

---

# Visual Odometry

---

**Jai Prakash**  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
jprakash@andrew.cmu.edu

**Utkarsh Sinha**  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
usinha@andrew.cmu.edu

## Abstract

In this project, trying to localize the camera using visual odometry. The major component of the project is to generate keyframes according to pre-defined heuristics and triangulate the points to create 3D reconstruction of the scene. The intermediate frames can be found using Perspective-n-point algorithm. In addition, we also perform local bundle adjustment over last few frames so that the localization is locally consistent. We also plan to exploit the onboard inertial sensors to get prior for the localization.

## 1 Introduction

Augmented reality has been around for years, yet not all problems are solved in the domain. One of the challenges is precise localization of the device in world coordinates. Several augmented reality applications on smartphones are based on markers. One good example of marker-based AR is Vuforia. On the other hand, there are standalone devices like Hololens, which has number of sensors to understand the scene and localize the head mounted display in the scene.

In many such applications, understanding the scene is not important. Localizing the camera in the world is enough to solve certain problems. In this project we focus on localizing the phone camera using the camera and the inertial sensors.

## 2 Background

**Visual SLAM vs. Visual Odometry:** The focus in the visual SLAM techniques is both in reconstructing the scene and also localizing the camera in the scene. However, our main focus is just in localization of the camera. For the scope of the project, we focus only to be locally consistent. So, our system's camera position might drift over time - however, we are only interested in accurate localization in a short timespan. We do not explore ideas like loop closure in this project.

## 3 Feature Extraction and Matching

We have experimented with OpenCV KLT features, AKAZE features and ORB features. The KLT features can also be used for tracking the features in the subsequent frames using optical flow. For AKAZE? and ORB features, the correspondences are found using feature matching. The outliers are removed using the epipolar constraints and finding unique matches i.e. feature from first image matches to a unique feature in second image and vice-versa.

### 3.1 3D reconstruction

Using the feature matching, we can triangulate the points. The fundamental matrix gives the relationship between the feature points. We used RANSAC based 8-point algorithm to find the fundamental matrix.

This gives rise to four possible camera configurations (with  $W/W^T$  and  $\pm t$ ). The correct location can be found using the camera configuration in which all the points are in front of both the cameras.

### 3.2 Camera pose recovery

Once the scene is reconstructed using the keyframes, the camera pose can be recovered using Perspective-n-point (PnP) algorithms. By knowing the 3D points from the reconstruction, and its corresponding feature location in any image the camera pose can be recovered using PnP.

### 3.3 Bundle Adjustment

In visual odometry, the current camera pose is obtained by adding the last observed motion to the current detection change. This leads to a superlinear increase in pose error over time. In this section, we look at the techniques we intend to use to correct this pose drift.

One solution is to use bundle adjustment to impose geometrical constraints over multiple frames. The computational cost increases with the cube of the number of frames used for computation. Thus, we limit the number of frames to a small window from the previously captured frames. This approach is called local bundle adjustment.

## 4 Results

So far, we have been working with the datasets available online. The results are illustrated on the Middlebury Temple dataset ?. We first find the feature correspondences and then remove the outliers using Epipolar constraints. The outliers can be found by using a threshold on distance from epipolar line and along the epipolar line.

We are able to reconstruct the temple structure using the two keyframes (handpicked for now). The results are shown in figure ?. Once the reconstruction is done, we are able to recover the camera poses using PnP algorithm as illustrated in the figure ?. We are using OpenCV for feature matching and visualization.

## 5 Future Work

Until now, we have evaluated the various landmark tracking techniques on the Temple dataset. We intend to test the different techniques across multiple datasets and pick an algorithm that works best.

We will integrate Ceres solver by Google for local bundle adjustment. This would improve the localization of the camera and generate a better tracking of the scene.

If time permits, we also intend to integrate inertial sensors readily available on mobile devices to improve the estimate?.

## References

References follow the acknowledgments. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. **Remember that you can use a ninth page as long as it contains only cited references.**

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

- [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.