

# **Visual Inertial Odometry using Focal Plane Binary Features (BIT-VIO)**

Matthew Lisondra <sup>1,\*</sup>, Junseo Kim <sup>1,\*</sup>, Riku Murai <sup>2</sup>, Kourosh Zareinia <sup>1</sup> and Sajad Saeedi <sup>1</sup>

<sup>1</sup>Toronto Metropolitan University (TMU).

<sup>2</sup>Imperial College London, Department of Computing.

\*Both authors contributed equally to this research at TMU.



# **Table of Contents**

I. Background

II. Problem Statement

III. Critical Design

IV. Results and Progress

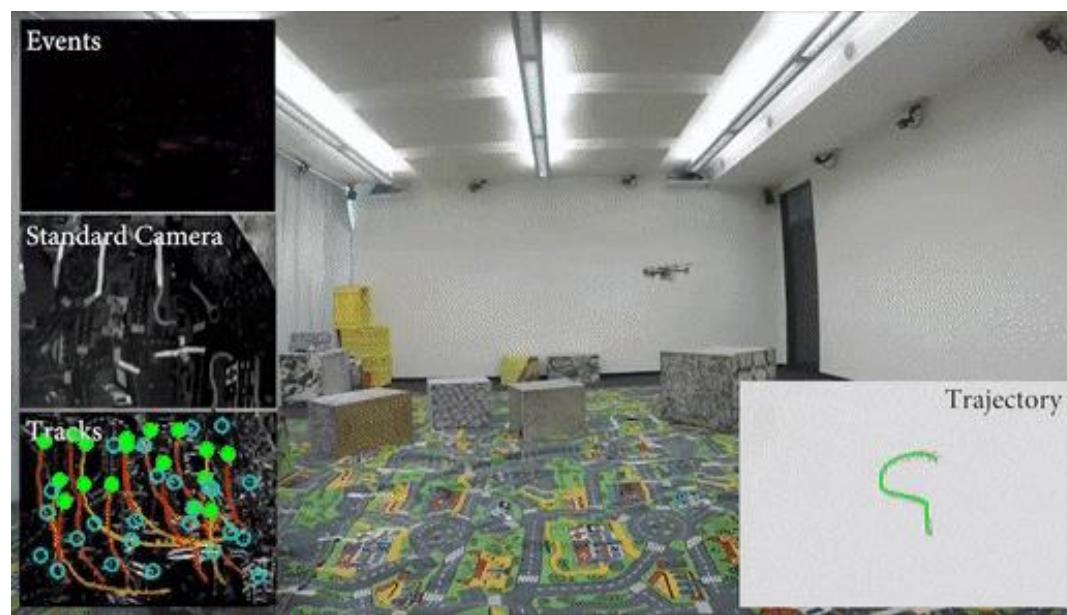
V. Conclusion and Next Steps

# I. Background

# I. Background – Introduction

Camera technology has advanced and developed since its inception since the 1800s to 1900s and is still continuing to do so today. Camera technology, as we know it, is the basis of so many different fields, responsible for so many different applications.

[1]



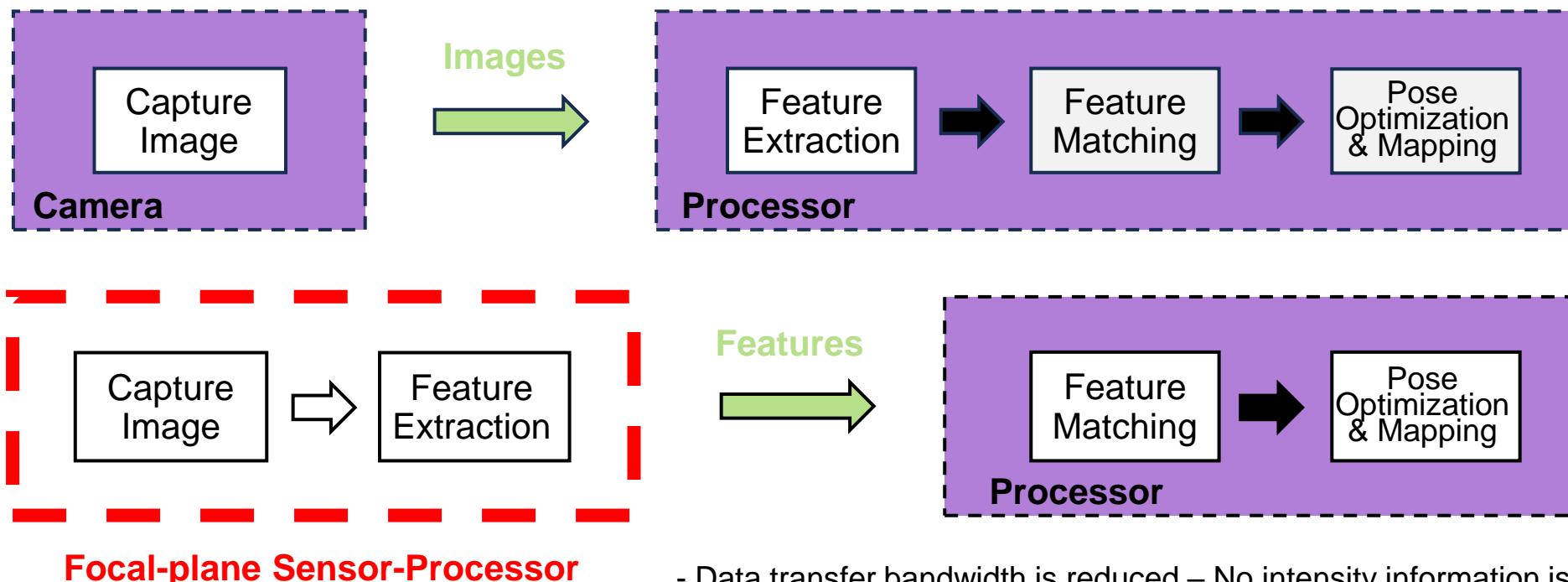
There are exciting types that are still being properly and fully developed, one of which is dealing with the event camera system.

Of particular interest to us today is not the event camera, but the FPSP camera technology, its advantages, and its potential use in robotic localization and mapping.

# FPSP vs. Conventional Camera

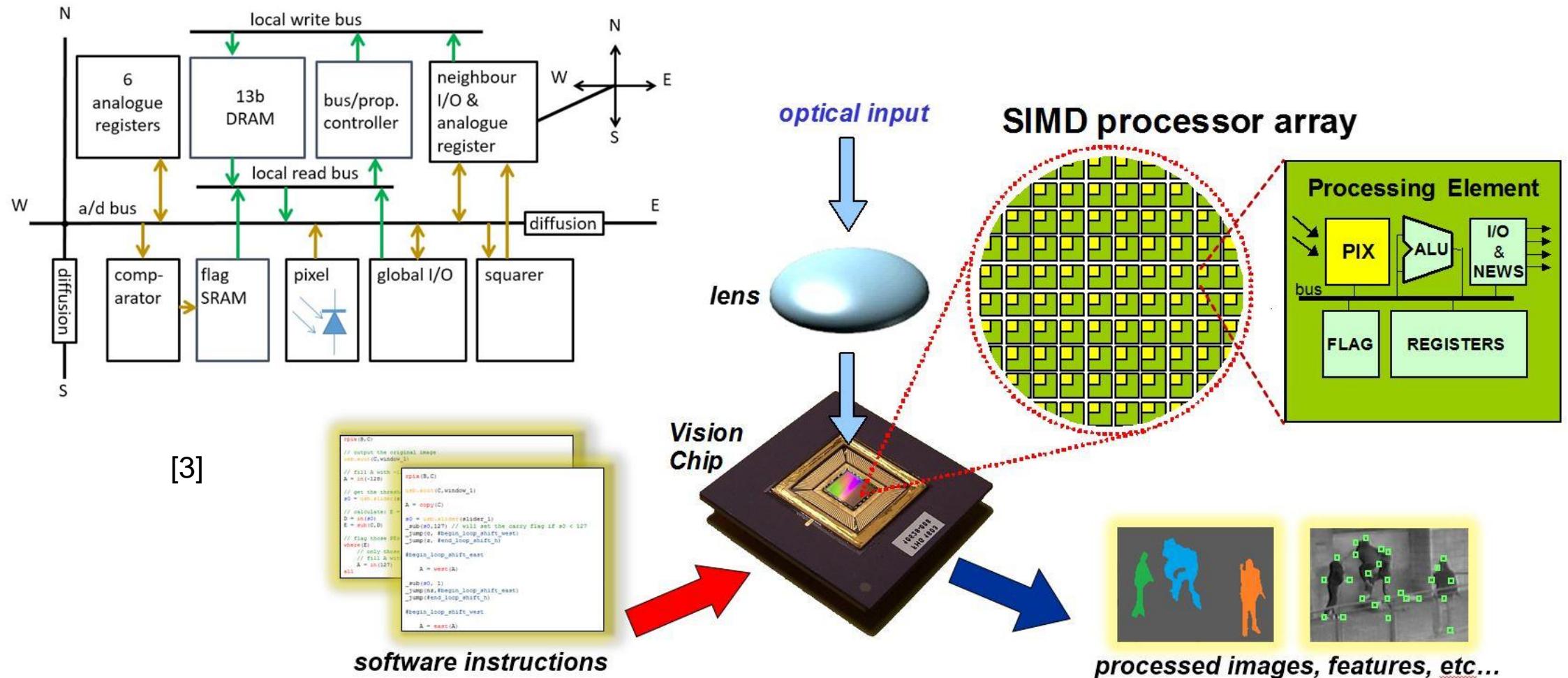
## Moving the Computation Towards the Focal-Plane

[1]



- Data transfer bandwidth is reduced – No intensity information is transferred
- Reduces computation costs on the processor side
- Only meaningful data is transmitted

# I. Background – FPSP SCAMP-5 Architecture



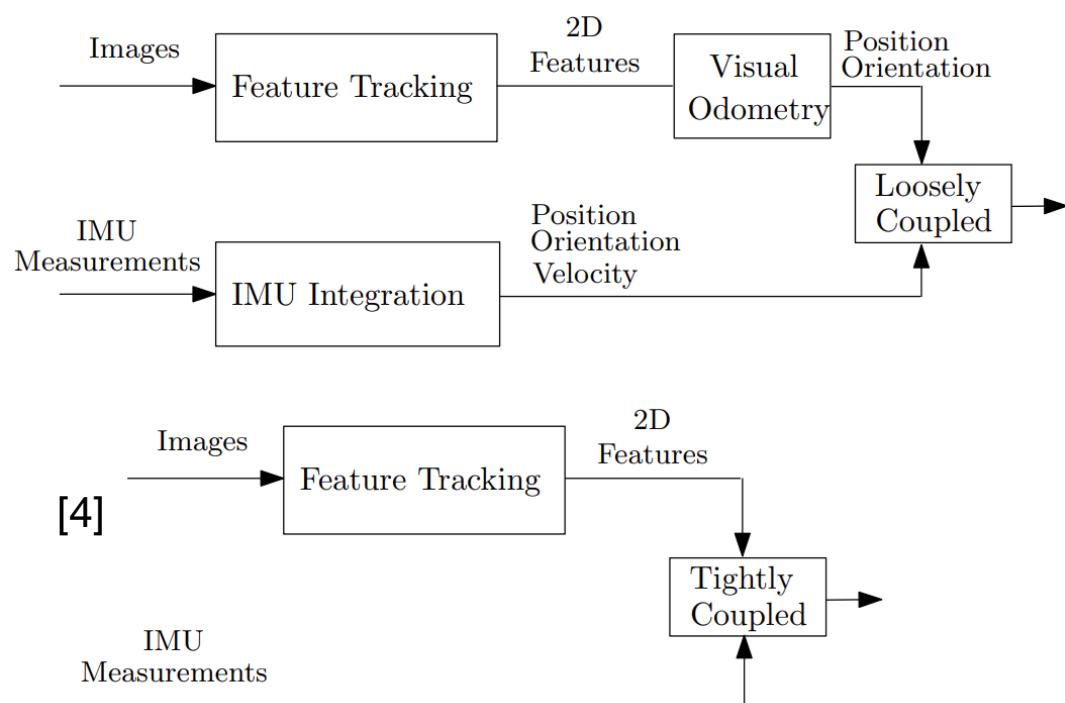
The FPSP SCAMP-5 camera technology is designed to provide lower latency with lower power consumption.

# I. Background – Visual Inertial Odometry

“Visual-Inertial odometry (VIO) is the process of estimating the state (pose and velocity) of an agent (e.g., an aerial robot) by using only the input of one or more cameras plus one or more Inertial Measurement Units (IMUs) attached to it.”

VIO is the only viable alternative to GPS and lidar-based odometry to achieve accurate state estimation.

[4]



[4]

Filtering	Fixed-lag Smoothing	Full smoothing
Only updates the most recent states • (e.g., extended Kalman filter)	Optimizes window of states • Marginalization • Nonlinear least squares optimization	Optimize all states • Nonlinear Least squares optimization
✗1 Linearization	✓ Re-Linearize	✓ Re-Linearize
✗Accumulation of linearization errors	✗Accumulation of linearization errors	✓ Sparse Matrices
✗Gaussian approximation of marginalized states	✗Gaussian approximation of marginalized states	✓ Highest Accuracy
✓ Fastest	✓ Fast	✗Slow (but fast with GTSAM)

Both VO and VIO would benefit greatly from higher framerates and lower power consumption vision systems. We need FPSPs!

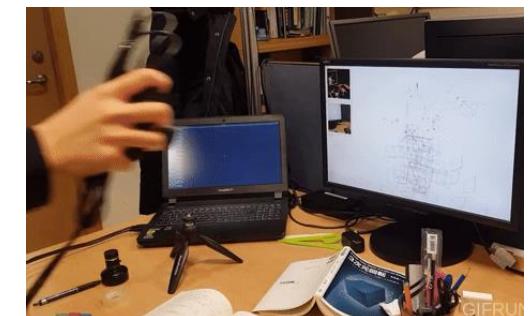
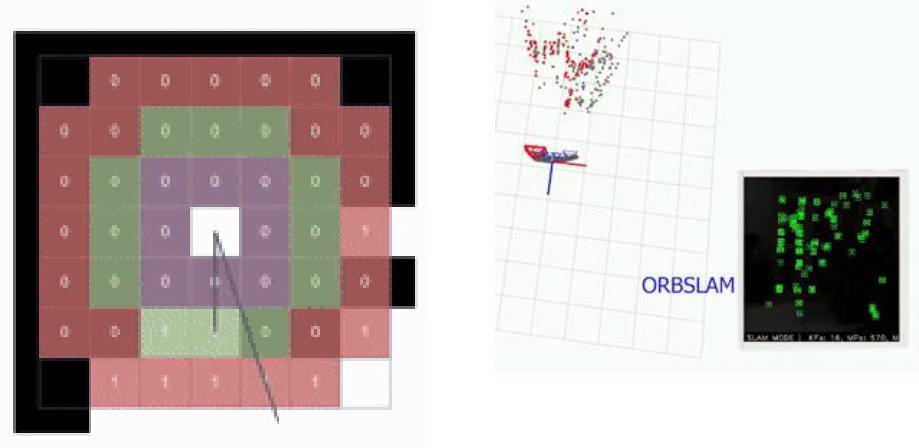
# I. Background – Prior Work 6-DOF VO (BIT-VO) on an FPSP

A prior work in 2020 titled BIT-VO: Visual Odometry at 300 FPS by Riku Murai, Sajad Saeedi and Paul H. J. Kelly makes full use of the Focal-plane Sensor-processor camera technology.

[2]

The contributions of this previous work consisted of:

1. An efficient Visual Odometry which operates at over 300 FPS. Using no intensity information, the proposed method can accurately track the pose, even under difficult situations where the state-of-the-art monocular SLAM fails.
2. A robust feature matching scheme, which uses the novel binary-edge based descriptor. Using a small, 44-bit descriptor, the system can track the noisy feature data computed on the focal plane in the SCAMP-5 FPSP image sensor itself.
3. Extensive evaluation of the system against measurements from a motion capture system, including difficult scenarios such as violently shaking the device 4-5 times in a second.



So, work has been done on an all 6-DOF algorithm with the FPSP camera technology, but what about more complex estimation, those of which are the more accurate 6-DOF algorithms SLAM, or VIO?

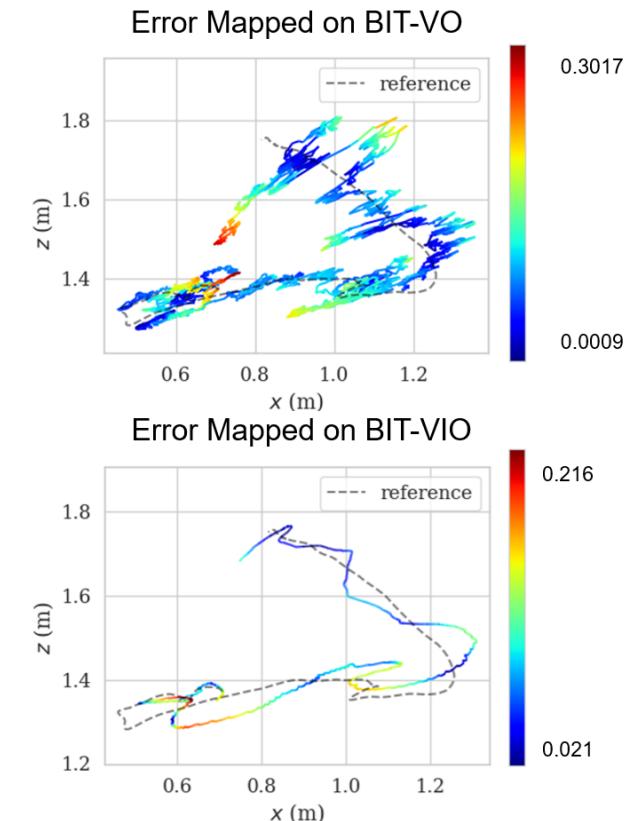
## II. Problem Statement

Does there exist a 6-DOF SLAM or VIO algorithm using the advantages of the FPSP camera technology?

Extending on the previous work BIT-VO, we present BIT-VIO, the first 6-DOF VIO algorithm to utilize the advantages of the FPSP for vision-IMU-fused state estimation.

