

# 16-720B Computer Vision: Homework 4

## 3D Reconstruction

**Instructors:** Srinivasa Narasimhan

**TAs:** N Dinesh Reddy, Alankar Kotwal, Martin Li,  
Shumian Xin, Ye Yuan, Zhou Xian

See course website for deadline: <https://canvas.cmu.edu/courses/13727>

Make sure to start early!

### Part I

### Theory

Before implementing our own 3D reconstruction, let's take a look at some simple theory questions that may arise. The answers to the below questions should be relatively short, consisting of a few lines of math and text (maybe a diagram if it helps your understanding).

**Q1.1 (5 points)** Suppose two cameras fixate on a point  $\mathbf{x}$  (see Figure 1) in space such that their principal axes intersect at that point. Show that if the image coordinates are normalized so that the coordinate origin  $(0, 0)$  coincides with the principal point, the  $F_{33}$  element of the fundamental matrix is zero.

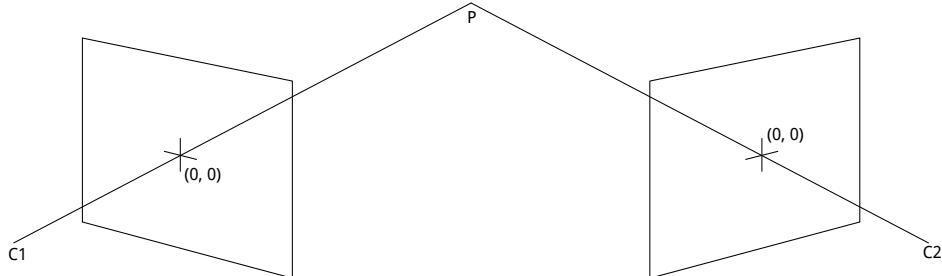


Figure 1: Figure for Q1.1.  $C_1$  and  $C_2$  are the optical centers. The principal axes intersect at point  $w$  ( $P$  in the figure).

**Q1.2 (5 points)** Consider the case of two cameras viewing an object such that the second camera differs from the first by a *pure translation* that is parallel to the  $x$ -axis. Show that the epipolar lines in the two cameras are also parallel to the  $x$ -axis. Backup your argument with relevant equations.

**Q1.3 (5 points)** Suppose we have an inertial sensor which gives us the accurate positions ( $\mathbf{R}_i$  and  $\mathbf{t}_i$ , the rotation matrix and translation vector) of the robot at time  $i$ . What will be the effective rotation ( $\mathbf{R}_{rel}$ ) and translation ( $\mathbf{t}_{rel}$ ) between two frames at different time stamps? Suppose the camera intrinsics ( $\mathbf{K}$ ) are known, express the essential matrix ( $\mathbf{E}$ ) and the fundamental matrix ( $\mathbf{F}$ ) in terms of  $\mathbf{K}$ ,  $\mathbf{R}_{rel}$  and  $\mathbf{t}_{rel}$ .

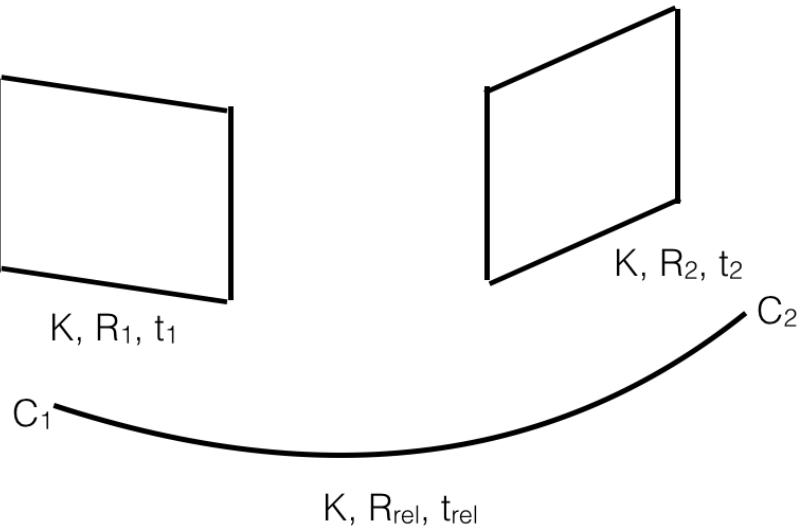


Figure 2: Figure for Q1.3.  $C_1$  and  $C_2$  are the optical centers. The rotation and the translation is obtained using inertial sensors.  $\mathbf{R}_{rel}$  and  $\mathbf{t}_{rel}$  are the relative rotation and translation between two frames.

