

# Fusion of IMU and Vision for Absolute Scale Estimation in Monocular SLAM

## ENPM667 Midterm Project Report

Harsh Kakashaniya · Rohitkrishna Nambiar

Submitted: 19th November 2018

**Abstract** The fusion of inertial and visual data is widely used to improve an objects pose estimation. However, this type of fusion is rarely used to estimate further unknowns in the visual framework. In this paper we present and compare two different approaches to estimate the unknown scale parameter in a monocular SLAM framework. Directly linked to the scale is the estimation of the objects absolute velocity and position in 3D. The first approach is a spline fitting task adapted from Jung and Taylor and the second is an extended Kalman filter. Both methods have been simulated offline on arbitrary camera paths to analyze their behavior and the quality of the resulting scale estimation. We then embedded an online multi rate extended Kalman filter in the Parallel Tracking and Mapping (PTAM) algorithm of Klein and Murray together with an inertial sensor. In this inertial/monocular SLAM framework, we show a real time, robust and fast converging scale estimation. Our approach does not depend on known patterns in the vision part nor a complex temporal synchronization between the visual and inertial sensor.

**Keywords** IMU vision fusion · Absolute scale · Monocular SLAM · Kalman filter

---

Harsh Kakashaniya  
UID: 116311236  
E-mail: harshbk@umd.edu

Rohitkrishna Nambiar  
UID: 115507944  
E-mail: rohit517@umd.edu

## Contents

1	Introduction . . . . .	4
2	Hardware Setup . . . . .	4
3	Camera and Image formation process . . . . .	5
4	Inertial Measurement Units . . . . .	7
4.1	Introduction . . . . .	7
4.2	Working . . . . .	7
5	What's wrong with monocular camera? . . . . .	8
6	Spline Fitting Method . . . . .	9
6.1	Types of curve fitting . . . . .	10
6.1.1	Maximum error . . . . .	10
6.1.2	Average error . . . . .	10
6.1.3	RMS (Least square method) . . . . .	11
6.2	IMU Reading . . . . .	12
6.3	Results . . . . .	14
7	Extended Kalman Filter . . . . .	15
8	Simultaneous Localization and Mapping . . . . .	20
8.1	Visual Odometry . . . . .	21
8.2	MonoSLAM . . . . .	23
8.3	Parallel Tracking and Mapping (PTAM) . . . . .	23
8.4	Results . . . . .	24
8.4.1	Simulated and Real Data . . . . .	24
8.4.2	Online Implementation . . . . .	25

## List of Figures

1	<b>a.</b> Camera/IMU setup. <b>b.</b> Frame transformations between IMU, Camera and World frame. . . . .	4
2	Pinhole Image formation [25] . . . . .	5
3	World co-ordinates to pixel co-ordinates [17] . . . . .	6
4	IMU Physical model [7] . . . . .	7
5	Working of IMU model [7] . . . . .	8
6	Monocular camera scale problem . . . . .	9
7	Curve Fitting with Maximum error . . . . .	11
8	Curve Fitting with Average error . . . . .	11
9	Curve Fitting with RMS error . . . . .	12
10	In this case Blue is actual curve traced by IMU this are points with noise. And green is plotted curve. With least square method.	15
11	Multiple spline fit . . . . .	16
12	<b>a.</b> Red curve for method with one curve and least square method. <b>b.</b> Green curve is using the given method by Jung and Taylor.	17
13	Actual 3D path with noise vs Optimized spline fit . . . . .	18
14	Kalman filter state estimation [1] . . . . .	18
15	Non-Linear system [1] . . . . .	19
16	SLAM architecture [4] . . . . .	21
17	Generalized Visual Odometry Pipeline [20] . . . . .	21
18	Offline simulation results with simulated and real data. Fewer directions (X,Y,Z) included in Kalman filter gives us a better estimate of $\lambda$ . [19] . . . . .	25
19	Online implementation with PTAM work flow. Blue color boxes are original PTAM threads. Yellow boxes are added threads for scale estimation $\lambda$ . [19] . . . . .	26

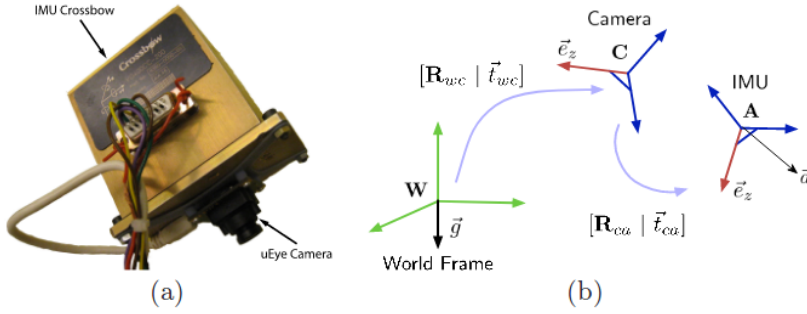
## 1 Introduction

This is a report of the paper titled *Fusion of IMU and Vision for Absolute Scale Estimation in Monocular SLAM* [19]. On-board pose estimation is very useful for a variety of applications ranging from autonomous robotics to augmented reality. Pose estimation can be done using different sensors such as cameras and inertial measurement units (IMU). Simultaneous localization and mapping (SLAM) using a monocular camera can be used to get the trajectory of the camera and a map of the environment. But monocular cameras suffer from scale ambiguity causing the overall trajectory to drift making it unusable in real time. Using a stereo camera helps solve the scale ambiguity. IMUs on the other hand can be used to get the trajectory traveled by integrating the acceleration measurements over time. But this leads to highly inaccurate trajectory estimates. The estimate of the scale factor is essential to fuse both these measurements. This fusion helps us determine the unknown scale factor  $\lambda$ .

In this paper [19], two methods are presented for scale estimation. The first one is a spline fitting method by Jung and Taylor [11]. The second is a multi rate Extended Kalman Filter(EKF). Both the approach have been simulated in MATLAB.

This report is organized as follows. Section 2 goes over the Camera and IMU setup, Section 3 covers the camera and image formation process, Section 5 goes over the scale estimation problem with monocular camera, Section 4 explains working of inertial measurement units, Section 6 outlines the spline fitting method with results, Section 7 covers Extended Kalman Filter and Section 8 explains visual odometry, SLAM along with filter based and key-frame based methods.

## 2 Hardware Setup

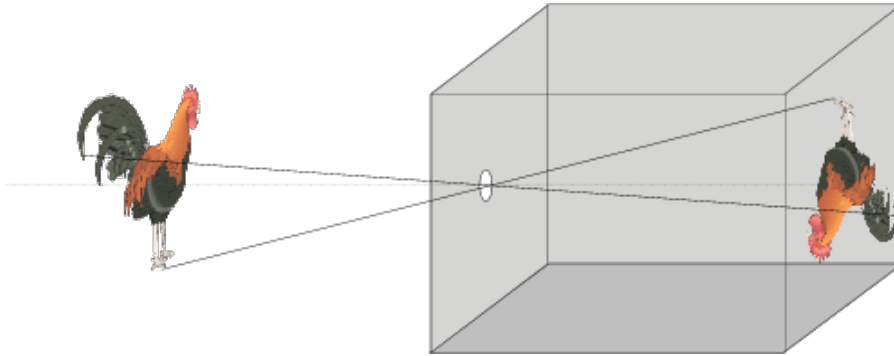


**Fig. 1** a. Camera/IMU setup. b. Frame transformations between IMU, Camera and World frame.

This section goes through the hardware setup used in the research paper. USB uEye UI-122xLE fisheye camera was used as the vision input. The camera has a resolution of 752x480 and a frame rate up to 87fps. Motion blur was minimized due to the high dynamic range and global shutter of the camera. VG400CC-200 solid state gyro which includes a tri-axial gyroscope and tri-axial accelerometer was used. It has an output frequency of 75 Hz with an input range of  $\pm 10g$  and  $< 1.25mg$  resolution. The IMU outputs acceleration around its 3-axis along with yaw, pitch and roll.

### 3 Camera and Image formation process

To understand how a camera perceives the environment, it is necessary we understand the image formation process. For us humans, we see things when a light originating from a light source is reflected on an object and enters our eyes. A camera acts very much similar to the human eye. The earliest and first model of an optical camera is the pinhole camera which is a simple and highly accurate representation of our eye model. This is the simplest device to form an image of a 3D scene on a 2D surface. As seen in figure 1 rays of light enter the pinhole and forms an inverted image of the object. This is called perspective projection. As the image formed in the image plane is inverted, we consider a virtual plane in front of the pinhole that acts as the image plane.



