# Visual-Inertial SLAM: A Project Review

Hao-Yuan Tang
*Department of Electrical Engineering*
*University of California, San Diego*
h6tang@ucsd.edu

## I. INTRODUCTION

In this project, Extended Kalman filter algorithm (EKF) was used to tackle simultaneous localization and mapping (SLAM)problem in visual-inertial sensor configuration. SLAM can be separated into two parts: localization given a surrounding map of the environment for estimating robot trajectory by odometry reading, and mapping surrounding with robot state simultaneously by two stereo cameras.

EKF, which is a form of Bayesian filter, can effectively model a robot's predicted locations in the world as Gaussian distribution based on past input controls and current observations. The EKF was used when we try to model non-linear observation and motion model, which provided more widely applicable scenario than particle filter.

In problem formulation session, we will discuss on essential parts of EKF SLAM in mathematical sense. In technical approach session, detailed implementation of EKF SLAM was given, such as transformation of coordinate systems, car model localization using odometry and surrounding map using stereo cameras.

## II. PROBLEM FORMULATION

### A. SLAM

The visual SLAM problem can be separated into localization by odometry and mapping by stereo cameras parts. Markov assumption was given for model simplification, i.e., the robots state $x_{t+1}$ only depends on the previous input $u_t$ and state $x_t$, and the map state $m_{t+1}$ only depends on the previous map state $m_t$. Given robot's pose $x_t$ and control input $u_t$ at discrete time steps t, the pose in next time step can be described as probabilistic function:

$$x_{t+1} = f(x_t, u_t, w_t) \sim p_f(\cdot \,|x_t, u_t) \qquad (1)$$

where $w_t$ is motion noise and this model was called motion model.

For mapping, the robot can sense environment via camera observations at each time step. In a similar way, observation model, which described camera observations within an environment, can be defined in probabilistic:

$$z_{t+1} = h(z_t, v_t) \sim p_f(\cdot \,|x_t, m) \qquad (2)$$

where $v_t$ is motion noise and this model was called observation model.

With motion model and observation model, SLAM is the problem to determine the environment and robot poses from observations and control inputs at each time step t. The objective to compute can be written in probabilistic form:

$$p(x_t, m | z_{0:t}, u_{0:t-1}) \qquad (3)$$

### B. Bayes Filtering

Bayes filtering is a probabilistic inference technique for estimating the state $x_t$ of dynamic systems combining observations and control inputs with Markov assumptions and Bayes rule. The Bayes filter relies on two steps to keep track of $p_{t|t}(x_t)$ and $p_{t|t+1}(x_{t+1})$.

For prediction step, given a prior density of robot state $p_{t|t}(x_t)$ and the control input $u_t$, we use the motion model Eq. (1) to compute the predicted density $p_{t|t+1}(x_{t+1})$ as the following equation:

$$p_{t+1|t}(x_t) = \int p_f(x|s, u_t) p_{t|t}(s)\, ds \qquad (4)$$

For update step, given the predicted density $p_{t+1|t}(x_{t+1})$ and the measurement $z_{t+1}$, we use the observation model Eq. (2) to incorporate the measurement information and obtain the posterior $p_{t+1|t+1}(x_{t+1})$ as following:

$$p_{t+1|t+1}(x) = \frac{p_h(z_{t+1}|x)p_{t+1|t}(x)}{\int p_h(z_{t+1}|s)p_{t+1|t}(s)ds} \qquad (5)$$

### C. Mapping of Landmarks

The landmark-based mapping aims to generate a map of the surrounding environment from sensor observations z with the assumption that the poses of the robot x are known. Also, uncertainty of system was modeled by adding noise into observation model.

The environment is represented by *M* static landmarks, and each of them is characterized by its location in the space denoted as $m_i$. These landmarks are considered as points in the 3D space and can be specified by three numerical values where $m_i \in \mathbb{R}^3$ and $m \in \mathbb{R}^{3 \times M}$.