**Q1.4 (10 points)** Suppose that a camera views an object and its reflection in a plane mirror. Show that this situation is equivalent to having two images of the object which are related by a skew-symmetric fundamental matrix. You may assume that the object is flat, meaning that all points on the object are of equal distance to the mirror (*Hint:* draw the relevant vectors to understand the relationships between the camera, the object and its reflected image.)

## Part II

### Practice

#### 1 Overview

In this part you will begin by implementing the two different methods seen in class to estimate the fundamental matrix from corresponding points in two images (Section 2). Next, given the fundamental matrix and calibrated intrinsics (which will be provided) you will compute the essential matrix and use this to compute a 3D metric reconstruction from 2D correspondences using triangulation (Section 3). Then, you will implement a method to automatically match points taking advantage of epipolar constraints and make a 3D visualization of the results (Section 4). Finally, you will implement RANSAC and bundle adjustment to further improve your algorithm (Section 5).

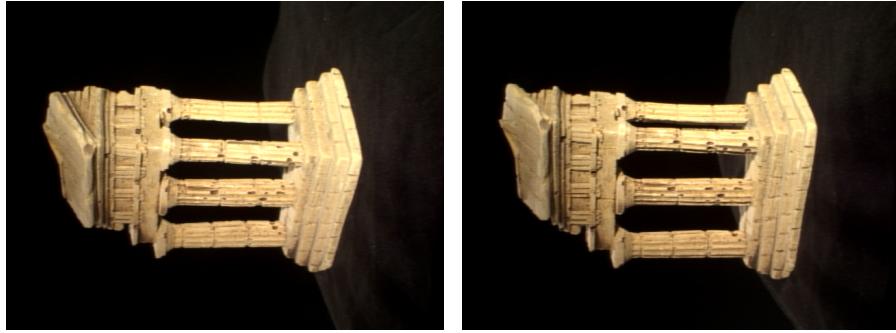


Figure 3: Temple images for this assignment

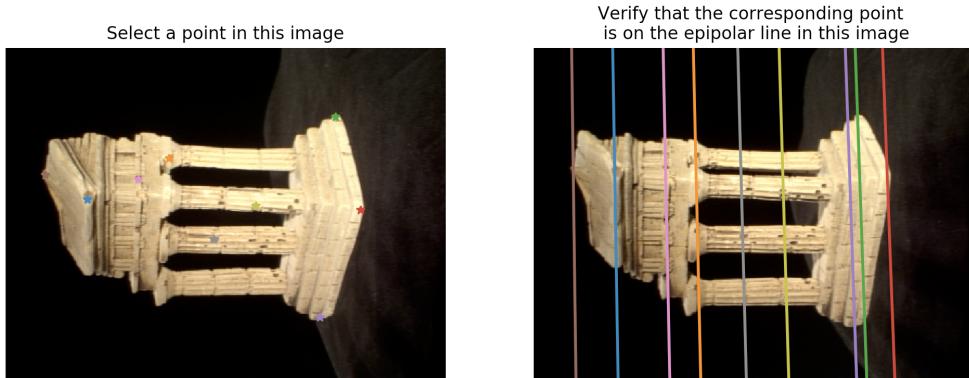


Figure 4: `displayEpipolarF` in `helper.py` creates a GUI for visualizing epipolar lines

## 2 Fundamental matrix estimation

In this section you will explore different methods of estimating the fundamental matrix given a pair of images. In the `data/` directory, you will find two images (see Figure 3) from the Middlebury multi-view dataset<sup>1</sup>, which is used to evaluate the performance of modern 3D reconstruction algorithms.

### The Eight Point Algorithm

The 8-point algorithm (discussed in class, and outlined in Section 10.1 of Forsyth & Ponce) is arguably the simplest method for estimating the fundamental matrix. For this section, you can use provided correspondences you can find in `data/some_corresp.npz`.