**Fig. 2** Pinhole Image formation [25]

Although the pinhole model is quite accurate, modern cameras have lenses. They help gather more light from a source leading to higher quality sharper images. Ignoring the internal diffraction, we assume thin lens equations for the imaging process. Every camera has a region of depths over which the scene is sharp. Modern cameras have variable apertures which help in focusing on objects at varying distances in the scene.

Every camera has an intrinsic, extrinsic and distortion parameters that are important to understand for every application. These parameters are used to

correct for lens distortion, measure objects in physical world and to determine the location of the camera in the real world. Camera calibration is the process of estimating these parameters specific to a given camera and application setup. Figure 2 demonstrates how the a real world 3D coordinate is related to a 2D pixel coordinate using the parameters mentioned above. The intrinsic camera matrix  $\mathbf{K}$  is given by

$$K = \begin{pmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \quad (1)$$

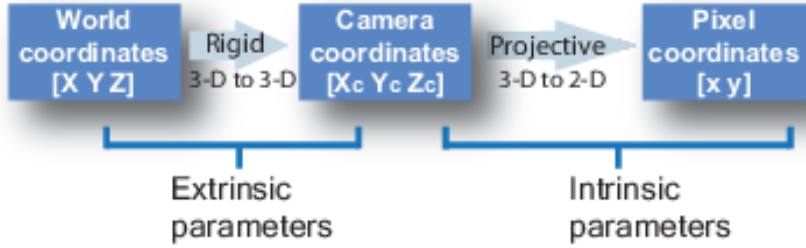
where  $(f_x, f_y)$  is the focal length,  $(c_x, c_y)$  is the optical center and  $s$  is skew factor. The extrinsic camera matrix  $[\mathbf{R} \mid \mathbf{t}]$  is given by

$$R = \begin{pmatrix} r_1 & r_2 & r_3 \\ r_4 & r_5 & r_6 \\ r_7 & r_8 & r_9 \end{pmatrix} \quad t = \begin{pmatrix} t_x \\ t_y \\ t_z \end{pmatrix} \quad (2)$$

where  $R$  is the rotation matrix and  $t$  is the translation vector between the camera center and world co-ordinates.

The camera matrix  $P$  is given by

$$P = K[Rt] \quad (3)$$



**Fig. 3** World co-ordinates to pixel co-ordinates [17]

The extrinsic parameters relate the rotation  $R$  and translation  $t$  of the camera origin at the optical center to the world frame. The intrinsic parameters consists of the focal length given by  $f_x$  and  $f_y$ , the optical center given by  $c_x$  and  $c_y$  and the skew coefficient  $s$ . The distortion parameters consists of radial and tangential distortion along the  $x$  and  $y$  direction respectively. The number of coefficients are decided based on the lens in consideration. Camera calibration techniques presented in [26, 2, 9] are widely used in commercial and open source camera calibration tool boxes.

## 4 Inertial Measurement Units

### 4.1 Introduction

Inertial measurement units (IMU) are self-contained systems that measure linear and angular motion of an object or vehicle. Measurements are summed over a time period to determine the instantaneous position, velocity, orientation, and direction of movement. The IMU is comprised of at least two dedicated sensors, one or more linear accelerometers and one or more gyroscopes or angular accelerometers. An optional magnetometer may be integrated into the unit to calibrate against orientation drift.



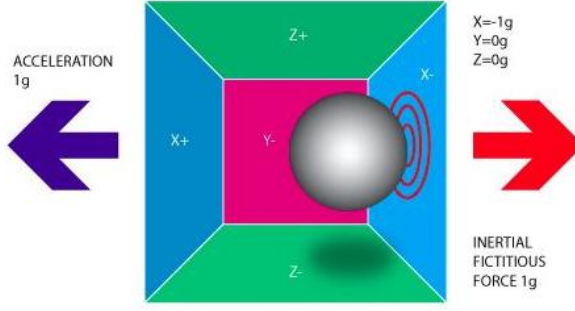
**Fig. 4** IMU Physical model [7]

### 4.2 Working

IMUs operate by use of reference data, bias values from an initial starting point, and calculate changes to these values using its integrated sensors. A central processing unit calculates directional information; position, speed, orientation, and direction of movement, at a given time in space using the IMU. The sensors suffer from orientation drift as they calculate these variables using a process known as dead-reckoning and are subject to accumulative errors. Calibration parameters can be stored in the memory of an IMU and are automatically reflected in the measurement data. In-situ calibrations may also be accomplished by use of a magnetometer or GPS unit to correct for orientation drift and improve the accuracy of the directional information.

Here we can see there is a weight less ball. And we are considering forward acceleration of 1 g hence ball is exerting inertial force on the box in -X direction as it is principle of action and reaction. So here we can learn that

$$Acceleration = \frac{\text{force exerted}}{\text{mass of ball}}$$



**Fig. 5** Working of IMU model [7]

Hence computing load on -X face of wall will tell us about the acceleration. similarly if we have a ground on one side it will exert a force of  $g$  in lower direction which will be constant. There will also be conditions where forces will be applied on two sides combined. So its magnitude will give is the angle of tilt and accelerations in any angular directions.

So if we consider enclosed box and a ball we can measure all 3 acceleration in X,Y,Z due to all side enclosed box. We can use sine and cosine functions to represent angles with different axis. So if we take Angular acceleration is differentiation of velocity where

$$Angularvelocity_x = \frac{\theta_{x2} - \theta_{x1}}{T_2 - T_1}$$

So with equations and forces in all direction we can compute all 3 axis linear accelerations and all 3 angular accelerations(Gyro accelerations).

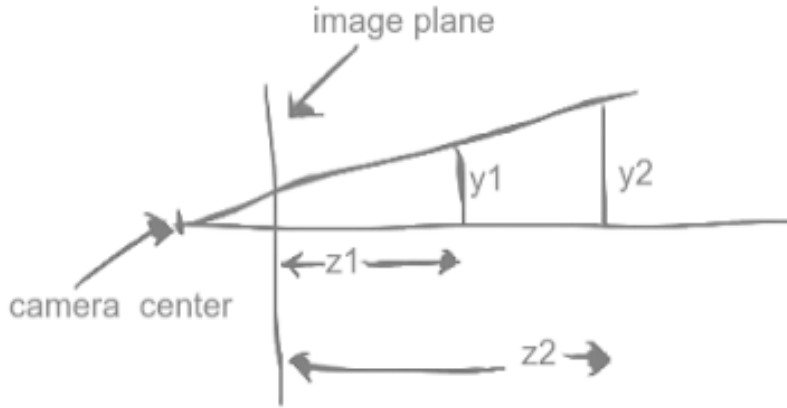
## 5 What's wrong with monocular camera?

As the name suggests a monocular camera is a camera with a single imaging sensor. In section X. we say that given a perfectly calibrated camera ie. with the intrinsic and extrinsic parameters known, we can measure size of objects in the real world. This however is not true for a single image taken from a monocular camera. When 3D objects are projected onto the image plane (2D), the depth information stored in the z-axis is lost. This can be seen in the example below.

As seen in the Fig 6, irrespective of the location of the object in the real world, the size of the object in the image plane remains the same. Mathematically, this can be shown as

$$\frac{y}{f} = \frac{y_1}{z_1 + f} = \frac{y_2}{z_2 + f} \quad (4)$$





**Fig. 6** Monocular camera scale problem

Where  $f$  is the focal length of the camera. The title of the paper mentions scale estimation. In the context of SLAM which will be covered in the next section, scale is a term used to denote the factor by which the computed trajectory by a SLAM algorithm needs to be multiplied/scaled to make it equal to the ground truth. For example, if the measured distance is 2m, whereas the ground truth is 4m, the scale value would be 2. This is an area of active research and the community has come up with different techniques to solve this problem. As the title suggests, we would be talking about the fusion of IMU and Vision for solving this challenge.

## 6 Spline Fitting Method

In this paper there are 2 methods of fusing IMU data and Monocular camera. First of which is spline fitting using scale factor with the data of IMU and Monocular camera. Depending on the reliability of the data from both sources to estimate distance. So here we did simulation of whole system with arbitrary value of data considering linear moment of robot with noisy data. These simulations can be seen in result section. With the matlab code. So initially theory behind curve fitting is explained to have a algebraic background of types of curve fitting and their applications. We use spline in curve fitting because we get discrete data from sensors and camera but tracing spline helps for interpolation and extrapolation to have continuous high probability data of path.