The robot can sense the landmarks at each time step t, with observation $z_t$, and since the robot can sense multiple landmarks at a single time step, $z_t$ is a general notation for composed observation from multiple landmarks. The goal of the mapping problem is then to estimate the locations of landmarks based on the pose of robot $x_t$ and the observation $z_t$. Therefore, we can define the observation model:

$$p(z_t|x_t, m, n_t) \qquad (6)$$

where $n_t$ is the index map with $|n_t| = N_t$ and $N_t$ is the number of observed landmarks at time *t*.

### D. Sensors Configuration

Our proposed solution aims to solve the SLAM problem with observations from an IMU and a stereo camera installed in a car. The IMU observations contain the linear velocity $v_t \in \mathbb{R}^3$ and angular velocity $\omega_t \in \mathbb{R}^3$ in IMU frame.

To collect landmarks, two calibrated stereo camera were install on the top of car and data are pre-computed to extract

the visual features and find the correspondence between left and right camera frames across time steps. Fig.1 shows the visual features extracted in this configuration. The visual features at time t is denoted as $z_t \in \mathbb{R}^{4 \times M}$, where the $i_{th}$ column contains the pixel coordinates of landmark $i$ in the left and right camera images. One thing to note is which landmarks showed up in specific time step, in other words, there might be some landmarks that are not observable at different time step.

In our configuration, we assume that the transformation from the IMU to the camera optical frame belongs to special Euclidean group SE(3) (extrinsic parameters) and the stereo camera calibration matrix M (intrinsic parameters) are known. The camera calibration matrix is defined as the following equation:

$$M = \begin{bmatrix} fs_u & 0 & c_u & 0 \\ 0 & fs_v & c_v & 0 \\ fs_u & 0 & c_u & -fs_u b \\ 0 & fs_v & c_v & 0 \end{bmatrix} \tag{7}$$

where f is the focal length, $s_u$, $s_v$ are pixel scaling, $c_u$ and $c_v$ are the principle points, and b is the stereo baseline.

## III. TECHNICAL APPROACH

### A. Extended Kalman Filter

The Extended Kalman Filter (EKF) is a nonlinear version of the Kalman Filter, which linearizes about an estimate of the current mean and covariance with a moment matching approach. The nonlinear Kalman Filter is a Bayes filter with the following assumptions:

1. The prior pdf is Gaussian distribution
2. The state $x_{t+1}$ is affected by Gaussian noise, that is:
   $$x_{t+1} = f(x_t, u_t, w_t), \ w_t \sim \mathcal{N}(0, W) \tag{8}$$
3. The observation $z_t$ is affected by Gaussian noise, that is:
   $$z_{t+1} = h(z_t, v_t), v_t \sim \mathcal{N}(0, V) \tag{9}$$
4. The process noise $w_t$ and measurement noise $v_t$ are independent of each other
5. The posterior pdf is forced to be Gaussian via approximation

The challenge of the nonlinear Kalman Filter is that the predicted and updated pdfs are not Gaussian and cannot be evaluated in closed form. Using moment matching, we can force the predicted and updated pdfs to be Gaussian by evaluating their first and second moments and modeling them with Gaussians with the same moments.

The Extended Kalman Filter uses a first-order Taylor series to approximate the integrals required to implement the nonlinear Kalman Filter. Thus, the approximation of the motion model would be as follow:

$$f(x_t, u_t, w_t) \approx f(\mu_{t|t}, u_t, 0) + \left[\frac{df}{dx}(\mu_{t|t}, u_t, 0)\right](x_t - \mu_{t|t})$$
$$+ \left[\frac{df}{dw}(\mu_{t|t}, u_t, 0)\right](w_t - 0)$$
$$\approx f(\mu_{t|t}, u_t, 0) + F_t(x_t - \mu_{t|t}) + Q_t w_t \tag{10}$$

where $F_t$ and $Q_t$ are Jacobians w.r.t. $x$ and w. Since $x_t$ and $w_t$ are independent of each other, and because the approximation given in Eq. (12) is linear, we know the distribution is Gaussian, i.e., $x_{t+1} \sim \mathcal{N}(\mu_{t+1|t}, \Sigma_{t+1|t})$:

$$\mu_{t+1|t} = f(\mu_{t|t}, u_t, 0) \tag{11}$$

$$\Sigma_{t+1|t} = F_t \Sigma_{t|t} F_t^T + Q_t W Q_t^T \tag{12}$$

We begin our localization problem in any given timestep t using our IMU readings to get some idea of our new position in space, after a control input $u_t$ has been applied. Our motion model, in this system, is a discretized version of nominal and perturbed kinematics. This allows us, with some time discretization, to employ the exponential map:

$$\mu_{t+1|t} = \exp(-\tau \hat{u}_t) \mu_{t|t} \tag{13}$$

where $\hat{u}_t \in \mathbb{R}^{4 \times 4}$ represents the hat map of the velocity $u_t = [\lambda_t, \omega_t]^T \in \mathbb{R}^6$ composed of the linear velocity vector $\lambda_t \in \mathbb{R}^3$ and rotational velocity vector $\omega_t \in \mathbb{R}^3$:

$$\hat{u}_t = \begin{bmatrix} \hat{\omega}_t & \hat{\lambda}_t \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{4 \times 4} \tag{14}$$

Here, for any vector $s \in \mathbb{R}^3$, $\hat{s}$ is the corresponding skew symmetric matrix. Similarly, the approximation of the observation model would be described as follow:

$$h(x_{t+1}, v_{t+1}) \approx h(\mu_{t+1|t}, 0) + \left[\frac{dh}{dx}(\mu_{t|t}, u_t, 0)\right](x_{t+1} - \mu_{t+1|t})$$
$$+ \left[\frac{df}{dv}(\mu_{t+1|t}, 0)\right](v_{t+1} - 0)$$
$$\approx h(\mu_{t+1|t}, 0) + H_{t+1}(x_{t+1} - \mu_{t+1|t}) + R_{t+1} v_{t+1} \tag{15}$$

where $H_{t+1}$ and $R_{t+1}$ is the Jacobian w.r.t. $x$ and $v$. Based on the equations above, the models of Extended Kalman Filter will be the following equations:

1. Prior:

$$x_t | z_{0:t}, u_{0:t-1} \sim \mathcal{N}(\mu_{t|t}, \Sigma_{t|t}) \tag{16}$$

2. Motion model:

$$x_{t+1} = f(x_t, u_t, w_t), w_t \sim \mathcal{N}(0, W) \tag{17}$$

$$F_t := \frac{df}{dx}(\mu_{t|t}, u_t, 0)$$

$$Q_t := \frac{df}{dw}(\mu_{t|t}, u_t, 0)$$

3. Observation Model:

$$z_t = h(z_t, v_t), v_t \sim \mathcal{N}(0, V) \tag{18}$$

$$H_t := \frac{dh}{dx}(\mu_{t|t-1}, 0)$$

$$R_t := \frac{dh}{dv}(\mu_{t|t-1}, 0)$$

4. Prediction:

$$\mu_{t+1|t} = f(\mu_{t|t}, \boldsymbol{u}_t, 0) \quad (19)$$

$$\Sigma_{t+1|t} = F_t \Sigma_{t|t} F_t^T + Q_t W Q_t^T \quad (20)$$

5. Update:

$$\mu_{t+1|t+1} = \mu_{t+1|t} + \mathbf{K}_{t+1|t}(\mathbf{z}_{t+1} - h(\mu_{t+1|t}, 0)) $$

$$\Sigma_{t+1|t+1} = (\mathbf{I} + \mathbf{K}_{t+1|t}\mathbf{H}_{t+1})\Sigma_{t+1|t} \quad (21)$$

6. Kalman gain:

$$\mathbf{K}_{t+1|t} := \Sigma_{t+1|t}\mathbf{H}_{t+1}^T(\mathbf{H}_{t+1}\Sigma_t\mathbf{H}_{t+1}^T + \mathbf{R_V}\mathbf{R}_{t+1}^T)^{-1} \quad (22)$$