**Q2.1 (10 points)** Finish the function `eightpoint` in `submission.py`. Make sure you follow the signature for this portion of the assignment:

```
F = eightpoint(pts1, pts2, M)
```

where `pts1` and `pts2` are  $N \times 2$  matrices corresponding to the  $(x, y)$  coordinates of the  $N$  points in the first and second image respectively. `M` is a scale parameter.

---

<sup>1</sup><http://vision.middlebury.edu/mview/data/>

- You should scale the data as was discussed in class, by dividing each coordinate by  $M$  (the maximum of the image's width and height). After computing  $\mathbf{F}$ , you will have to “unscale” the fundamental matrix.

*Hint:* If  $\mathbf{x}_{\text{normalized}} = \mathbf{T}\mathbf{x}$ , then  $\mathbf{F}_{\text{unnormalized}} = \mathbf{T}^T \mathbf{F} \mathbf{T}$ .

You must enforce the singularity condition of the  $\mathbf{F}$  before unscaling.

- You may find it helpful to refine the solution by using local minimization. This probably won't fix a completely broken solution, but may make a good solution better by locally minimizing a geometric cost function.

For this we have provided a helper function `refineF` in `helper.py` taking in  $\mathbf{F}$  and the two sets of points, which you can call from `eightpoint` before unscaling  $\mathbf{F}$ .

- Remember that the  $x$ -coordinate of a point in the image is its column entry, and  $y$ -coordinate is the row entry. Also note that eight-point is just a figurative name, it just means that you need at least 8 points; your algorithm should use an over-determined system ( $N > 8$  points).

- To visualize the correctness of your estimated  $\mathbf{F}$ , use the supplied function `displayEpipolarF` in `helper.py`, which takes in  $\mathbf{F}$ , and the two images. This GUI lets you select a point in one of the images and visualize the corresponding epipolar line in the other image (Figure 4).

- **Output:** Save your matrix  $\mathbf{F}$ , scale  $M$  to the file `q2_1.npz`.

**In your write-up:** Write your recovered  $\mathbf{F}$  and include an image of some example output of `displayEpipolarF`.

### 3 Metric Reconstruction

You will compute the camera matrices and triangulate the 2D points to obtain the 3D scene structure. To obtain the Euclidean scene structure, first convert the fundamental matrix  $\mathbf{F}$  to an essential matrix  $\mathbf{E}$ . Examine the lecture notes and the textbook to find out how to do this when the internal camera calibration matrices  $\mathbf{K}_1$  and  $\mathbf{K}_2$  are known; these are provided in `data/intrinsics.npz`.

**Q3.1 (5 points)** Write a function to compute the essential matrix  $\mathbf{E}$  given  $\mathbf{F}$ ,  $\mathbf{K}_1$  and  $\mathbf{K}_2$  with the signature:

```
E = essentialMatrix(F, K1, K2)
```

**In your write-up:** Write your estimated  $\mathbf{E}$  using  $\mathbf{F}$  from the eight-point algorithm.

Given an essential matrix, it is possible to retrieve the projective camera matrices  $\mathbf{M}_1$  and  $\mathbf{M}_2$  from it. Assuming  $\mathbf{M}_1$  is fixed at  $[\mathbf{I}, \mathbf{0}]$ ,  $\mathbf{M}_2$  can be retrieved up to a scale and four-fold rotation ambiguity. For details on recovering  $\mathbf{M}_2$ , see section 7.2 in Szeliski. We have provided you with the function `camera2` in `python/helper.py` to recover the four possible  $\mathbf{M}_2$  matrices given  $\mathbf{E}$ .

**Note:** The  $\mathbf{M}_1$  and  $\mathbf{M}_2$  here are projection matrices of the form:  $\mathbf{M}_1 = [\mathbf{I}|\mathbf{0}]$  and  $\mathbf{M}_2 = [\mathbf{R}|\mathbf{t}]$ .