Spline is fitted on the data points which is output of sensor in our case. Spline fitting also helps in reducing error due to noise. For Improving accuracy we can increase number of spline but this will also result into increase in computation. So we have to trade off between accuracy and computation.

Spline fitting is useful to give smooth curve .N degree spline may has N-1 curves in its shape. So if spline has X3 term that means it can have 2 curves in its spline. Ideal method to cover n points in a plan is by tracing nth order spline which will result into zero error. So spline will pass through all the points and will have a general equation as

$$A_0 + A_1X + A_2X^2 + A_3X^3 + \dots + A_nX^n = Y \quad (5)$$

This problem then converts in simple  $Ax = B$  form. we can get values of n terms by plugging in n points values.

Where

$$A = \begin{pmatrix} 1 & X_1 & X_1^2 & \dots & X_1^n \\ 1 & X_2 & X_2^2 & \dots & X_2^n \\ 1 & X_3 & X_3^2 & \dots & X_3^n \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^n \end{pmatrix} B = \begin{pmatrix} A_0 \\ A_1 \\ A_2 \\ \vdots \\ A_n \end{pmatrix} C = \begin{pmatrix} Y_0 \\ Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \quad (6)$$

But for practical applications we use other methods and do not plot n degree spline as it is computationally expensive .Secondly if there is more noise in system it will give wrong results. For practical application there are 3 kinds of spline fitting according to the application:-

## 6.1 Types of curve fitting

### 6.1.1 Maximum error

This is the method where spline is fit in a way where the point with maximum error is reduced and curve move towards the outlier point. If there exist in order to reduce maximum error of a point. This on the other hand results into

$$E(f) = |\max(f(x_i) - y_i)| \text{ where } i \text{ goes from } 1 \text{ to } n \quad (7)$$

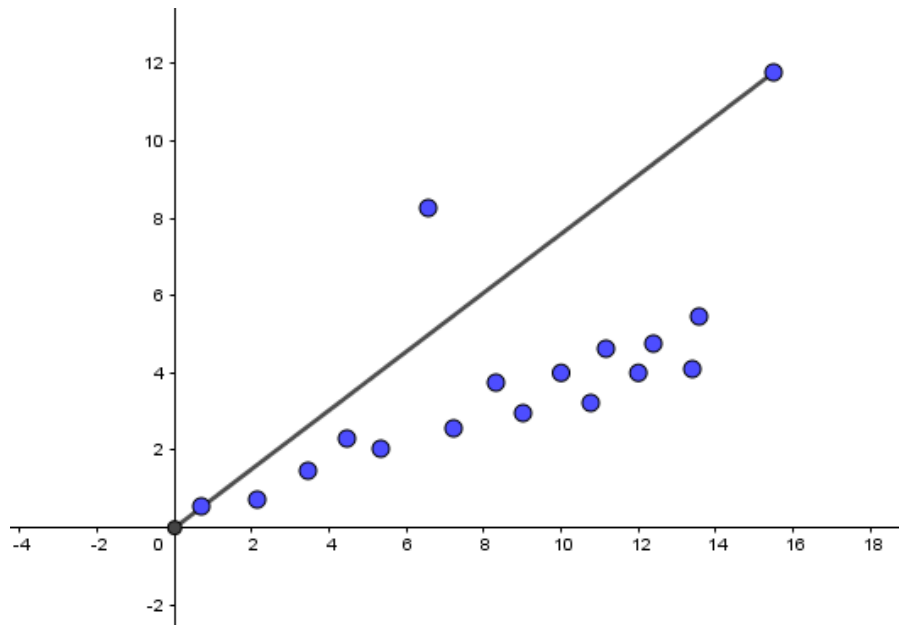
Actually due to one outlier here line shifted in order to reduce maximum error. So this is not reliable process but is easy to compute.

### 6.1.2 Average error

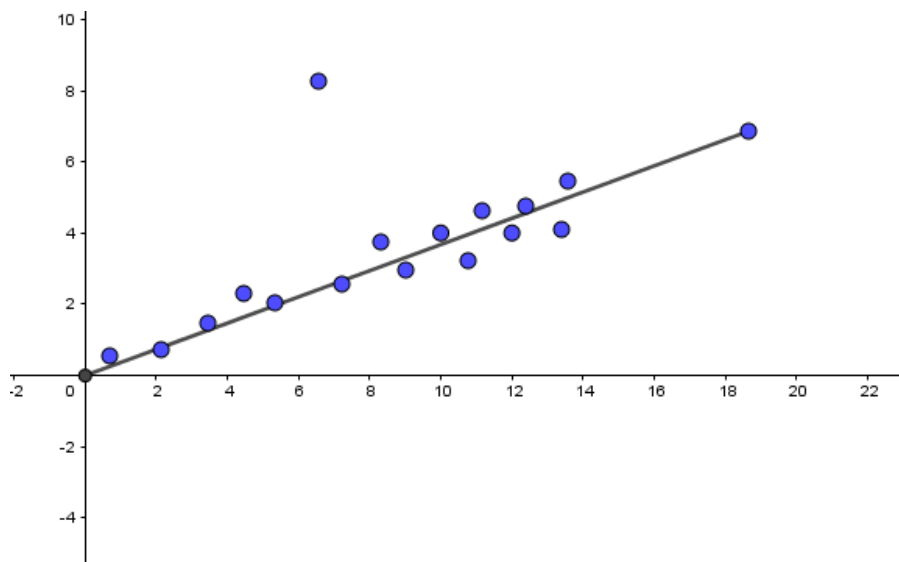
This method is used to minimize average and fit the spline with minimum error conditions.

$$E(f) = \frac{1}{n} * \sum |f(x_k) - Y_k| \quad (8)$$

Actually due to one outlier here line shifted in order to reduce maximum error. So this is not reliable process but is easy to compute.



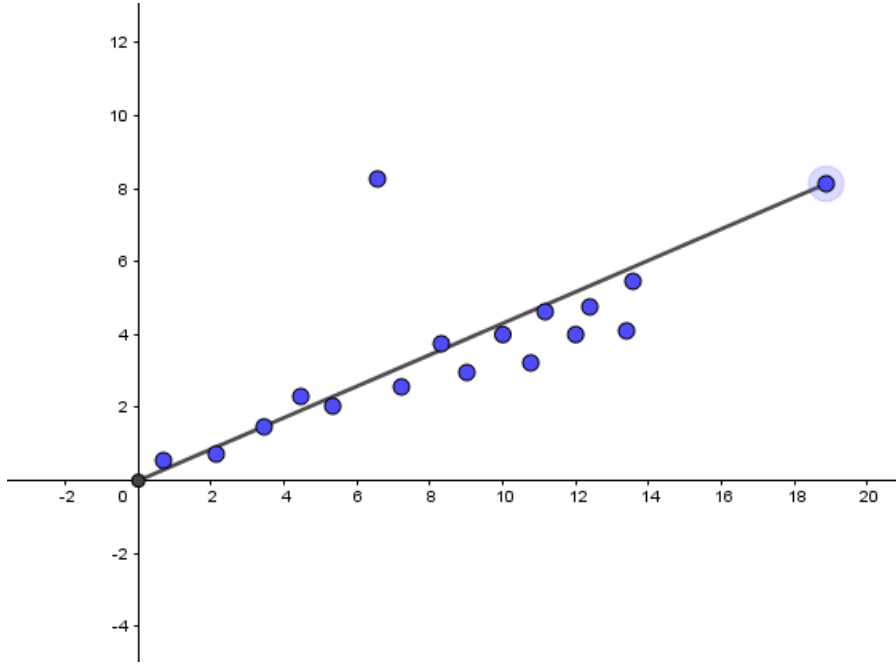
**Fig. 7** Curve Fitting with Maximum error



**Fig. 8** Curve Fitting with Average error

### 6.1.3 RMS (Least square method)

This is the most efficient method to minimize the error this is also called as least square method. Hence this is mostly used in curve fitting. We have also created an algorithm with the help of same type of curve fitting.



