

# Visual Inertial Odometry based on Image Lines

Duarte Correia, José Gaspar

Instituto Superior Técnico / UTL, Lisbon, Portugal

duartecorreiascs@hotmail.com, jag@isr.tecnico.ulisboa.pt

## ABSTRACT

This work proposes to give the needed building blocks to achieve a Simultaneous Localization and Mapping (SLAM) solution using Visual Inertial Odometry (VIO) techniques with Lines as features. It was motivated by new advances in machine learning that facilitated the detection of lines in images that could provide a richer mapping visualization tool than classical point features. The main obstacle encountered was the matching of lines through multiple frames to be initialized as 3D structures. An Iterative Closest Point (ICP) variation was introduced to remedy this problem, where the main focus of the algorithm was to match the features instead of finding the transformation between different views. They are introduced and analyzed two methods for line initialization from 2D images, one that utilizes Plucker coordinates to represent the lines and another that minimizes 3D the distance of a general point to back-projected rays. A Square Root Unscented Kalman Filter (SR-UKF) is proposed based on Lie Groups as a base for the filtering process, which can integrate measurements from an Inertial Measurement Unit (IMU) and visual information. As a proof of concept, the proposed approach provided a sound localization of the drone through a complex trajectory and provided line initializations that could potentially maintain the filter alive. It was shown that lines provide a rich descriptor of an environment compared to a point cloud and that the line algorithms proposed could give a first approximation of the map without the need for filtering processes in cases where the localization is provided.

## I. INTRODUCTION

Visual Odometry systems have been used successfully to provide robot's pose estimations for over 16 years, with inertial sensors being added in recent years. Most of these systems utilize 3D points as features that are then propagated across multiple types of Kalman filters, usually an EKF or UKF, but since visual odometry should be run online, the system cannot create detailed maps of the environment traversed. A solution to this problem comes in utilizing lines as features since even though they can be described from six parameters, they provide a rich map of the environment, even with few lines. Since the main cost of updating a state through a filter is usually bound to the dimension of the observation, lines could reduce the execution time of filters, given that they can provide more visual information than a single point, and as such, provide a richer estimation of the robot's pose.

One of the main obstacles in using lines as features has been the difficulty in Line detection using classical methods, but with the appearance of neural networks trained explicitly for

this purpose, this process has become faster, more reliable and consistent. As such, the perfect opportunity arises to explore methods that take advantage of these new detectors.

Another weakness of using lines is that usually, they do not provide good descriptors that make it challenging to match features across multiple images so that mapping the environment is possible. With this problem in mind, line matching methods should be researched.

### A. Related Work

One important precedent of SLAM is the Structure from Motion (SfM) approach, which has been studied and used with impressive results for the past three decades [7]. Monocular Simultaneous Localization and Mapping (MonoSLAM [6]) was the first approach in which a real-time filter achieved good results, which leverages an Extended Kalman filter (EKF) [10] to propagate the state. New parameterizations were proposed to worth features in the infinite in Inverse Depth Parameterization [13]. The Unscented Kalman Filter [25] was introduced as an alternative to the EKF, leveraging the Unscented transform [9] to better propagate the state and Covariance without the need for the linearization of the system. Since the UKF increased the computational complexity to the filter, a Square Root-UKF (SR-UKF [24]) proposes propagating the squared root of the Covariance, reducing the complexity closer to the one of the EKF. Novel approaches like [3] take advantage of Lie Groups and algebra which present a natural state representation and permit Inertial Measurement Units (IMUs) for state estimation.

SLAM variations that use Lines as features are first presented in [19], where a UKF is used to propagate 6-dimension Plücker Coordinates directly and proposes a transformation of the 2D lines to a subset of  $\mathbb{R}^3$  paired with a representative patch of the Line for matching purposes. Another approach using an EKF filter is shown in [22], where it is given a method to initialize lines using plane intersections until their endpoints can represent them; until then, the lines are propagated using an exterior EKF filter. In [23], different approaches of landmark parameterizations that can be used in EKF-SLAM with points and lines are shown. However, these line approaches do not use inertial measurements to update the state. As such, the mapping and localization have an inherent scale to the actual map associated with them.

### B. Problem Formulation

This thesis aims to provide the necessary building blocks to implement a Simultaneous Localization and Mapping (SLAM) solution using Visual Inertial Odometry with lines as features.

It is proposed to take advantage of a line detection neural network to provide feature detection. For matching, it is proposed a variation the Iterative Closest Point (ICP) where the objective is not to obtain the transformation between the points, but the line matches between two sets of lines.

A method that performs line initializations from the obtained matches and benefits from using Plucker line coordinates is analyzed. It is also introduced an initialization technique based on the 3D distance to back-projected rays.

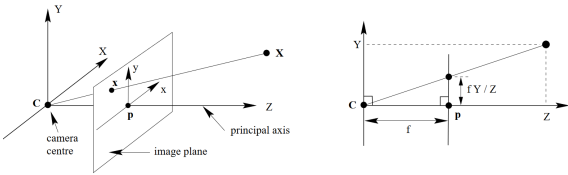
At last, The filtering strategy takes advantage of the lines obtained and inertial information given by an Inertial Measurement Unit (IMU) to estimate a drone's pose and a map of the environment using a Square Root Unscented Kalman Filter based on Lie Groups.

## II. BACKGROUND THEORY

### A. Pinhole Camera Model

The pinhole camera model describes the mathematical relation between the 3D coordinates of a point in space and its 2D projection on an image. It only serves as an excellent first-order approximation of the mapping from a 3D scene to a 2D image since it assumes no lenses or radial distortion.

The reference frame is set on the camera using its camera center of projection  $\mathbf{C}$  as its origin in the pinhole model. The image frame is orthogonal to the  $Z$ -axis direction of the frame at a focal distance,  $f$ , the distance between the lens and the image sensor when the subject is in focus.



**Figure 1: Pinhole Camera Geometry.**  $\mathbf{C}$  is the camera center and  $\mathbf{p}$  the principal point. The camera center is here placed at the coordinate origin. Note the image plane is placed in front of the camera center. From [7]

Using the geometric properties of the frames, the normalized pinhole image projection is given by:

$$\mathbf{x}_n = \begin{bmatrix} x \\ y \end{bmatrix} = f \begin{bmatrix} X/Z \\ Y/Z \end{bmatrix} \quad (1)$$

From [8] comes that after including lens distortion, the new normalized point coordinate  $\mathbf{x}_d$  can then be calculated as follows:

$$\mathbf{x}_d = \begin{bmatrix} \mathbf{x}_{d1} \\ \mathbf{x}_{d2} \end{bmatrix} = (1 + k_{c1}r^2 + k_{c2}r^4 + k_{c5}r^6)\mathbf{x}_n + d\mathbf{x} \quad , \quad (2)$$

where  $r^2 = x^2 + y^2$  and  $d\mathbf{x}$  is the tangential distortion vector:

$$d\mathbf{x} = \begin{bmatrix} 2k_{c3}xy + k_{c4}(r^2 + 2x^2) \\ k_{c4}xy + k_{c3}(r^2 + 2y^2) \end{bmatrix}. \quad (3)$$

Finally, the points can be converted into homogenous pixel coordinates in the image plane using the matrix known as the intrinsic camera matrix  $\mathbf{K}$ :

$$\begin{bmatrix} x_p \\ y_p \\ 1 \end{bmatrix} = \mathbf{K} \begin{bmatrix} \mathbf{x}_{d1} \\ \mathbf{x}_{d2} \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{f}_{c1} & \alpha_c \mathbf{f}_{c1} & \mathbf{cc}_1 \\ 0 & \mathbf{f}_{c2} & \mathbf{cc}_2 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{x}_{d1} \\ \mathbf{x}_{d2} \\ 1 \end{bmatrix} \quad , \quad (4)$$

where  $\mathbf{f}_c$  are the focal distances (usually expressed in mm),  $\mathbf{cc}$  are the coordinates of the principal point of the camera, and  $s = \alpha_c \mathbf{f}_{c1}$  encodes the skew of the angle between the  $x$  and  $y$  sensor axes, in digital cameras it is usually 0.