**Q3.2 (10 points)** Using the above, write a function to triangulate a set of 2D coordinates in the image to a set of 3D points with the signature:

```
[w, err] = triangulate(C1, pts1, C2, pts2)
```

where  $\text{pts1}$  and  $\text{pts2}$  are the  $N \times 2$  matrices with the 2D image coordinates and  $\mathbf{w}$  is an  $N \times 3$  matrix with the corresponding 3D points per row.  $\mathbf{C1}$  and  $\mathbf{C2}$  are the  $3 \times 4$  camera matrices. Remember that you will need to multiply the given intrinsics matrices with your solution for the canonical camera matrices to obtain the final camera matrices. Various methods exist for triangulation - probably the most familiar for you is based on least squares (see Szeliski Chapter 7 if you want to learn about other methods):

For each point  $i$ , we want to solve for 3D coordinates  $\mathbf{w}_i = [x_i, y_i, z_i]^T$ , such that when they are projected back to the two images, they are close to the original 2D points. To project the 3D coordinates back to 2D images, we first write  $\mathbf{w}_i$  in homogeneous coordinates, and compute  $\mathbf{C1}\tilde{\mathbf{w}}_i$  and  $\mathbf{C2}\tilde{\mathbf{w}}_i$  to obtain the 2D homogeneous coordinates projected to camera 1 and camera 2, respectively.

For each point  $i$ , we can write this problem in the following form:

$$\mathbf{A}_i \mathbf{w}_i = 0,$$

where  $\mathbf{A}_i$  is a  $4 \times 4$  matrix, and  $\tilde{\mathbf{w}}_i$  is a  $4 \times 1$  vector of the 3D coordinates in the homogeneous form. Then, you can obtain the homogeneous least-squares solution (discussed in class) to solve for each  $\mathbf{w}_i$ .

**In your write-up:** Write down the expression for the matrix  $\mathbf{A}_i$ .

Once you have implemented triangulation, check the performance by looking at the reprojection error:

$$\text{err} = \sum_i \|\mathbf{x}_{1i}, \hat{\mathbf{x}}_{1i}\|^2 + \|\mathbf{x}_{2i}, \hat{\mathbf{x}}_{2i}\|^2$$

where  $\hat{\mathbf{x}}_{1i} = \text{Proj}(\mathbf{C}_1, \mathbf{w}_i)$  and  $\hat{\mathbf{x}}_{2i} = \text{Proj}(\mathbf{C}_2, \mathbf{w}_i)$ .

**Note:**  $\mathbf{C1}$  and  $\mathbf{C2}$  here are projection matrices of the form:  $\mathbf{C}_1 = \mathbf{K}_1 \mathbf{M}_1 = \mathbf{K}_1 [\mathbf{I}|0]$  and  $\mathbf{C}_2 = \mathbf{K}_2 \mathbf{M}_2 = \mathbf{K}_2 [\mathbf{R}|\mathbf{t}]$ .

**Q3.3 (10 points)** Write a script `findM2.py` to obtain the correct  $\mathbf{M2}$  from  $\mathbf{M2s}$  by testing the four solutions through triangulations. Use the correspondences from `data/some_corresp.npz`.

**Output:** Save the correct  $\mathbf{M2}$ , the corresponding  $\mathbf{C2}$ , and 3D points  $\mathbf{P}$  to `q3_3.npz`.

## 4 3D Visualization

You will now create a 3D visualization of the temple images. By treating our two images as a stereo-pair, we can triangulate corresponding points in each image, and render their 3D locations.

**Q4.1 (15 points)** Implement a function with the signature:

```
[x2, y2] = epipolarCorrespondence(im1, im2, F, x1, y1)
```

This function takes in the  $x$  and  $y$  coordinates of a pixel on  $\text{im1}$  and your fundamental matrix  $\mathbf{F}$ , and returns the coordinates of the pixel on  $\text{im2}$  which correspond to the input point. The match is obtained by computing the similarity of a small window around the  $(x_1, y_1)$  coordinates in  $\text{im1}$  to various windows around possible matches in the  $\text{im2}$  and returning the closest.

Instead of searching for the matching point at every possible location in  $\text{im2}$ , we can use  $\mathbf{F}$  and simply search over the set of pixels that lie along the epipolar line (recall that the epipolar line passes through a single point in  $\text{im2}$  which corresponds to the point  $(x_1, y_1)$  in  $\text{im1}$ ).

There are various possible ways to compute the window similarity. For this assignment, simple methods such as the Euclidean or Manhattan distances between the intensity of the pixels should suffice. See Szeliski Chapter 11, on stereo matching, for a brief overview of these and other methods.  
*Implementation hints:*

- Experiment with various window sizes.
- It may help to use a Gaussian weighting of the window, so that the center has greater influence than the periphery.
- Since the two images only differ by a small amount, it might be beneficial to consider matches for which the distance from  $(x_1, y_1)$  to  $(x_2, y_2)$  is small.