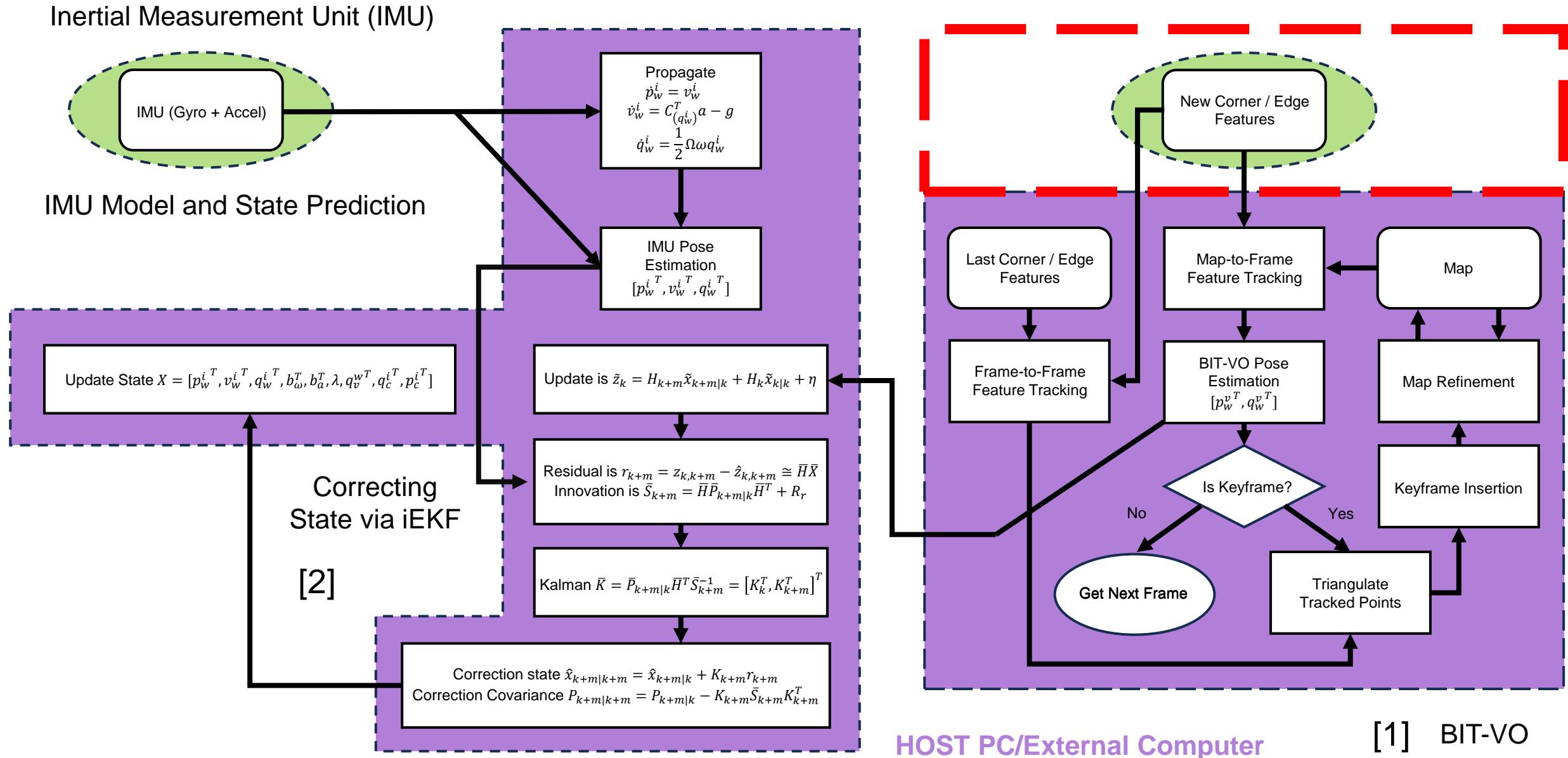
The contributions of this work are:

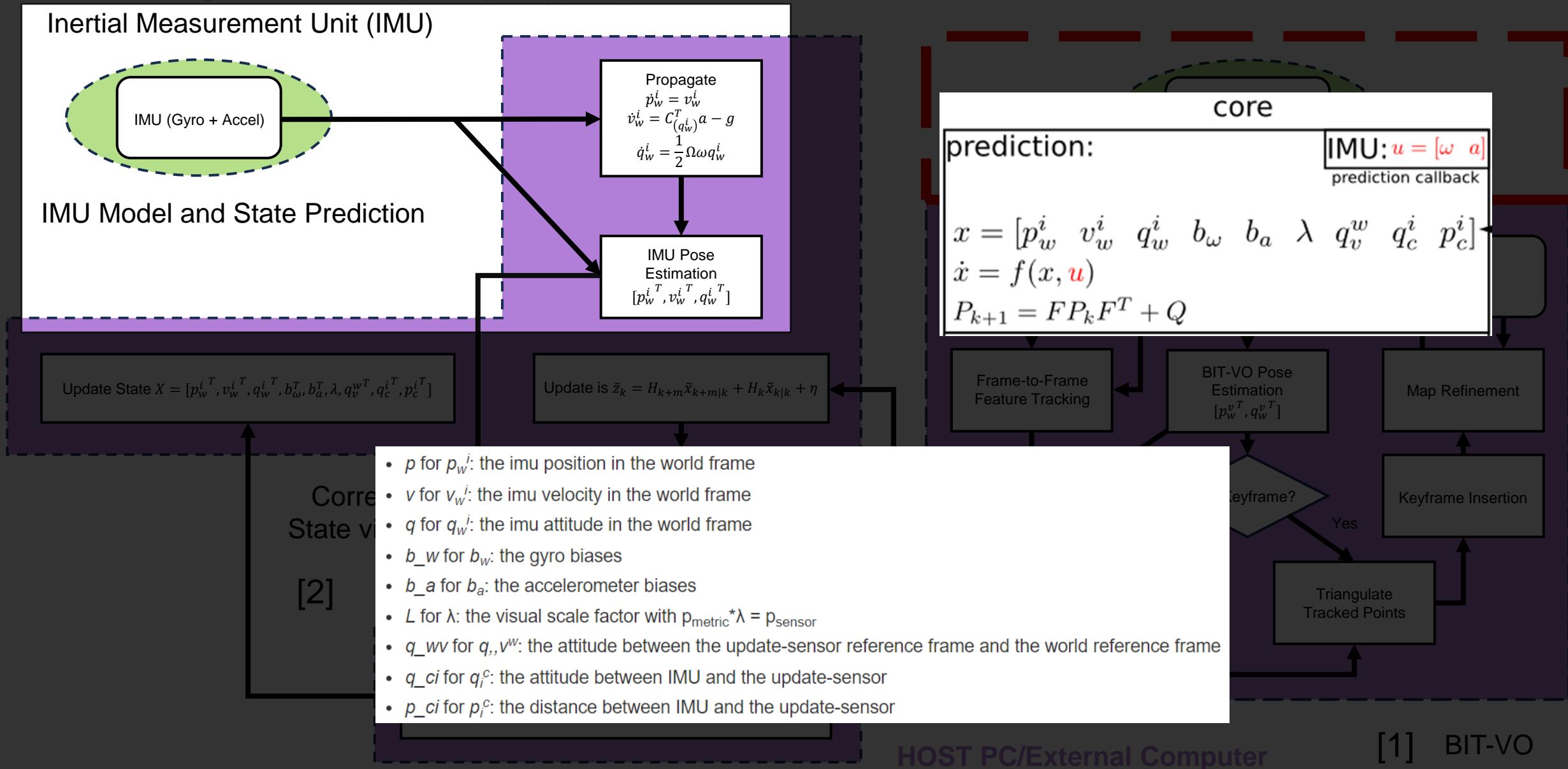
- (I) Efficient Visual Inertial Odometry operating and correcting by loosely-coupled sensor-fusion iterated Extended Kalman Filter (iEKF) at 300 FPS with an IMU at 400 Hz.
- (II) Uncertainty propagation for BIT-VO's pose as it is based on binary-edge-based descriptor extraction, 2D to 3D re-projection, and transform from pose to pose.
- (III) Comparison with BIT-VO and ground-truth, tested on a moving robotic system against a motion capture system.



This work is the first steps toward a 6-DOF SLAM or VIO algorithm using the advantages of the FPSP camera technology.

# III. Critical Design





### Inertial Measurement Unit (IMU)



### IMU Model and State Prediction

$$\text{Propagate}$$

$$\dot{p}_w^i = v_w^i$$

$$\dot{v}_w^i = C_{(q_w^i)}^T a - g$$

$$\dot{q}_w^i = \frac{1}{2} \Omega \omega q_w^i$$

$$\text{IMU Pose Estimation}$$

$$[p_w^{i^T}, v_w^{i^T}, q_w^{i^T}]$$

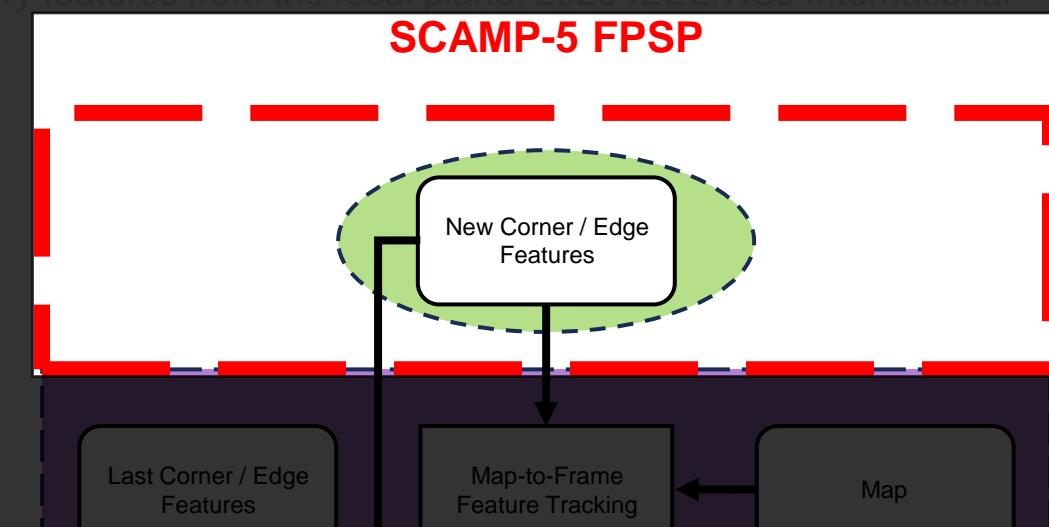
Update State  $X = [p_w^{i^T}, v_w^{i^T}, q_w^{i^T}, b_a^T, b_a^T, \lambda, q_v^{w^T}, q_c^{i^T}, p_c^{i^T}]$

Update is  $\hat{z}_k = H_{k+m} \hat{x}_k$

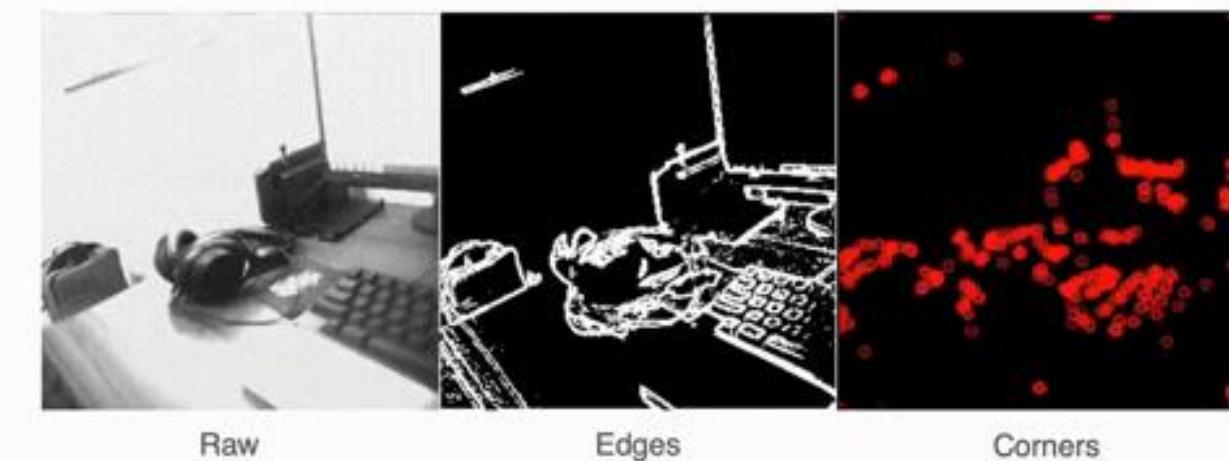
Feature extractions are performed on the FPSP SCAMP-5, while feature tracking and VO operates on the host device which is, for example, a consumer grade laptop. The system operates only on the binary edge image and corner coordinates; thus, **no pixel intensity information is ever transferred**.

Correction state  $\hat{x}_{k+m|k+m} = \hat{x}_{k+m|k}$   
 Correction Covariance  $P_{k+m|k+m} = P_{k+m|k}$

### SCAMP-5 FPSP



### Focal-plane Sensor-processor



## Inertial Measurement Unit (IMU)

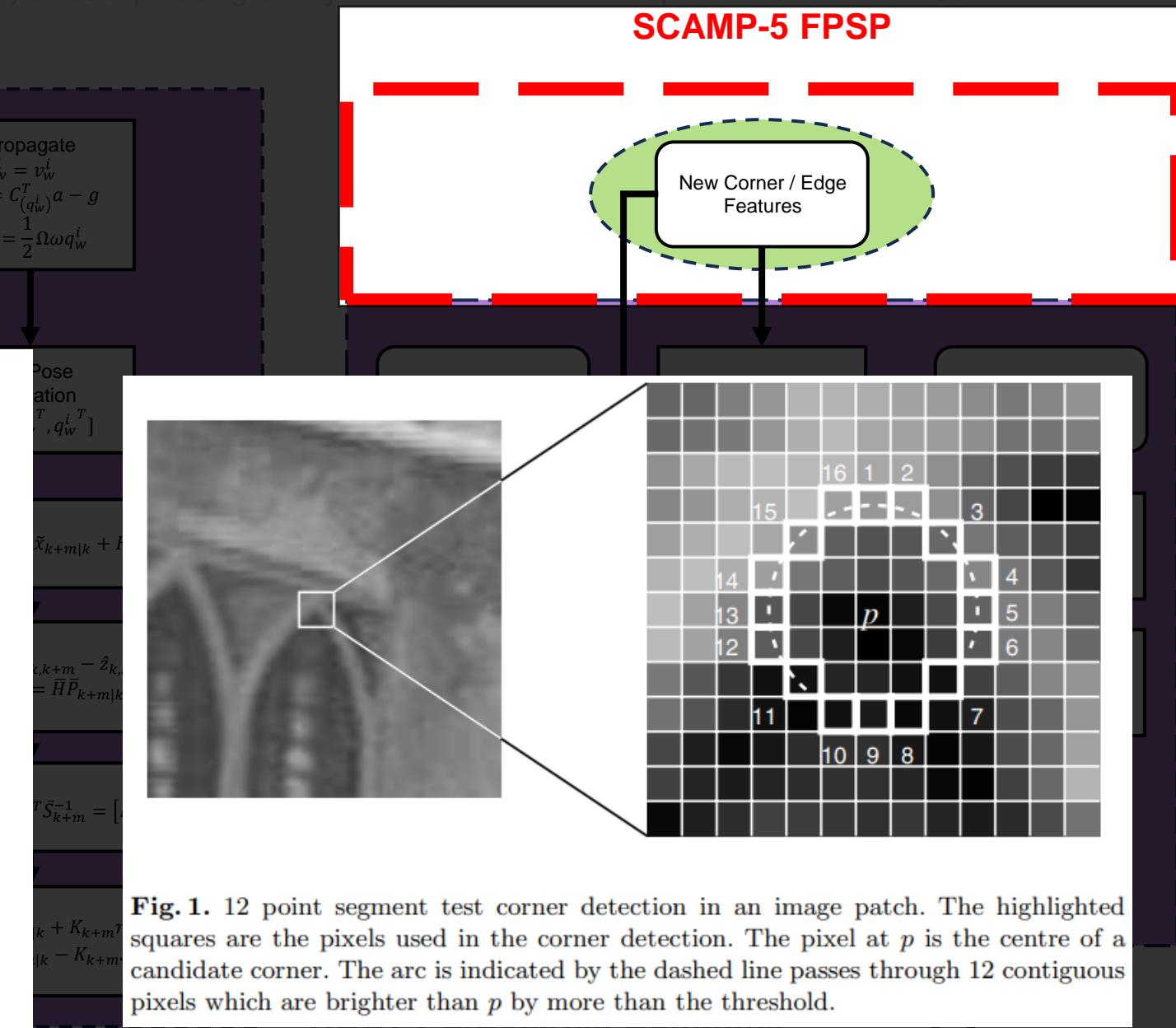


## IMU Model and State Prediction

Corner and edge features are computed on the FPSP, and it operates at a high frame-rate of **330 FPS**. FAST Keypoint Detector is used for the corner detection.

Unfortunately, performing repeatable suppression techniques such as non-maximal suppression is difficult due to the noisy analog computation on SCAMP-5.

The noisy computation leads to not only incorrect inequality comparisons but also to incorrect computation of the compared values.



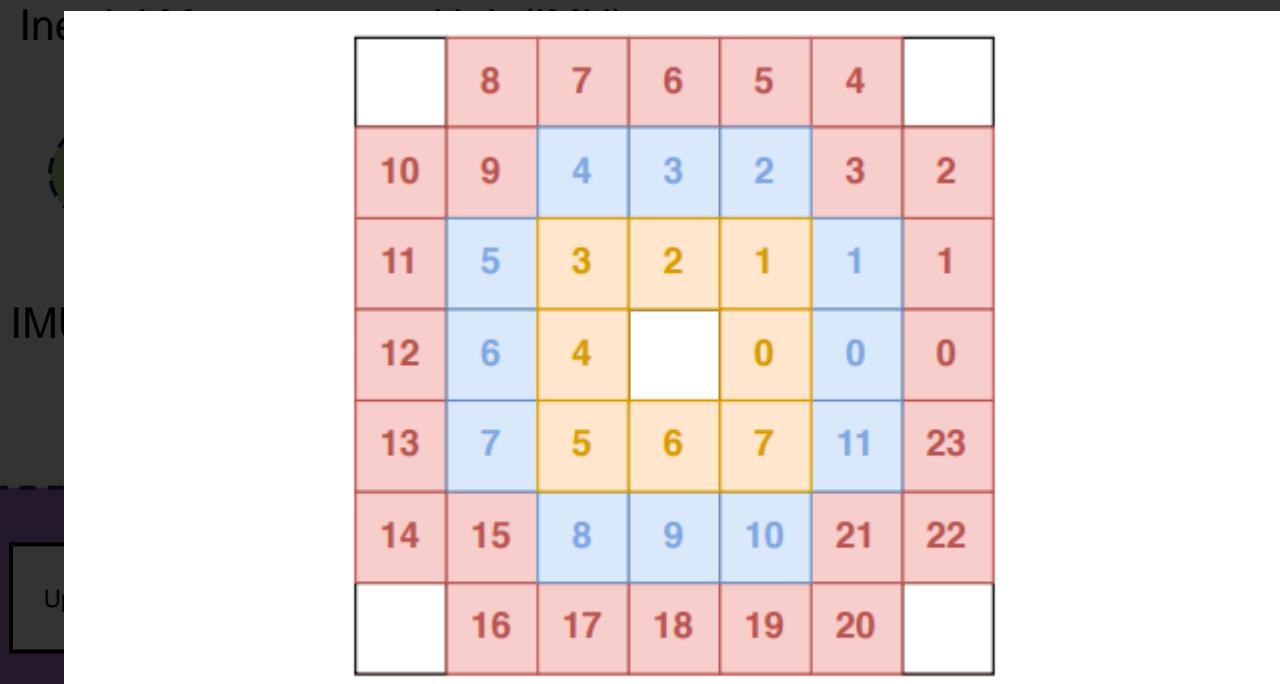


Fig. 4: Descriptor sampling pattern. Different colours denote a different ring, and indices correspond to the bit index.