Since the camera is not usually centered at the origin of the world frame, points need to be mapped to the camera frame. Supposing a general camera rotation,  $\mathbf{R} \in \text{SO}(3)$ , and translation,  $\mathbf{T} \in \mathbb{R}^{3 \times 1}$ , the full projection matrix becomes

$$\mathbf{P} = \mathbf{K} [\mathbf{R} \quad \mathbf{T}] \in \mathbb{R}^{3 \times 4}. \quad (5)$$

1) *Back-Projection of Points to Rays:* Given a point  $\mathbf{x}$  in an image, the 3-dimensional points that map to that point constitute a ray in space that passes through the camera center. The ray can be represented as the connection of two points, one being the camera center and the other  $\mathbf{P}^+\mathbf{x}$ , where  $\mathbf{P}^+$  is the pseudo-inverse of  $\mathbf{P}$ . The line can then be formally described as

$$\mathbf{X}(\lambda) = \mathbf{P}^+\mathbf{x} + \lambda \mathbf{C} \quad . \quad (6)$$

### B. Line Detection

One of the most critical aspects of computer vision and image processing is to detect image features that should be distinct from others so that matching between images is possible. Since the type of features used in this thesis is lines, detecting them in images is crucial. A neural network will be employed to detect lines; this approach will be more detailed explained in Section IV-A.

### C. Line Matching

Since the final objective is to have a 3D reconstruction of the world, it is vital to have correct matches between lines over different camera frames. Lines will also be matched to reprojections of already seen lines to improve the estimations of those features.

**Sum of Squared Distances (SSD)** is a naive approach where the cost function is simply the squared distance of the endpoints of the line between the frames. This method expects a high image frequency or a slow-moving drone, so that lines have limited movement between poses keeping them in similar pixel positions.

Using only SSD proved unsuccessful to match lines but gave an initialization to the Procrustes Orthogonal problem [20], which can then calculate a transformation between both sets of matched lines. This transformation can then be applied to all lines making them closer to other lines to be matched, providing that the first set of matches is correct. Iteratively repeating this method is known as the **Iterative Closest Point (ICP)** first proposed by [5].

#### D. Plücker Coordinates

Introduced by Julius Plücker in the 19th century Plücker coordinates present a convenient representation for directed lines in affine 3-Dimensional space. There all multiple conventions to represent Plücker coordinates, in this thesis [2] is used as a guide. Given two 3D points  $\mathbf{M}^T \sim (\bar{\mathbf{M}}^T|m)$  and  $\mathbf{N}^T \sim (\bar{\mathbf{N}}^T|n)$ , one can represent the line joining them by a homogeneous "Plücker" 6-vector  $\mathbf{L}^T \sim (\mathbf{a}^T|\mathbf{b}^T)$ , where:

$$\begin{cases} \mathbf{a} = \bar{\mathbf{M}} \times \bar{\mathbf{N}} \\ \mathbf{b} = m\bar{\mathbf{N}} - n\bar{\mathbf{M}} \end{cases} \quad (7)$$

Other conventions for Plücker 6-vectors are also possible, each with a bilinear constraint that the 6-vector must satisfy in order to represent valid line coordinates. For the definition chosen the constraint comes as

$$\mathcal{C}(\mathbf{L}) = 0 \quad \text{where} \quad \mathcal{C}(\mathbf{L}) = \mathbf{a}^T \mathbf{b} \quad . \quad (8)$$

Given a standard  $(3 \times 4)$  perspective projection matrix  $\mathbf{P} \sim (\bar{\mathbf{P}}|p)$ , as defined in Equation 5, a  $(3 \times 6)$  matrix projecting Plücker line coordinates is given by

$$\tilde{\mathbf{P}} \sim (\det(\bar{\mathbf{P}})\bar{\mathbf{P}}^{-T}||_{\mathbf{p}}\bar{\mathbf{P}}) \quad . \quad (9)$$

#### E. Unscented Kalman Filter (UKF)

The Extended Kalman Filter [11] has been widely used to do state-estimation on non-linear systems by approximating the state distribution to a Gaussian distribution using a first-order linearization of the non-linear system. This recurrent first-order approximation introduces significant errors in the true posterior mean and covariance of the system.

The Unscented Kalman Filter [25] proposes that the state's Gaussian random variable should be approximated by a minimal set of carefully chosen sample points denoted by sigma points. After selection, the points are propagated through the non-linear system.

1) *The Square Root UKF (SR-UKF)*: The most computationally expensive operation of the UKF is on the square root of the covariance needed whenever sigma points are calculated. An efficient implementation using Cholesky factorization has  $\mathcal{O}(L^3/6)$  time complexity. The Square-Root UKF proposed by [24] reduces this complexity to  $\mathcal{O}(LM^2)$  by propagating the covariance directly, where  $M$  is the observation dimension.

In order to achieve this performance, it takes advantage of three powerful linear algebra techniques, QR decomposition, Cholesky factor updating and efficient least squares.

- **QR decomposition**: The QR decomposition or factorization of a matrix  $\mathbf{A} \in \mathbb{R}^{L \times N}$  is given by,  $\mathbf{A}^T = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q} \in \mathbb{R}^{N \times N}$  is orthogonal,  $\mathbf{R} \in \mathbb{R}^{N \times L}$  is upper triangular and  $N \geq L$ . The upper triangular part of  $\mathbf{R}$ ,  $\tilde{\mathbf{R}}$ , is the transpose of the Cholesky factor of  $\Sigma = \mathbf{A}\mathbf{A}^T$ , i.e.,  $\tilde{\mathbf{R}}$  is the square root of the covariance such that  $\tilde{\mathbf{R}}^T\tilde{\mathbf{R}} = \mathbf{A}\mathbf{A}^T$ .
- **Cholesky factor updating**: If  $\mathbf{S}$  is the original Cholesky factor of  $\Sigma = \mathbf{A}\mathbf{A}^T$ , then the Cholesky factor of the rank-1 update  $\Sigma \pm \sqrt{v}\mathbf{u}\mathbf{u}^T$  is denoted as  $\mathbf{S} = \text{cholupdate}(\mathbf{S}, \mathbf{u}, \pm v)$ . If  $\mathbf{u}$  is a matrix and not a vector,

then the result is  $M$  consecutive updates of the Cholesky factor using the  $M$  columns of  $\mathbf{u}$ .

- **Efficient Least Squares**: The solution to the equation  $(\mathbf{A}\mathbf{A}^T)\mathbf{x} = \mathbf{A}^T\mathbf{b}$  also corresponds to the solution of the over determined least squares problem  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . This can be solved efficiently using a QR decomposition with pivoting.

According to [24] even though the computational complexity of the filters are of the same order the SR-UKF is about 20% faster than the UKF and about 10% faster than the EKF.

