

Project 4: Deep and Un-Deep Visual Inertial Odometry

RBE 549: Computer Vision

Riley Blair

School of Engineering
Worcester Polytechnic Institute
Email: rpblair@wpi.edu

Scott Pena

School of Engineering
Worcester Polytechnic Institute
Email: sb pena@wpi.edu

Abstract—Visual Inertial Odometry (VIO) aims to provide accurate odometry of a moving object (like a drone or car) in the wake of imperfect sensors. An Inertial measurement Unit (IMU) is noisy, and prone to drifting over long periods of integration, whereas a camera is data-hungry, and susceptible to short-term errors (motion blur, data-loss), when polled infrequently. By combining the two in a Multi-State Constraint Kalman Filter (MSCKF), we aim to create a high-precision robust tracking software stack for any mobile system.

I. INTRODUCTION

A classic problem in computer vision and robotics is detecting scale and depth via onboard sensors. Using a single RGB camera, it is impossible to obtain depth without any prior knowledge of the environment without using a deep learning model. Even then, the model is highly constrained to the dataset it was trained on. An alternative solution is to use a stereo camera, but obtaining depth and scale using stereo is computationally expensive and can be very slow. The best option is to use an IMU onboard the robot. An IMU can be fused with the power of an RGB camera to estimate the camera's pose and depth within the camera's image. The process to accomplish localization using the combination of computer vision and an IMU sensor is known as Visual-Inertial Odometry, or VIO. This paper highlights our implementation of VIO on a drone flying in an unknown environment, using the implementation of a Multi-State Constraint Kalman Filter provided by a paper written by Anastasios I. Mourikis and Stergios I. Roumeliotis. [2]

II. DATASET

The dataset provided for this assignment was the Machine Hall 01 Easy, or MH_01_easy subset of the EuRoC dataset. The dataset was generated via a VI sensor integrated within a drone flying a trajectory in an environment. This recorded data regarding the IMU readings for linear acceleration and angular velocity, as well as captured a video onboard the drone. The ground truth of the drone's pose was recorded using the Vicon Motion capture system. The IMU readings and RGB images were used to perform our implementation of VIO, where we compared our output with the ground truth trajectory.

III. IMPLEMENTATION

A. Initialize Gravity and Bias

A strong disadvantage to using an IMU is that the readings for linear acceleration and angular velocity are prone to high amounts of error due to unpreventable mechanical inconsistencies. To combat the readings from drifting too far when the drone first launches, we initialize the IMU's gravity vector and the IMU's bias for both the linear acceleration and angular velocity.

We assume the drone is stationary at the beginning of the flight, and therefore can calculate the biases before any movement has occurred. For the first buffer of IMU messages (200 messages), we calculate the average linear acceleration and the average angular velocity. We set the gyroscope's bias to be equal to the average angular velocity over the 200 messages. To find the gravity vector, we calculate the norm of the average linear acceleration. Since the IMU detects the acceleration of gravity, we can create the gravity vector to be the calculated norm pointed in the negative z axis. Now that we have estimated the bias of the gyroscope and the gravity vector, we can also initialize the IMU's orientation with respect to the world.

B. Batch IMU Processing

During execution, this function processes the data recorded in IMU messages in batches. It iterates through each message in the batch, extracting the timestamp, angular velocity, and linear acceleration to be processed through the model. The function also has a time limit, preventing the function from lagging behind the sensor readings. If the current message's timestamp is too far behind the current IMU's timestamp, then the function moves onto the next batch. The function also keeps track of the current and next IMU IDs to keep track of the IMU state.

C. Process Model

This function calculates the dynamics of the error IMU state, and assists in predicting the new state of the IMU. We first calculate the error of the angular velocity and linear

acceleration by subtracting the currently tracked velocity and acceleration by the new IMU message data. We also get a time difference subtracting the previous and current timestamps. The IMU state error can be calculated using Equation 1.

$$\dot{\tilde{X}}_I = F\tilde{X}_I + Gn_I \quad (1)$$

where F and G are defined as per [2]:

$$F = \begin{bmatrix} -\hat{\omega}_x & -\mathbf{I}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ -C(\hat{\mathbf{q}}^I)^T[\hat{\mathbf{a}}_x] & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & -C(\hat{\mathbf{q}}^I)^T & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{I}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \end{bmatrix}$$

$$G = \begin{bmatrix} -\mathbf{I}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{I}_3 & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & -C(\hat{\mathbf{q}}^I)^T & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{I}_3 \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \end{bmatrix}$$

After computing F and G , we approximate the matrix exponentially to the third order using Equation 2 of Sun et. Al. [4].

$$\Phi = \mathbf{I}_{21} + F * dt + 0.5 * (F * dt)^2 + \frac{1}{6} * (F * dt)^3 \quad (2)$$

After this approximation, we predict the new state of the system before continuing this method (see Section 3D).