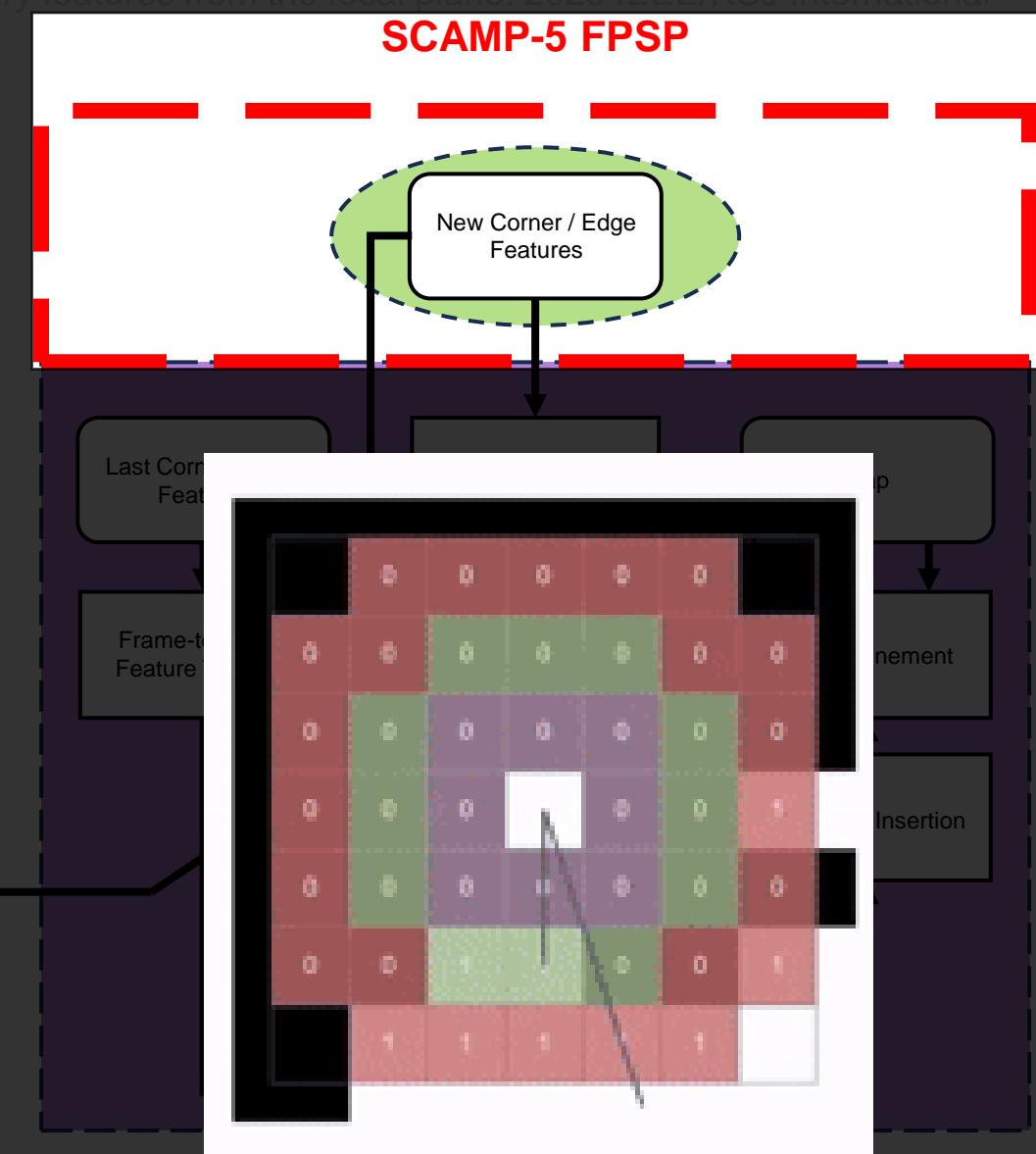
State via iEKF

$$\theta = \tan^{-1} \frac{\sum_{x,y} y G(x, y)}{\sum_{x,y} x G(x, y)} \quad (1)$$

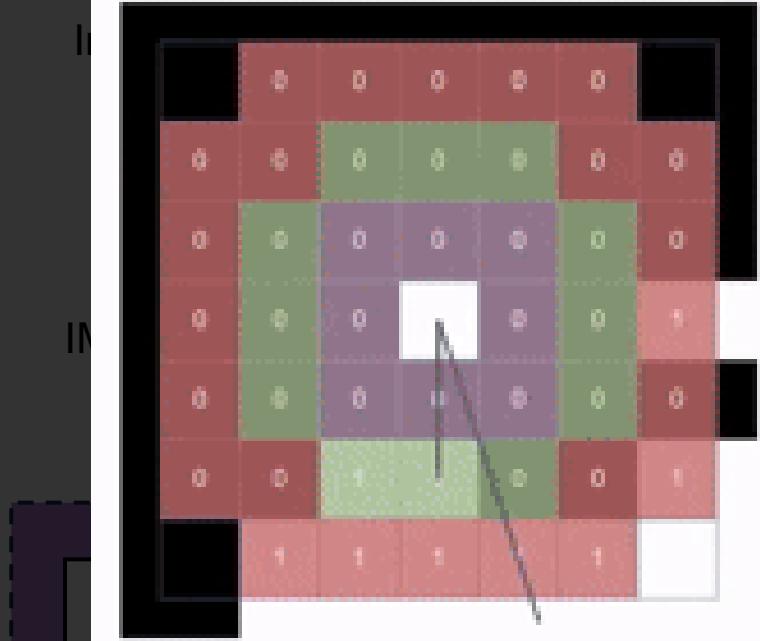
$$\theta = \tan^{-1} \frac{\sum_{x,y} y B(x, y)}{\sum_{x,y} x B(x, y)} \quad (2)$$

$\bar{H}^T \bar{S}_{k+m}^{-1} = [K_k^T, K_{k+m}^T]^T$

$\bar{m}|k + K_{k+m} r_{k+m}$   
 $\bar{m}|k - K_{k+m} \bar{S}_{k+m} K_{k+m}^T$



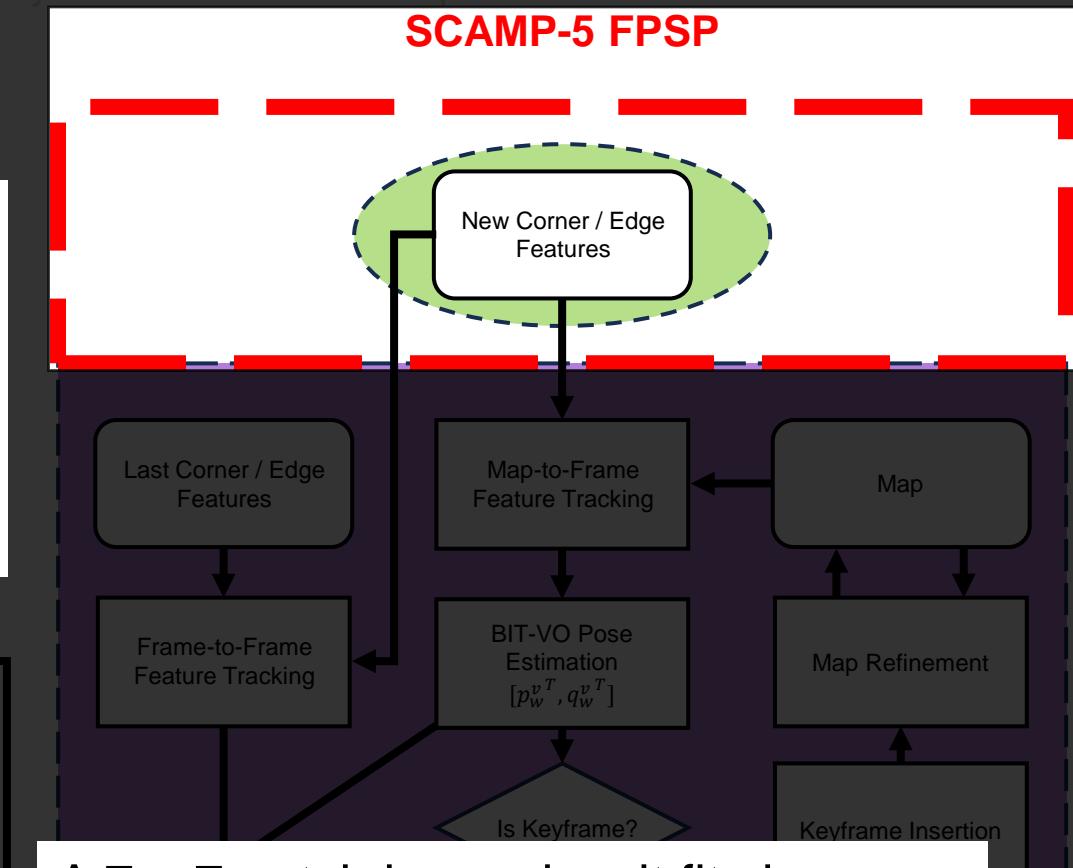
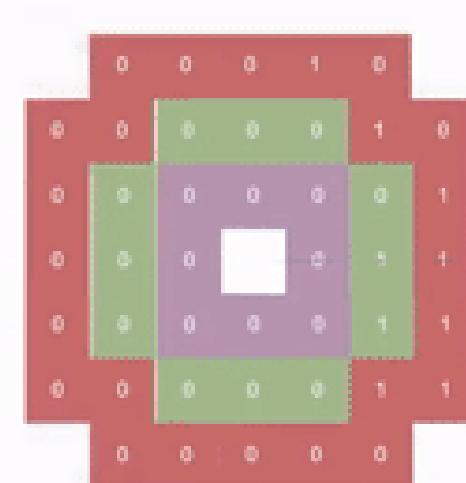
The rotation invariance is achieved by bit-rotations of the rings independently, based on the orientation  $\theta$ .



$$\theta = \tan^{-1} \frac{\sum_{x,y} yG(x,y)}{\sum_{x,y} xG(x,y)} \quad (1)$$

$$\theta = \tan^{-1} \frac{\sum_{x,y} yB(x,y)}{\sum_{x,y} xB(x,y)} \quad (2)$$

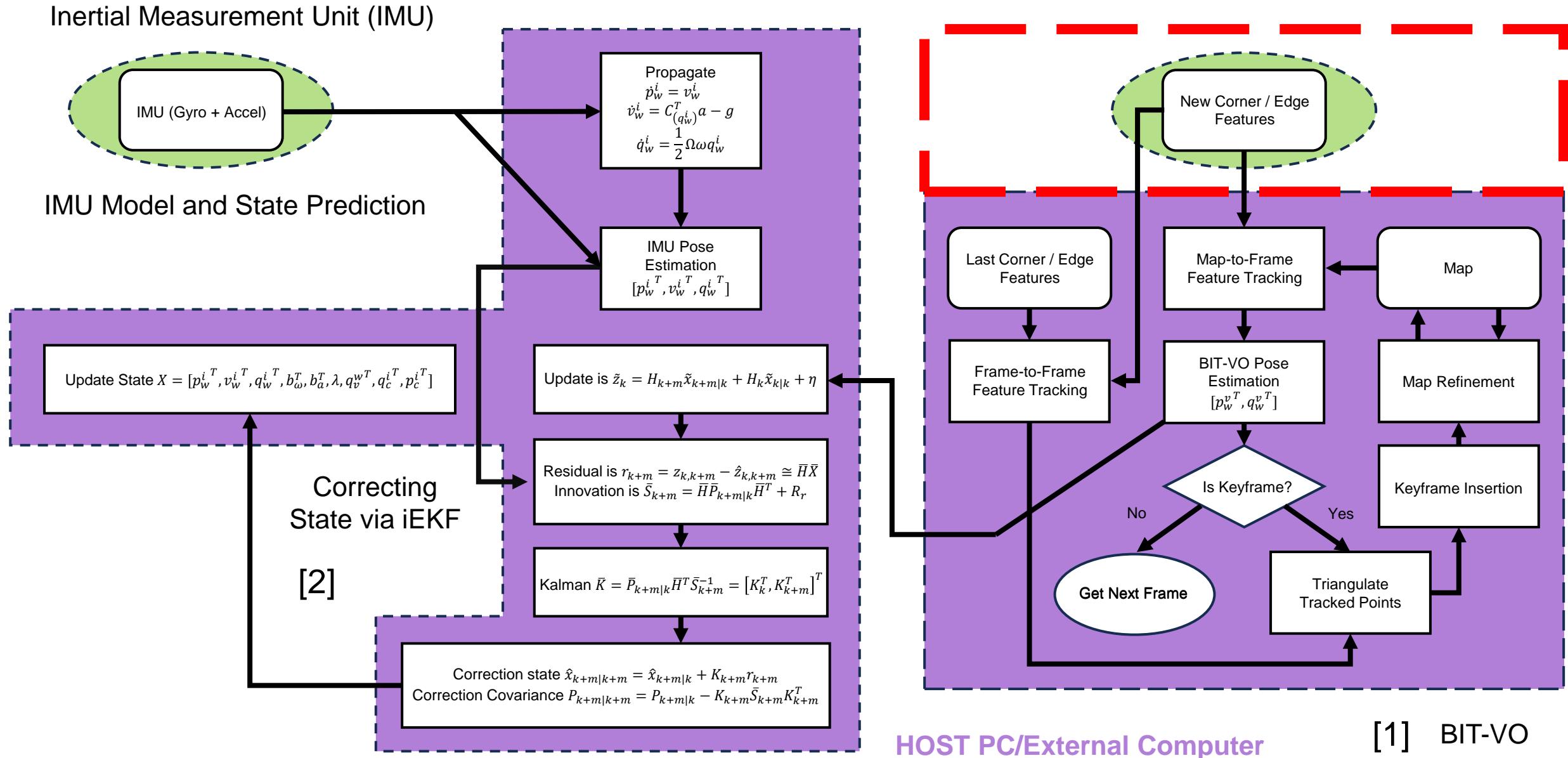
Update is  $\tilde{z}_k = H_{k+m}\tilde{x}_{k+m|k} + H_k\tilde{x}_{k|k} + \eta$



A  $7 \times 7$  patch is used as it fits in a single **64-bit unsigned integer**. This allows the patch data to be converted into rings efficiently using bitwise manipulation only.

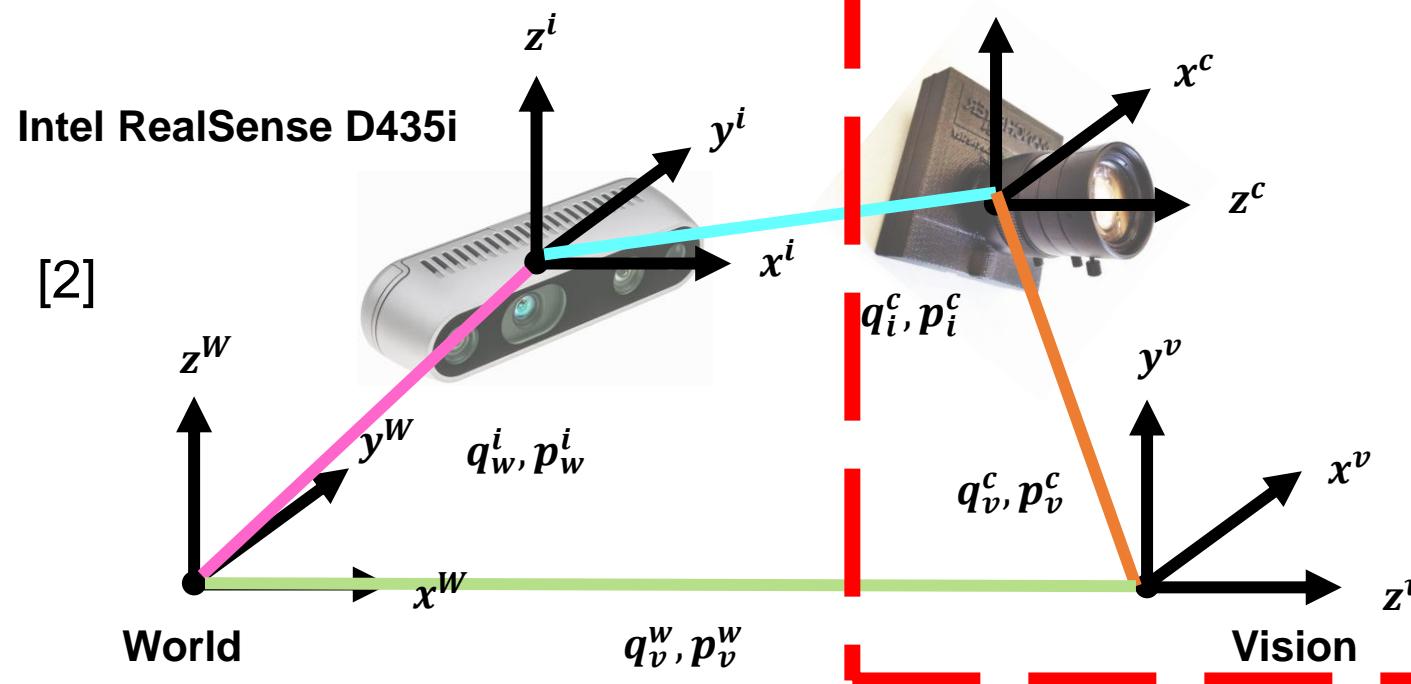
HOST PC/External Computer

[1] BIT-VO



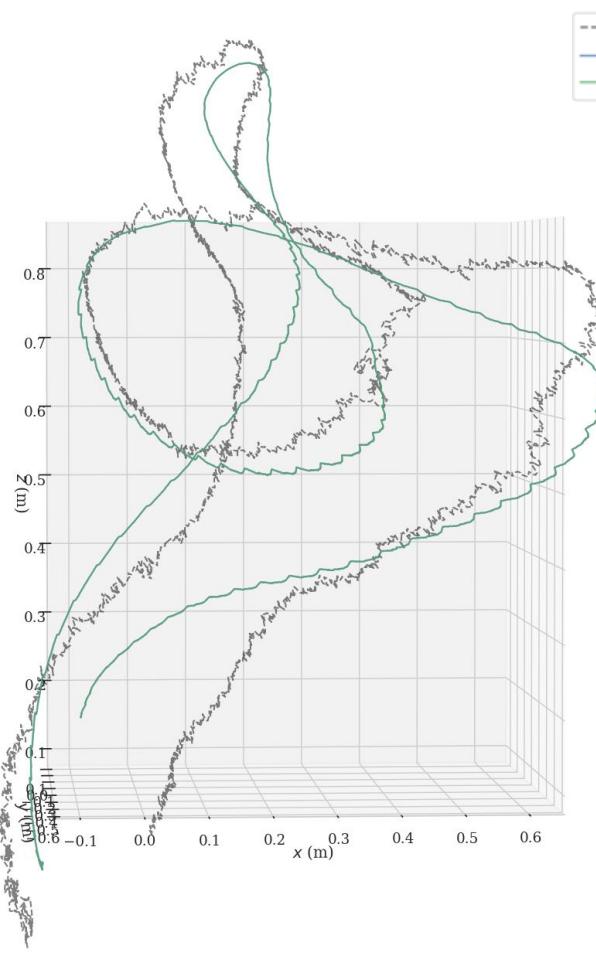
# IMU and SCAMP-5 FPSP Frames

IMU and SCAMP-5 FPSP frames. Intel RealSense D435i IMU is used in this work, and SCAMP-5 is the FPSP used as a camera sensor. Four coordinate frames, with two being a part of SCAMP-5 FPSP (camera and vision frames).