### III. VISUAL INERTIAL FEATURE-BASED SIMULTANEOUS LOCALIZATION AND MAPPING

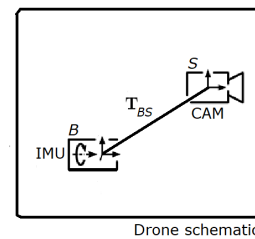
#### A. Problem Formulation

Odometry comes from the Greek words *odos* (meaning "route") and *metron* (meaning "measure") and, as such, odometry is the problem of measuring a robot route through an environment. In robots that use wheels as a way of moving, rotary encoders can estimate the movement, but for aerial drones, this is not possible (since they do not use wheels). From this problem came the approach of using a camera to estimate the robot's pose by adjusting the position of specific features that it could see over time; this approach is called Visual Odometry, first proposed by [17].

An Inertial Measurement Unit (IMU) provides an inertial estimation of the movement that can help the camera measure the drone's pose; when this happens, it is then called Visual Inertial Odometry. IMUs report the angular velocity ( $\mathbf{w}$ ) and linear accelerations ( $\mathbf{a}$ ) in the sensor frame. A rigid transformation can be used to map one frame to another:

$$\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{bmatrix} \quad , \quad (10)$$

where  $\mathbf{T}$  is the transformation matrix,  $\mathbf{R} \in \text{SO}(3)$  the rotation matrix and  $\mathbf{t} \in \mathbb{R}^3$  the translation from one frame to the other. This set of transformations constitute the Special Euclidean Group ( $\text{SE}(3)$ ).



**Figure 2:** Example transformation from  $\text{SE}(3)$  where  $\mathbf{T}$  provides a transformation from the IMU frame to the camera frame. Adapted from [4].

One of the main advantages of using an IMU as an odometry sensor is the speed at which it can update the system, given its high frequency (higher than 100 Hz) and that the update on the pose is usually straightforward. However, an IMU cannot be used as the only sensor since the pose is obtained through continuous integration of the accelerations from the IMU with quantification noise. Since this noise is additive, it introduces a bias that drifts the system away from the absolute pose.

## B. Proposed Approach

The proposed approach closely follows the work done by [3] in which an SR-UKF is used to integrate IMU and visual information. One of the main innovations that it introduces is the use of Lie Groups as a structure for the state leading to a direct representation and propagation of the uncertainty in the Lie Group. The filter estimates the robot's pose (Rotation, velocity and location), the 3D location of the features detected, and the IMU biases, the latter being appended to the state estimation.

1) *State, Action and Measurements Model:* In [3], the Special Euclidean of multiple direct spatial isometries  $SE_{2+p}(3)$  group is presented to naturally represent the problem in a Lie Group structure as leveraged in [1]. This group extends the Special Euclidean group and provides the state representation and a closed-form solution to the Lie exponential of the uncertainties. The state can be then represented as

$$\chi = \begin{bmatrix} \mathbf{R} & \mathbf{v} & \mathbf{x} & \mathbf{p}_1 \cdots \mathbf{p}_p \\ 0_{(p+2) \times 3} & I_{(p+2) \times (p+2)} & & \end{bmatrix} \in SE_{2+p}(3) \quad (11)$$

where  $\mathbf{R} \in SO(3)$  is the rotation matrix,  $\mathbf{v} \in \mathbb{R}^3$  is the velocity of the robot,  $\mathbf{x} \in \mathbb{R}^3$  is the position of robot and finally  $\mathbf{p}_k \in \mathbb{R}^3$  is the position of the  $k$ th landmark. The uncertainties, defined as  $\xi = [\xi_{\mathbf{R}}^T \xi_{\mathbf{v}}^T \xi_{\mathbf{x}}^T \xi_{\mathbf{p}_1}^T \cdots \xi_{\mathbf{p}_p}^T]^T \in \mathbb{R}^{9+3p}$  are mapped to the Lie algebra through the transformation  $\xi \rightarrow \hat{\xi}$  defined as

$$\hat{\xi} = \begin{bmatrix} [\xi_{\mathbf{R}}]_{\times} & \xi_{\mathbf{v}} & \xi_{\mathbf{x}} & \xi_{\mathbf{p}_1} \cdots \xi_{\mathbf{p}_p} \\ & & 0_{(2+p) \times (5+p)} & \end{bmatrix} . \quad (12)$$

The biases of the IMU are appended to the state and are represented as

$$\mathbf{b} = \begin{bmatrix} \mathbf{b}_{\omega} \\ \mathbf{b}_a \end{bmatrix} \in \mathbb{R}^6 , \quad (13)$$

where  $\mathbf{b}_{\omega}$  is the bias of the angular velocity and  $\mathbf{b}_a$  the bias of the linear velocity reported by the IMU.

a) *Action model:* The action model represents how the system can change in time. It uses the information given by the IMU to change the robot's pose and assumes that the landmarks do not move. The dynamics of the system read

$$\text{body state} \quad \begin{cases} \dot{\mathbf{R}} = \mathbf{R} [\boldsymbol{\omega} - \mathbf{b}_{\omega} + \mathbf{n}_{\omega}]_{\times} \\ \dot{\mathbf{v}} = \mathbf{R} (\mathbf{a} - \mathbf{b}_a + \mathbf{n}_a) - g \\ \dot{\mathbf{x}} = \mathbf{v} \end{cases} \quad (14)$$

$$\text{IMU biases} \quad \begin{cases} \dot{\mathbf{b}}_{\omega} = \mathbf{n}_{\mathbf{b}_{\omega}} \\ \dot{\mathbf{b}}_a = \mathbf{n}_{\mathbf{b}_a} \end{cases} \quad (15)$$

$$\text{landmarks} \quad \dot{\mathbf{p}}_i = 0, i = 1, \dots, p , \quad (16)$$

where the various white Gaussian noises can be stacked as

$$\mathbf{n} = [\mathbf{n}_{\omega}^T \quad \mathbf{n}_a^T \quad \mathbf{n}_{\mathbf{b}_{\omega}}^T \quad \mathbf{n}_{\mathbf{b}_a}^T]^T \sim \mathcal{N}(0, \mathbf{Q}) . \quad (17)$$

b) *Measurement Model:* The measurements are given by the calibrated monocular camera that observes the  $p$  landmarks through a standard perspective projection model, like the one described in Section II-A. Each landmark  $\mathbf{p}_i$  is observed through the camera as

$$\mathbf{y}_i = \begin{bmatrix} y_u^i \\ y_v^i \end{bmatrix} + \mathbf{n}_y^i , \quad (18)$$

where  $\mathbf{y}_i$  is the measured pixel location of the landmark in the camera, that is,

$$\lambda \begin{bmatrix} y_u^i \\ y_v^i \\ 1 \end{bmatrix} = \mathbf{P} \mathbf{p}_i , \quad (19)$$

with  $\lambda$  being the scale parameter,  $\mathbf{P}$  the camera projection matrix from Equation 5 and  $\mathbf{p}_i$  the 3-dimensional position of the landmark in homogeneous coordinates. Finally,  $\mathbf{n}_y^i \sim \mathcal{N}(0, N)$  represents the pixel image noise.

2) *Filter Implementation:* Since the filter needs a discrete model to propagate the state through time, To integrate the action model on the filter, it was discretized using the Euler method except for the Rotation.