To help you test your `epipolarCorrespondence`, we have included a helper function `epipolarMatchGUI` in `python/helper.py`, which takes in two images the fundamental matrix. This GUI allows you to click on a point in `im1`, and will use your function to display the corresponding point in `im2`. See Figure 5.

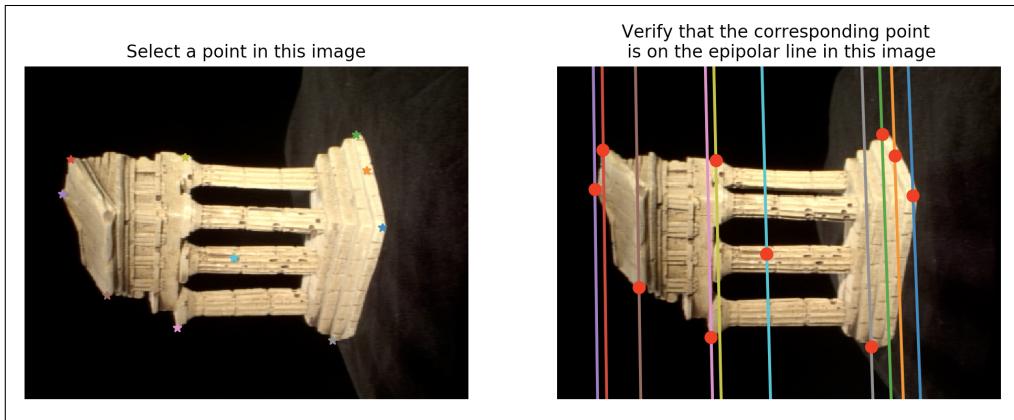


Figure 5: `epipolarMatchGUI` shows the corresponding point found by calling `epipolarCorrespondence`

It's not necessary for your matcher to get *every* possible point right, but it should get easy points (such as those with distinctive, corner-like windows). It should also be good enough to render an intelligible representation in the next question.

**Output:** Save the matrix  $\mathbf{F}$ , points  $\mathbf{pts1}$  and  $\mathbf{pts2}$  which you used to generate the screenshot to the file `q4.1.npz`.

**In your write-up:** Include a screenshot of `epipolarMatchGUI` with some detected correspondences.

**Q4.2 (10 points)** Included in this homework is a file `data/templeCoords.npz` which contains 288 hand-selected points from `im1` saved in the variables `x1` and `y1`.

Now, we can determine the 3D location of these point correspondences using the `triangulate` function. These 3D point locations can then plotted using the Matplotlib or plotly package. Write a script `visualize.py`, which loads the necessary files from `../data/` to generate the 3D reconstruction using `scatter`. An example is shown in Figure 6.

**Output:** Again, save the matrix  $\mathbf{F}$ , matrices  $\mathbf{M1}, \mathbf{M2}, \mathbf{C1}, \mathbf{C2}$  which you used to generate the screenshots to the file `q4.2.npz`.