**Fig. 9** Curve Fitting with RMS error

$$E(f) = \left\{ \frac{1}{n} * \sum |f(x_k) - Y_k| \right\}^{\frac{1}{2}} \quad (9)$$

## 6.2 IMU Reading

As we dont have real IMU lets simulate virtual readings for this we took X, Y, Z as a function on time

$$X_1 = A_y t^2 + B_y t + C_x \quad (10)$$

$$Y_1 = A_y t^2 + B_y t + C_y \quad (11)$$

$$Z_1 = A_z t^2 + B_z t + C_z \quad (12)$$

With this treatment we will add some random noise to all the terms so that the data can be similar to what we obtain form IMU.

$$noise = A * (rand(1, n) - 0.5) \quad (13)$$

So our noise will range from  $-A/2$  to  $A/2$ . As function is  $rand()$  so it will change for every point.

$$X = X_1 + noise \quad (14)$$

$$Y = Y_1 + noise \quad (15)$$

$$Z = Z_1 + noise \quad (16)$$

We also took some random points of camera pose. And provided our system with data set. Proof of Least square method in our case: Lets consider optimization of X first.

From equation (10) we have

Here  $A$ ,  $B$ ,  $C$  are known We have data set of X with respect to t To calculate error in RMS we have

$$Error_x = \left\{ \sum_{i=1}^n (|X_i - (A_x t_i^2 + B_x t_i + C_x)|)^2 \right\}^{\frac{1}{2}} \quad (17)$$

When we differentiate some quantity and equate it to zero it goes to either minima or maxima And if double derivative is positive it goes to minima.

So in this case if we take partial derivatives of error with respect to  $A_x$ ,  $B_x$ ,  $C_x$

$$\frac{\partial Error_x}{\partial A_x} = \left\{ \sum_{i=1}^n 2(|X_i - (A_x t_i^2 + B_x t_i + C_x)|) \right\} (-t_i^2) \quad (18)$$

$$\sum_{i=1}^n (|-t_i^2 X_i + (A_x t_i^4 + B_x t_i^3 + C_x t_i^2)|) = 0 \quad (19)$$

$$\frac{\partial Error_x}{\partial B_x} = \left\{ \sum_{i=1}^n 2(|X_i - (A_x t_i^2 + B_x t_i + C_x)|) \right\} (-t_i) \quad (20)$$

$$\sum_{i=1}^n (|-t_i X_i + (A_x t_i^3 + B_x t_i^2 + C_x t_i^1)|) = 0 \quad (21)$$

$$\frac{\partial Error_x}{\partial C_x} = \left\{ \sum_{i=1}^n 2(|X_i - (A_x t_i^2 + B_x t_i + C_x)|) \right\} (-1) \quad (22)$$

$$\sum_{i=1}^n (|-1 X_i + (A_x t_i^2 + B_x t_i + C_x)|) = 0 \quad (23)$$

So we can compute matrix in such a way that the problem changes to

$Ax=B$

Where,

$$A = \begin{pmatrix} \sum t^4 & \sum t^3 & \sum t^2 \\ \sum t^3 & \sum t^2 & \sum t \\ \sum t^2 & \sum t & \sum 1 \end{pmatrix} B = \begin{pmatrix} A_x \\ B_x \\ C_x \end{pmatrix} C = \begin{pmatrix} \sum X_i t^2 \\ \sum X_i t \\ \sum X_i \end{pmatrix} \quad (24)$$

From this equation we can calculate X matrix. Similarly, We can get X matrix of Y and Z axis.

$$A = \begin{pmatrix} \Sigma t^4 & \Sigma t^3 & \Sigma t^2 \\ \Sigma t^3 & \Sigma t^2 & \Sigma t \\ \Sigma t^2 & \Sigma t & \Sigma 1 \end{pmatrix} B = \begin{pmatrix} A_y \\ B_y \\ C_y \end{pmatrix} C = \begin{pmatrix} \Sigma Y_i t^2 \\ \Sigma Y_i t \\ \Sigma Y_i \end{pmatrix} \quad (25)$$

$$A = \begin{pmatrix} \Sigma t^4 & \Sigma t^3 & \Sigma t^2 \\ \Sigma t^3 & \Sigma t^2 & \Sigma t \\ \Sigma t^2 & \Sigma t & \Sigma 1 \end{pmatrix} B = \begin{pmatrix} A_z \\ B_z \\ C_z \end{pmatrix} C = \begin{pmatrix} \Sigma Z_i t^2 \\ \Sigma Z_i t \\ \Sigma Z_i \end{pmatrix} \quad (26)$$

By this treatment we converted given points into equation with least square method here we used condition of

$$\min \begin{pmatrix} A_x t^2 + B_x t + C_x - \lambda_i X_c \\ A_y t^2 + B_y t + C_y - \lambda_i Y_c \\ A_z t^2 + B_z t + C_z - \lambda_i Z_c \end{pmatrix}^2 \quad (27)$$

Here in this formula we have different value of for different spline according to accuracy of camera data and IMU reading.

### 6.3 Results

In our case we simulated the results and found following outputs. Plotting the second order curve for given data of IMU with just least square method. We get

Then according to the paper we fused data of camera and IMU so we get following graph.

In Figure 11 we have 3 curved fitted as shown as red which gives better results and it also uses data of camera shown with black dots. So here according to the above formula scale factor of different section is different it is according to the quality of reading by camera. For example, in our case

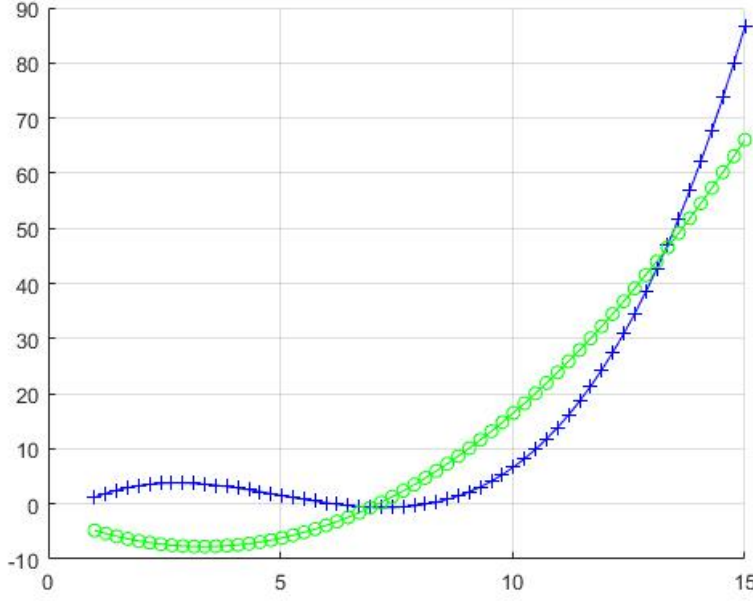
$$\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 0.1 \quad (28)$$

In 3rd section camera readings were not accurate so taken into smaller scalar value. This camera data further improved the result. And we get minimum error and good fit. With the given spline fitting method.

So if we look into different of error with traditional Least square method with Jung and Taylor method. It gives good results. Almost we get the same curve.

The graph in Figure 13 has error comparison with and without Jung and Taylor method of split curve fitting and scaling camera input. So it is better to follow the given first method in the paper. It results into considerable reduction in error which is clearly seen in graph.

With the same approach we can optimize and plot spline in 3 Dimension so when we use this algorithm for 3D space we get the following result where we get an second order optimized curve. we also included considerable noise and tested our algorithm.



**Fig. 10** In this case Blue is actual curve traced by IMU this are points with noise. And green is plotted curve. With least square method.

## 7 Extended Kalman Filter

In estimation theory, the extended Kalman filter (EKF) is the nonlinear version of the Kalman filter which linearizes about an estimate of the current mean and covariance. In the case of well defined transition models, the EKF has been considered to estimate the state. In this system we have a state and an observer these are given by equation.  $\dot{X} = AX + BU + W_d$  and  $Y = CX + W_n$ . In this  $W_d$  is disturbance and  $W_n$  is the noise by sensor. So Kalman filter is used to judge between data as disturbance and noise.