Having the measurement from the IMU at time step  $k$  given as  $\mathbf{u}_k = [\omega_k^T \quad a_k^T]^T$  the discrete action model equations correspond to the  $f$  function that serves as the action model for the filter

$$\chi_{k+1} = f(\chi_k, \mathbf{u}_k, \omega_k) , \quad (20)$$

where the state  $\chi_k$  lives in  $G$ , and  $\omega_k \sim \mathcal{N}(0, \mathbf{Q}_k)$  is a white Gaussian noise, associated with generic discrete measurements of the form

$$\mathbf{y}_k = h(\chi_k, \mathbf{v}_k) , \quad (21)$$

where  $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{R}_k)$  is a white Gaussian noise. The final full model associated to the Right-SR-UKF-LG (Right Square Root Unscented Kalman Filter on Lie Groups), denominated as Fusion2018 from now on, becomes

$$\text{state} \quad \begin{cases} \chi_k = \exp(\xi) \bar{\chi}_k \\ \mathbf{b}_k = \bar{\mathbf{b}}_k + \tilde{\mathbf{b}} \end{cases} , \quad \begin{bmatrix} \xi \\ \tilde{\mathbf{b}} \end{bmatrix} \sim \mathcal{N}(0, \mathbf{P}_k), \quad (22)$$

$$\text{dynamics} \quad \chi_k, \mathbf{b}_k = f(\chi_{k-1}, \mathbf{u}_k - \mathbf{b}_{k-1}, \mathbf{n}_k) , \quad (23)$$

$$\text{observations} \quad \begin{cases} \mathbf{Y}_k = [\mathbf{y}_1^T \quad \cdots \quad \mathbf{y}_p^T]^T := \mathbf{Y}(\chi_k, \omega_k) \\ \mathbf{y}_k \text{ given in 18, } i = 1, \dots, p \end{cases} , \quad (24)$$

such that  $(\bar{\chi}_k, \bar{\mathbf{b}}_k) \in \mathbb{R}^{(15+3p)}$  represents the mean (estimated) state at time  $k$ ,  $\mathbf{P}_k \in \mathbb{R}^{(15+3p) \times (15+3p)}$  is the covariance matrix that defines the state uncertainties  $(\xi, \tilde{\mathbf{b}})$ , and the vector  $\mathbf{Y}_k$  contains the observations of the  $p$  landmarks with associated Gaussian noise  $\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{W})$ .

The complete filter operation is shown in Figure 6, where it is described the introduction of line features. The Propagation and Update steps are the ones earlier defined, and steps to

detect features use the method described in [21] and match features using the Matlab built-in functions. Feature Initialization is done by triangulation using the external map.

It should be noted that the Fusion2018 filter expects a good initialization, where the state is initialized with the ground truth values from the dataset for the pose and 30 3D landmark positions given by [15]. The external map also uses 200 initial features given by [15]. Landmark initializations require at least 8 different views of the same feature.

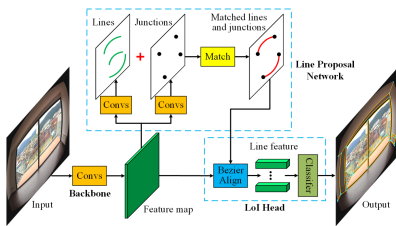
#### IV. LINES BASED VISUAL INERTIAL ODOMETRY

##### A. Line Detection using Unified Line Segment Detection (ULSD)

Following the approach in Section II-B, where lines are detected using the Canny Edge detector and the Hough Transform, a more robust way of extracting 2D lines from images was needed. A solution is presented in [12] by taking advantage of a stacked hourglass network ([16]) to extract a feature map used to propose and validate line predictions.

The complete network can be separated into two modules; the first module, denominated by "Line Proposal Network", tries to predict the junctions of lines using another network outputting confidence and a junction offset map; in this first step, duplicated junctions are also removed from the map. Lines are predicted by trying to find line segments represented by a Bezier curve. Finally, line segments that can match the endpoints to two junction proposals based on the Euclidean distance are kept. If multiple line proposals are matched with the same junctions, only those with the shortest cost are kept.

The second module, "Line of Interest", takes the candidate lines from the first module and feature map from the first network to validate the line segment prediction. The candidate is first passed through a Bezier Align function that samples uniformly the line to form the final prediction, which is then fed into another classifier network to get a confidence level for each line. In Figure 3 is shown an overview of the network.



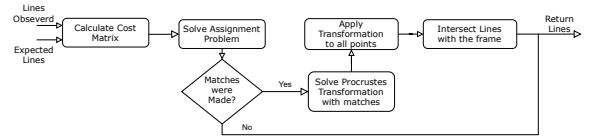
**Figure 3:** Unified Line Segment Detection (ULSD) - Network Overview. From [12]

##### B. Line Matching Algorithm

Following Section II-C the chosen algorithm to pair lines between images or the projection of 3D lines is a variation of the Iterative Closest Point. The first step is to match lines using the Sum of Squared Distances (SSD) between the endpoints of two lines; minimizing the sum of this error is known as the Assignment Problem. In order to correctly match endpoints, the endpoints of the second line are swapped, and the distance is recalculated; It is then chosen the orientation which provides

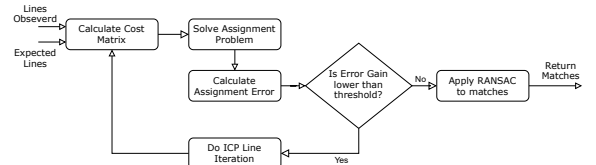
the smaller distance. The minimization problem is solved using the Hungarian method proposed by James Munkres in [14]. Matches are also filtered by a maximum error cost between any two lines.

After the first match, Procrustes is applied to the endpoints of the matched lines. The obtained transformation is then applied to all the points in the image so that lines that should be matched to each other get closer together. As such, this method requires that at least a match so that a transformation can be calculated. An overview of this iteration is shown in Figure 4.



**Figure 4:** ICP Line Iteration Block Diagram This algorithm matches lines using the SSD and then adjusts the endpoints positions of all points using the transformation found by applying Procrustes to the matched endpoints.

Since the objective is to minimize the error given by SSD, the method described in Figure 4 is iteratively run until the error no longer decreases above a fixed threshold. Since the ICP is used with the objective of matching in mind, it simply returns the matches obtained. After the ICP algorithm is done, RANSAC is used to filter some outliers that may appear from overfitting to some matches. The entire algorithm is shown in Figure 5.



**Figure 5:** ICP Line Block Diagram This algorithm performs iterative matches and lines transformations until the error gain is lower than a fixed threshold ( $10^{-6}$  in the case of this thesis). The first solve of the assignment problem tolerates a higher error in matches than subsequent iterations.