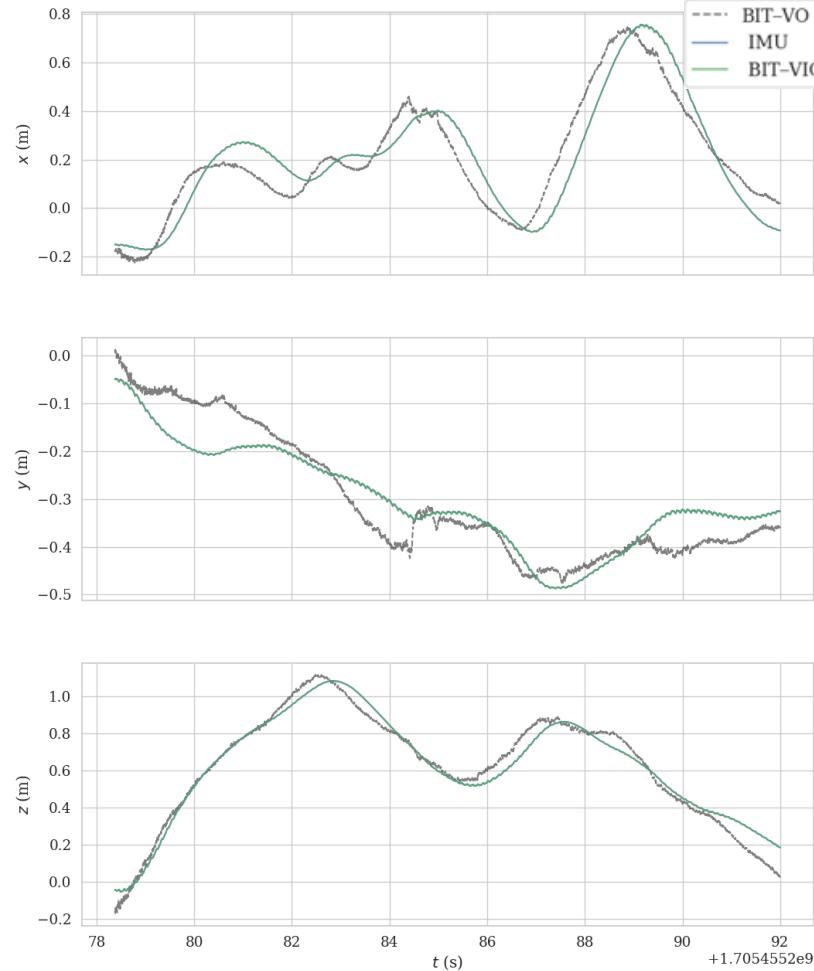
# IV. Results and Progress

# Compared just BIT-VO with BIT-VIO and IMU

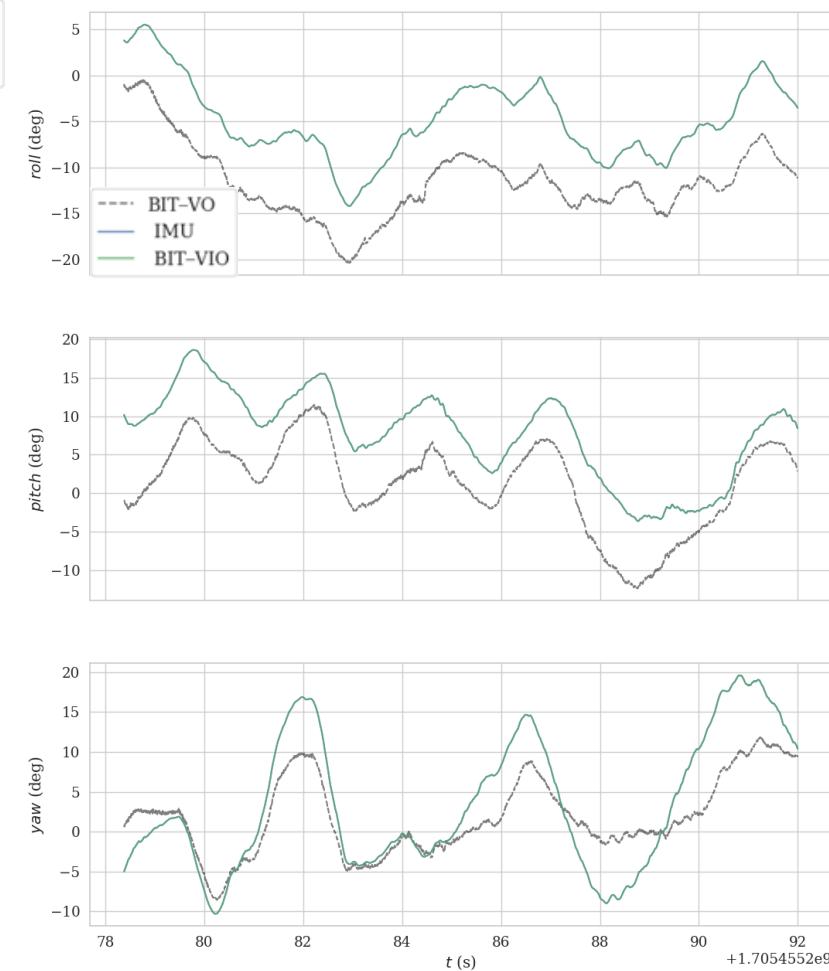


Traj 9/20 ~ 3.4 m

Full 6-DOF long trajs., fast motions  
Is fully achievable as we solved for both  
Now time to run experiments with GT.



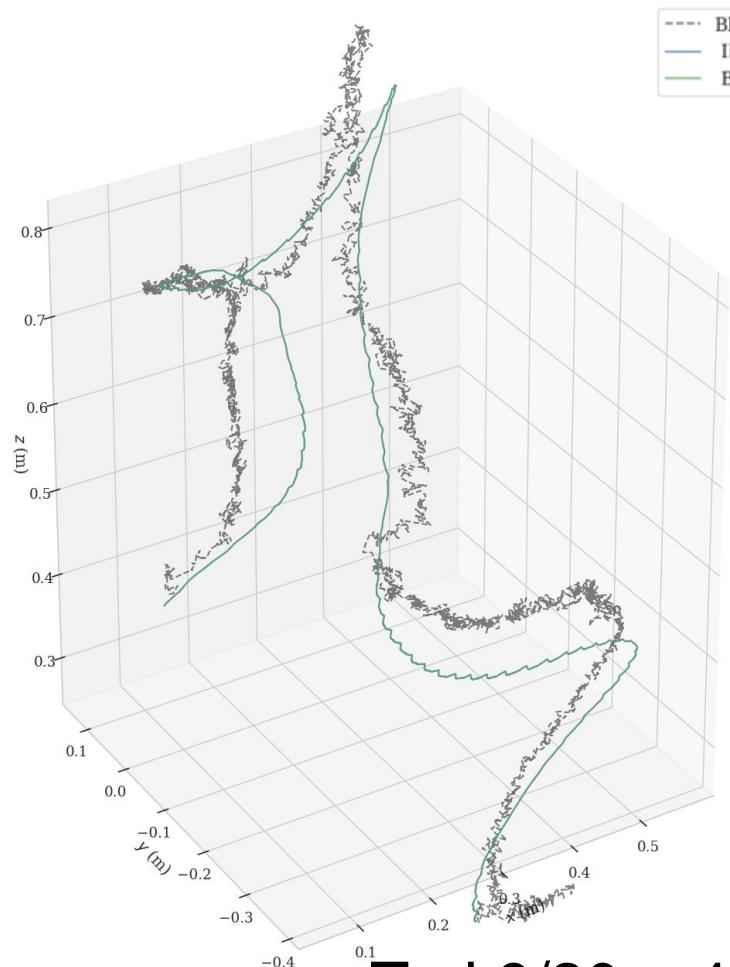
XYZ strong



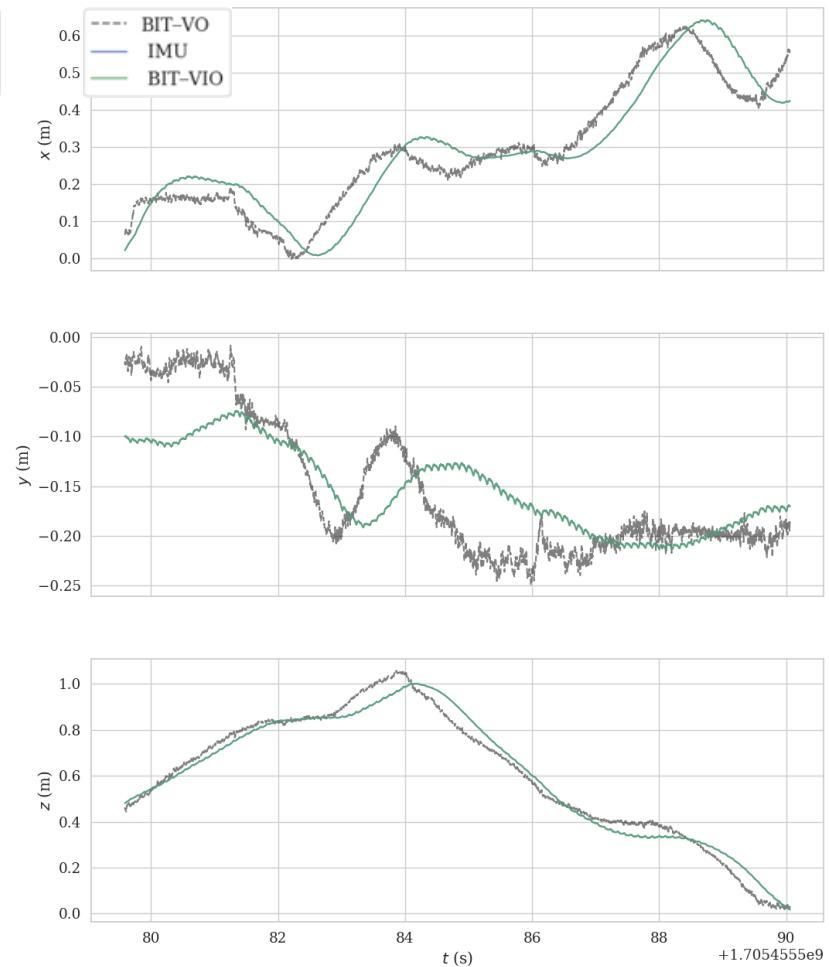
RPY trends (~5deg)

Our BIT-VIO silences the high frequency noise of prior BIT-VO.

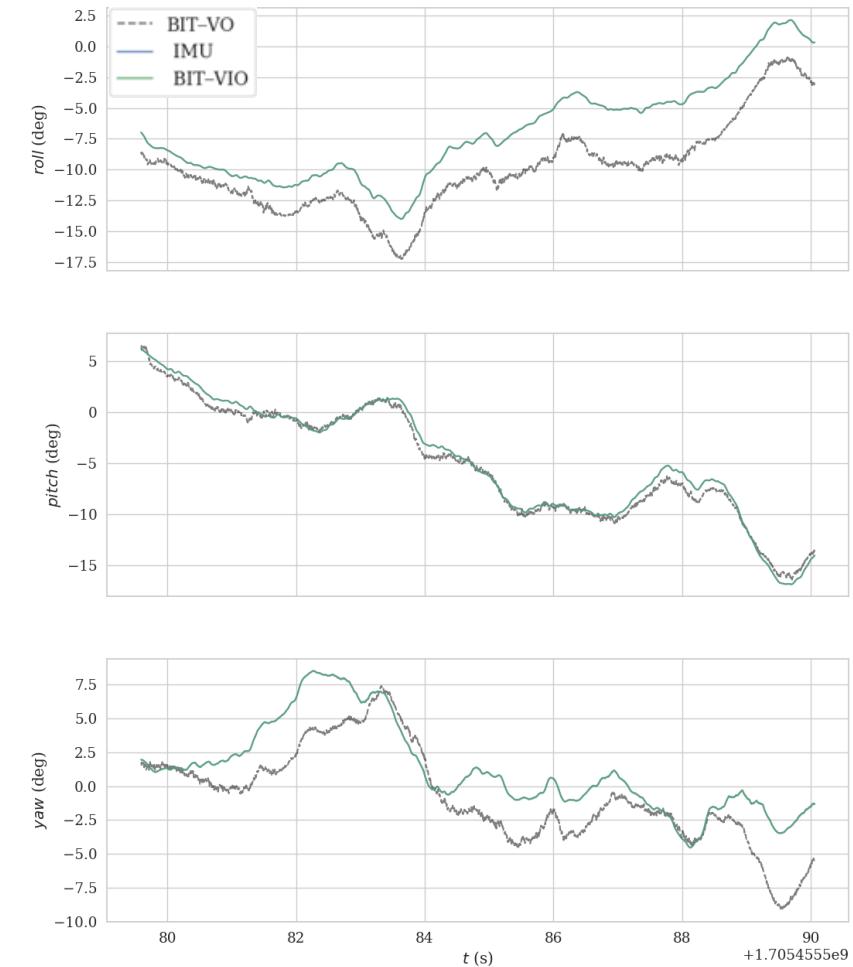
# Compared just BIT-VO with BIT-VIO and IMU



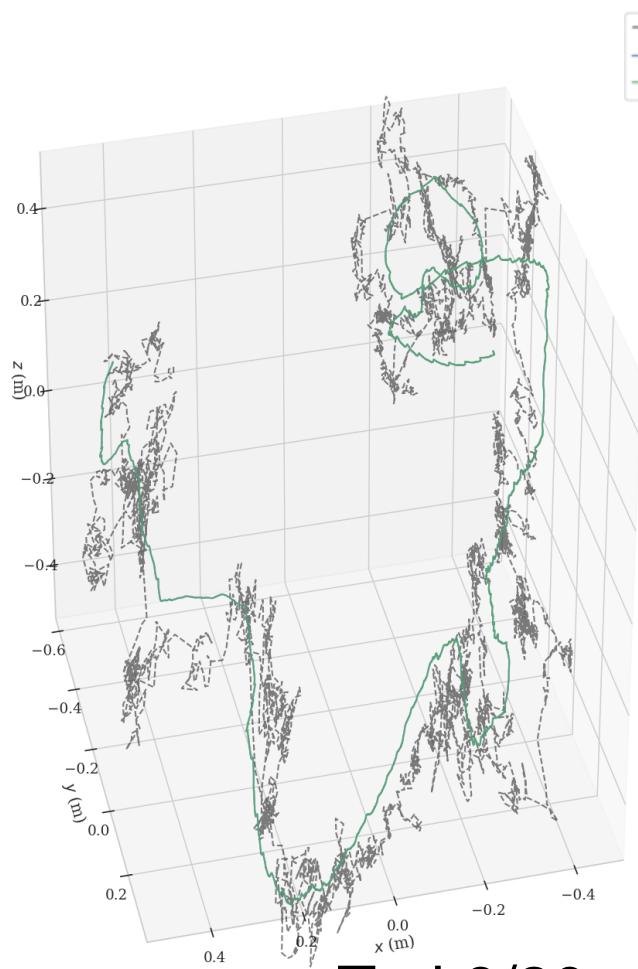
Full 6-DOF long trajs., fast motions  
Is fully achievable as we solved for both  
Now time to run experiments with GT.



Our BIT-VIO silences the high frequency noise of prior BIT-VO.

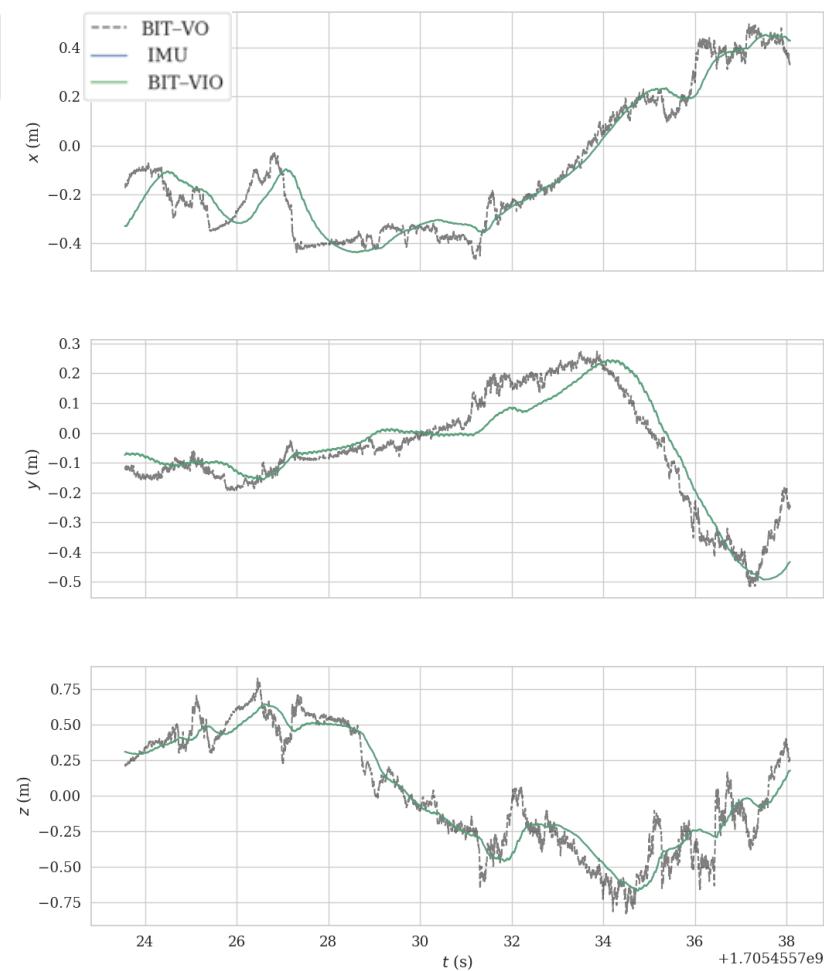


# Compared just BIT-VO with BIT-VIO and IMU

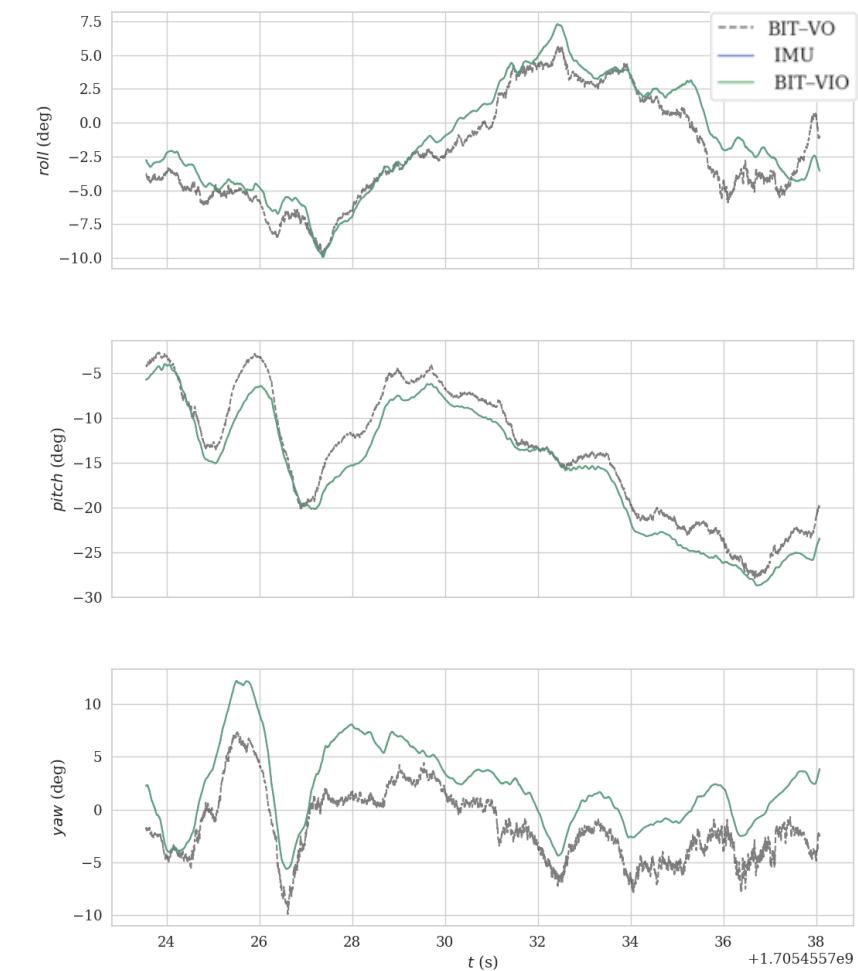


Traj 9/20 ~ 4.2 m

Full 6-DOF long trajs., fast motions  
Is fully achievable as we solved for both  
Now time to run experiments with GT.