*B. EKF-based Visual Mapping*

For the landmark-based visual mapping problem, we assume that the inverse IMU pose is known. Moreover, as described above, we assume that the landmarks are static and the data association $\pi_t = \{1, ..., M\} \to \{1, ..., N_t\}$ stipulating which landmarks were observed at each time t is pre-computed by an external algorithm. The observation model can be written with the Gaussian prior and observation noise:

Prior:

$$m|z_{0:t} \sim \mathcal{N}(\mu_t, \Sigma_t) \quad (23)$$

$$\mu_t \in \mathbb{R}^{3M}$$

$$\Sigma_t \in \mathbb{R}^{3M \times 3M}$$

Observation model:

$$z_{t,i} = h(U_t, m_j) + v_{t,i}$$

$$= M\pi(_oT_iU_t, \underline{m}_j) + v_{t,i} \quad (24)$$

$$v_{t,i} \sim \mathcal{N}(0, V)$$

where $\mu_t \in \mathbb{R}^{3M}$ is the expectation of locations of all landmarks, $\Sigma_t \in \mathbb{R}^{3M \times 3M}$ is the covariance of the estimate, and $M$ is the calibration matrix in Eq. 7. And $\pi$ refers to the projection function:

$$\pi(\mathbf{q}) := \frac{1}{q_3}\mathbf{q} \in \mathbb{R}^4 \quad (25)$$

and whose derivative is given by:

$$\frac{d\pi}{d\mathbf{q}} = \frac{1}{q_3}\begin{bmatrix} 1 & 0 & -\frac{q_1}{q_3} & 0 \\ 0 & 1 & -\frac{q_2}{q_3} & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -\frac{q_4}{q_3} & 1 \end{bmatrix} \quad (23)$$

We can stack all the observation $z_{t,i}$ as a single vector $z_t \in \mathbb{R}^{4N_t}$ and rewrite the observation model equation:

$$z_t = M\pi(_oT_iU_t\underline{m}) + v_t \quad (27)$$

$$v_t \sim \mathcal{N}(0, I \otimes V) \quad (28)$$

where $\otimes$ is the Kronecker product.
For Extended Kalman Filter, we need to derive the differentiation of the observation model with respect to m evaluated with $\mu_t$. Consider a small perturbation $\delta\mu_{t,j}$ for the location of landmark $j$:

$$m_j = \mu_{t,j} + \delta\mu_{t,j} \quad (29)$$

The first-order Taylor series approximation to observation time t can be then written as:

$$\begin{aligned} z_{t,i} &= M\pi(_oT_IU_t(\mu_{t,j} + \delta\mu_{t,j})) + v_{t,i} \\ &= M\pi(_oT_IU_t(\underline{\mu}_{t,j} + P^T\delta\mu_{t,j})) + v_{t,i} \\ &\approx M\pi(_oT_IU_t\underline{\mu}_{t,j}) + \\ &\quad M\frac{d\pi}{d\mathbf{q}}(_oT_IU_t\underline{\mu}_{t,j})_oT_IU_tP^T\delta\mu_{t,j} + v_{t,i} \end{aligned} \quad (30)$$

where $P = [\mathbf{I}, \mathbf{0}] \in \mathbb{R}^{3 \times 4}$ is the projection matrix.
The stereo camera Jacobian is given by:

$$H_{t,i,j} = M\frac{d\pi}{d\mathbf{q}}(_oT_{iw}T_i^{-1}\underline{\mu}_{t,j})_oT_{iw}T_i^{-1}P^T \quad (31)$$

and the EKF update:

$$K_t = \Sigma_t H_t^T(H_t\Sigma_t H_t^T + I \otimes V)^{-1} \quad (32)$$

$$\mu_{t+1} = \mu_t + K_t(z_t - \tilde{z}_t) \quad (33)$$

$$\Sigma_{t+1} = (\mathbf{I} - K_tH_t)\Sigma_t \quad (34)$$

Since we assume all the landmarks are static, there is no need to perform prediction steps in EKF-based visual mapping. For each landmark, we initialize its position at the first timestep t we observe it. The following equation with the observation $z_{t,i}$ is used to compute the landmark initialization:

$$z_{t,i} = M\pi(_oT_iU_t, \underline{m}_j) \quad (35)$$

*C. EKF-based Visual-Inertial Odometry*

The localization problem aims to estimate the inverse IMU pose of the robot given the IMU measurements $u_t = [v_t^T, \omega_t^T]$, the visual feature observation $z_{0:t}$, and the landmark coordinates $\underline{m}$ in the world frame. With the same assumption in the visual mapping problem, data association between observations and landmarks is computed by an external algorithm. The motion model can be written with the Gaussian prior and process noise:

Prior:

$$U_t | z_{0:t}, u_{0:t-1} \sim \mathcal{N}(\mu_{t|t}, \Sigma_{t|t}) \quad (36)$$

$$\mu_{t|t} \in SE(3)$$

$$\Sigma_{t|t} \in \mathbb{R}^{6 \times 6}$$

Motion model:

$$U_{t+1} = \exp(-\tau((u_t + w_t))^\wedge) U_t \quad (37)$$

$$w_{t,i} \sim \mathcal{N}(0, W)$$

where $\tau$ is time difference in seconds between two contiguous frames. In order to separate the effect of the noise $w_t$ from the motion of the deterministic part of $wTi, t = U_t^{-1}$, we take advantages of the discrete-time perturbation idea in and rewrite the motion model in terms of nominal kinematics and zero-mean perturbation kinematics:

$$\boldsymbol{\mu}_{t+1|t} = \exp(-\tau \mathbf{u}_t^\wedge) \boldsymbol{\mu}_{t|t}$$
$$\boldsymbol{\Sigma}_{t+1|t} = \mathbb{E}[\delta \boldsymbol{\mu}_{t+1|t} \delta \boldsymbol{\mu}_{t+1|t}^T]$$
$$= \exp(-\tau \mathbf{u}_t^\wedge) \boldsymbol{\Sigma}_{t|t} \exp(-\tau \mathbf{u}_t^\wedge)^T + \mathbf{W} \quad (38)$$

The observation model is the same as the definition in the visual mapping problem. In order to derive the update steps for Extended Kalman Filter, we need the observation model Jacobian $H_{t+1|t} \in \mathbb{R}^{4N_t \times 6}$ with respect to the inverse IMU pose $U_t$ evaluated at $\mu_{t|t}$. Similar to the update steps in EKF-based visual mapping, the update steps for visual-inertial odometry can be composed of the following equations:
Prior:

$$U_{t+1} | z_{0:t}, u_{0:t} \sim \mathcal{N}(\mu_{t+1|t}, \Sigma_{t+1|t}) \quad (39)$$

$$\mu_{t|t} \in SE(3)$$

$$\Sigma_{t|t} \in \mathbb{R}^{6 \times 6}$$

Predicted observations:

$$\tilde{z}_{t+1,i} = M\pi({}_o T_i \mu_{t+1|t}, \underline{m}_j) \quad (40)$$

Observation matrix:

$$\mathbf{H}_{i,t+1|t} = \mathbf{M} \frac{d\pi}{d\mathbf{q}} ({}_O \mathbf{T}_I \mu_{t+1|t} \mathbf{m}_j)_O \mathbf{T}_I (\mu_{t+1|t} \mathbf{m}_j)^\odot$$

$$\mathbf{H}_{t+1|t} = \begin{bmatrix} \mathbf{H}_{1,t+1|t} \\ \mathbf{H}_{2,t+1|t} \\ \vdots \\ \mathbf{H}_{N_{t+1},t+1|t} \end{bmatrix} \in \mathbb{R}^{4N_t \times 6} \quad (41)$$