D. Predict New State

Using the IMU readings from the IMU message, we can predict the next state of the drone. Using the calculated state error estimated in Section 3C, we can propagate the state using 4th order Runge-Kutta via the MSCKF paper. We first calculate the Ω matrix, containing the time evolution of the IMU state. The Ω matrix is defined as:

$$\Omega = \begin{bmatrix} [-\hat{\omega}] & \omega \\ -\omega^T & 0 \end{bmatrix} \quad (3)$$

Afterwards, we get the current orientation, position, and velocity of the IMU state, we can re-estimate the IMU's angular velocity and linear acceleration using 4th order Runge-Kutta:

$$k_1 = f(t_n, y_n) \quad (4)$$

$$k_2 = f(t_n + 0.5dt, y_n + 0.5k_1dt) \quad (5)$$

$$k_3 = f(t_n + 0.5dt, y_n + 0.5k_2dt) \quad (6)$$

$$k_4 = f(t_n + 0.5dt, y_n + 0.5k_3dt) \quad (7)$$

We update the IMU state by converting the estimated orientation into a quaternion and assigned the estimated orientation, velocity, and acceleration to the current IMU state.

E. State Augmentation

This function calculates the state covariance matrix using Equation 3 in the MSCKF paper. First, it gets the current IMU state, as well as the rotation and translation from the IMU frame to the camera frame. We then build the Jacobian matrix J to estimate the covariance matrix P of the EKF.

$$J = \begin{bmatrix} C(\hat{\mathbf{q}}^I) & \mathbf{0}_{3 \times 9} & \mathbf{0}_{3 \times 3} & \mathbf{I}_3 & \mathbf{0}_{3 \times 3} \\ -C(\hat{\mathbf{q}}^I)^T[\hat{\mathbf{p}}_{cx}] & \mathbf{0}_{3 \times 9} & \mathbf{I}_3 & \mathbf{0}_{3 \times 3} & \mathbf{I}_3 \end{bmatrix} \quad (8)$$

$$P_{k+1|k} = \begin{bmatrix} P_{II_{k+1|k}} & \phi_k P_{IC_{k|k}} \\ P_{IC_{k|k}}^T \phi_k^T & P_{CC_{k|k}} \end{bmatrix} \quad (9)$$

F. Add Feature Observations

The camera produces a set of messages based on the features it is tracking within an image. That message gets sent to this function, where for each feature, we append its position to the set with the same feature id. The set contains every frame for which this feature is visible, and is used in the measurement step to calculate the residuals in the measurement step.

G. Measurement Update

The measurement step reduces uncertainty in the system by using the tracked features as a reference for a duration of uncertain movement. Update steps happen once a feature is no longer tracked (falls outside of the image frame), or when there is enough tracked features in the map server for which we can prune them.

We use three parameters for the measurement model: the residual r , the measurement Jacobian H and the state error \tilde{X} . They are related via this equations

$$r = H * \tilde{X} + noise \quad (10)$$

Where the noise is a zero-mean gaussian. We reduce the size of the Jacobian to increase speed via QR decomposition, to find the parameters Q_1 , Q_2 , and T_H .

$$H\tilde{X} = [Q_1 \quad Q_2] \begin{bmatrix} T_H \\ 0 \end{bmatrix} \quad (11)$$

We then compute the Kalman gain K with the covariance of the system P , covariance of the noise $R_N = Q_1^T R_o Q_1$, and the upper triangle matrix from the previous equation T_H .

$$K = PT_H^T(T_H P T_H^T + R_N)^{-1} \quad (12)$$

Then using the residual and our Kalman gain, we update the current state $\Delta X = Kr_n$. Finally we update our state covariance based on the Kalman gain, and fix the covariance to be symmetric.

$$P_{k+1|k+1} = (I - KH)P_{k|k} + KR_N K^T \quad (13)$$

IV. RESULTS

We tested our methods on the EuRoC dataset, and used the rpg trajectory evaluation package [5] for generating our metric.

The RMSE error from the Machine Hall 01 easy flight are shown below.

