

# Object Recognition, Pose Estimation and Tracking of a Rigid Object with Planar Surface

Ramkumar Natarajan Siddharthan Rajasekaran Sri Ramana Sekharan Vishaal Kabilan

## Abstract

In this project we recognize and estimate the pose of a Rubik's Pyramid. We model the object using a 3D-mesh, recognize it using OpponentSURF, estimate the pose using EPnP and then track the object using Linear Kalman filter. We discuss the theory behind our approach, its performance and its limitations. In this project our focus is not general object recognition, instead we exploit the structure of our object and tune our algorithm based on it.

## 1. Introduction

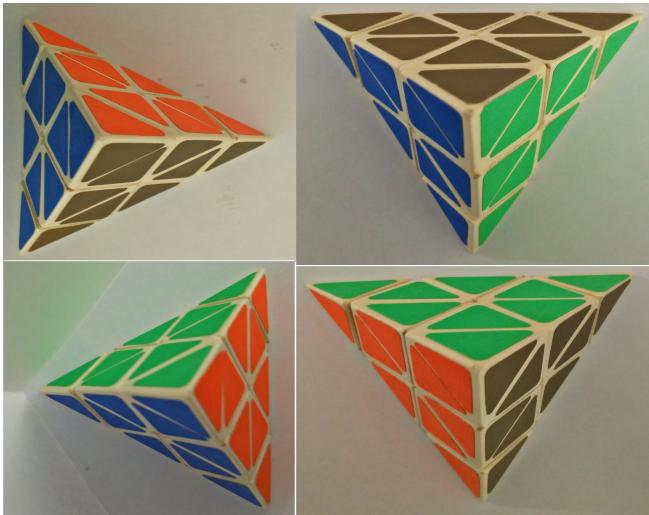
The task of object recognition from an image or motion picture can be accomplished for many applications by identifying an unknown object as a member of a set of well known objects. In our project we deal with rigid Rubik's pyramid object whose model is known beforehand.

Many object recognition methods use a set of abstract characterizations in order to efficiently compare and represent the object. Such characterizations are canonically defined by measurable object feature extracted from various types of images and any previous knowledge available. Similarity between characterizations is interpreted as similarity between the objects themselves. Thus the ability of a given algorithm to describe the object uniquely from the available information determines the effectiveness of the for the given application.

## 2. Problem Statement

Our object is a pyramid with four faces each of which is shown in Fig. 1. Each of these faces has a single

RBE/CS 549: Computer Vision. Project work in partial fulfillment of the requirements for course evaluation.



*Figure 1. Different Poses of Object:* The object to be detected in different poses shown here. From top-left clockwise: pose1, pose2, pose3, pose4

color namely green, blue, orange, gold. Each face of the pyramid is joint by a number of smaller triangles. So the smooth color on every side has interleaved triangular patterns in white. Since this is a well-defined regular object, we can accurately model the object. We use the 3D modelling approach followed by recognition using feature detection. We then estimate the pose using EPnP and the track using kalman filter.

## 3. Our Approach

We explain the theory behind each component of our approach here and how it applies to our object.

### 3.1. Model Registration

Model registration is the essential part of our method. We first register the model of the given object, then compare the image frame captured online with our model to classify the features observed as inliers and outliers. (Pi, 2015) To do this we do model

registration offline before the detection stage. We have two major steps in model registration. First, we define the 3-D mesh of the object. A 3-D mesh is a set of triangles that lay on the surface of the object. We choose triangles over other polygons because 1) we are directed towards recognizing a rigid body with planar surface 2) computing area of a triangle is easier than computing area of other polygons. The latter is required for area based sampling which we'll explain in the following sections.

Apart from the 3-D mesh, the model consists of  $n$ -unique feature points that characterize the object. We will be referring to these as correspondence points. In our case we pick the corners of pyramid and the middle points of each edge as correspondence points. This is because the corner points have 3 colors around it in a specific order. For example, point  $P4$  in Fig.3.1.2 has blue,green and gold colors in anticlockwise direction. Also, we need four non-coplanar control points on the boundary of the object to span through any interior point of the object. This is important to preserve the constraints (5.  $\sum_{j=1}^4 \alpha_{ij} = 1$ ) of the barycentric coordinates of all the other correspondence points(Pi, 2015). However we need 6 unique points apart from the control points to find the camera parameters. We choose 6 points on the mid point of the 6 edges of the pyramid.