EKF Update:

$$K_{t+1|t} = \Sigma_{t+1|t} H_{t+1|t}^T (H_{t+1|t} \Sigma_{t+1|t} H_{t+1|t}^T + I \otimes V)^{-1}$$

$$\mu_{t+1|t+1} = \exp(K_{t+1|t}(\widehat{\boldsymbol{z_{t+1}}} - \tilde{z}_{t+1})) \mu_{t+1|t} \quad (42)$$

$$\Sigma_{t+1|t+1} = (I - K_{t+1|t} H_{t+1|t}) \Sigma_{t+1|t} \quad (43)$$

### D. EKF SLAM

To achieve the goal of estimating position of landmarks and the pose of the robot simultaneously, the proposed idea is to merge the predict and update steps of Extended Kalman Filter based visual mapping and visual-inertial odometry. First of all, the joint estimated state and covariance under the Gaussian assumption are defined as follows:

$$\mu = \begin{bmatrix} \mu_m \\ \mu_p \end{bmatrix} \in \mathbb{R}^{6 \times 6} \quad (44)$$

$$\Sigma_{t|t} \in \mathbb{R}^{(6+3M) \times (6+3M)}$$

where $\mu_m$ is the estimated landmark position in Eq. 23. and $\mu_p$ is the estimated six degrees of freedom of inverse IMU pose in Eq. 36.

The Extended Kalman Filter predict step on the joint estimated state and covariance is derived from the predict step of visual-inertial odometry only because all the landmarks are assumed static. The equations given the IMU measurement $u_t$ are listed as follows:

$$\mu_{t+1|t} = \begin{bmatrix} \mu_{m,t+1|t} \\ \mu_{p,t+1|t} \end{bmatrix} = \begin{bmatrix} \mu_{m,t+1|t} \\ \exp(-\tau \hat{u}_t) \mu_{p,t|t} \end{bmatrix} \quad (45)$$

$$\Sigma_{t|t} = F_t \Sigma_{t|t} F_t^T + W$$

$$F_t = \begin{bmatrix} I & 0 \\ 0 & \exp(-\tau \hat{u}_t) \end{bmatrix}$$

$$W = \begin{bmatrix} 0 & 0 \\ 0 & W_p \end{bmatrix}$$

The update step is formed by combining the update step of both visual mapping and visual-inertial odometry. The update step equations are listed as follows:
Predicted observations:

$$\tilde{z}_{t,i} = M\pi({}_o T_i U_t, \underline{\mu}_{m,t,j}) \in \mathbb{R}^4 \quad (46)$$

Observation matrix:

$$H_{t+1|t} = [H_{m,t+1|t} \quad H_{p,t+1|t}] \in \mathbb{R}^{4N_t \times (3M+6)} \quad (47)$$

EKF update:

$$K_{t+1|t} = \Sigma_{t+1|t} H_{t+1|t}^T (H_{t+1|t} \Sigma_{t+1|t} H_{t+1|t}^T + I \otimes V)^{-1} \quad (48)$$

$$\mu_{t+1|t} = \begin{bmatrix} \mu_{m,t+1|t+1} \\ \mu_{p,t+1|t+1} \end{bmatrix} = \begin{bmatrix} \mu_{m,t+1|t} + K_{t+1|t}(z_t - \tilde{z}_t) \\ \exp(K_{t+1|t}(\widehat{z_{t+1}} - \tilde{z}_{t+1})) \mu_{p,t+1|t} \end{bmatrix} \quad (49)$$

## IV. RESULTS

The proposed EKF based SLAM algorithm is tested with 2 data set, and all of them are collected in real driving scenarios. The hyper-parameters used in our algorithm are same for all the data set except for the number of landmarks. The results are shown in Fig.2., which show the estimated trajectory and the 2D position of the visual features. The red line in the figure is the estimate robot trajectory, and the green dots are the position of landmarks. The estimate position of the visual features also distributes in a reasonable way.

On the other hand, the IMU-based landmark mapping does not provide excellent results. The reason is that the trajectory is not corrected based on the observations of visual features. The quality of the estimate position of the visual features also relies on the correct pose of the robot. Thus, the landmark estimates are also not accurate.