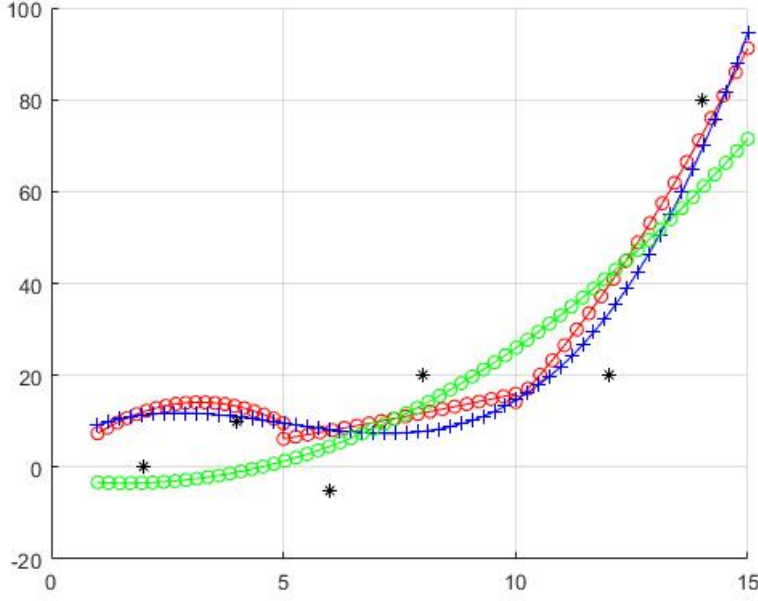
If  $W_d \gg W_n$  that means disturbance amount is higher than noise hence state estimation can rely more on sensor data for mapping than on previous state. And sensor is more reliable.

If  $W_n \gg W_d$  that means noise amount is higher than disturbance hence state estimation can rely more on previous plot data for mapping than on sensor. And state is more reliable than sensor.

So according to the paper we learned that non-linear system is defined as.

So now let's take non-linear state equation as follows

$$\vec{z}_{k+1} = \vec{f}_k(\vec{z}_k) + v_k \quad (29)$$



**Fig. 11** Multiple spline fit

---

**Algorithm 1** Kalman filter algorithm

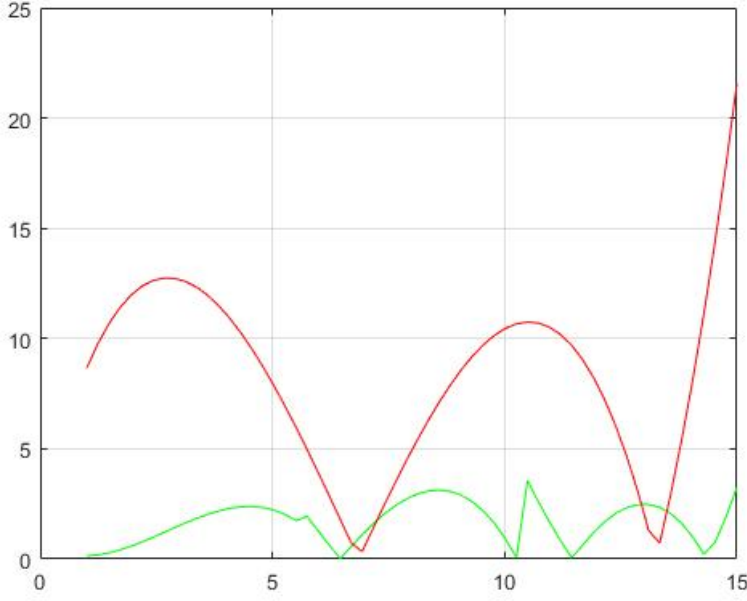
---

- 1: Kalman Filter ( $\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$ )
  - 2:  $\mu_t = A_t \mu_{t-1} + B_t u_t$
  - 3:  $\Sigma_t = A_t \Sigma_{t-1} A_t^T + R_t$
  - 4:  $K_t = \Sigma_t C_t^T (C_t^T \Sigma_t C_t^T + Q_t)^{-1}$
  - 5:  $\mu_t = \mu_t + K_t (z_t - C_t \mu_t)$
  - 6:  $\Sigma_t = (I - K_t C_t) \Sigma_t$
  - 7: return  $\mu_t, \Sigma_t$
- 

$$\begin{pmatrix} \vec{x}_{k+1} \\ \vec{v}_{k+1} \\ \vec{a}_{k+1} \\ \vec{\lambda}_{k+1} \end{pmatrix} = \begin{pmatrix} I_3 & \frac{T}{\lambda} I_3 & \frac{T^2}{2\lambda} I_3 & 0 \\ 0 & I_3 & T I_3 & 0 \\ 0 & 0 & I_3 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \vec{x}_k \\ \vec{v}_k \\ \vec{a}_k \\ \vec{\lambda}_k \end{pmatrix} \quad (30)$$

where  $x_k + 1$  is the position without scale of the IMU/Camera and  $v_{k+1}$ ,  $a_k + 1$  are the velocity and acceleration of the IMU/Camera in metric unit [m].  $k$  is the gaussian process noise. Main reason behind extended kalman filter working is when ever transition takes places between matrics resultant output remains Gaussian and have defined mean and co-variance so in out treatment we can only do those transformation where we dont loose condition of system being Gaussian. For example, Linear curve when transformed is Gaussian but if the function is non linear it is resultant does not result into Gaussian so



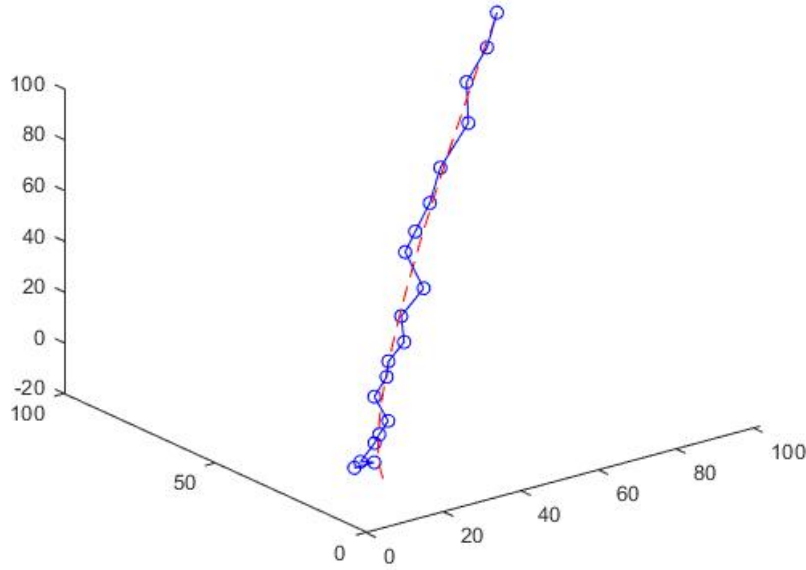


**Fig. 12** **a.** Red curve for method with one curve and least square method. **b.** Green curve is using the given method by Jung and Taylor.

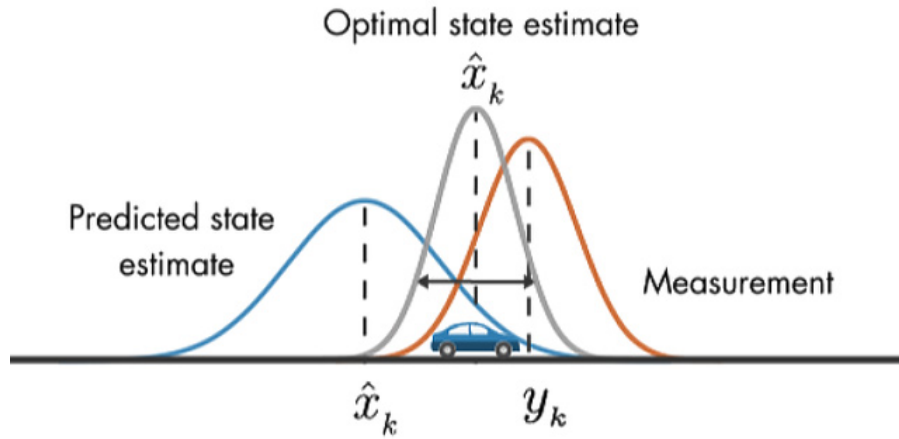
whole point of using Kalman filter is killed. So in Extended Kalman filter we linearize the system. Every vector in  $z_k$  is resolved in the world frame  $W$ . Note that we do not include the orientation information in the model nor use it as a measurement in order to keep the algorithm simple and fast. On each acceleration measurement we do the conversion from the inertial to the world frame by using a zero order hold of the unfiltered attitude measurement returned by the visual SLAM framework. As we work in a middle size environment with enough loops we assume negligible drift in the SLAM map and assume thus highly accurate attitude estimation from the visual SLAM framework. The model in its linearized form yields,

$$F_k = \begin{pmatrix} I_3 & \frac{T}{\lambda} I_3 & \frac{T^2}{2\lambda} I_3 & -\frac{T}{\lambda^2} I_3 & -\frac{T^2}{2\lambda^2} I_3 \\ 0 & I_3 & T I_3 & 0 & 0 \\ 0 & 0 & I_3 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix} \quad (31)$$