We hand code the 3-D coordinates of these points and store it alongside the 3-D mesh. We find the points that are visible for different poses of the object. Our object has only four such poses (with different set of correspondence points visible) and hence this method works well with few definitions of these different poses. For example, in pose 1 3.1.2 we see all the points i.e. P1,P2,P3,P4,P5,P6,P7,P8,P9,P10 given the perspective. Once we have such definition for the vertices and edges we have sufficient information for associating 2D correspondence points detected using features to 3D coordinates and hence find the pose of the camera.

### 3.1.1. BARYCENTRIC COORDINATES

The 3-D mesh of an object is a set of triangles. We store the mesh as a collection of three vertices (for each triangle) which stay connected to each other in the mesh. To convert the mesh of an object to point cloud we need to uniformly sample points from the mesh all over its surface area. We randomly sample a triangle  $\Delta_i$  from the mesh with a weight proportional to its area  $A_i$  calculated from its vertices  $\{V_1^i, V_2^i, V_3^i\}$  (stored alongside mesh definition).  $A_i = \frac{1}{2} \cdot \| (V_1^i - V_3^i) \times (V_1^i - V_2^i) \|$ . We now have to

sample a point from this triangle with uniform distribution. Barycentric coordinates (?) provide a perfect framework for uniform sampling. Barycentric coordinates of any  $n$  dimensional point can be represented as an  $m$  dimensional point when it is on an  $m$  dimensional simplex. Since we use triangles ( $m = 3$ ) to represent the mesh, the barycentric coordinates of the point ( $\vec{X}^i$ ) sampled from this triangle will be three dimensional.

$$\vec{X}' = \begin{bmatrix} x_1^i \\ x_2^i \\ x_3^i \end{bmatrix} \xleftarrow{\text{sample}} \begin{bmatrix} u \\ v \\ w \end{bmatrix} \mid u+v+w = 1, 0 \leq u, v, w \leq 1 \quad (1)$$

$u, v, w$  is all the space spanned by  $\vec{X}'$ . To get a uniform distribution over this space we choose two random numbers,  $r_1$  and  $r_2$  between 0 and 1. We then assign,

$x_1^i, x_2^i, \& x_3^i$  using

$$x_1^i = 1 - \sqrt{r_1}$$

$$x_2^i = \sqrt{r_1}(1 - r_2)$$

$$x_3^i = 1 - x_1^i - x_2^i$$

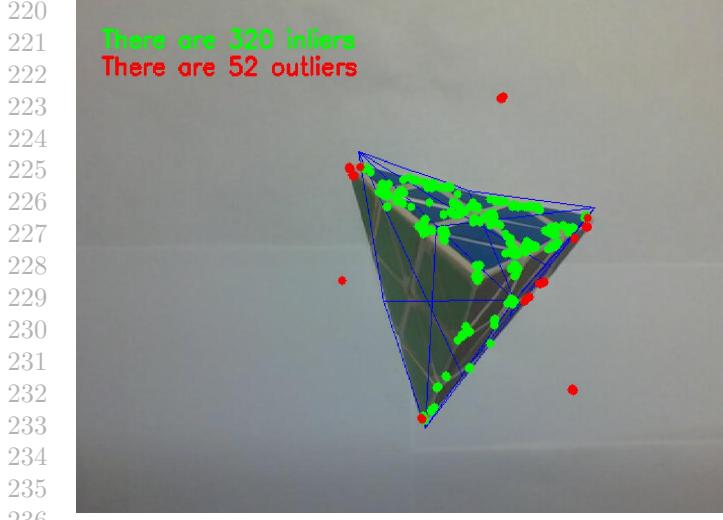
### 3.1.2. MOLLER-TRUMBORÉ ALGORITHM

Any point  $P$  inside triangle  $\Delta_i$  can be represented as an intersection of a parametric ray given by  $P = O + t.D$ , where  $O$  is the origin of the ray,  $D$  is the direction and  $t$  is the distance of the point  $P$  from the origin  $O$ . Inversely, the barycentric coordinates of the intersection of the ray with the triangle can be found using

$$[-D(V_2^i - V_1^i)(V_3^i - V_1^i)] \begin{bmatrix} t \\ u \\ v \end{bmatrix} = (O - A)$$