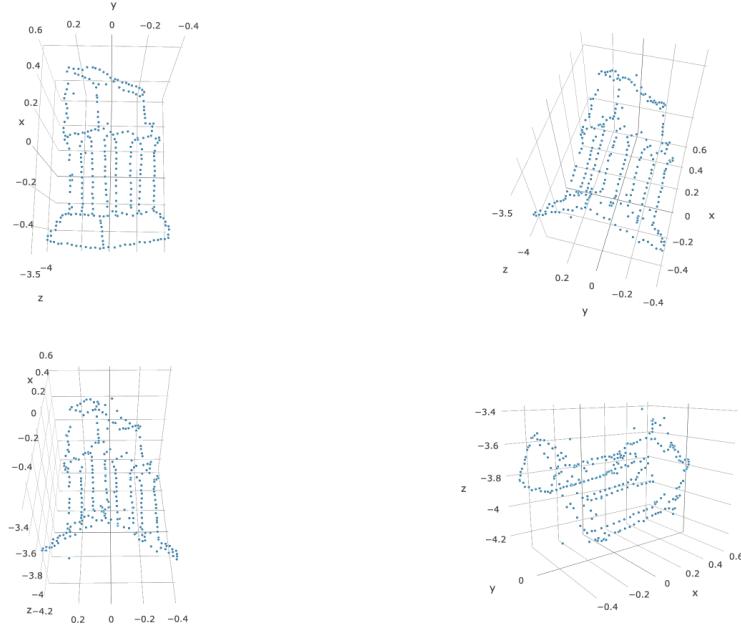


Figure 6: An example point cloud

**In your write-up:** Take a few screenshots of the 3D visualization so that the outline of the temple is clearly visible, and include them with your homework submission.

## 5 Bundle Adjustment

**Q5.1 (15 points)** In some real world applications, manually determining correspondences is infeasible and often there will be noisy coorespondances. Fortunately, the RANSAC method seen in class can be applied to the problem of fundamental matrix estimation.

Implement the above algorithm with the signature:

```
[F, inliers] = ransacF(pts1, pts2, M, nIters, tol)
```

where  $M$  is defined in the same way as in Section 2 and  $\text{inliers}$  is a boolean vector of size equivalent to the number of points. Here  $\text{inliers}$  is set to true only for the points that satisfy the threshold defined for the given fundamental matrix  $F$ .

We have provided some noisy coorespondances in `some_corresp_noisy.npz` in which around 75% of the points are inliers. Compare the result of RANSAC with the result of the eightpoint when ran on the noisy coorespondances. Briefly explain the error metrics you used, how you decided which points were inliers and any other optimizations you may have made. `nIters` is the maximum number of iterations of The RANSAC and `tol` is the tolerance of the error to be considered as inliers. Discuss the effect on the Fundamental matrix by varying these values.

- *Hints:* Use the seven point to compute the fundamental matrix from the minimal set of points. Then compute the inliers, and refine your estimate using all the inliers.

### Q5.2 (15 points)

So far we have independently solved for camera matrix,  $\mathbf{M}_j$  and 3D projections,  $\mathbf{w}_i$ . In bundle adjustment, we will jointly optimize the reprojection error with respect to the points  $\mathbf{w}_i$  and the camera matrix  $\mathbf{w}_j$ .

$$err = \sum_{ij} \|\mathbf{x}_{ij} - Proj(\mathbf{C}_j, \mathbf{w}_i)\|^2,$$

where  $\mathbf{w}_j = \mathbf{K}_j \mathbf{M}_j$ , same as in Q3.2.

For this homework we are going to only look at optimizing the extrinsic matrix. To do this we will be parametrizing the rotation matrix  $\mathbf{R}$  using Rodrigues formula to produce vector  $\mathbf{r} \in \mathbb{R}^3$ . Write a function that converts a Rodrigues vector  $\mathbf{r}$  to a rotation matrix  $\mathbf{R}$

$$\mathbf{R} = \text{rodrigues}(\mathbf{r})$$

as well as the inverse function that converts a rotation matrix  $\mathbf{R}$  to a Rodrigues vector  $\mathbf{r}$

$$\mathbf{r} = \text{invRodrigues}(\mathbf{R})$$

### Q5.3 (10 points)

Using this parameterization, write an optimization function

$$\text{residuals} = \text{rodriguesResidual}(K1, M1, p1, K2, p2, w)$$

where  $x$  is the flattened concatenation of  $\mathbf{x}$ ,  $\mathbf{r}_2$ , and  $\mathbf{t}_2$ .  $w$  are the 3D points;  $\mathbf{r}_2$  and  $\mathbf{t}_2$  are the rotation (in the Rodrigues vector form) and translation vectors associated with the projection matrix  $\mathbf{M}_2$ . The `residuals` are the difference between original image projections and estimated projections (the square of 2-norm of this vector corresponds to the error we computed in Q3.2):

$$\text{residuals} = \text{numpy.concatenate}([(p1-p1_hat).reshape([-1]), (p2-p2_hat).reshape([-1])])$$

Use this error function and Scipy's nonlinear least square optimizer `leastsq` write a function to optimize for the best extrinsic matrix and 3D points using the inlier correspondences from `some_corresp_noisy.npz` and the RANSAC estimate of the extrinsics and 3D points as an initialization.

$$[M2, w] = \text{bundleAdjustment}(K1, M1, p1, K2, M2_init, p2, w_init)$$

**In your write-up:** include an image of the original 3D points and the optimized points as well as the reprojection error with your initial  $\mathbf{M}_2$  and  $\mathbf{w}$ , and with the optimized matrices.