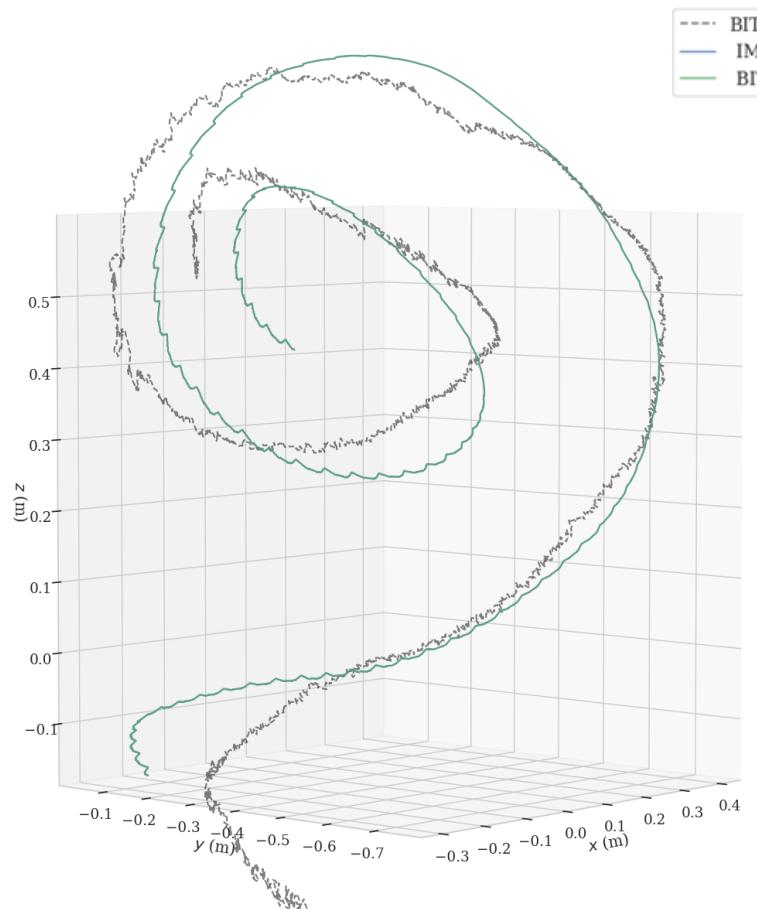
XYZ strong



RPY trends (~5deg)

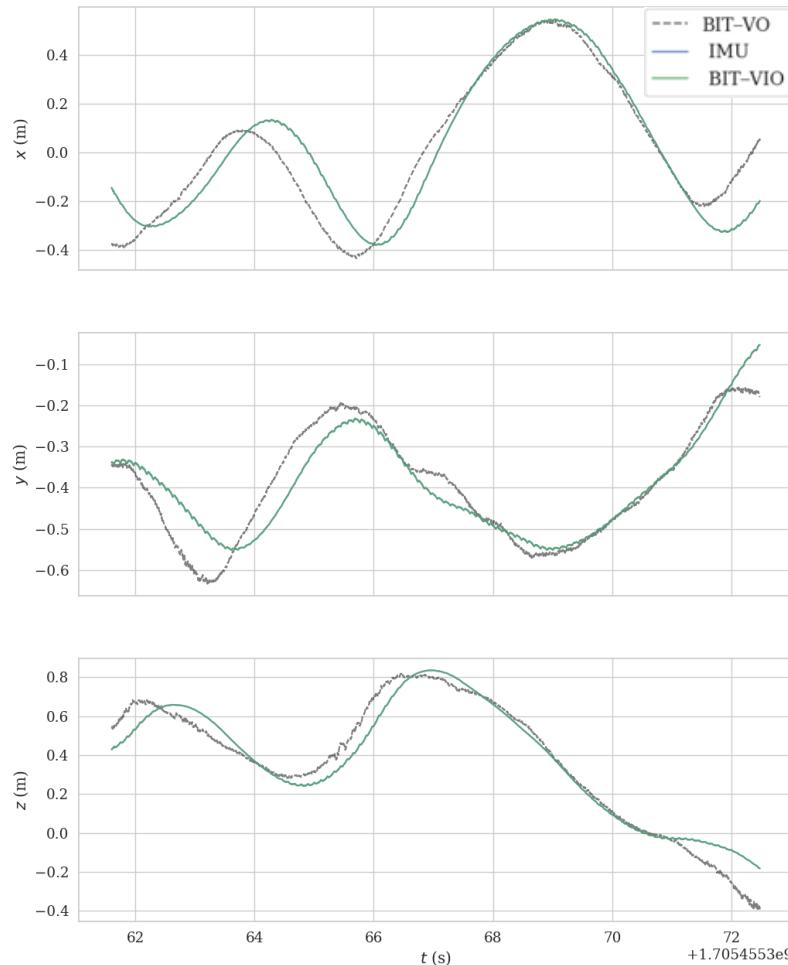
Our BIT-VIO silences the high frequency noise of prior BIT-VO.

# Compared just BIT-VO with BIT-VIO and IMU



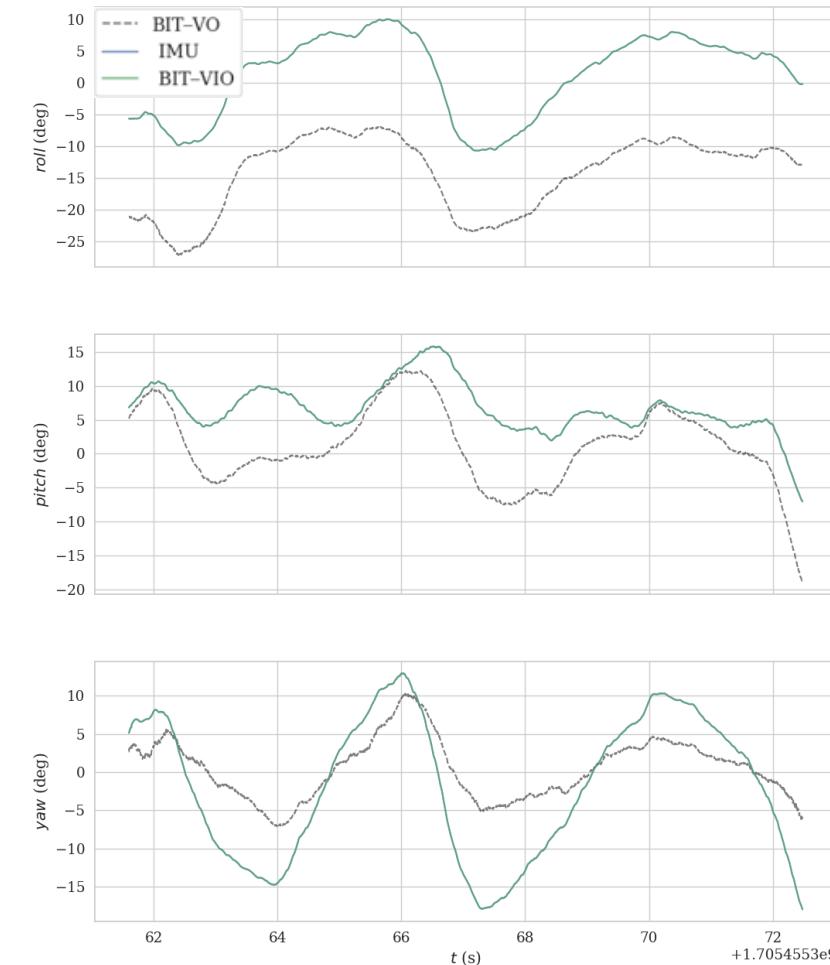
Traj 6/20 ~ 5.4 m

Full 6-DOF long trajs., fast motions  
Is fully achievable as we solved for both  
Now time to run experiments with GT.



XYZ strong

Our BIT-VIO silences the high frequency noise of prior BIT-VO.



RPY trends (~5deg)

# BIT-VO:

Absolute Trajectory

Error (ATE)

RMSE: 0.005504

median: 0.005067

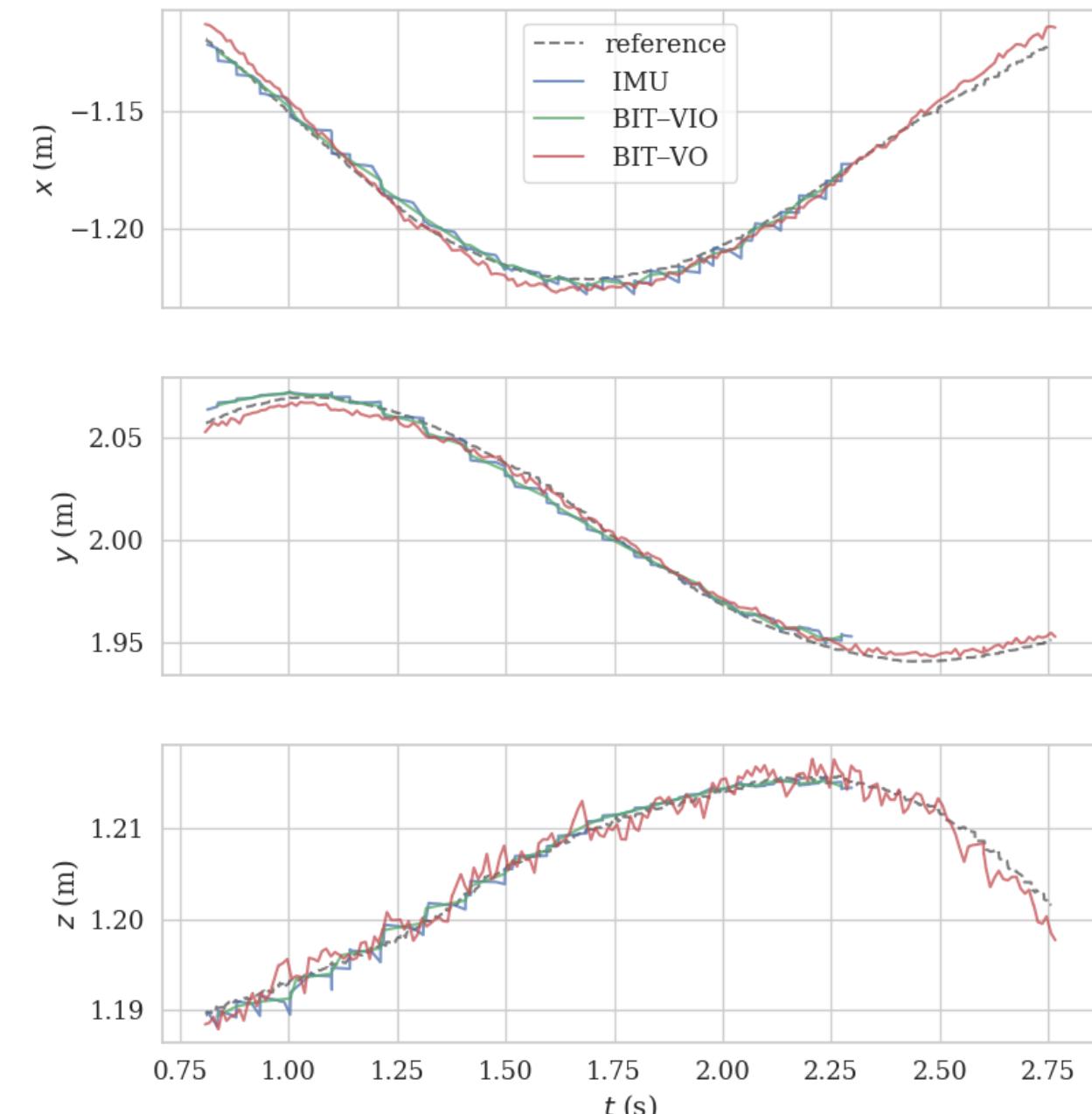
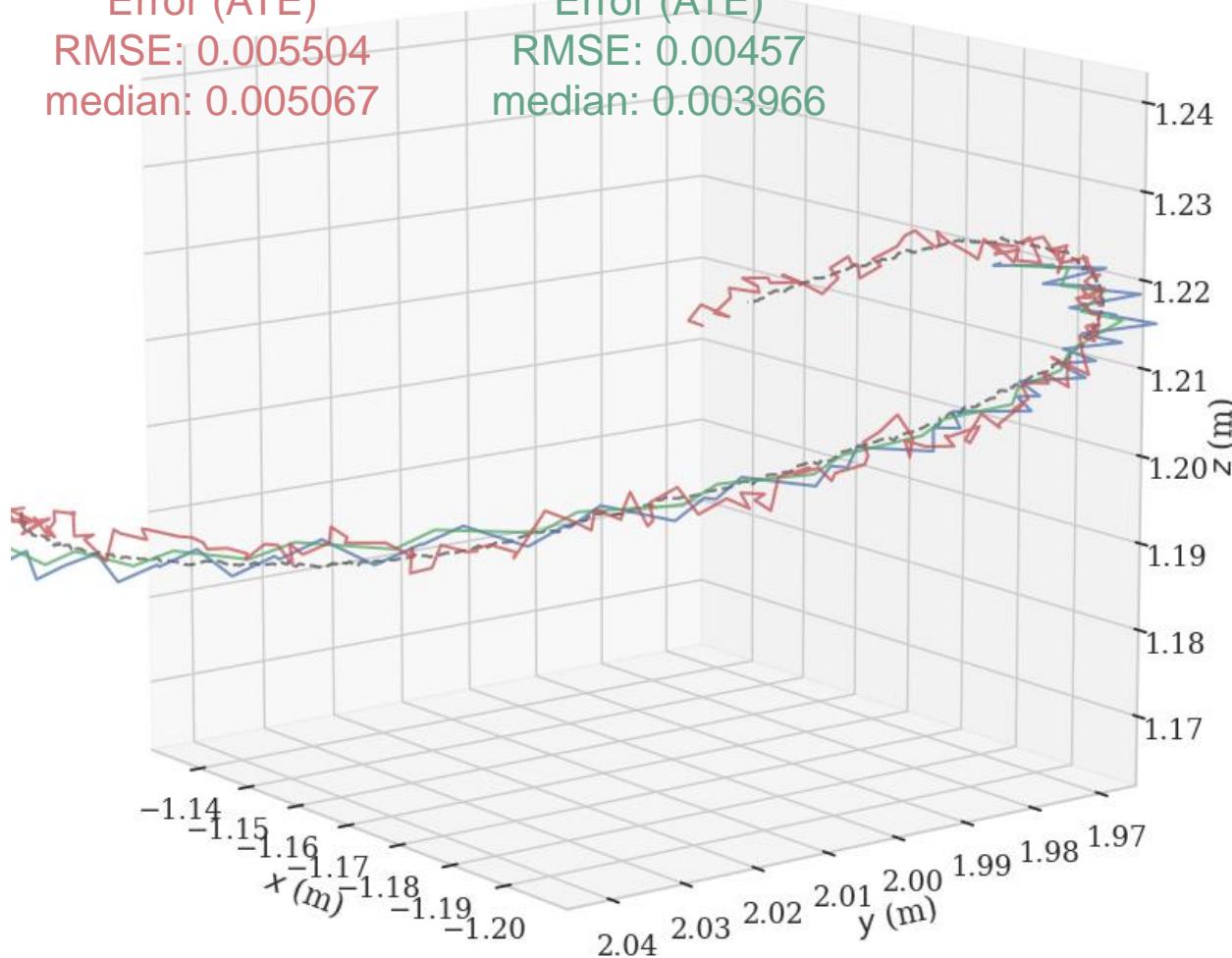
# BIT-VIO:

Absolute Trajectory

Error (ATE)

RMSE: 0.00457

median: 0.003966



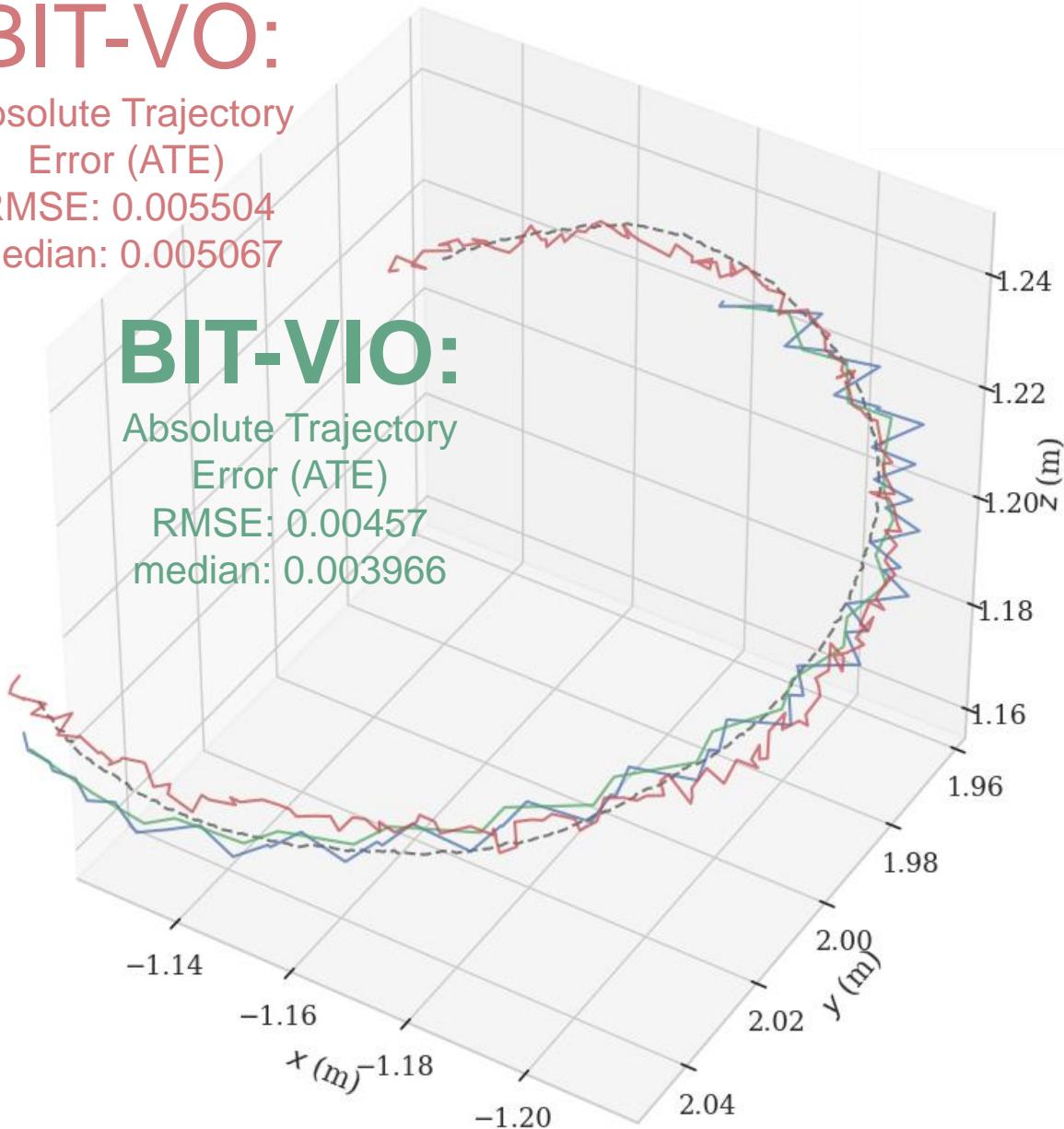
Comparison of the proposed BIT-VIO with inertial (IMU) and visual odometry (BIT-VO) overlaid on the reference ground-truth trajectory. Our BIT-VIO estimates are closer to the ground-truth compared to predictions from IMU or BIT-VO. Notice that our BIT-VIO effectively removes the high frequency noise visible in BIT-VO's trajectory.