It is easy to find the mesh 3-D of rigid bodies with planar surface. Given our object (triangular faces), which is a collection of only four faces and can be exactly represented using triangles without any approximation, we pick the farthest three points on the surface of the object such that the triangle obtained by joining these three points is exactly the surface area of the object lying under this triangular section. In Fig. 3.1.2 we have shown the mesh using blue lines. We choose triangles (for mesh) on a face such that all the correspondence points on that face get selected at least once. For any detected feature point  $f$  we use k-nearest neighbors to classify the point as an inlier or

165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219



220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274

Figure 2. 3D Mesh: The mesh is shown using blue lines. The inliers are shown in green and the outliers in red

outlier(Pi, 2015). If the point lies near the point-cloud sampled from the mesh, we classify it as an inlier. Fig. 3.1.2 shows the number of inliers (green) and the number of outliers (red).

In Fig.3.1.2 we can see the manually registered correspondence points marked in red while that detected by our feature matching algorithm are marked in green.

After we find the inliers we still have some features that distort the orientation calculations. We filter these features using left/bottom-right/top constraint and the neighbor ratio check. That in the Fig.registration if point  $P_2 \Rightarrow P_3$  then  $P_4 \Rightarrow P_5$  (' $\Rightarrow$ ' is read 'is right/top side of'). Mathematically,

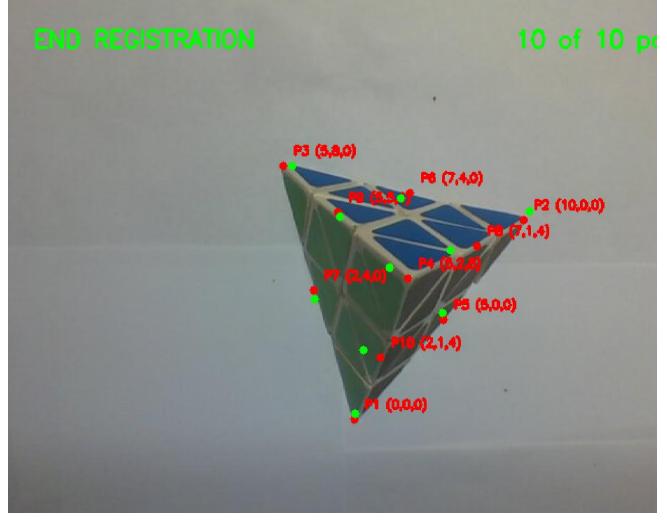
$$(\vec{P}_2 - \vec{P}_3) \cdot (\vec{P}_4 - \vec{P}_5) \geq 0$$

We check the distance between matches in order to remove repeated elements. We keep matches with ratio of distance of first nearest neighbor to the second less than a threshold  $\delta$ .

### 3.2. Camera Calibration

Object recognition is a tract of the broader problem of pose estimation. To get the position and orientation of the object, we formalize the problem into a *PnP* problem. This is because using Direct Linear Transform (DLT) was computationally expensive and slow. Thus we calibrate the intrinsic parameters. The camera projection matrix is given as:

$$M = K[RT] \quad (2)$$



275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
Figure 3. Correspondence points: The user marked correspondence points using mouse is given in red color. The correspondence points detected by feature matching algorithm is given in green.

where  $K$  is the intrinsic(inner) and  $R, T$  are the extrinsic(outer) camera parameters.