For fusion implementation we consider the measurements in different observation vectors. For a multi rate filter, as it is in our case, the literature suggests two solutions. One would be using a (higher order) hold to synchronize the different measurements. Another is to weight the uncertainty of the measurement according to its temporal occurrence. We claim no certainty at



**Fig. 13** Actual 3D path with noise vs Optimized spline fit



**Fig. 14** Kalman filter state estimation [1]

all if no measurement is available (i.e. the measurement noise variance is infinite). Thus the update equations simplify to improve results. A more complex weighting function (i.e. exponential decay in time) could also be applied, how-

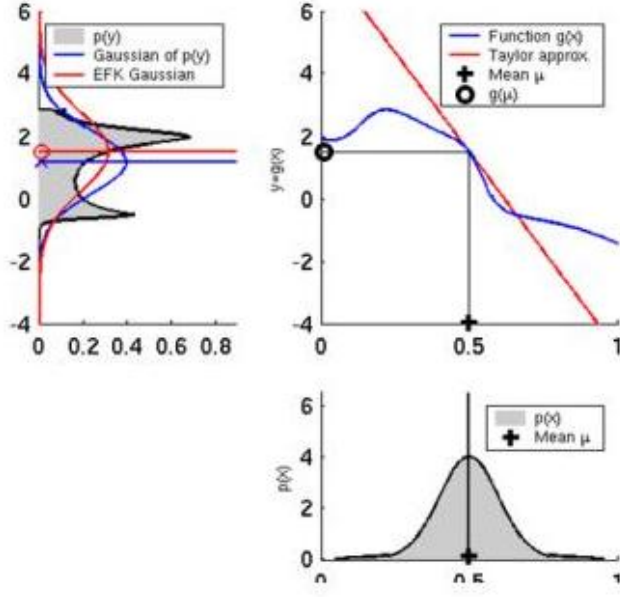


Fig. 15 Non-Linear system [1]

ever, at the cost of speed. The measurement updates for the vision and the IMU yields (V and I denotes Vision and IMU)

The innovation done by authors in for the vision part is,

$$K_{V,k} = P_k^- H_{V,k}^T (H_{V,k} P_k H_{V,k}^T + R_v)^{-1} \quad (32)$$

$$\vec{z}_k = \vec{z}_k + K_{V,k} (\vec{x}_{SLAM} - H_{V,k} \vec{z}_k) \quad (33)$$

$$P_k = (I - K_{V,k} H_{V,k}) P_k^- \quad (34)$$

The innovation for by authors in the IMU part is,

$$K_{I,k} = P_k^- H_{I,k}^T (H_{I,k} P_k H_{I,k}^T + R_I)^{-1} \quad (35)$$

$$\vec{z}_k = \vec{z}_k + K_{I,k} (\vec{x}_{IMU} - H_{I,k} \vec{z}_k) \quad (36)$$

$$P_k = (I - K_{I,k} H_{I,k}) P_k^- \quad (37)$$

The two matrices  $R_I$ ,  $R_V$  are the noise covariance matrices for the vision and IMU measurement inputs  $x_{SLAM}$ ,  $a_{IMU}$  which are resolved in the world frame W. The vector  $x_{SLAM}$  is the position without scale obtained from the vision algorithm (SLAM). The IMU measurement  $a_{IMU}$  needs special attention, because significant errors arise in the conversion from the raw IMU output.

$$\vec{a}_w = R_{wc} R_{ca} (\vec{a}_a - \vec{b}) - \vec{g}_w \quad (38)$$

## 8 Simultaneous Localization and Mapping

In this section, we cover SLAM, and the different concepts related to SLAM such as Visual Odometry. We then explore the two different approaches to SLAM mentioned in the paper. We briefly cover filter based SLAM called MonoSLAM [6] and a key-frame based SLAM called Parallel Tracking and Mapping (PTAM) [13]. A comprehensive analysis of filtering based SLAM is given in [22]. Simultaneous Localization and Mapping (SLAM) is a technique for estimating the motion of the robot and reconstructing the map/structure of the unknown environment. SLAM using only visual information only is specifically referred to as visual SLAM (vSLAM). The SLAM problem can be stated as follows:

*How can a body navigate in a previously unknown environment while constantly building and updating a map of its workspace using onboard sensors only?* [16]

Mathematically, this can be expressed as [21]

Input:

$$\text{Control input : } u_{1:T} = u_1, u_2, \dots, u_T \quad (39)$$

$$\text{Observations : } z_{1:T} = z_1, z_2, \dots, z_T \quad (40)$$

Output:

$$\text{Map of Environment : } m \quad (41)$$

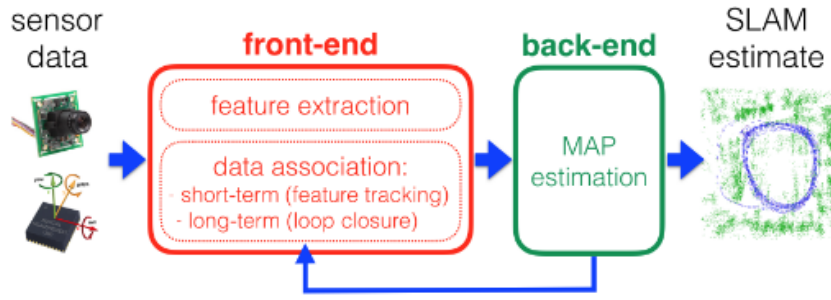
$$\text{Path of robot : } x_{0:T} = x_0, x_1, \dots, x_T \quad (42)$$

From the problem statement above, we notice that the robot has no a priori knowledge of the workspace or environment that it is in. This makes SLAM a very challenging problem in probabilistic robotics. This is also referred to as a chicken-egg problem in that we need to map the environment to get an accurate pose, but at the same time we also need an accurate pose to build a correct map. Therefore, it is an iterative process of estimating pose and building a map simultaneously. Accurate pose estimation is critical for many applications in computer vision, autonomous robotics and augmented reality. A survey of SLAM implementations is given in [3, 4, 23]. The basic modules for a vSLAM system are

- Initialization
- Tracking
- Mapping
- Relocalization
- Global Map optimization

According to [20], the relationship between vSLAM and VO can be represented as follows

$$vSLAM = VO + \text{global map optimization}$$

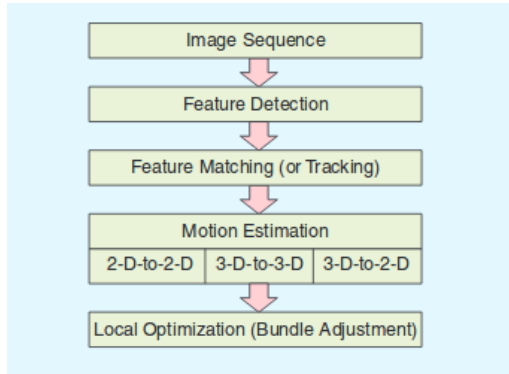


**Fig. 16** SLAM architecture [4]

From Fig. 16 we see that a typical SLAM system consists of a front-end and a back-end. The front-end is responsible for feature extraction while the back-end is responsible for MAP estimation. Both the methods (MonoSLAM & PTAM) are feature based implementations.

### 8.1 Visual Odometry

Odometry is the technique of estimating the position of a robot over time using sensors such as cameras, wheel encoders or any sensor measuring relative movement. Compared to SLAM which maintains a global consistent map, visual odometry (VO) maintains a local consistent map optimized over the last  $n$  frames. A generalized VO pipeline can be summarized as follows



**Fig. 17** Generalized Visual Odometry Pipeline [20]

Different algorithms exist due to the different type of cameras available such as Stereo, Mono, RGB-D [20,8]. As we use a monocular setup in our research paper, we would be going through the visual odometry algorithm from 2D to

2D correspondences for a monocular sensor. The algorithm is summarized as follows [20].