# BIT-VO:

Absolute Trajectory  
Error (ATE)

RMSE: 0.005504

median: 0.005067

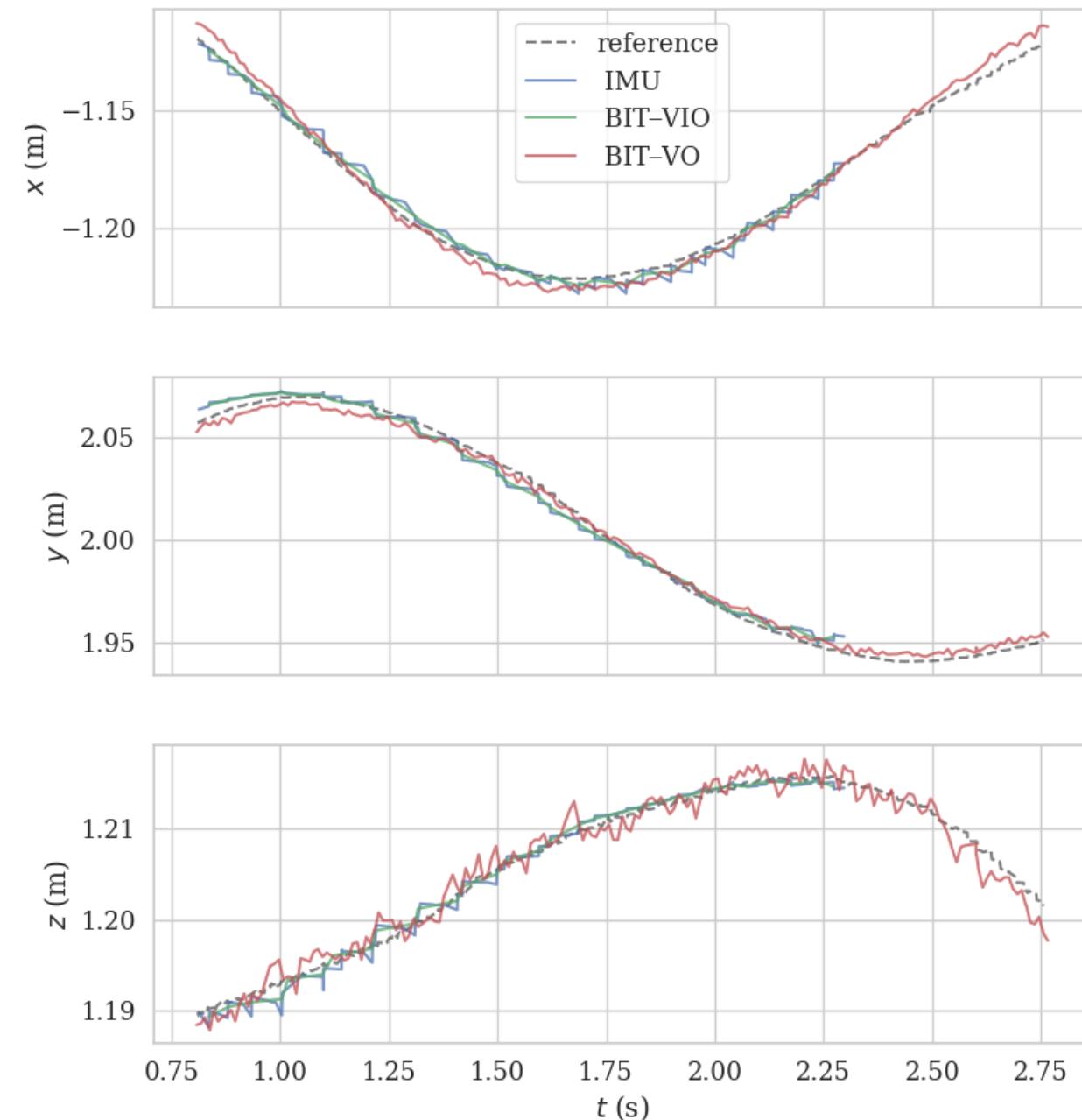


# BIT-VIO:

Absolute Trajectory  
Error (ATE)

RMSE: 0.00457

median: 0.003966



Comparison of the proposed BIT-VIO with inertial (IMU) and visual odometry (BIT-VO) overlaid on the reference ground-truth trajectory. Our BIT-VIO estimates are closer to the ground-truth compared to predictions from IMU or BIT-VO. Notice that our BIT-VIO effectively removes the high frequency noise visible in BIT-VO's trajectory.

# BIT-VO:

Absolute Trajectory  
Error (ATE)

RMSE: 0.037553

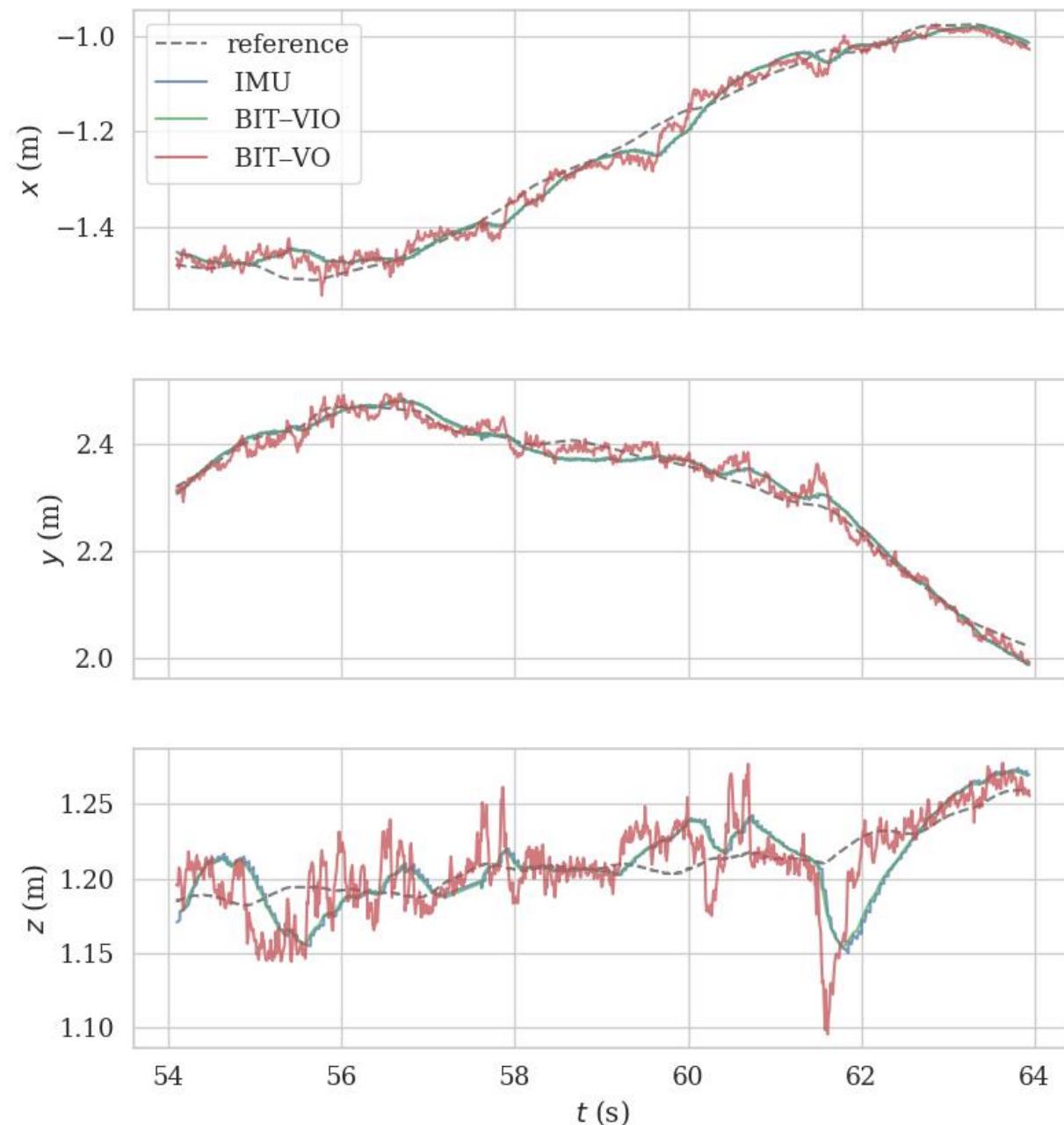
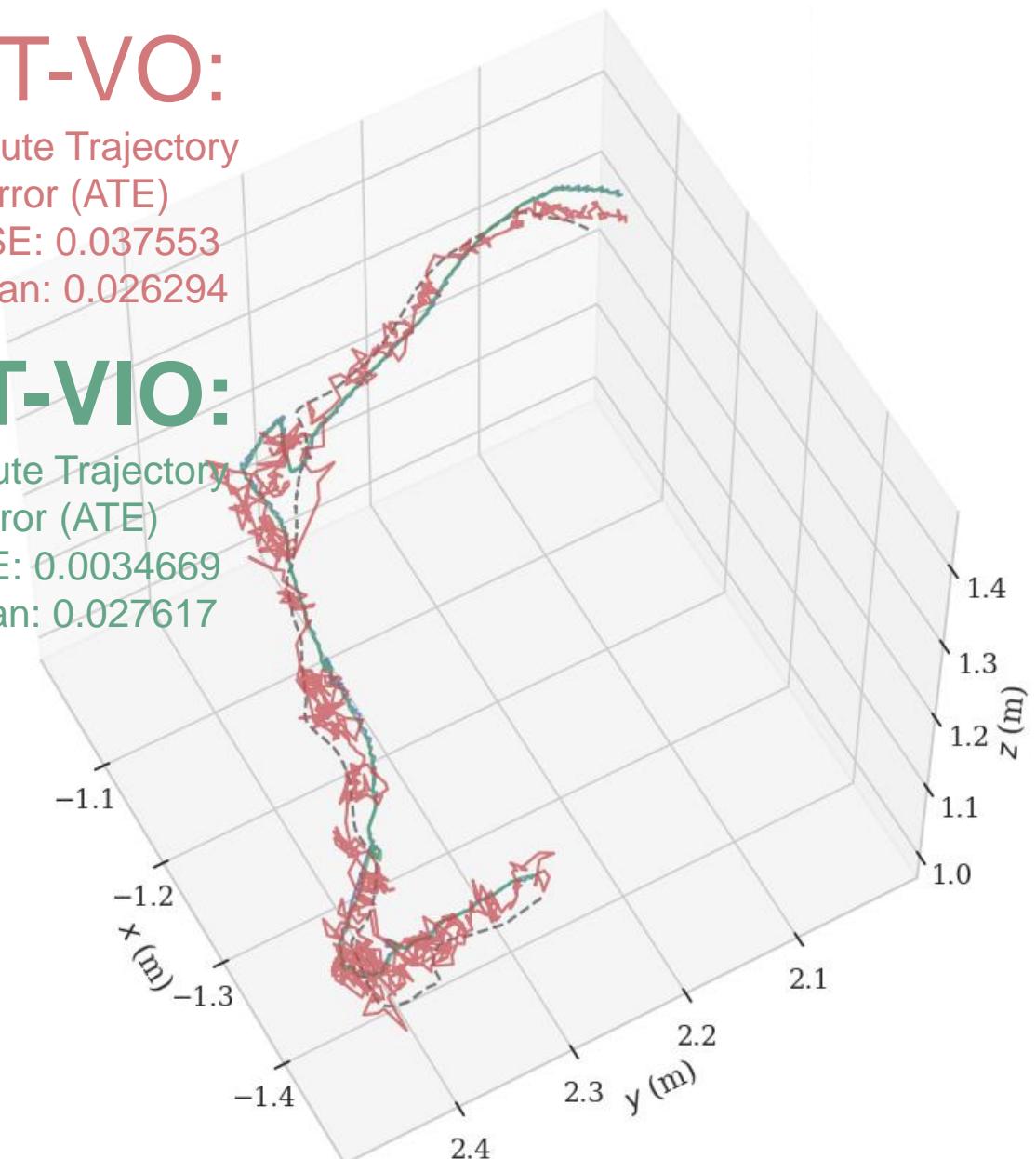
median: 0.026294

# BIT-VIO:

Absolute Trajectory  
Error (ATE)

RMSE: 0.0034669

median: 0.027617



Comparison of the proposed BIT-VIO with inertial (IMU) and visual odometry (BIT-VO) overlaid on the reference ground-truth trajectory. Our BIT-VIO estimates are closer to the ground-truth compared to predictions from IMU or BIT-VO. Notice that our BIT-VIO effectively removes the high frequency noise visible in BIT-VO's trajectory.

# BIT-VO:

Absolute Trajectory  
Error (ATE)

RMSE: 0.037553

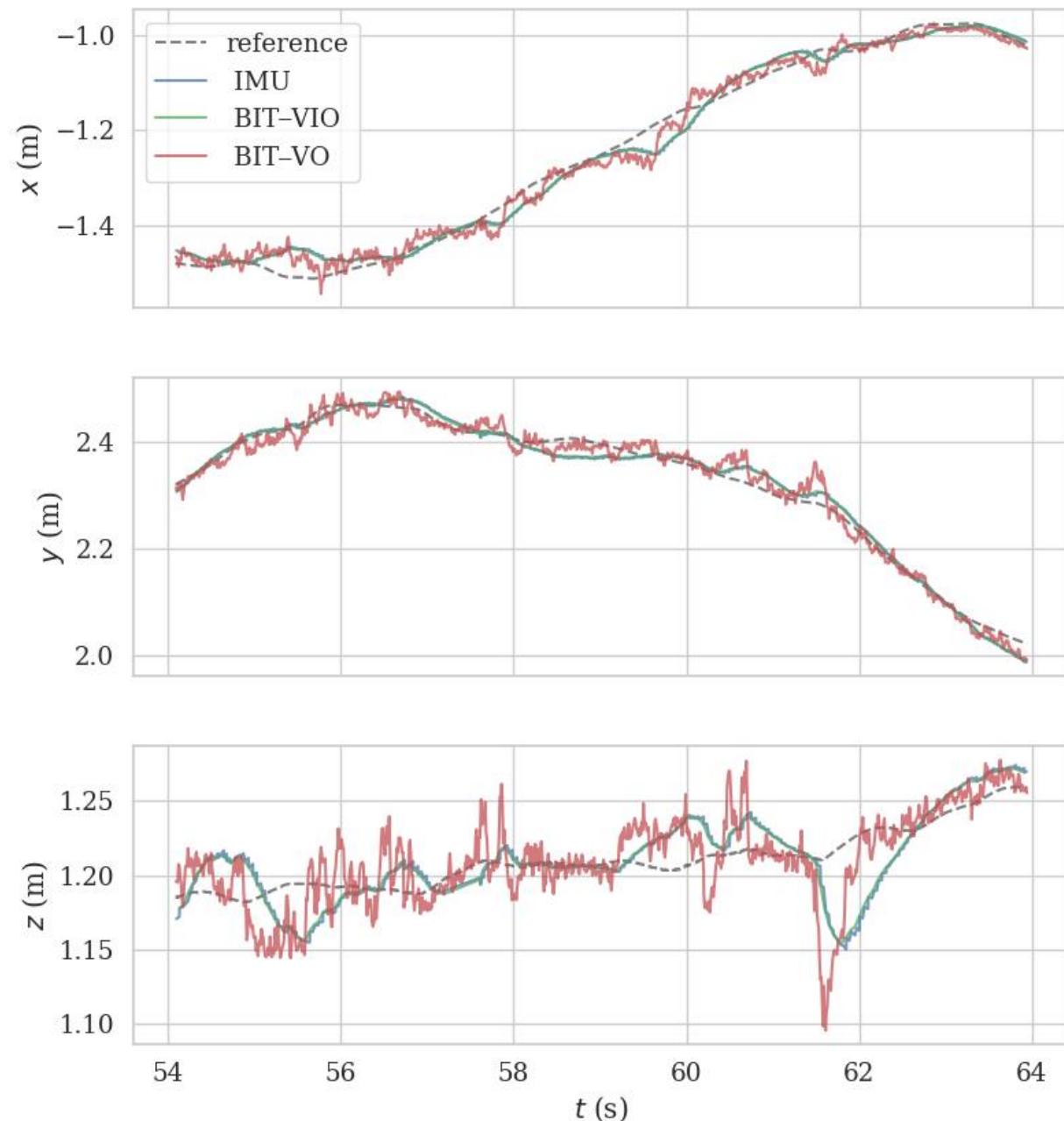
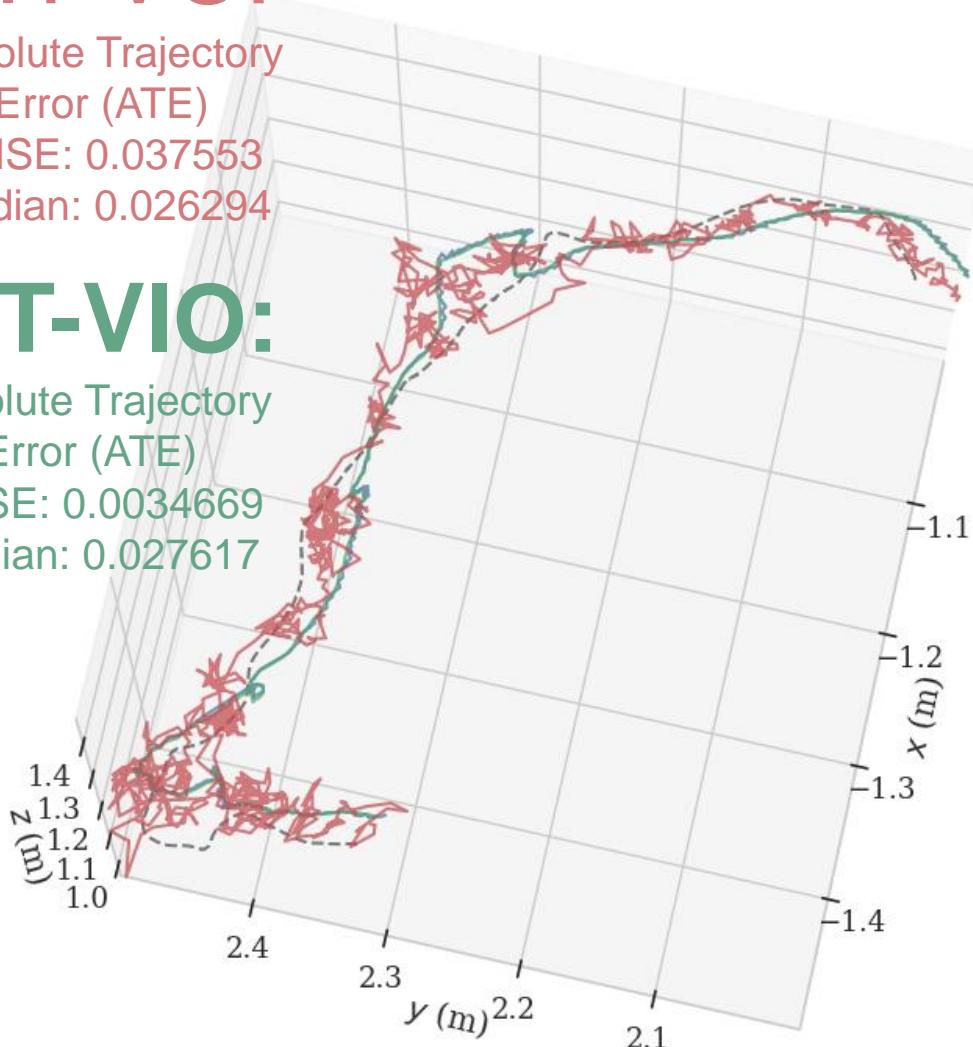
median: 0.026294

