SLAM: A semantic visual SLAM towards dynamic environments”. In: *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2018, pp. 1168–1174.
- [8] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. “Segnet: A deep convolutional encoder-decoder architecture for image segmentation”. In: *IEEE transactions on pattern analysis and machine intelligence* 39.12 (2017), pp. 2481–2495.
- [9] Irene Ballester et al. “DOT: dynamic object tracking for visual SLAM”. In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2021, pp. 11705–11711.
- [10] Yuxin Wu et al. “Detectron2. 2019”. In: *URL <https://github.com/facebookresearch/detectron2>* 2.3 (2019).