---

**Algorithm 2** Monocular Visual Odometry
 

---

- 1: Capture a new image frame  $I_k$
  - 2: Extract and match features between frames  $I_k$  and  $I_{k-1}$
  - 3: Compute the essential matrix  $E$  from the matched feature points
  - 4: Decompose  $E$  into  $R_k$  and  $t_k$  to form  $T_k$
  - 5: Choose the correct  $T_k$  matrix and scale  $t_k$  accordingly from scale factor
  - 6: Compute  $C_k = C_{k-1} * T_k$
  - 7: Goto step 1.
- 

Since we are using a monocular sensor, we always compare the image  $I_k$  captured at time instant  $k$  with the image  $I_{k-1}$  capture at time instant  $k - 1$ . In a stereo setup, we would compare images from the left and right part of the stereo. Once an image is captured, we compute the features for image  $I_k$  and  $I_{k-1}$ . In the feature detection step, the image is searched for key-points called features which are likely to be present in successive images. Features can be defined as image patterns that differ from its immediate neighborhood in terms of intensity, color and texture. There exist a variety of feature detectors in literature which vary in performance and properties. Some of the properties of feature detectors are rotation, scale and affine invariant, repeatability, localization accuracy, robustness and efficiency. Heitanen et al. present a comparison of feature detectors and descriptors for object class matching in [10]. Once we detect features, we then need to compute a descriptor such that features detected in two images can be compared with each other. The easiest approach would be to form a patch of pixels and compare using a sum of squared distances (SSD) or normalized cross correlation (NCC) metric. Such descriptors are not invariant to any of the above mentioned properties and perform poorly in practical applications. One of the popular descriptors called SIFT [15] uses gradient orientations as its descriptors. This forms a 128-element descriptor that is invariant to most of the above mentioned properties such as rotation, scale and illumination which makes it widely applicable in practical real-time applications. For extracting features from two images there are two paths with one being to detect and match features independently in two frames and the other being to track features in subsequent frames an example of which is a KLT Tracker [24]. There are various methods employed to match features accurately such a RANSAC which stands for Random Sample Consensus and is used to remove outliers among feature matches. Once we detect and match features in  $I_k$  and  $I_{k-1}$ , we calculate the essential matrix  $E$  which describes the geometric relationship between two images. We can use Nisters five point algorithm [18] or Longuet-Higgins eight point algorithm [14] to get the essential matrix  $E$ . Once we get the essential matrix  $E$ , we decompose it to extract the rotation and translation parts. Four different solutions are obtained of which the correct pair can be found out by triangulation. The solutions are given as

$$R = U(\pm W^T)V^T \quad (43)$$

$$t = U(\pm W^T)SV^T \quad (44)$$

where  $R$  is the rotation matrix and  $t$  is the translation vector.

## 8.2 MonoSLAM

MonoSLAM was developed by Davison et al. [5,6]. It uses EKF as an estimator. In EKFSLAM, the state  $\mathbf{x}$  space is represented by the robot state  $\hat{x}_\nu$  along with the position of the different landmarks  $\hat{y}_i$  that the robot observes. The state vector along with the covariance matrix can be represented as follows

$$\hat{x} = \begin{pmatrix} \hat{x}_\nu \\ \hat{y}_1 \\ \hat{y}_2 \\ \vdots \end{pmatrix} P = \begin{pmatrix} P_{xx} & P_{xy1} & P_{xy2} & \cdots \\ P_{y1x} & P_{y1y1} & P_{y1y2} & \cdots \\ P_{y2x} & P_{y2y1} & P_{y2y2} & \cdots \\ \vdots & \vdots & \vdots & \cdots \end{pmatrix} \quad (45)$$

The covariance matrix is used to represent the first order uncertainty. As the robot moves around, new features are added to the state vector. It is important to note that the robot initialization is done by using a known object. The prediction and correction steps are performed as mentioned in the EKF algorithm in section 7. One of the problems with this method is that the state vector becomes large when the number of land marks or key-points tracked increases. This increases the computation time.

## 8.3 Parallel Tracking and Mapping (PTAM)

PTAM is a SLAM system specifically designed to track a hand-held camera in a small AR workspace. For real-time operation, PTAM splits the tracking and mapping threads into two separate tasks running in parallel. The tracking thread robustly tracks hand-held motion and the mapping thread produces a 3D map of the features. Although developed specific to AR applications, the concept of running operations on multiple threads was incorporated into future SLAM algorithms.

The map consists of a collection of key-points/features that are located in a world frame  $W$ . They are represented in the homogeneous form. The map also contains key-frames which are instances of the frames taken from the hand-held camera. Each point is stored with a source key-frame. Unlike earlier methods which process each and every frame, this method only processes frames when there is sufficient information. Thus incremental mapping is replaced with a more computationally expensive batch method called bundle adjustment.

The tracking algorithm can be summarized as follows.

PTAM is initialized by a lateral offset movement of the hand-held camera. This can be assumed as a stereo image for starting the mapping. The initial

---

**Algorithm 3** PTAM - Tracking Algorithm

---

- 1: A new frame is captured and a prior pose estimate is generated from motion model.
  - 2: Map points are projected into the image based on prior pose.
  - 3: Small number of features (50) are searched in the image.
  - 4: The camera pose is updated from the feature matches.
  - 5: Large number of points is re-projected and searched in the image.
  - 6: Final pose estimate is computed from all matched found.
  - 7: Goto step 1.
- 

map is constructed using the 5 point algorithm [18] with arbitrary scale. As the camera moves, the key-frames increase from an initial two. Here, we only add key-frames if the tracking quality is good and a minimum of twenty frames has passed from the last keyframe. Bundle adjustemnt is iteratively performed to adjust the map based on a cost function. The mapping algorithm can then be summarized as.

---

**Algorithm 4** PTAM - Mapping Algorithm

---

**Require:** Stereo Initialization

- 1: **if** New Keyframe **then**
  - 2:     Update key-frame
  - 3:     Integrate key-frame
  - 4:     Add new features
  - 5: **else if** Locally Converged **then**
  - 6:     **if** Globally Converged **then**
  - 7:         Update data association
  - 8:     **else**
  - 9:         Global bundle adjust
  - 10:     **end if**
  - 11: **else**
  - 12:     Locally bundle adjust
  - 13: **end if**
  - 14: Sleep 5ms
  - 15: Goto step 1.
- 

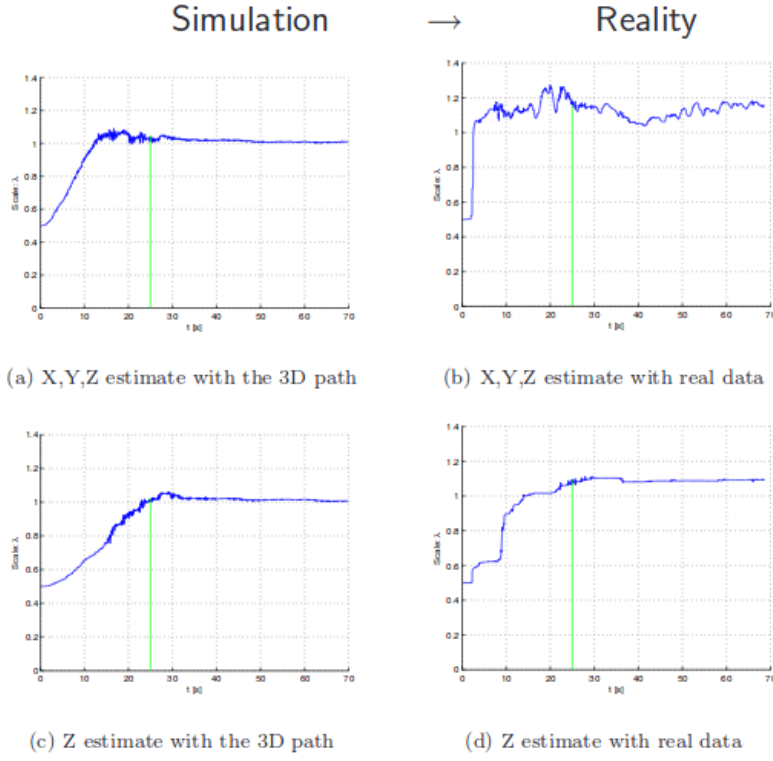
Finally when compared with Filter based SLAM, PTAM can handle thousands of features by splitting the tracking and mapping onto different threads on the CPU. We see that SLAM algorithms built now-a-days have used this philosophy to great success.