# BIT-VIO:

Absolute Trajectory  
Error (ATE)

RMSE: 0.0034669

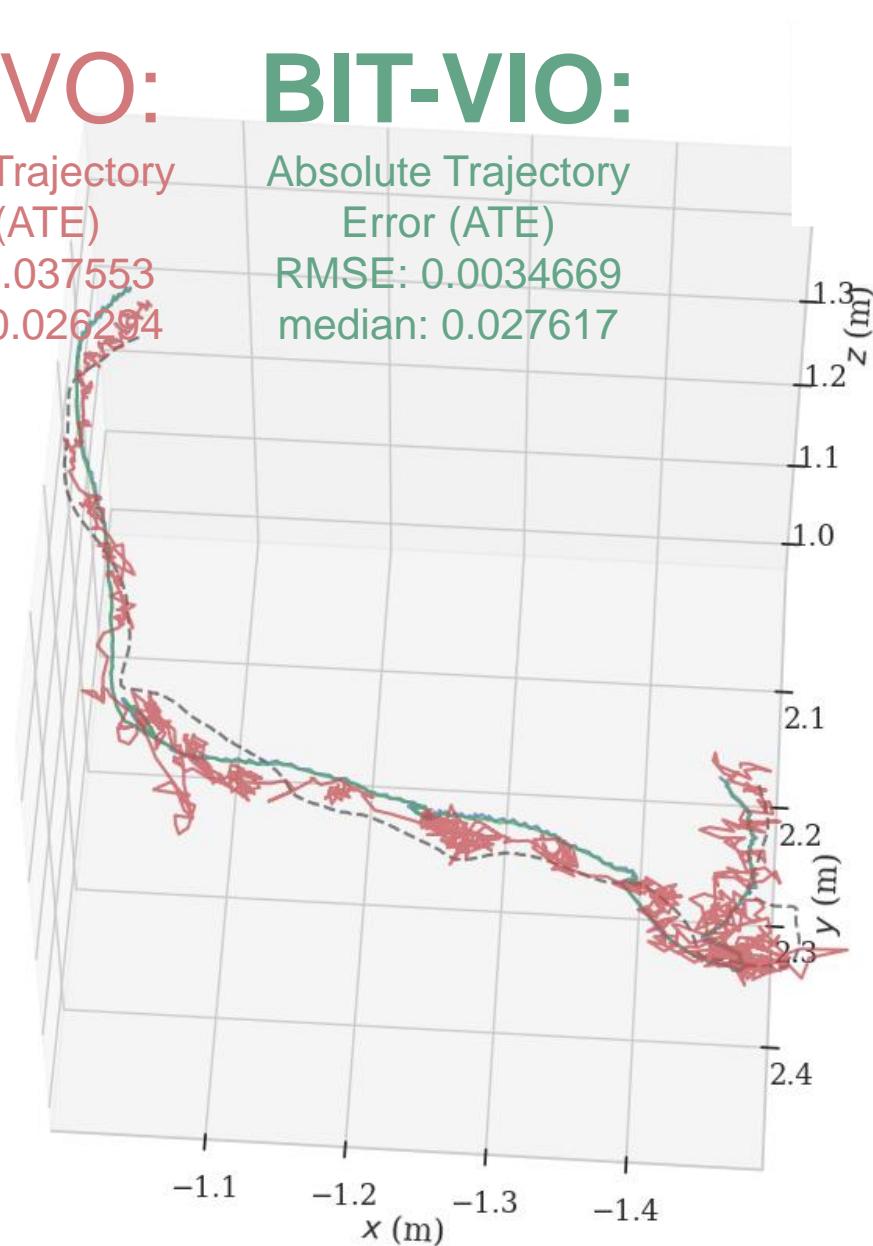
median: 0.027617



Comparison of the proposed BIT-VIO with inertial (IMU) and visual odometry (BIT-VO) overlaid on the reference ground-truth trajectory. Our BIT-VIO estimates are closer to the ground-truth compared to predictions from IMU or BIT-VO. Notice that our BIT-VIO effectively removes the high frequency noise visible in BIT-VO's trajectory.

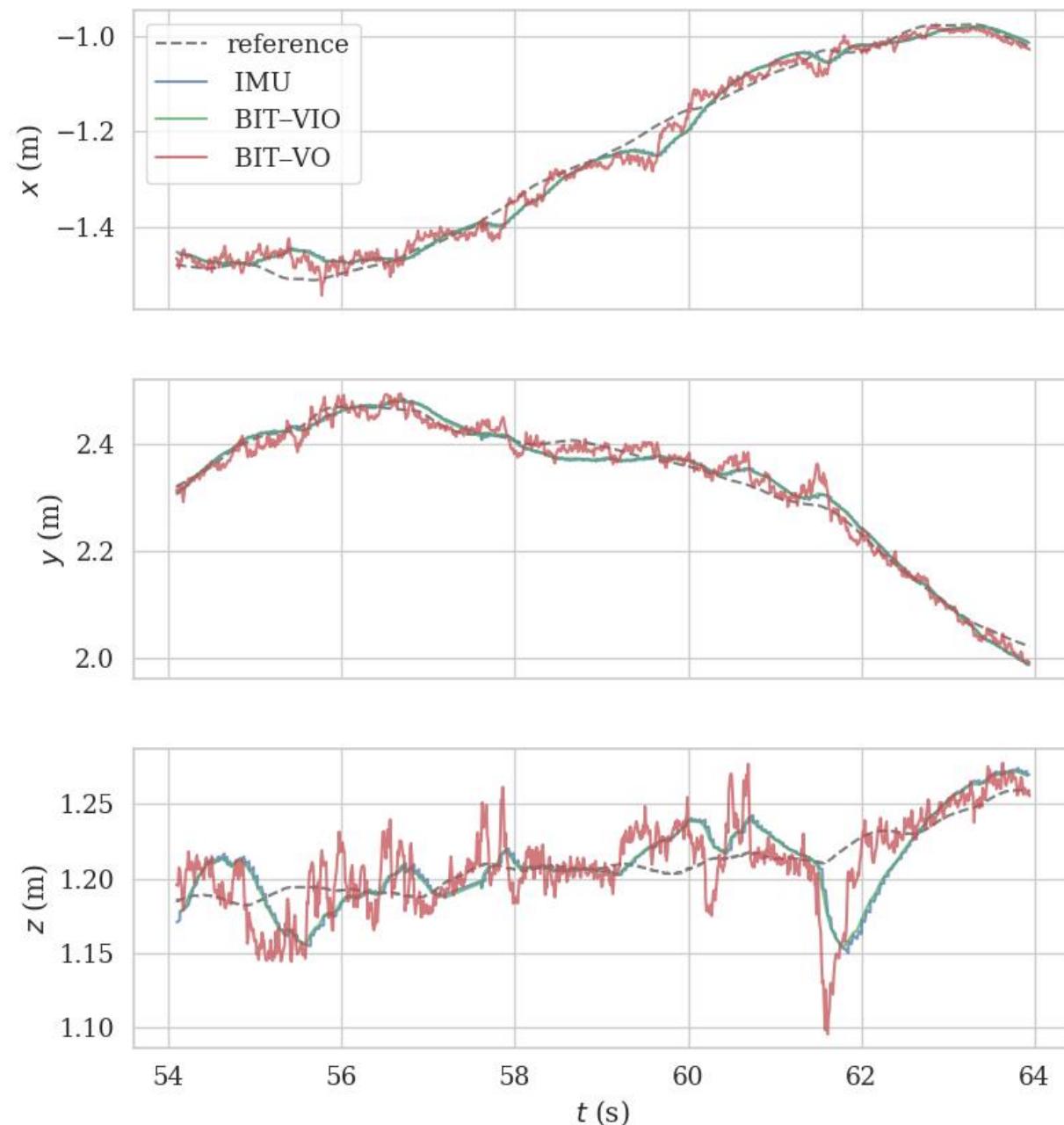
# BIT-VO:

Absolute Trajectory  
Error (ATE)  
RMSE: 0.037553  
median: 0.026294

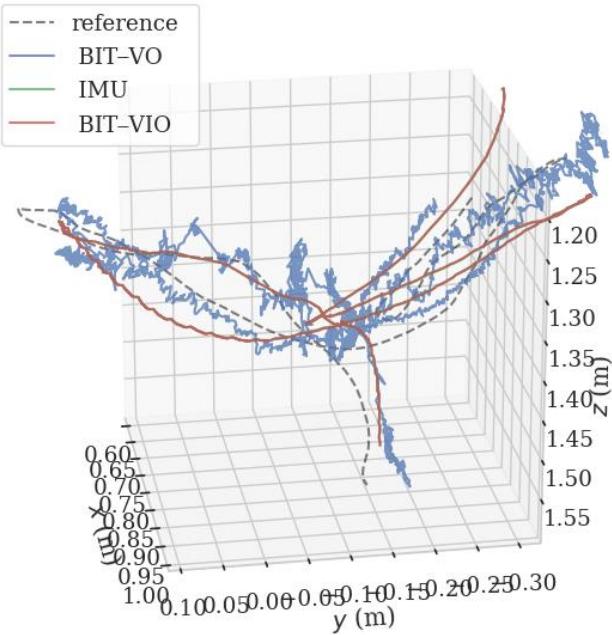


# BIT-VIO:

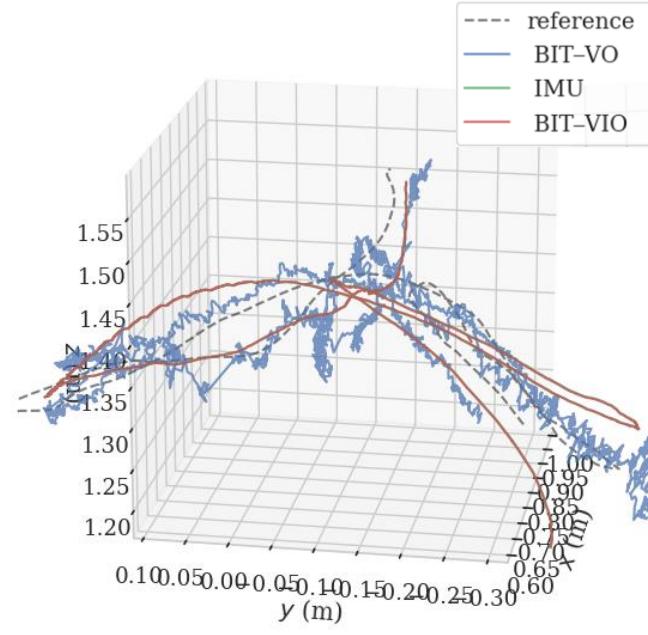
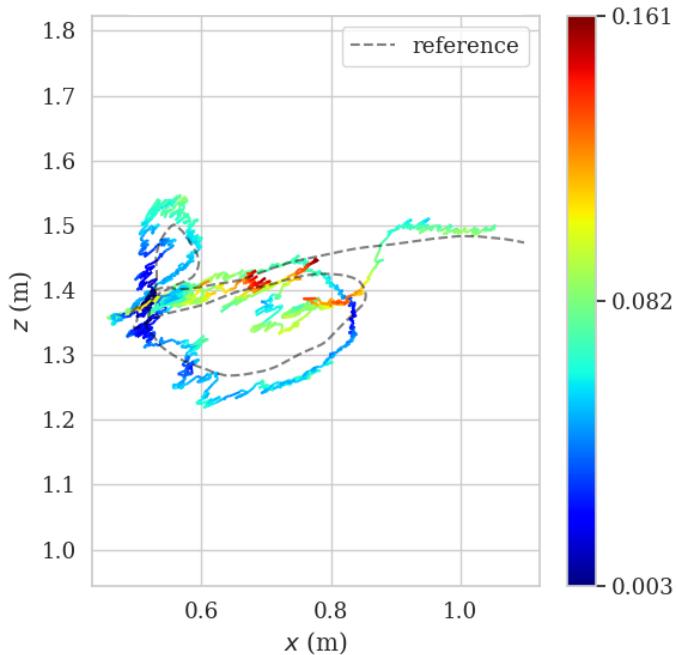
Absolute Trajectory  
Error (ATE)  
RMSE: 0.0034669  
median: 0.027617



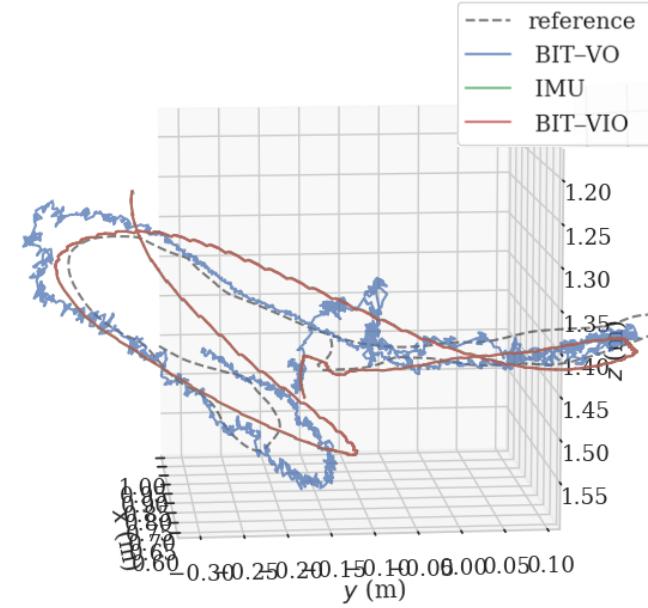
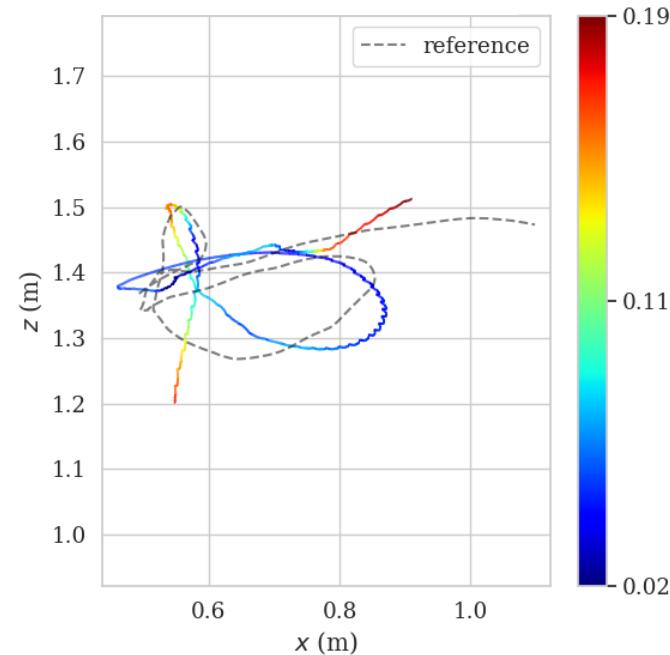
Comparison of the proposed BIT-VIO with inertial (IMU) and visual odometry (BIT-VO) overlaid on the reference ground-truth trajectory. Our BIT-VIO estimates are closer to the ground-truth compared to predictions from IMU or BIT-VO. Notice that our BIT-VIO effectively removes the high frequency noise visible in BIT-VO's trajectory.