Out of the various techniques used for camera calibration, we used the 2D plane based calibration. Techniques in this method requires to observe a planar pattern shown at a few different orientations. Unlike Tsai's technique[v1], the knowledge of the plane motion is not necessary. In our experiment we used the checkerboard pattern. The pattern was first printed and attached to planar wall 3.2. Then a few images of the model plane were taken under different orientations by moving the camera. Then the feature points in the images are detected. Then the five intrinsic parameters are estimated by the closed form solution discussed in [v2]. This was followed by estimating the coefficients of radial distortion by solving the least squares

$$k = (D^T D)^{-1} D^T d \quad (3)$$

Finally all the parameters are refined, including lens distortion parameters by minimizing

$$\sum_{i=1}^n \sum_{j=1}^m \|m_{ij} - m(A, k_1, k_2, R_i, t_i, M_j)\|^2 \quad (4)$$

Once the intrinsic parameters like the focal length and principal point were obtained, they were then used for pose estimation.

```

330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439

```

The screenshot shows a camera calibration interface. On the left is a grayscale image of a Rubik's cube with a grid pattern. On the right is a terminal window displaying the camera matrix. The matrix includes values for 'R' (rotation), 'P' (projection), and 'K' (camera matrix). Below the terminal are four buttons: 'CALIBRATE', 'SAVE', and 'COMMIT'.

```

scale (000/100) □
[x=720,y=217] - R80 C155 B155

```

```

rscorer http://raam-SV... x raam@raam-SV15117F... x raam@raam-SV15117F... x
'R = ', [1.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 1.0]
'P = ', [486, 54693251953125, 0.0, 348, 4607733742656, 0.0, 0.0, 576, 10
09521484375, 207, 47569953969752, 0.0, 0.0, 0.0, 1.0, 0.0]
None
# oST version 5.0 parameters

[Image]
width
640
height
360
[narrow_stereo]
camera_matrix
006,638219 0,000000 356,752160
0,000000 611,309745 206,372199
0,000000 0,000000 1,000000
distortion
-0,598284 0,267622 -0,697477 -0,024298 0,000000
rectification
1,000000 0,000000 0,000000
0,000000 1,000000 0,000000
0,000000 0,000000 1,000000
projection
486,546933 0,000000 348,460773 0,000000
0,000000 576,10952 207,475699 0,000000
0,000000 0,000000 1,000000 0,000000

```

Figure 4. Camera calibration : On the right (output) the 'camera matrix' shows the scaling factor and camera offset during calibration

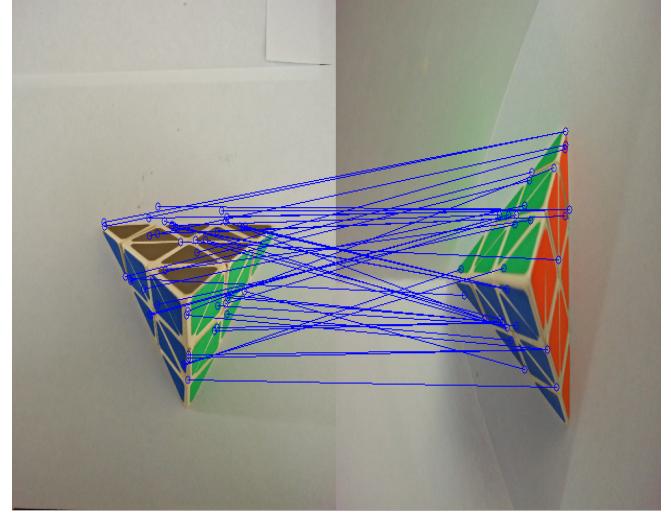


Figure 5. Failure of SURF: The image shows color side mismatch. Keypoints from blue side are matched erroneously with orange/green and vice versa.

We used SURF algorithm for feature matching on our pyramid object. It resulted in very bad performance 3.3.1. The primary reason for this behavior, is the lack of distinguishing features for different colored sides in the gray scale image of the object. Gray scale image does not contain enough information to distinguish between different colored regions. Taking a look at the object, we see that there are a small set of features in the object, that repeats itself on every colored side. Since gray scale image is not able to distinguish colors, there is a high probability of occurrence of wrong matching pairs from one colored side to a differently colored side. Failure to incorporate color information results in poor performance.