TABLE I: RMSE error

Relative Type	Value
Relative Translation	38.43
Relative Rotation	7.68
Relative Rotational Velocity	15.37

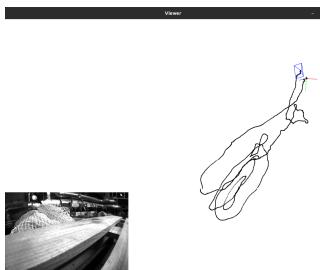


Fig. 1: Final Path graphed in Pangolin window

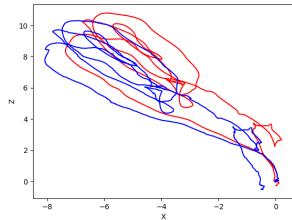


Fig. 2: Estimated vs Ground-Truth xy coordinates over full flight

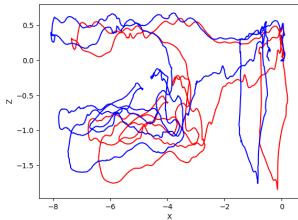


Fig. 3: Estimated vs Ground-Truth xz coordinates over full flight

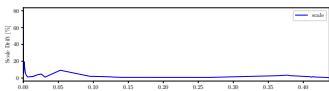


Fig. 4: Scale Error over distance. Note the beginning has extremely high percentages because of division of very small deltas in position.

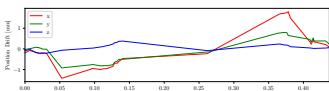


Fig. 5: Translation error over distance

V. PHASE TWO: DEEP VIO

The second phase of Project 4 is to use a deep-learning approach to combine the IMU readings and images to perform Inertial Odometry, Visual Odometry, and Visual-Inertial Odometry. With the assistance of Blender, we simulated the linear acceleration and angular velocity readings from an IMU with estimated noise and simulated images from a camera on the bottom of the drone. The camera would output images that were directly beneath the drone, rotating with the orientation of the drone, as well.

A. Dataset Generation

* Took several previously tested trajectories and turned them into paths Our dataset was generated in three parts - path generation, IMU+Gyro data from positions, image data from blender simulation.

1) *Paths*: A select set of trajectories for the quadcopter to follow were modeled numerically. These paths include easy to learn paths such as a straight line, and circular path, but also paths mentioned in [3]. We created 5 generic paths in total, with adjustable hyperparameters to differentiate the training, testing, and validation paths. Each class contained two methods: `get_position()` and `get_orientation()`, that took a time t as a parameter. By modeling each path as a function of time, we were able to pass in a uniform density list of times from $[0, 10]$ with a spacing of 0.005. This spacing modeled our IMUs 200hz operating frequency. `get_position()` generated a set of x,y,z values by evaluating the function $f(t)$, while our orientation function was designed to better model the realistic orientation of a drone while flying.

We thought one way to increase the realism in the simulated data was to estimate the orientation at any given position based on the direction vector. Instead of creating a complex dynamics model we made the assumption that the quadrotor will always be rolled and pitched proportional to the direction it is heading. To find this angle, we first found the direction vector $\dot{x}(t)$ by taking the difference between two adjacent poses.

$$\dot{x}(t) = x(t) - x(t - 1) \quad (14)$$

Using the norm of this direction, we found the quaternion between the direction vector and the "hover" vector $[0, 0, 1]^t$ (direction normal when the quadrotor is stationary). From there were used Spherical-Linear-Interpolation to generate an orientation in between the direction and hover quaternions.

$$q_d = (q_v \bigotimes (q_h)^{-1})^{0.2} \bigotimes q_h \quad (15)$$

This quaternion was used as is for the ground truth data, and translated into euler angles for use in Gyro data generation using the standard transformation

$$\begin{aligned} \text{roll } (\phi) &= \tan^{-1} \left(\frac{2(wx + yz)}{1 - 2(x^2 + y^2)} \right) \\ \text{pitch } (\theta) &= \sin^{-1} (2(wy - zx)) \\ \text{yaw } (\psi) &= \tan^{-1} \left(\frac{2(wz + xy)}{1 - 2(y^2 + z^2)} \right) \end{aligned}$$

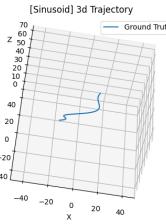


Fig. 6: Example path generated from our Sinusoid function

2) *IMU+Gyro Data*: Using the array of position and euler angles, we could generate IMU and Gyro data. We utilized an existing 3rd party library for doing this called OysterSim [1]. Using this library, we generated data with realistic noise added. We exported this data along with an image filename under each Train, Test, and Val directories called "input_data.csv".