##### C. Line Initialization in Plücker Coordinates

Line initialization is crucial to obtain 3D features that the filter can propagate to update the system's state. With this objective in mind, it was proposed by [2] a maximum likelihood estimator using the squared 2D orthogonal (Euclidean) distance from the line  $\mathbf{L}$  to the projected endpoints weighted by their distance to the camera, which is formally written as the sum of the distance to both points and defined by

$$d_{\perp}^2(\mathbf{q}, \mathbf{l}) = (\mathbf{q}^T \mathbf{l})^2 / (l_1^2 + l_2^2) \quad (25)$$

Having 3D lines  $\mathcal{S} = \{\mathbf{L}^1, \dots, \mathbf{L}^m\}$  and cameras  $\mathcal{M} = \{\mathbf{P}^1, \dots, \mathbf{P}^n\}$  the negative log likelihood function  $\mathcal{E}(\mathcal{S}, \mathcal{M})$  for the reconstruction, corresponding to the total reprojection error, can be written in terms of individual reprojection errors  $\mathcal{E}(\mathbf{L}^j, \mathcal{M})$  for each line  $j$  but since all projections should

correspond to the same 3D line, the final functional to be minimized becomes

$$\underset{\mathbf{L}}{\text{minimize}} \quad \mathcal{E}(\mathbf{L}, \mathcal{M}) = \sum_{i=1}^n (d_{\perp}^2(\mathbf{x}^i, \mathbf{I}^i) + d_{\perp}^2(\mathbf{y}^i, \mathbf{I}^i)) \quad (26)$$

1) *Linear Algorithm*: Ignoring the denominator in equation (25) leads to an algebraic distance denoted  $d_a$ , biased compared to the orthogonal distance. It is linear in the predicted line and defined by  $d_a^2(\mathbf{q}, \mathbf{l}) = d_{\perp}^2(\mathbf{q}, \mathbf{l})w^2 = (\mathbf{q}^T \mathbf{l})^2$ , where the scalar factor  $w$  encapsulates the bias as  $w^2 = l_1^2 + l_2^2$ :

$$(w^i)^2 = \left( (\tilde{\mathbf{P}}^i \mathbf{l})_1 \right)^2 + \left( (\tilde{\mathbf{P}}^i \mathbf{l})_2 \right)^2 \quad (27)$$

The biased linear least squares error function is then defined by

$$\begin{aligned} \mathcal{B}(\mathbf{L}, \mathcal{M}) &= \sum_{i=1}^n \left( \left( x^{iT} \tilde{\mathbf{P}}^i \mathbf{l} \right)^2 + \left( y^{iT} \tilde{\mathbf{P}}^i \mathbf{l} \right)^2 \right) \\ &= \|\mathbf{A}_{(2n \times 6)} \mathbf{l}\|^2 \quad \text{with} \quad \mathbf{A} = \begin{pmatrix} \dots \\ x^{iT} \tilde{\mathbf{P}}^i \\ y^{iT} \tilde{\mathbf{P}}^i \\ \dots \end{pmatrix} \end{aligned} \quad (28)$$

Since  $\mathbf{l}$  is a homogeneous vector, the constraint  $\|\mathbf{l}\|^2 = 1$  must be added. The  $\mathbf{l}$  that minimizes  $\mathcal{B}(\mathbf{L}, \mathcal{M})$  is then given by the singular vector of  $\mathbf{A}$  associated to its smallest singular value, that can be computed using SVD.

The linear algorithm provides a biased estimate of the plucker line  $\mathbf{l}$  that can estimate the weight factors that contain the bias of the linear least-squares error function. This process is then used iteratively until convergence, using the difference between consecutive errors as a threshold. A correction should be made between each iteration to ensure that the plucker constraint is satisfied.

2) *Quasi-linear Algorithm*: According to [2], the linear algorithm misbehaves since the plucker constraint is not considered while solving the linear least-squares problem.

A new approach is suggested by rewriting the constraint as  $\mathcal{C}(\mathbf{l}) = \mathbf{l}^T \mathbf{G} \mathbf{l} = 0$  where  $\mathbf{G} = \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix}$ . Computing the null space of  $\mathbf{l}^T \mathbf{G}$ , using SVD, gives a base of all vectors that satisfy the plucker constraint; this base can then be plugged into the linear estimator (Equation. 28) to give a line  $\mathbf{l}$  that satisfies the plucker constraint. The estimate is then used iteratively, as before, to compute the bias that reweighs the least-squares problem (Equation. 27) until convergence.

3) *Endpoints Extraction from the Plucker Coordinates*: Since the endpoints will be propagated through the SR-UKF filter, they must be extracted from the Plucker Coordinates. First, the Plucker line is transformed back into Euclidean Coordinates, then the endpoints are estimated by finding the closest line points to the back-projected rays of the endpoints observed in each image. These values are then averaged to give a final prediction of the 3D endpoints. The points chosen are always part of the line estimated using the quasi-linear algorithm.

D. *Line Initialization Using 3D Distance to Back Projected Rays*.

In a noiseless world, all back-projected rays from a single 3D point should intersect in that same point, but this is not always true due to noises in the image or the camera's pose. In reality, the rays will pass close to the real point but probably not over, as shown in Figure ?? ; as such, the real point should be close to every back-projected ray from every view that can see him. The 3D distance from a generic point  $\mathbf{P}$  to a line defined by  $\mathbf{L} = \mathbf{A} + \lambda \mathbf{u}$  can be calculated as:

$$d(\mathbf{L}, \mathbf{P}) = \frac{\|\overrightarrow{\mathbf{A}\mathbf{P}} \times \mathbf{u}\|}{\|\mathbf{u}\|} \Rightarrow \|\overrightarrow{\mathbf{A}\mathbf{P}} \times \mathbf{u}\|, \quad \text{if} \quad \|\mathbf{u}\| = 1. \quad (29)$$

Having the back-projected rays  $\{\mathbf{L}^1 = \mathbf{A}^1 + \lambda \mathbf{u}^1, \dots, \mathbf{L}^n = \mathbf{A}^n + \lambda \mathbf{u}^n\}$ , a mean squared error formulation can be written as

$$\underset{\mathbf{P}}{\text{minimize}} \quad F(\mathcal{L}, \mathbf{P}) = \frac{1}{2} \sum_{k=1}^n \|d(\mathbf{L}^k, \mathbf{P})\|^2 \quad (30)$$

Since the minimum of the functional is found when the gradient is equal to zero, a closed form solution is given by:

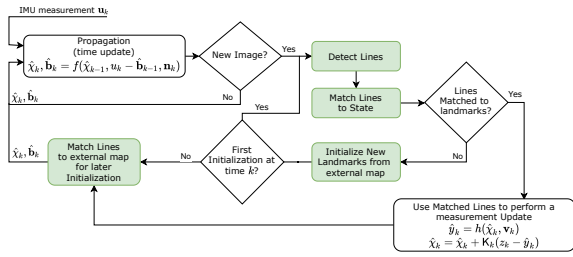
$$\begin{aligned} \nabla F(\mathcal{L}, \mathbf{P}) &= 0 \\ \Leftrightarrow \sum_{k=1}^n \begin{bmatrix} \mathbf{u}_2^{k2} + \mathbf{u}_3^{k2} & -\mathbf{u}_1^k \mathbf{u}_2^k & -\mathbf{u}_1^k \mathbf{u}_3^k \\ -\mathbf{u}_1^k \mathbf{u}_2^k & \mathbf{u}_1^{k2} + \mathbf{u}_3^{k2} & -\mathbf{u}_2^k \mathbf{u}_3^k \\ -\mathbf{u}_1^k \mathbf{u}_3^k & -\mathbf{u}_2^k \mathbf{u}_3^k & \mathbf{u}_1^{k2} + \mathbf{u}_2^{k2} \end{bmatrix} \begin{bmatrix} \overrightarrow{\mathbf{A}^k \mathbf{P}}_x \\ \overrightarrow{\mathbf{A}^k \mathbf{P}}_y \\ \overrightarrow{\mathbf{A}^k \mathbf{P}}_z \end{bmatrix} &= 0 \\ \Leftrightarrow \sum_{k=1}^n \mathbf{U}^k (\mathbf{P} - \mathbf{A}^k) \Leftrightarrow \left( \sum_{k=1}^n \mathbf{U}^k \right) \mathbf{P} &= \sum_{k=1}^n \mathbf{U}^k \mathbf{A}^k \end{aligned} \quad (31)$$

where  $\overrightarrow{\mathbf{A}^k \mathbf{P}}_i$  is  $i$ th coordinate of the vector defined by  $\overrightarrow{\mathbf{A}^k \mathbf{P}} = \mathbf{P} - \mathbf{A}^k$ , and  $\mathbf{u}_i$  is the  $i$ th coordinate of the vector  $\mathbf{u}$ . The vectors  $\mathbf{u}$  that define the direction of the line should be normalized. It should be noted that this system may not be determined for  $n = 2$ . For  $n = 1 \Rightarrow \mathbf{P} = \mathbf{A}$ .