### 3.3.2. OPPONENT SURF

In order to incorporate color information, we use a new variant of SURF known as OpponentSURF (Van De Sande et al., 2010). The opponent color space ( $O_1, O_2, O_3$ ) of an RGB image is defined by the following expression.

$$\begin{bmatrix} O_1 \\ O_2 \\ O_3 \end{bmatrix} = \begin{bmatrix} \frac{R-G}{\sqrt{2}} \\ \frac{R+G-2B}{\sqrt{6}} \\ \frac{R+G+B}{\sqrt{3}} \end{bmatrix}$$

The  $O_1$  and  $O_2$  channels carry color information and the intensity information will cancel out due to the subtraction in their calculation. The  $O_3$  channel carry intensity information, since we average the three

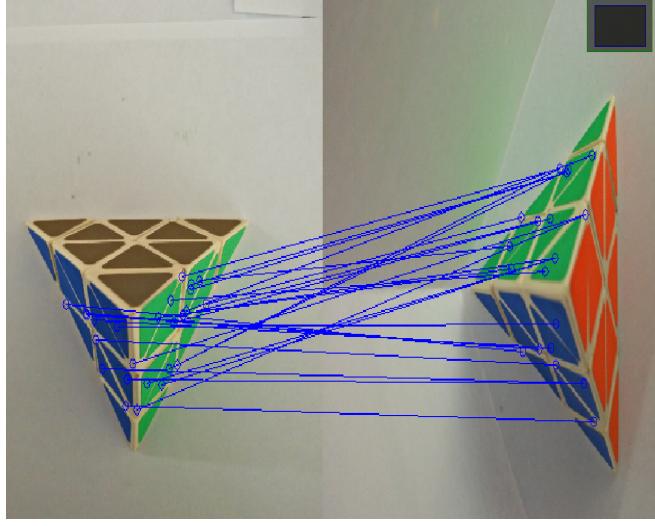


Figure 6. OpponentSURF: The image shows that same colored side matches on both images due to the incorporation of color information in OpponentSURF

channels. In OpponentSURF method, the keypoints are obtained from the same stage as the SURF algorithm. The difference is that instead of taking SURF descriptor around the keypoints in the gray scale image, OpponentSURF uses each of the opponent color channels. So for a given key-point, we obtain three SURF descriptors of length 128 corresponding to three color channels. The OpponentSURF descriptor is the concatenation of the three SURF descriptors and so it is a  $3 \times 128$  length vector.

OpponentSURF was able to find robust representations for the object, that gives good matching performance due to the incorporation of color information

### 3.3.2.

#### 3.3.3. FLANN MATCHING

We have a database of images with stored feature descriptors. When we have a new image, we extract features from that. We need to match these features on the incoming image with that of the features stored in the database. This process is called feature matching. We use FLANN matcher for this operation.

FLANN stands for Fast Library for Approximate Nearest Neighbors. It contains a collection of algorithms optimized for fast nearest neighbor search in large datasets and for high dimensional features(Muja & Lowe, 2014). It uses kNN algorithm to find the closest match. There is a problem in this approach, in that it results in poor matches. To alleviate that problem, ratio test is done. In this method, the first

and the second best match is taken and if the ratio between them is less than a threshold it is taken as an authentic match. We took the ratio to be .85 as suggested by Lowe(Lowe, 2004).

$$(F_A^i, F_B^j) : j = \min_k d(\text{opp}(F_A^i), \text{opp}(F_B^j)) < \delta \min_{k \neq j} d(\text{opp}(F_A^i), \text{opp}(F_B^k))$$

In the above expression,  $\text{opp}(F_A^i)$  indicates the opponent SIFT descriptors of feature point  $F_A$ .

### 3.4. Pose Estimation

Given the internal parameters, we find the external parameters for estimating the pose of the camera. Without the internal parameters or the correspondence points we found the camera parameters (both intrinsic and extrinsic) using Direct Linear Transform (DLT)(Hartley & Zisserman, 2003). DLT requires only one world coordinate (equivalent to one correspondence point) to find these parameters. The algorithm turned out to be slow and inaccurate for a fairly simple 4 dimensional simplex (object). Hence we choose calibrating the intrinsic parameters and hand-coding the coordinates of 6 correspondence and 4 control points over DLT.