## 8.4 Results

### 8.4.1 Simulated and Real Data

In the paper, the Kalman filter was simulated offline with simulated data from the spline fitting section. The two inputs to the Kalman filter were the position from the vision sensor and acceleration from the IMU. Three different approaches were taken. The first method is same as Eq. 7, the second method





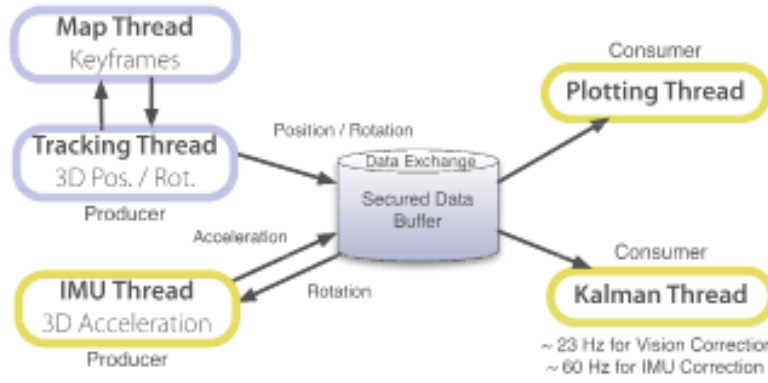
**Fig. 18** Offline simulation results with simulated and real data. Fewer directions (X,Y,Z) included in Kalman filter gives us a better estimate of  $\lambda$ . [19]

only uses the Z-axis which gives the states as  $[\vec{x}_z, \vec{v}_z, \vec{a}_z, \lambda]$  and the third uses only the X and Y-axis.

Figure 18 shows the scale estimation  $\lambda(t)$  for a simulated 3D path on left and actual path on right. The standard deviation for the acceleration noise for simulation was chosen same as real data ( $\sigma_{SLAM} = 0.01, \sigma_{IMU} = 0.2m/s^2$ ) with initial velocity and acceleration set to zero. Plots in Fig. 18a and Fig. 18c do not differ due to simulating the acceleration from the spline ideally in world frame. The estimate becomes more accurate when we use less number of directions (X,Y,Z). Wrong measurements in acceleration influence  $\lambda$  and make the scale estimation sensitive as seen in Fig. 18b. Hence only using a single axis (Z-axis) gives us the best result.

#### 8.4.2 Online Implementation

For the online implementation, the third setup was incorporated into the PTAM code [12]. As mentioned in section 8.3 PTAM employs two parallel threads called Tracker for tracking and MapMaker for mapping. Two additional threads for IMU and Kalman we added. The IMU thread provides the



**Fig. 19** Online implementation with PTAM work flow. Blue color boxes are original PTAM threads. Yellow boxes are added threads for scale estimation  $\lambda$ . [19]

acceleration measurements. The Kalman thread starts with  $\lambda$  calculated from integrating acceleration values, position from SLAM algorithm and acceleration and velocity set to 0. Values for covariance matrix  $\mathbf{Q}$  is time-varied which provides control over the sensitivity of the Kalman filter. For our report we have not incorporated the online implementation and have reported results from the paper.

## References

1. Byron Boots. Statistical techniques in robotics, 2015. URL: [https://www.cc.gatech.edu/~bboots3/STR-Spring2015/Lectures/Lecture12/Lecture12\\_notes.pdf](https://www.cc.gatech.edu/~bboots3/STR-Spring2015/Lectures/Lecture12/Lecture12_notes.pdf). Last visited on 11/18/2018.
2. J-Y Bouguet. Camera calibration toolbox for matlab. [http://www.vision.caltech.edu/bouguetj/calib\\_doc/index.html](http://www.vision.caltech.edu/bouguetj/calib_doc/index.html), 2004.
3. Guillaume Bresson, Zayed Alsayed, Li Yu, and Sébastien Glaser. Simultaneous localization and mapping: A survey of current trends in autonomous driving. *IEEE Transactions on Intelligent Vehicles*, 20:1–1, 2017.
4. Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif, Davide Scaramuzza, José Neira, Ian Reid, and John J Leonard. Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Transactions on Robotics*, 32(6):1309–1332, 2016.
5. Andrew J Davison. Real-time simultaneous localisation and mapping with a single camera. page 1403. IEEE, 2003.
6. Andrew J Davison, Ian D Reid, Nicholas D Molton, and Olivier Stasse. Monoslam: Real-time single camera slam. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (6):1052–1067, 2007.
7. Starlino Electronics. Imu guide. URL: [http://www.starlino.com/imu\\_guide.html](http://www.starlino.com/imu_guide.html). Last visited on 11/18/2018.
8. Friedrich Fraundorfer and Davide Scaramuzza. Visual odometry: Part ii: Matching, robustness, optimization, and applications. *IEEE Robotics & Automation Magazine*, 19(2):78–90, 2012.
9. Janne Heikkila and Olli Silven. A four-step camera calibration procedure with implicit image correction. In *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*, pages 1106–1112. IEEE, 1997.
10. Antti Hietanen, Jukka Lankinen, Joni-Kristian Kämäräinen, Anders Glent Buch, and Norbert Krüger. A comparison of feature detectors and descriptors for object class matching. *Neurocomputing*, 184:3–12, 2016.
11. S-H Jung and Camillo J Taylor. Camera trajectory estimation using inertial sensor measurements and structure from motion results. In *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, volume 2, pages II–II. IEEE, 2001.
12. Georg Klein and David Murray. Source code of ptam (parallel tracking and mapping). URL: <http://www.robots.ox.ac.uk/~gk/PTAM/>. Last visited on 11/18/2018.
13. Georg Klein and David Murray. Parallel tracking and mapping for small ar workspaces. In *Mixed and Augmented Reality, 2007. ISMAR 2007. 6th IEEE and ACM International Symposium on*, pages 225–234. IEEE, 2007.
14. H Christopher Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293(5828):133, 1981.
15. David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.
16. Roland Siegwart Margarita Chli, Martin Rufli. Lecture notes in autonomous mobile robots, 2017. URL: [https://www.ethz.ch/content/dam/ethz/special-interest/mavt/robotics-n-intelligent-systems/asl-dam/documents/lectures/autonomous\\_mobile\\_robots/spring-2017/2017%20-%20SLAM%20I%20-%20after%20lecture.pdf](https://www.ethz.ch/content/dam/ethz/special-interest/mavt/robotics-n-intelligent-systems/asl-dam/documents/lectures/autonomous_mobile_robots/spring-2017/2017%20-%20SLAM%20I%20-%20after%20lecture.pdf). Last visited on 11/18/2018.
17. MATLAB. What is camera calibration. URL: <https://www.mathworks.com/help/vision/ug/camera-calibration.html>. Last visited on 11/18/2018.
18. David Nistér. An efficient solution to the five-point relative pose problem. *IEEE transactions on pattern analysis and machine intelligence*, 26(6):756–770, 2004.
19. Gabriel Nützi, Stephan Weiss, Davide Scaramuzza, and Roland Siegwart. Fusion of imu and vision for absolute scale estimation in monocular slam. *Journal of intelligent & robotic systems*, 61(1-4):287–299, 2011.
20. Davide Scaramuzza and Friedrich Fraundorfer. Visual odometry [tutorial]. *IEEE robotics & automation magazine*, 18(4):80–92, 2011.

21. Cyrill Stachniss. Lecture notes in robot mapping, 2013. URL: <http://ais.informatik.uni-freiburg.de/teaching/ws12/mapping/pdf/slam04-ekf-slam.pdf>. Last visited on 11/18/2018.
22. Hauke Strasdat, JMM Montiel, and Andrew J Davison. Visual slam: Why filter?
23. Takafumi Taketomi, Hideaki Uchiyama, and Sei Ikeda. Visual slam algorithms: a survey from 2010 to 2016. *IPSJ Transactions on Computer Vision and Applications*, 9(1):16, 2017.
24. Carlo Tomasi and Takeo Kanade. Detection and tracking of point features. 1991.
25. Rein van den Boomgaard. Lecture notes image processing and computer vision. URL: <https://staff.fnwi.uva.nl/r.vandenboomgaard/IPCV20162017/LectureNotes/IP/Images/ImageFormation.html>. Last visited on 11/18/2018.
26. Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22, 2000.