Also, one important observance is that noise covariances W are quite small, and it would lead poor results if such hyperparameters getting to large. In comparison, another noise V was quite robust to change, with maps changed only slightly between 1000 to 4000.

One failure in this project review is the EKF SLAM part. The problem researcher of this paper faced is singular matrix when computing matrix inversion in observation matrix. Since the whole part of $H_{t+1|t}\Sigma_{t+1|t}H_{t+1|t}^T$ in Eq. 48. Is positive semi-definite, this mistake should not happen. Another reason related to this is inappropriate variance value, which too large may cause instability in the algorithm. The expected result of visual-inertial SLAM should be similar to the one in Fig. 2. with little bias.

### REFERENCES

[1] A. Filipescu, V. Minzu, B. Dumitrascu, A. Filipescu and E. Minca, "Trajectory-tracking and discrete-time sliding-mode control of wheeled mobile robots," 2011 IEEE International Conference on Information and Automation, 2011, pp. 27-32, doi: 10.1109/ICINFA.2011.5948958.

[2] https://docs.opencv.org/3.4/d9/dba/classcv_1_1StereoBM.html

[3] UCSD ECE276A: Sensing & Estimation in Robotics (Winter 2022) https://natanaso.github.io/ece276a

| Parameter | Description | Value |
|---|---|---|
| $\Sigma_m$ | Prior landmark estimate covariance | $0.001 * I \in \mathbb{R}^{3\times3}$ |
| $\Sigma_p$ | Prior pose estimate covariance | $0.001 * I \in \mathbb{R}^{6\times6}$ |
| V | Observation | $1000 * I \in \mathbb{R}^{4\times4}$ |
| $W_p$ | Process noise covariance | $0.001 * I \in \mathbb{R}^{6\times6}$ |

Table 1. Hyperparameters for EKF visual inertial SLAM
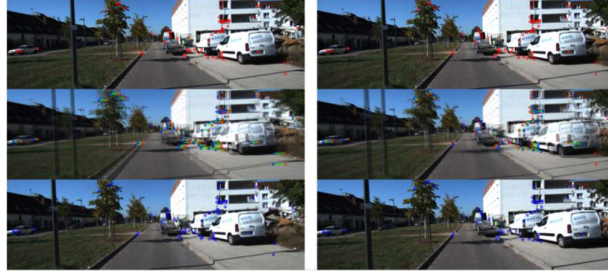


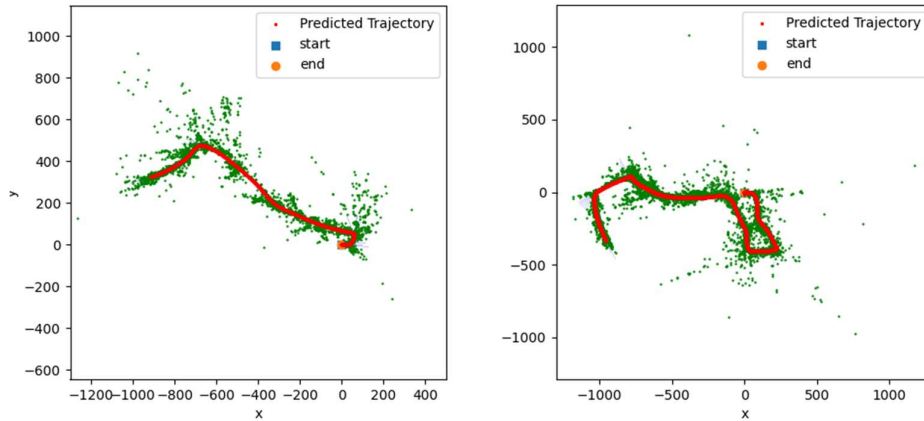Fig.1. Visual features matched across the left-right camera frames (left) and across time (right).



Fig. 2. Results of Experiments (a) 10.npz (b) 03.npz