We use EPnP algorithm to reduce the time complexity of iterating over n-correspondence points. We can represent any point inside the object with respect to camera ( $p_i^c$ ) as a linear combination ( $\alpha_{ij}$  barycentric coordinates) of the four control points  $c_j$  available to us.

$$p_i^c = \sum_{j=1}^4 \alpha_{ij} * c_j^C \quad (5)$$

The coordinates in image plane are given by,

$$p_i^I = \begin{bmatrix} f_x & 0 & C_x \\ 0 & f_y & C_y \\ 0 & 0 & 1 \end{bmatrix} * \sum_{j=1}^4 \alpha_{ij} * c_j^I \quad (6)$$

We can rewrite the above equation as  $Mx = 0$  where  $M$  is  $2n+12$  parametric matrix. Assuming that the distance between two control points in the world coordinates and the camera coordinates are same, we find the constants of linear combination ( $\mathbf{x} = \sum_i \beta_i v_i$ ) of eigen vectors of  $12 \times 12$  parametric matrix ( $M^T M$ ) for finding the coordinates of the correspondence points as

$$\beta = \frac{\sum_{i,j \in 1:4} \|\vec{v}_i - \vec{v}_j\| \cdot \|c_i - c_j\|}{\sum_{i,j \in 1:4} \|\vec{v}_i - \vec{v}_j\|^2}$$

where  $\vec{v}$ 's are the eigen vectors of  $M$ . And hence we find  $x = \sum_i \beta_i$ . The advantage of this method is that we convert any n- correspondence points into

an unknown of constant 12- dimensions. Though we have (including the control points) only 10 correspondence points, this method gives user the flexibility to choose more correspondence points. This was helpful in forming the mesh around the detected object in case of extreme environments(extremely cluttered or ill-illuminated i.e strong deviation, too high or too low, from the reference image). One such instance is shown in Fig.cluttered.

To find the pose of the camera we calculate  $p_i^I$  from the image. Then we pick one of the control points as 0. With this as reference we compute  $p_j^C$  for the remaining points. We compare the  $p_j^W$  with  $p_j^C$  and hence find R and T by comparing the initial and final frames of the object. We use the screw motions formula to estimate the amount of rotation.

$$\begin{bmatrix} R & T \\ 0 & 1 \end{bmatrix} = e^{[S] \cdot \theta}$$

where  $S$  is the screw (a vector) around which we rotate and  $\theta$  is the amount of rotation. Using  $S$  and  $\theta$  we can estimate  $R$  and  $T$  ??

### 3.5. Object Tracking

The objective of object tracking is estimate the trajectory of an object in the image plane as it moves around a scene. (Roach, 1980a) and (Roach, 1980b) explain few of the basic approaches of tracking using several points of object correspondence from image sequence. However object tracking in our case boils down to tracking the 3-D mesh of the object alone. We use Kalman Filter to estimate the position and orientation of the 3-D mesh. Another key advantage of Kalman filtering is to keep track of the object pose in those cases where the number of inliers is lower than a given threshold.

### 3.5.1. KALMAN FILTER AS A OBJECT TRACKER.

The Kalman filter is a recursive method to estimate the state of the object in the process of tracking (Broida & Chellappa, 1986). We use linear kalman filter to estimate the transnational and rotational position, velocity and acceleration. The successive states  $s_t \in \mathbb{R}^n$  of a discrete-time controlled process are assumed to evolve according to a dynamics model written as follows,

$$s_t = F s_{t-1} + w_t$$

where  $F$  is the state transition matrix. For the purpose of tracking the object the state vector consists of

six defining parameters of the 3-D mesh's pose along with their translational and rotational velocities and acceleration.

The state transition model  $F$  is given by the following matrix. The model computes the values of state components at time  $t$  given the estimate at time  $t-1$ , where the duration between successive time steps is  $\Delta t$ .