3) *Image Data*: We utilized blender along with a high-quality texture map to generate images at each location given by the first step. We move the camera's location using the x,y,z offsets, and rotate the camera using the provided Euler angles. We choose an image of a child's cityscape rug as it has many detectable features for our vision model to detect. We used AI-based up-scaling to increase the fidelity of the image, and resized it again in blender to ensure we were able to move without detecting the edges, and projected the image onto a plane centered at (0,0,0)



Fig. 7: Image Rendering

B. Inertial Odometry (IO)

We used a very simple architecture for to perform IO, adapted from previous works in the class. The architecture utilizes two Long Short-Term Memory layers as well as two Fully Connected layers separated by a PReLU function. The inputs to the model are (batch size x 10 x 6), 10 for the number of IMU readings between updates and 6 for the number of states provided by the IMU. We trained this model for 200 epochs using the Adam optimizer with a learning rate of 0.001 and a batch size of 32.



Fig. 8: Deep-IO Architecture

C. Visual Odometry (VO)

VO only uses the images provided by the onboard camera to perform odometry, working very well in environments containing many high-quality image features.

To extract and associate the features in each image inputted into the model, we utilize the feature descriptor used in GMFlowNet, FNet. We used this pre-trained model as it was proven to be robust and has publicly accessible pre-trained weights we were able to download onto our machines. We freeze the weights used in FNet to prevent the weights from being altered during training. We trained our model for 50 epochs using the Adam optimizer with a learning rate of 0.001 and a batch size of 16.

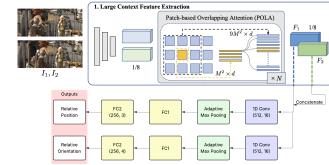


Fig. 9: Deep-VO Architecture

D. Visual-Inertial Odometry (VIO)

We combine both the VO and IO models to create our Deep-VIO model. The architectures for each VO and IO models remain mostly unchanged, with the exception of removing final steps from each model. Instead, we concatenate the outputs of each model and send the concatenation through an additional few layers to further improve the odometry estimation. The hyperparameters used to train the VIO model were identical to VO's hyperparameters.

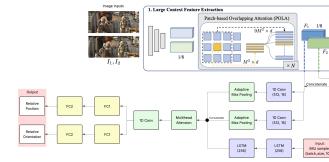


Fig. 10: Deep-VIO Architecture

E. Experiments and Results

We tested each model against a path with the same shape, but different hyperparameters to see how well it generalizes to paths it had not seen before. We do this by summatting the relative pose from the model at each timestep from $[0, t_f]$

We then graphed each w.r.t the ground truth path it was given. Here are the results of the estimation in 3D. Below in the appendix are the breakdowns per axis on each path.

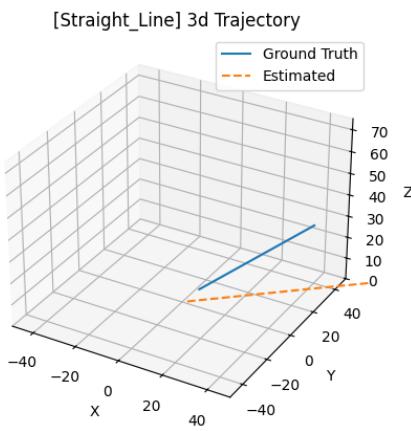


Fig. 11: Straight line path

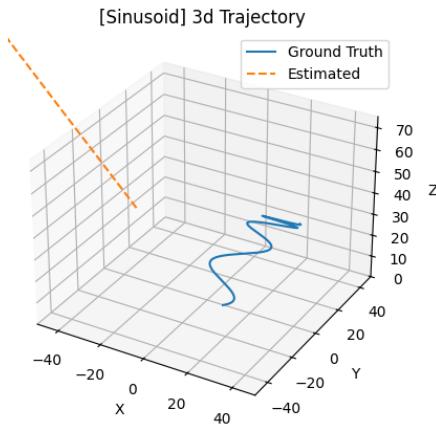


Fig. 12: Straight line path with sinusoidal variations in each axis

VI. RESEARCH PROBLEMS WITHIN THE FIELD

* Riley - Estimating Heading of simulated quadcopter data.

VII. CONCLUSION

In this phase, we implemented a method for performing state estimation via the combination of two ordinary sensors, a camera and an IMU. Although the IMU is great for short term - high frequency values, its propensity to drift makes it unstable long term. And while a camera provides robust features usable for long-term tracking, it struggles with sudden movements where features cannot be easily tracked. By fusing the two, you attain stable positional state estimation in real time, and can correct for the biases generated during a flight.

REFERENCES

- [1] Xiaomin Lin, Nitesh Jha, Mayank Joshi, Nare Karapetyan, Yiannis Aloimonos, and Miao Yu. Oystersim: Underwater simulation for enhancing

oyster reef monitoring. In *OCEANS 2022, Hampton Roads*, page 1–6. IEEE, October 2022.