## 6 Multi View Keypoint Reconstruction

You will use multi-view capture of moving vehicles and reconstruct the motion of a car. The first part of the problem will be using a single time instance capture from Three views (Figure 7(Top)) and reconstruct vehicle keypoints and render from multiple views(Figure 7(Bottom))

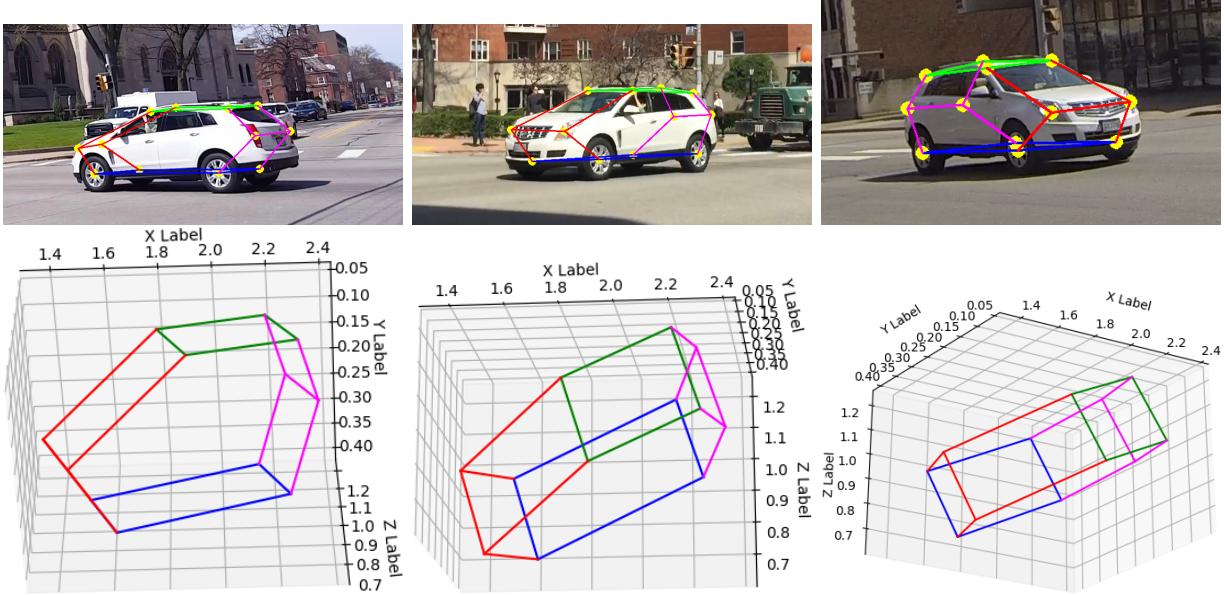


Figure 7: An example detections on the top and the reconstructions from multiple views

**Q6.1 (15 points)** Write a function to compute the 3D keypoint locations  $P$  given the 2d part detections  $\text{pts1}, \text{pts2}$  and  $\text{pts3}$  and the camera projection matrices  $C1, C2, C3$ .

```
[P, err] = MultiviewReconstruction(C1, pts1, C2, pts2, C3, pts3, Thres)
```

**In your write-up:** Describe the method you used to compute the 3D locations and include an image of the Reconstructed 3D points with the points connected using the helper function `plot_3d_keypoint(P)` with the reprojection error.

The 2D part detections( $\text{pts}$ ) are computed using a neural network <sup>2</sup> and correspond to different locations on a car like the wheels, headlights etc. The third column in  $\text{pts}$  is the confidence of localization of the keypoints. Higher confidence value represents more accurate localization of the keypoint in 2D. To visualize the 2D detections run the `visualize_keypoints(image, pts, Thres)` helper function.  $\text{Thres}$  is defined as the confidence threshold of the 2D detected keypoints. The camera matrices ( $C$ ) are computed by running an SfM from multiple views. By varying confidence Threshold  $\text{Thres}$ .(i.e. considering only the points above the threshold), We get different reconstruction and accuracy. Try varying the thresholds and analyze its effects on the accuracy of the reconstruction.