As we have already mentioned, measuring the camera pose is equivalent to estimating the real-time orientation of the object (Horn et al., 1988). Our implementation of linear kalman filter assumes that the measurement  $z_t$  of camera pose at any time  $t$  is related to state vector by a linear measurement model.

$$z_t = \mathcal{H} s_t + v_t$$

where  $\mathcal{H}$  is the linear measurement model.  $w_t$  and  $v_t$  are the process and measurement noise respectively which are assumed to be normally distributed with zero mean. The measurement model  $\mathcal{H}$  which can observe the 6 defining parameters of mesh from the state vector is given below.

$$\mathcal{H} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

The state of the kalman filter is represented as follows.

$$s_t = (x, y, z, \dot{x}, \dot{y}, \dot{z}, \ddot{x}, \ddot{y}, \ddot{z}, \psi, \theta, \phi, \dot{\psi}, \dot{\theta}, \dot{\phi}, \ddot{\psi}, \ddot{\theta}, \ddot{\phi})$$

where  $(.)$  represents the velocity component and  $(..)$  represents the acceleration component of any parameter. The six primary parameters of the camera pose are shown in fig 7,

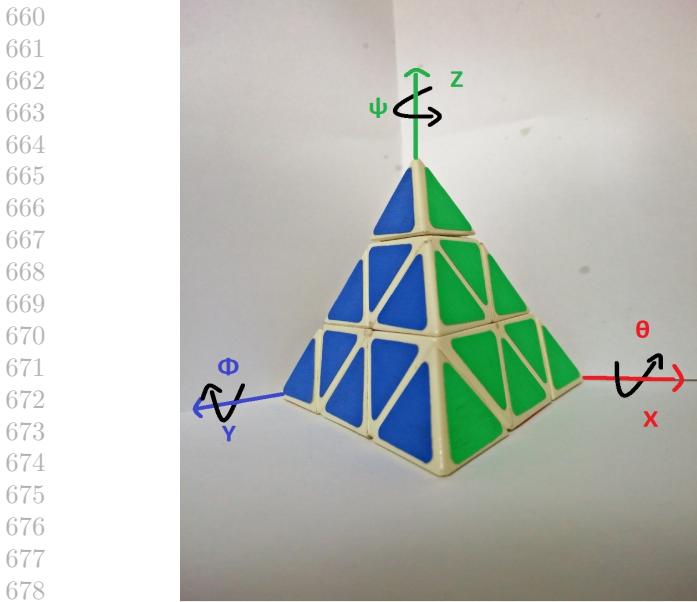


Figure 7. Image showing reference axis and pose angles used for implementing Kalman filter.

At each time step, the Kalman Filter makes an 'estimation' of the current 3-D mesh's pose parameters called the *apriori* state estimate  $s_t^-$ , which is refined by incorporating the measurements to yield the *aposteriori* estimate  $s_t$ .  $s_t^-$  and its covariance matrix  $S_t^-$ , are computed during the prediction stage and can be written as,

$$s_t^- = \mathbf{A}s_{t-1} \\ S_t^- = F S_{t-1} F^T + \mathbf{Q}$$

where  $S_{t-1}$  is the *aposteriori* estimate error covariance for the previous time step and  $\mathbf{Q}$  is the process noise covariance that measures quality of the motion model respect to the reality. Next, the Kalman Filter does a 'measurement update' or correction. The *aposteriori* state estimate  $s_t$  and its covariance matrix  $S_t$  are now generated by adding the measurements  $z_t$ .

$$s_t = s_t^- + \mathbf{K}(z_t - H s_t^-) \\ S_t = S_t^- - \mathbf{K} H S_t^-$$

In the above equations the term  $\mathbf{K}$  is called the Kalman Gain which controls a relation between the filter's use of predicting the parameters of mesh and measurement of camera's pose.  $\mathbf{K}$  can be computed from the following equation,

$$\mathcal{K}_t = S_t^- H^T (\mathbf{H} S_t^- \mathbf{H}^T + R)^{-1}$$

where  $\mathbf{Q}$  is the measurement covariance matrix. In our context of tracking, the parameters of pose of the mesh is given by the *apriori* estimate  $s_t$ . Hence the measurement prediction can be expressed as,

$$z_t^- = \mathbf{H} s_t^-$$

But this measurement prediction may be uncertain which is represented by covariance matrix  $\mathbf{P}$  by propagating the error.