- [2] Anastasios I. Mourikis and Stergios I. Roumeliotis. A multi-state constraint kalman filter for vision-aided inertial navigation. In *Proceedings 2007 IEEE International Conference on Robotics and Automation*, pages 3565–3572, 2007.
- [3] Nitin J. Sanket, Chahat Deep Singh, Cornelia Fermüller, and Yiannis Aloimonos. Prgflow: Benchmarking swap-aware unified deep visual inertial odometry, 2020.
- [4] Ke Sun, Kartik Mohta, Bernd Pfommer, Michael Watterson, Sikang Liu, Yash Mulgaonkar, Camillo J. Taylor, and Vijay Kumar. Robust stereo visual inertial odometry for fast autonomous flight, 2018.
- [5] Zichao Zhang and Davide Scaramuzza. A tutorial on quantitative trajectory evaluation for visual(-inertial) odometry. In *IEEE/RSJ Int. Conf. Intell. Robot. Syst. (IROS)*, 2018.

VIII. APPENDIX A: TESTING RESULTS PER AXIS

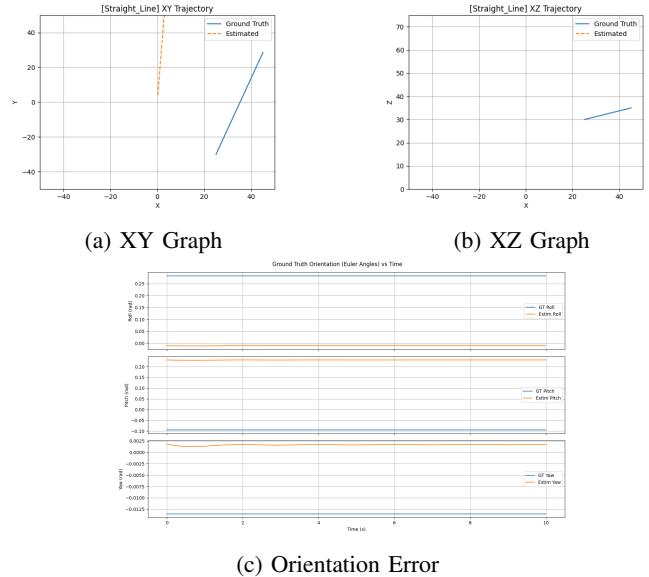
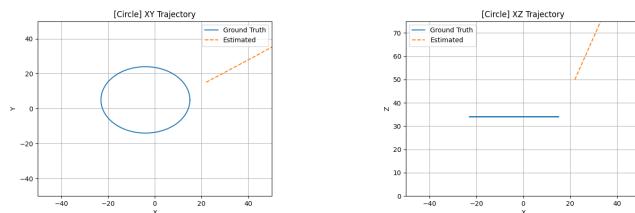
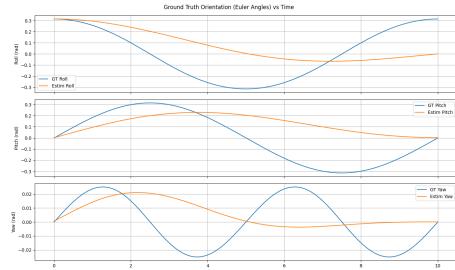


Fig. 13: Straight Line Path



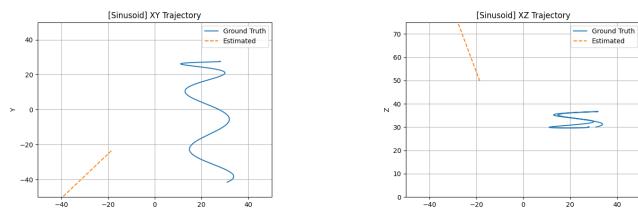
(a) XY Graph

(b) XZ Graph



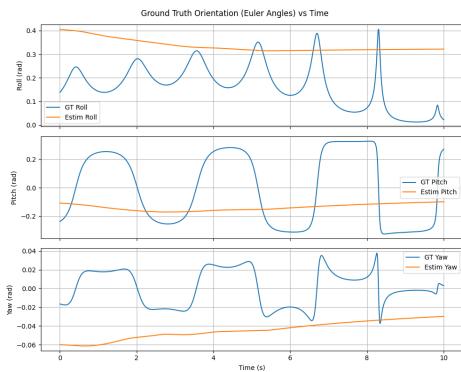
(c) Orientation Error

Fig. 14: Circle Path



(a) XY Graph

(b) XZ Graph



(c) Orientation Error

Fig. 15: Sinusoid Path