### Q6.2 (15 points)(Extra Credit)

From the previous question you have done a 3D reconstruction at a time instance. Now you are going to iteratively repeat the process over time and compute a spatio temporal reconstruction of the car. The images in the `data` folder shows the motion of the car at an intersection captured from multiple views. The images are given as (`cam1_time0.jpg, ..., cam1_time9.jpg`) for camera 1 and (`cam2_time0.jpg, ..., cam2_time9.jpg`) for camera2 and (`cam3_time0.jpg, ..., cam3_time9.jpg`) for camera3. The corresponding detections and camera matrices are given in (`time0.npz, ..., time9.npz`). Use the above details and compute the spatio temporal recon-

---

<sup>2</sup>Code Used For Detection and Reconstruction

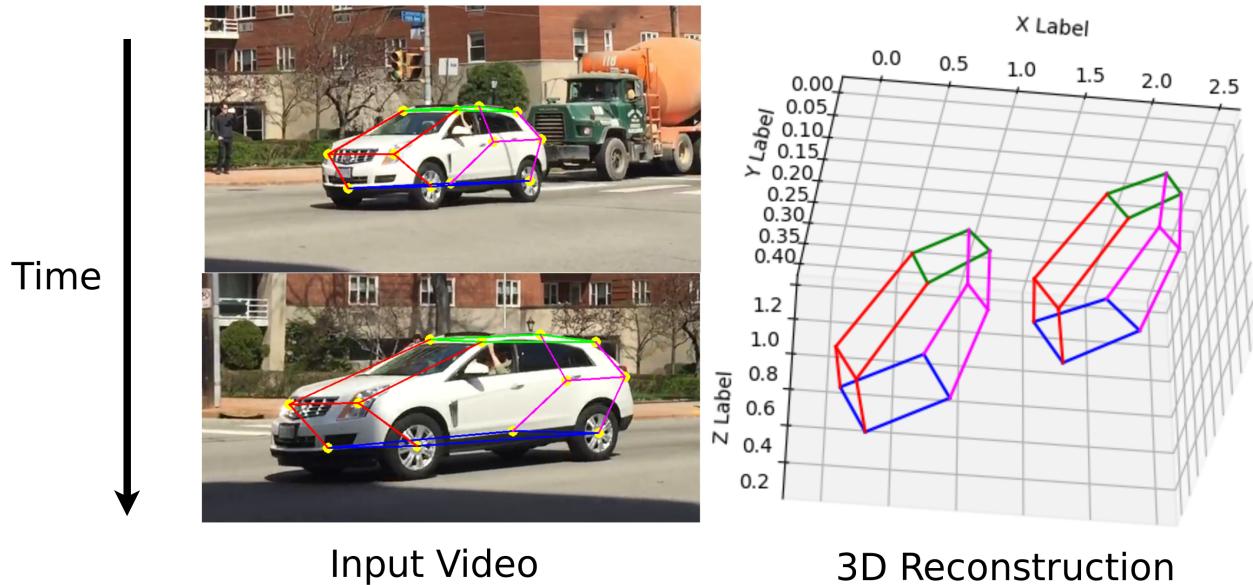


Figure 8: Spatiotemporal reconstruction of the car (right) with the projections at two different time instances in a single view(left)

struction of the car for all 10 time instances and plot them. A sample plot with the first and last time instance reconstruction of the car with the reprojections shown in the Figure.

## 7 Deliverables

If your andrew id is `bovik`, your submission should be the writeup `bovik.pdf` and a zip file `bovik.zip`. **Please submit the zip file to Canvas and the pdf file to Gradescope.**

The zip file should include the following directory structure:

- `submission.py`: your implementation of algorithms.
- `findM2.py`: script to compute the correct camera matrix.
- `visualize.py`: script to visualize the 3d points.
- `q2_1.npz`: file with output of Q2.1.
- `q2_2.npz`: file with output of Q2.2.
- `q3_3.npz`: file with output of Q3.3.
- `q4_1.npz`: file with output of Q4.1.
- `q4_2.npz`: file with output of Q4.2.
- `q6_1.npz`: file with output of Q4.2.