$$P = \mathbf{H} S_t^- \mathbf{H}^T + R$$

## 4. Experiments and Result

The algorithm was able to successfully track the Rubik's pyramid in real time. Pose estimation was also successful, however in partial occlusion if any of the control points are occluded there was substantial error in pose-estimation. Even without occlusion during motion, simultaneous determination of both pose and velocities was not accurate. We have not incorporated algorithms that heuristically search objects that recover from complete occlusion. Hence during recovery from complete occlusions we search the whole image. This leads to a delay of 0.59 seconds for the recovery of tracking. Because of the low frame rate of 30fps captured by Logitech C270 and also our sampling rate of 10fps, the kalman filter was not able to track approximately 10000 pix/sec in image. Kalman filter considers such velocity as above threshold. When the angle of incidence of light source and the viewing angle was equal and both have a high value, the reflection from the light source distorted the color information. Due to this OpponentSURF was not able to detect feature points.

## 5. Conclusion and Future Work

Our algorithm was able to successfully track the given object by exploiting its structure but with some limitations. One good direction for future work is to experiment with more robust tracking algorithm other than kalman filter. Also for pose estimation we could use UPnP instead of EPnP since the former doe snot need camera calibration.

## 6. Acknowledgement

We thank Prof. Mike Gennert for his constant motivation. We also extend our gratitude to Koushik Bal-

660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714

715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769

770	asubramiam.	825
771		826
772	<b>7. Video Demonstration</b>	827
773		828
774	Please visit the following links for the video showing	829
775	our implementation.	830
776		831
777	• Object recognition and Pose estimation of	832
778	a pyramid.	833
779		834
780		835
781	• Recognition of a Rubik's Pyramid	836
782		837
783		838
784	<b>References</b>	839
785		840
786	Bay, Herbert, Tuytelaars, Tinne, and Van Gool, Luc.	841
787	Surf: Speeded up robust features. In <i>Computer</i>	842
788	<i>vision-ECCV 2006</i> , pp. 404–417. Springer, 2006.	843
789		844
790	Broida, Ted J and Chellappa, Rama. Estimation of ob-	845
791	ject motion parameters from noisy images. <i>Pattern</i>	846
792	<i>Analysis and Machine Intelligence, IEEE Transac-</i>	847
793	<i>tions on</i> , (1):90–99, 1986.	848
794	Hartley, R. and Zisserman, A. <i>Multiple View Geometry</i>	849
795	<i>in Computer Vision</i> . Cambridge University Press,	850
796	2003. ISBN 9780511184512. URL <a href="https://books.google.com/books?id=eLbfQEACAAJ">https://books.</a>	851
797	<a href="https://books.google.com/books?id=eLbfQEACAAJ">google.com/books?id=eLbfQEACAAJ</a> .	852
798		853
799	Horn, Berthold KP, Hilden, Hugh M, and Negah-	854
800	daripour, Shahriar. Closed-form solution of absolute	855
801	orientation using orthonormal matrices. <i>JOSA A</i> , 5	856
802	(7):1127–1135, 1988.	857
803		858
804	Lowe, David G. Distinctive image features from scale-	859
805	invariant keypoints. <i>International journal of com-</i>	860
806	<i>puter vision</i> , 60(2):91–110, 2004.	861
807		862
808	Muja, Marius and Lowe, David G. Scalable near-	863
809	est neighbor algorithms for high dimensional data.	864
810	<i>Pattern Analysis and Machine Intelligence, IEEE</i>	865
811	<i>Transactions on</i> , 36(11):2227–2240, 2014.	866
812	Pi, Edgar Riba. Implementation of a 3d pose estima-	867
813	tion algorithm. <i>MasterâŽs Degree in Automatic</i>	868
814	<i>Control and Robotics</i> , 2015.	869
815		870
816	Roach, John Wilson. Determining the movement of	871
817	objects from a sequence of images. <i>Pattern Analysis</i>	872
818	<i>and Machine Intelligence, IEEE Transactions on</i> ,	873
819	(6):554–562, 1980a.	874
820		875
821	Roach, John Wilson. Determining the movement of	876
822	objects from a sequence of images. <i>Pattern Analysis</i>	877
823	<i>and Machine Intelligence, IEEE Transactions on</i> ,	878
824	(6):554–562, 1980b.	879