E. *Integration of Line Features into Fusion 2018*

The algorithms proposed in this chapter served as a replacement to the algorithms used to detect, match and initialize features to the Fusion2018. One main difference is that even though the filter propagates the endpoints of the lines as independent landmarks, they maintain the line structures that are then used to match the state landmarks and the images taken by the camera. The blocks that are replaced in the Fusion2018 filter are shown in green in Figure 6.

The external map saves the image coordinates and projection matrix of each matched view for every line ID that are still not initialized. If there is a low number of matches to the landmarks (they do not appear in frame) the initialization step can be called to try and create landmarks which can be seen by the system. The state initialization will be discussed in Section V-E.



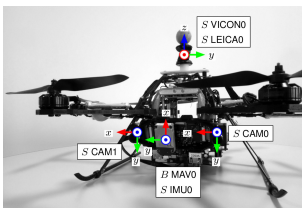
**Figure 6: Fusion2018 Block Diagram Filter Proposed by [3] to perform Visual Inertial Odometry with Lines as Features. Green blocks are the new blocks introduced to use line features.**

## V. EXPERIMENTS

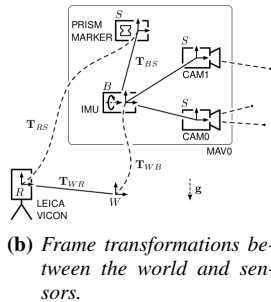
### A. The EuRoC MAV Dataset

The dataset used in this thesis was recorded by [4] in the context of the European Robotics Challenge (EuRoC) to assess the contestant’s visual-inertial simultaneous localization and mapping (SLAM) and 3D reconstruction capabilities on micro aerial vehicles (MAVs). More specifically it will be used the medium difficulty data set recorded on the Vicon MoCap Room.

The dataset features an accurate 3D point cloud of the environment captured with a laser scanner and also the 6D pose ground truth (position and rotation) from the flight of the AscTec ”Firefly” hex rotor helicopter, shown in Figure 7. The drone features an Inertial Measurement Unit (IMU) and two cameras (even though only one is used).



(a) The AscTec ”Firefly” helicopter.



(b) Frame transformations between the world and sensors.

**Figure 7: The drone used in the recorded dataset and the transformations between frames of the sensors in the drone and the world. From [4].**

The Fusion18 dataset consists of 16020 IMU measurements and 1602 images. The IMU values can be used directly since its frame is aligned with the body’s referential. On the other hand, when used for image reconstruction, the camera image coordinates need an extra transformation to ensure that the values are mapped to the correct frame. As such, the camera projection matrix from equation 5 becomes

$$P = K(R_w^b R_b^c)^T [I_3 \quad -(\mathbf{T}_w^b + R_w^b R_b^c \mathbf{T}_b^c)] \in \mathbb{R}^{3 \times 4}, \quad (32)$$

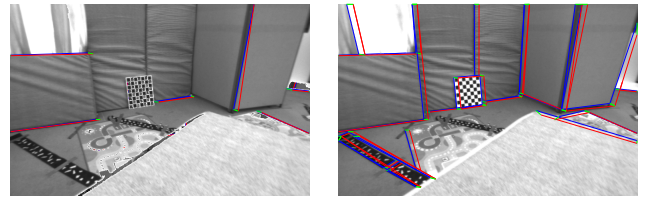
where  $R_a^b \in SO(3)$  is the rotation matrix from referential  $a$  to  $b$ , and  $\mathbf{T}_a^b \in \mathbb{R}^3$  is the translation vector from referential  $a$  to  $b$ ;  $w$ ,  $b$  and  $c$  are the referential of the world, the drone (IMU) and the camera, respectively.

### B. Line Detection with ULSD

The ULSD network can be trained separately, but since [12] provide a pre-trained model for each type of camera, the pinhole model is used for line detection. The confidence interval ( $CI$ ) which declines or accepts lines can be adjusted if there is a need to detect more or fewer lines. The  $CI$  set to 70% is used since there is an improvement in the number of detected lines without losing accuracy, being detected on average 67 lines per image. Lower thresholds identify repeated lines, increasing the complexity for the matching task to maintain lines alive through multiple frames.

### C. Line Matching Validation

The algorithm was validated using subsequent images from the Fusion18 dataset where lines were identified using the method in Section IV-A. It is possible to check how much the algorithm improved the number of matches obtained in Figure 8. This growth is vital so that lines are consistently matched across images making possible a posterior 3D initialization.



(a) Matches obtained by SSD (b) Matches obtained by ICP Line

**Figure 8: Example of ICP Line improvement on matching. Blue lines represent the lines identified on the image, and red lines correspond to lines found in other images that were matched to this image. Finally, green lines join the matched endpoints.**

The Fusion18 dataset was then used to better understand how long the algorithm would take to converge and how many matches were gained by utilizing the ICP. It was found that it took 6.3 iterations to converge, and more than 17 matches were gained per a pair of images. More importantly, the percentage of pairs of images with less than 20 matches went from 20% to less than 5%, demonstrating the power of this algorithm in increasing the matches found. The algorithm also showed promising results when matching lines from images with greater time steps between them.

### D. Simulation

A Matlab simulator was created to test the line initialization techniques, which consisted of a cube centered on the origin with random lines inside and a variable number of cameras that would take pictures of the cube and lines. The line initialization techniques were then tested and then compared to the ”real” lines. In the simulator, it is possible to modify the number of cameras, the positions of each camera and the option to use a camera for the projection of the reconstructed lines to assess the reprojection errors of the initialized lines.

The first validation was done with no noise added to the simulator to ensure the validity of the algorithms in a case where perfect measurements are made. The techniques resulted in euclidean distance errors lower than  $10^{-15}$ , appearing from Matlab functions' numerical precision. These results supported the validity of the initializations using the [2] and Section IV-D methods, at least in a noiseless environment.

1) *Effects of noise:* Since reality is not noiseless, the simulator provided three types of noise that could be tested separately or at the same time; the first type of noise appears from vibrations on the camera or for incorrect line detection, and it is manifested as pixel deviations on the line positions in the pictures taken by the camera. The other two errors come from the camera's pose localization and appear as location and rotation differences from the actual pose. The noise was modelled as a Gaussian distribution with zero average and covariance given by noise value selected.

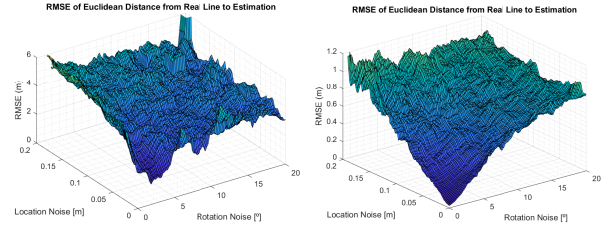
The cameras were placed near each other to understand the effect of noise on the lines' initialization. The reason was that the simulation could be closely related to the drone flight where there is little movement between each frame. As such, cameras were placed less than 20 cm apart and with less than 10 degrees of total rotation between each consecutive frame.

As a first validation, all noise types were tested separately by setting all types to zero minus the one to be tested. A common occurrence is that the method described in section IV-C using the Plücker coordinates misbehaves in the presence of noise and rapidly increases the reprojection error to values where the last frame would discard the lines. One intriguing aspect is that the error seems to reach a plateau, which may come from intersecting the infinite Plücker line with the back-projected rays to retrieve the endpoints.

On the other hand, the technique using the minimum 3D distance to the back-projected rays seemed to behave predictably when exposed to all kinds of noise. It presented worse results in the cases where the location noise would come close to the difference in location of each frame (around 7 cm and higher noise). Even though it presents a lower error than the first method, this latter algorithm seems to be more affected by rotation noises.

Since the pixel noise should suffer fewer variations while running the SR-UKF, it was set as fixed for the following two experiments, where location and rotation noise were analyzed in detail. First, the pixel noise was set to zero, the location noise varied from 0 to 20 cm, and the rotation noise from 0 to 20 degrees, which is higher than the difference in the camera's pose for consecutive frames, giving a complete picture on how the error of the initialization should vary for different combinations of noise. The error was measured by averaging the length from the actual endpoints and their initialization. The results are shown in Figure 9.

As seen previously, the Plücker coordinates method showed little resistance to noise, increasing almost instantly to values above 2 m which would place the lines outside 1 by 1-meter cube. A strange phenomenon is that it seemed to achieve a little better results for small values of rotation noise but higher than zero. This phenomenon was found across multiple experiences, and it may be a result of the representation



(a) RMSE for the Plücker Line Method. (b) RMSE for the Back-Projected Rays Method.

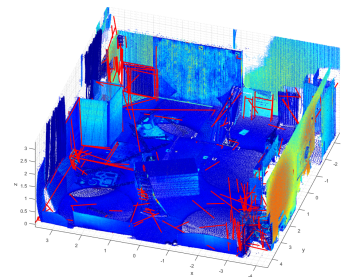
**Figure 9:** RMSE of the euclidean distance of the initialized lines to the real lines in function of the rotation and location noise with pixel noise fixed as zero.

used where it benefits from a slightly higher rotation between frames to converge to a result, even though this rotation came from noise.

Once again, the Back-Projected Rays method proved to behave with a consistent increase in error for similar increases in noise, one crucial aspect that should be noted is that for the values of noise 5 degrees and 5 cm, the method showed an error of 20 cm (on average 10 cm for each endpoint) in the initialization close to having a pixel noise of 4. For a reliable line initialization, the filter should maintain the location noise under 5 cm and rotation noise under 5 degrees, expecting that the filter can then correct the position of the lines.

#### E. Application of the Proposed Approaches on the EuRoC MAV Dataset

As a first naïve approach, the matching algorithm was run across all images of the EuRoC MAV Dataset, and then lines that survived for longer than 15 frames were initialized using the method of the back-projected rays described in Section IV-D. These results can be observed in Figure 10, and they show that many of the lines follow the environment correctly.



**Figure 10:** Line Initialization using the matching and Initialization methods proposed and using the ground truth for the robot's pose.

Two main zones were identified where the technique could not initialize many lines, one where the drone observes a curtain and the other where it sees a wall. It was recognized that this problem appears from the line detector. This theory was validated upon inspection of the pictures taken in those areas. The main problem from those frames is that most observed lines disappear; Since lines are only matched on



subsequent images, they are interpreted as new features when the system sees the same line again, with no relation to the past observations.

1) *Complete System Experiment*: Following the Fusion2018 filter, the robot's pose was initialized with the ground truth. The landmarks were initialized using ten lines selected from Figure 10; the proposed filter needed a higher initial landmark covariance since the method's precision was lower than the ORB SLAM filter. The external map was initialized without any feature and is constructed as the filter progresses.

One of the problems encountered was that lines exhibited more pixel error from the reprojections of the landmarks than initially expected. This error is again due to the not entirely perfect line initializations selected for the landmarks. One solution to this problem was to allow a higher error for the first ICP line iteration. This problem is significantly intensified in the zone where the drone looks to curtains displayed in Figure ??, where few lines could be used to update the system state. This technique was called "loose matching", and resulted in a increased number of features matched during the critical time steps between 3 and 4 seconds where the curtains occupy most of the frame.

The position error is obtained from the Euclidean distance between the robot's estimated pose and the ground truth. The rotation error is calculated using the metric proposed by [18], where the distance between two rotation matrices can be calculated using the logarithm map of the product of the inverse of one of the matrices by the other. This metric can be formally described as

$$d(R_1, R_2) = \|\log(R_1^{-1}R_2)\| = \left\| \cos^{-1} \left( \frac{\text{Tr}(R_1^{-1}R_2) - 1}{2} \right) \right\|, \quad (33)$$

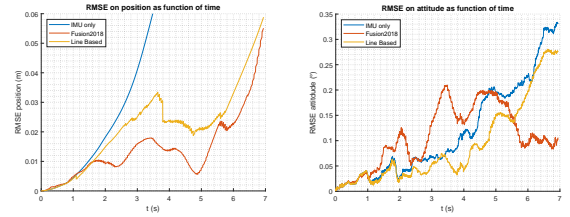
where  $\|\cdot\|$  denotes the Euclidean norm,  $R_1, R_2 \in \text{SO}(3)$  are rotation matrices and  $\log(\cdot)$  represents the logarithm map  $\log: \text{SO}(3) \rightarrow \mathfrak{so}(3)$ .

After the first six seconds of flight, the filter starts to deviate and is incapable of initializing lines soon enough due to loss of the combination of observing a new environment and passing through another zone with a low number of lines, shown in Figure ?? as the "wall". The filter tries to initialize the lines to replace the state landmarks but diverges before fixing its pose beyond irreparable damage. A mix of points and lines could probably be used to combat these zones where the environment is not very geometric; however, this approach was not explored.

In the following steps, the filter was compared to the Fusion2018, which works in the same circumstances since it does not initialize new features before seven seconds into the flight. The Fusion2018 filter is not lead to divergence since it can find features on the textured wall.

Since the used lines have lower precision than the landmarks used by the Fusion2018 filter, the proposed filter takes a longer time to start updating its position and trusts the IMU readings for the first three seconds. After the curtain section, the filter rapidly updates its position, almost obtaining the same error as the Fusion2018 filter. The filter excels in the rotation error while it has observable features, always maintaining a lower

error than the IMU and the Fusion2018 filter; this may be due to lines giving better rotational information than points. These results are shown in Figure 11

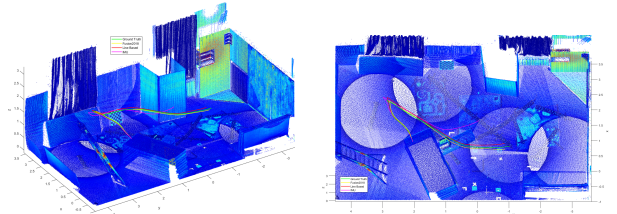


(a) RMSE of the position in function of time.

(b) RMSE of the rotation in function of time.

**Figure 11:** Error of the full proposed solution compared to the Fusion2018 filter and only using IMU.

The complete filter and ground truth trajectories are shown in detail in Figure 12, where it is possible to check that the filter responds to the aggressive change in direction faster than the IMU and maintains the actual trajectory path where IMU rapidly diverges.



(a) Side view of the flight estimated by the filters.

(b) Top view of the flight estimated by the filters.

**Figure 12:** Detailed view of the flight estimated by the filters and the ground truth trajectory.

A summary of the obtained errors is presented in Table I, where the improvement of using image-based features to update the IMU predictions is evident. An aspect that should be mentioned is that the tuning parameters of the proposed filter had to change so that the filter would not diverge due to noise of not only the initializations but also to the difficulty of accurately identifying the same endpoints of lines in the dataset.

**Table I:** Average RMSE of the filters tested on the EuRoC MAV dataset.

Filter	Mean Position error [cm]	Mean Rotation error [deg]
IMU Only	11.18	0.117
Fusion2018	1.42	0.108
Line Base	2.15	0.091

The proof of concept presented shows that lines can be used with success to perform Visual Inertial Odometry. An exciting observation is that the filter can maintain a low error value even with fewer landmarks since only 10 lines are used (20 points) compared to the 30 used by the Fusion2018 filter. This reduction in the state size may allow using a higher number of updates and a higher frequency camera since the update step scales poorly with the increase of landmarks,

having a complexity of  $\mathcal{O}(LM^2)$ , where  $M$  is the observation dimension and  $L$  one of the dimensions of the state.

Another promising result is observed by viewing the initialized lines in Figure 10 without the ground truth point cloud, where it is possible to get a general idea of the environment while using a limited number of lines when compared to the point cloud, that for reference, uses over 3 million points.

## VI. CONCLUSION AND FUTURE WORK

### A. Conclusion

This work proposed to give the needed building blocks to achieve a SLAM solution using VIO techniques with Lines as features. It was motivated by new advances in machine learning that facilitated the detection of lines in images that could provide a richer mapping visualization tool than classical point features.

The main obstacle encountered was the matching of lines through multiple frames to be initialized as 3D structures. An ICP variation was introduced to remedy this problem, where the main focus of the algorithm was to match the features instead of finding the transformation between different views. This solution proved to be very successful in increasing the number of matches by 17 for each image pair. It also reduced the percentage of images with less than 20 matches from 20% to less than 5%.

It was also shown that the method proposed to initialize lines using Plucker coordinates did not provide a great response to noise and rapidly diverged. As such, a method was proposed that minimizes the distance of the 3D point to back-projected rays; this method shows a robust and predictable response to the tested types of noise.

Finally, it was found that some zones of the environment present little to no lines, which leads to a potential loss of lines to be initialized and leads the filter to divergence in the worst-case scenario; it is shown that the filter can survive through brief moments through those zones, but some adaptations may be needed to surpass these obstacles.

As a proof of concept, the proposed approach provided a sound localization of the drone through a somewhat complex trajectory and provided some line initializations which could potentially maintain alive the filter. It was shown that lines provide a rich descriptor of an environment compared to a point cloud and that the line algorithms proposed could give a first approximation of the map without the need for filtering processes in cases where the localization is provided.

### B. Future Work

I think that that a mixed approach may solve the problem of zones where the filter cannot detect lines. A parallel map or even filter, which uses the robot's pose from the primary filter, can save or propagate the point features. In the problematic zones, the point features would substitute the line features while line initialization is impossible. Ways of directly propagating line coordinates inside the filtering process should improve the precision of the lines detected. However, this approach would need an external map where endpoints are saved so that line segments that may be part of the same infinite line are matched correctly.

## REFERENCES

- [1] Axel Barrau and Silvere Bonnabel. An EKF-SLAM algorithm with consistency properties. *CoRR*, abs/1510.06263, 2015.
- [2] Adrien Bartoli and Peter Sturm. Structure-From-Motion Using Lines: Representation, Triangulation and Bundle Adjustment. *Computer Vision and Image Understanding*, 100:416–441, 2005.
- [3] M. Brossard, S. Bonnabel, and J. Condomines. Unscented kalman filtering on lie groups. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2485–2491, 2017.
- [4] Michael Burri, Janosch Nikolic, Pascal Gohl, Thomas Schneider, Joern Rehder, Sammy Omari, Markus W Achtelik, and Roland Siegwart. The euroc micro aerial vehicle datasets. *The International Journal of Robotics Research*, 35(10):1157–1163, 2016.
- [5] Yang Chen and Gérard Medioni. Object modelling by registration of multiple range images. *Image and Vision Computing*, 10(3):145–155, 1992. Range Image Understanding.
- [6] Davison. Real-time simultaneous localisation and mapping with a single camera. In *Proceedings Ninth IEEE International Conference on Computer Vision*, pages 1403–1410 vol.2, 2003.
- [7] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, USA, 2000.
- [8] J. Heikkila and O. Silven. A four-step camera calibration procedure with implicit image correction. In *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1106–1112, 1997.
- [9] Simon J. Julier and Jeffrey K. Uhlmann. New extension of the Kalman filter to nonlinear systems. In Ivan Kadar, editor, *Signal Processing, Sensor Fusion, and Target Recognition VI*, volume 3068, pages 182 – 193. International Society for Optics and Photonics, SPIE, 1997.
- [10] R.E. Kalman. Contributions to the theory of optimal control, 1960.
- [11] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45, 1960.
- [12] H. Li, H. Yu, W. Yang, L. Yu, and S. Scherer. Ulsd: Unified line segment detection across pinhole, fisheye, and spherical cameras, 2020.
- [13] J. M. M. Montiel, Javier Civera, and Andrew J. Davison. Unified inverse depth parametrization for monocular slam. In *Robotics: Science and Systems*, 2006.
- [14] James R. Munkres. Algorithms for the assignment and transportation problems. *Journal of The Society for Industrial and Applied Mathematics*, 10:196–210, 1957.
- [15] Raul Mur-Artal and Juan D. Tardos. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, Oct 2017.
- [16] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation, 2016.
- [17] D. Nister, O. Naroditsky, and J. Bergen. Visual odometry. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 1, pages I–I, 2004.
- [18] F. C. Park. Distance Metrics on the Rigid-Body Motions with Applications to Mechanism Design. *Journal of Mechanical Design*, 117(1):48–54, 03 1995.
- [19] Eduardo Perdices, Luis Lopez-Ramos, and José Cañas Plaza. Lineslam: Visual real time localization using lines and ukf. *Advances in Intelligent Systems and Computing*, 252:663–678, 01 2014.
- [20] P. H. Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31:1–10, 1966.
- [21] Jianbo Shi and Tomasi. Good features to track. In *1994 Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 593–600, 1994.
- [22] Paul Smith, Ian Reid, and Andrew Davison. Real-time monocular slam with straight lines. pages 17–26, 01 2006.
- [23] Joan Solà, Teresa Vidal-Calleja, Javier Civera, and Jose Maria Martinez-Monti. Impact of landmark parametrization on monocular EKF-SLAM with points and lines. *International Journal of Computer Vision*, 97(3):pp.339–368, September 2011.
- [24] R. Van der Merwe and E.A. Wan. The square-root unscented kalman filter for state and parameter-estimation. In *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, volume 6, pages 3461–3464 vol.6, 2001.
- [25] E. A. Wan and R. Van Der Merwe. The unscented kalman filter for nonlinear estimation. In *Proceedings of the IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium (Cat. No.00EX373)*, pages 153–158, 2000.