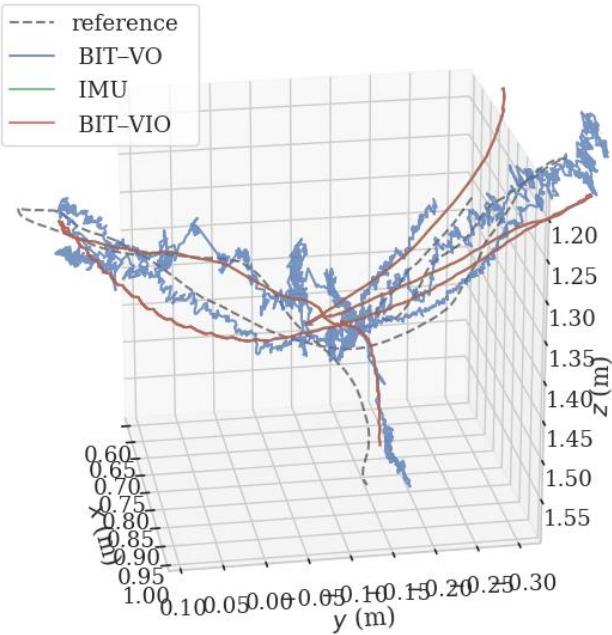
Error Mapped on BIT-VO



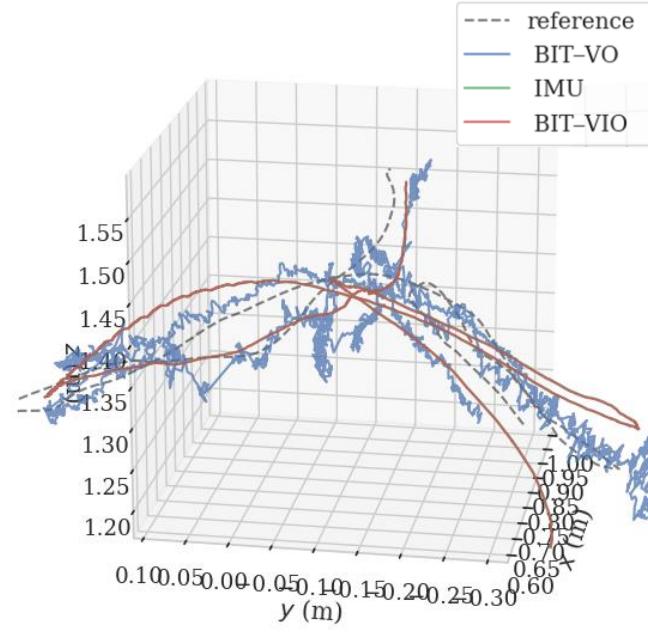
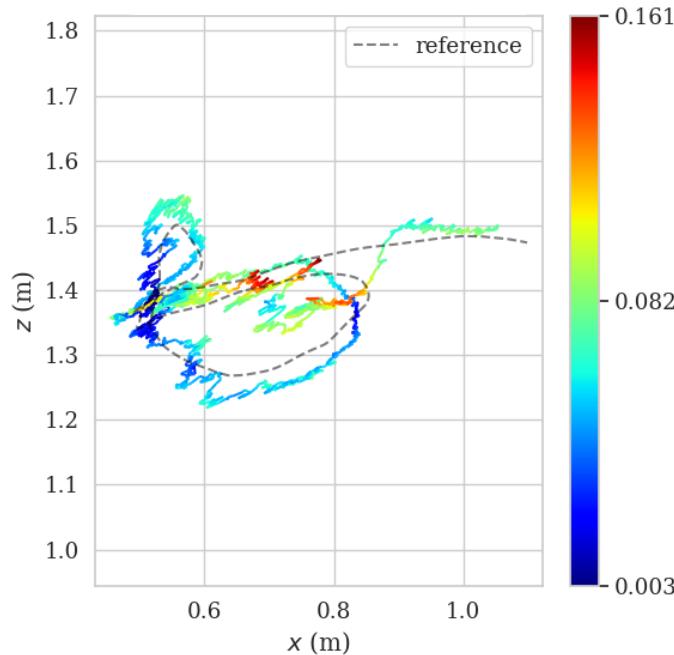
Error Mapped on BIT-VIO



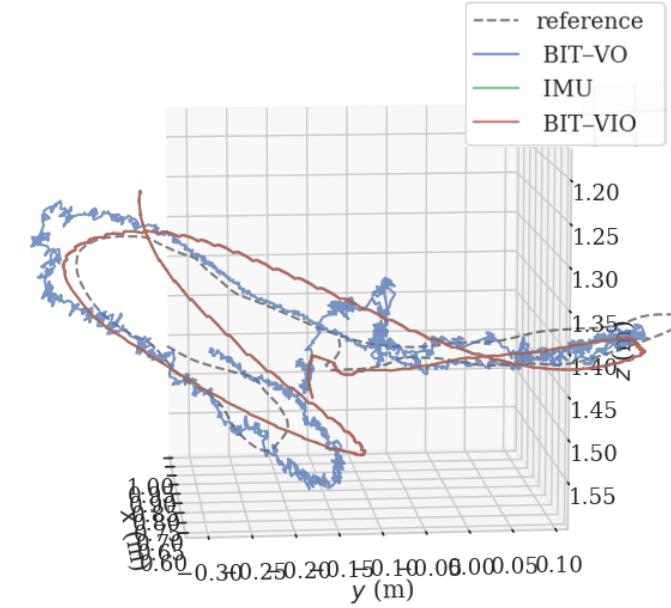
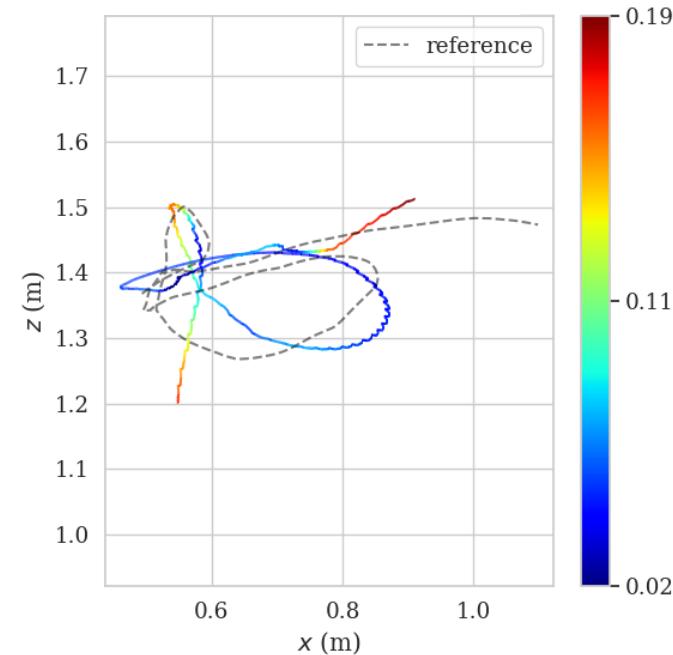
When projecting the error on BIT-VO (left) and BIT-VIO (right), we see that BIT-VO has high frequency noise with red tail-ends on its trajectory when compared to BIT-VIO's stabler closer to ground-truth trajectory with little regions of large ATE error (in red).



Error Mapped on BIT-VO



Error Mapped on BIT-VIO



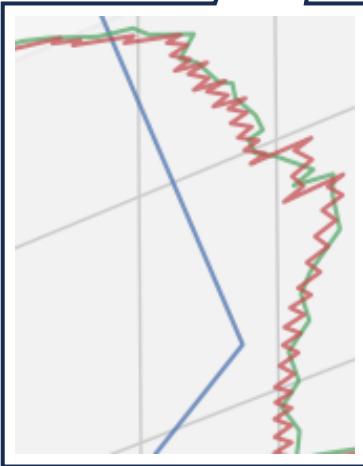
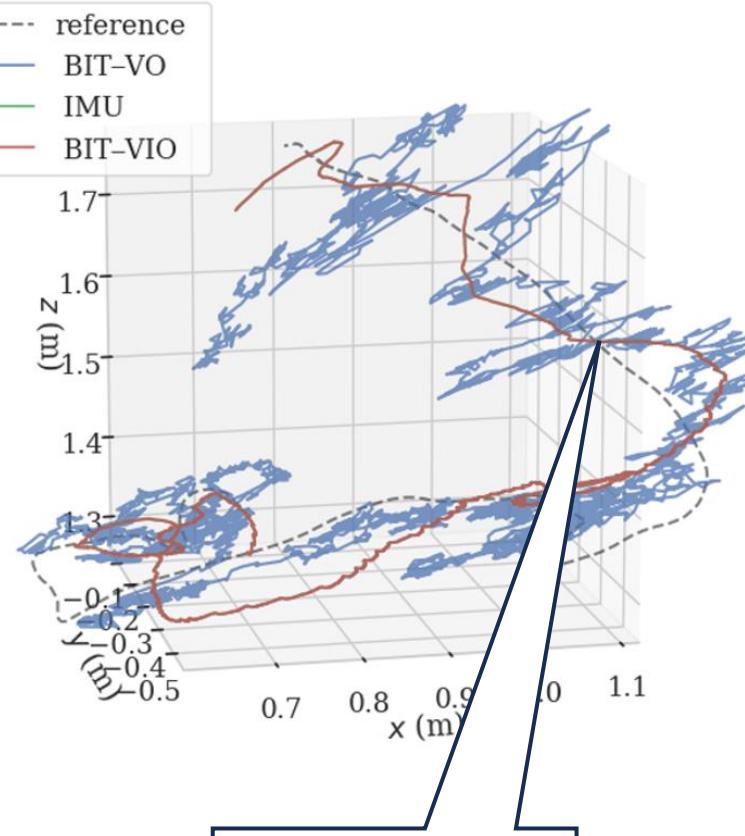
## TRAJ. K

Doing fast, hostile motions, BIT-VO still has such high frequency noise.

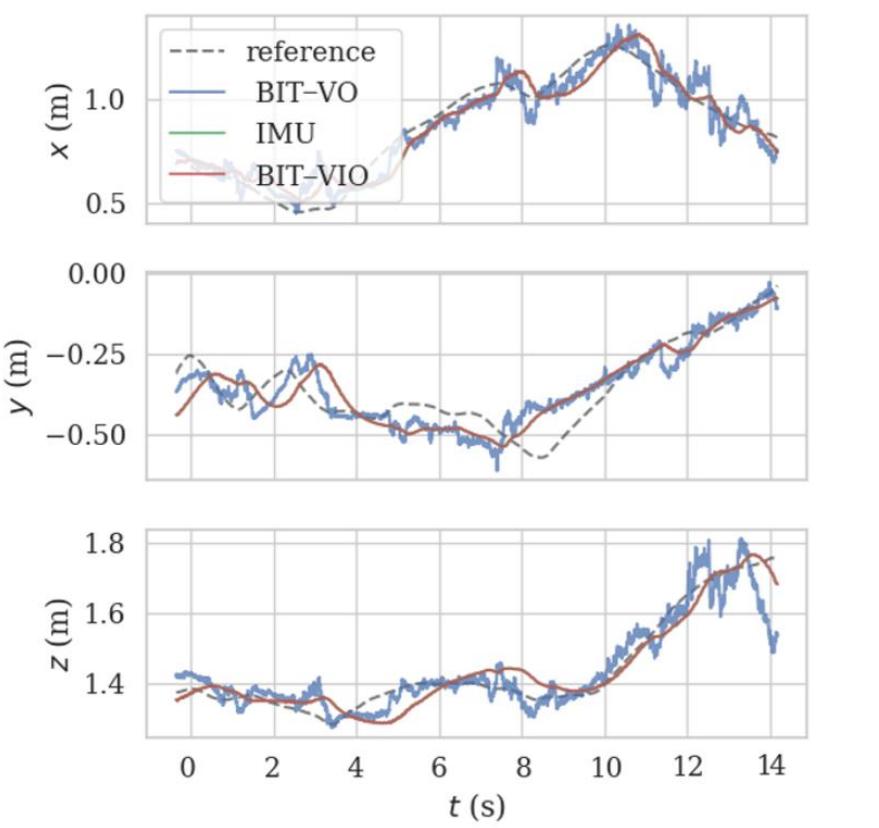
Our BIT-VIO algorithm can effectively remove this while having lower ATE error.

APE with IMU at 400 Hz, BIT-VO at 300 FPS

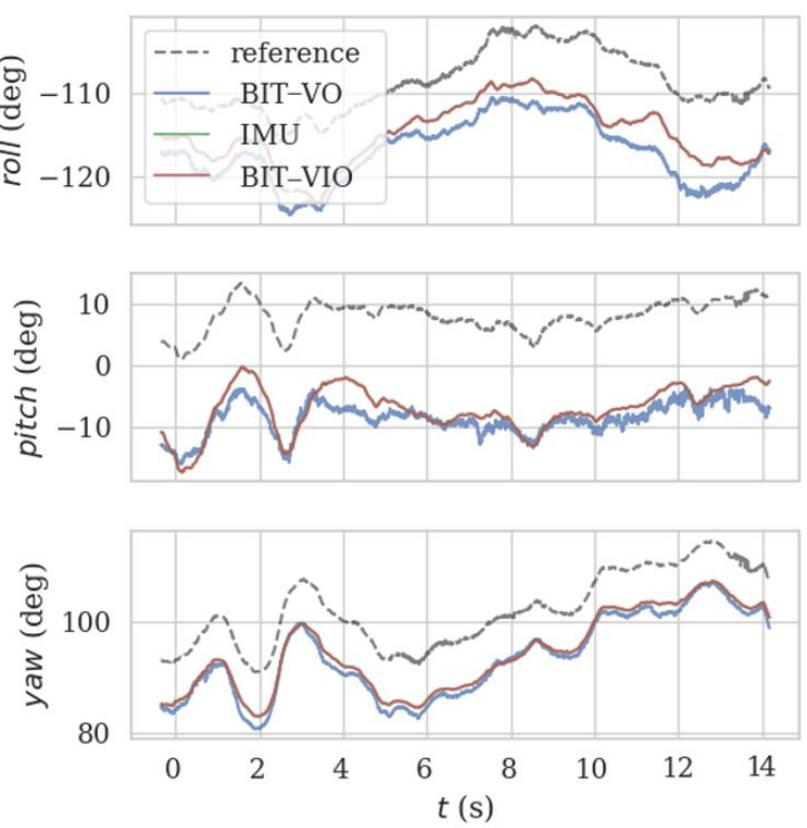
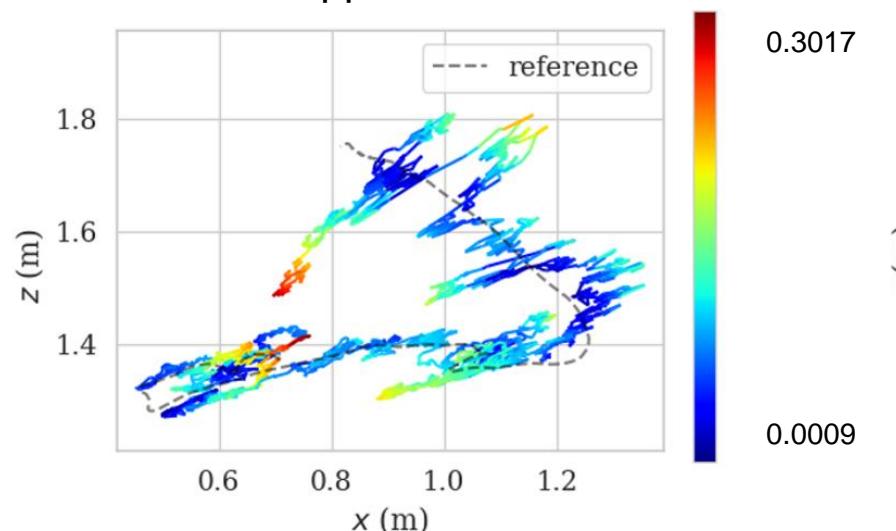
Traj.	Type	BIT-VO APE (m)	IMU APE (m)	BIT-VIO APE (m)	Length (m)
K	RMSE:	0.075214	0.092716	<b>0.092355</b>	2.48
	median:	0.067571	0.062949	<b>0.062812</b>	



IMU as prediction, BIT-VO as update, BIT-VIO correcting via iEKF



Error Mapped on BIT-VO



Error Mapped on BIT-VIO

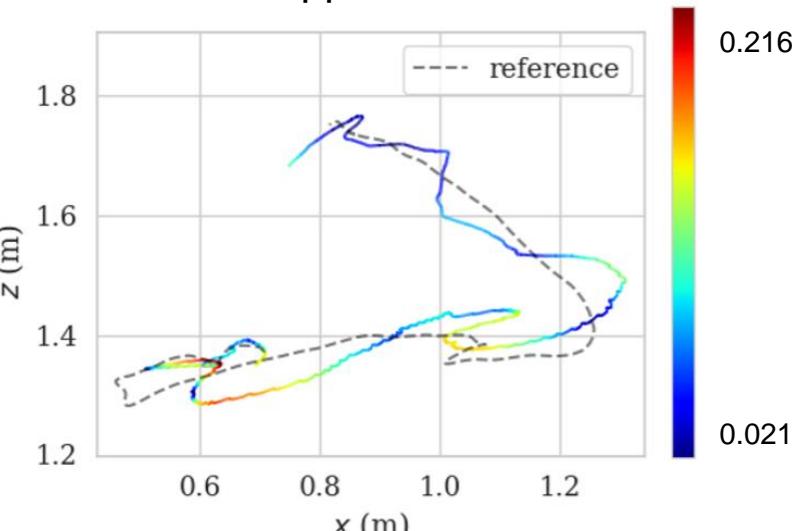
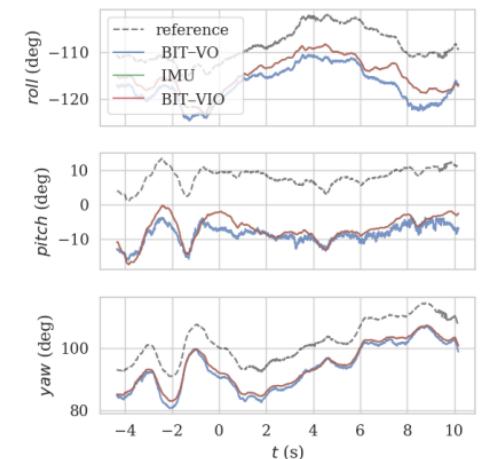
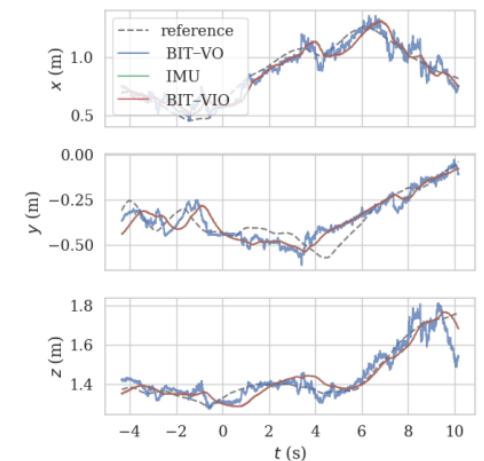
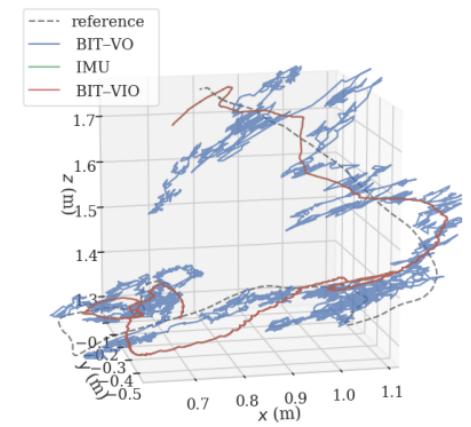
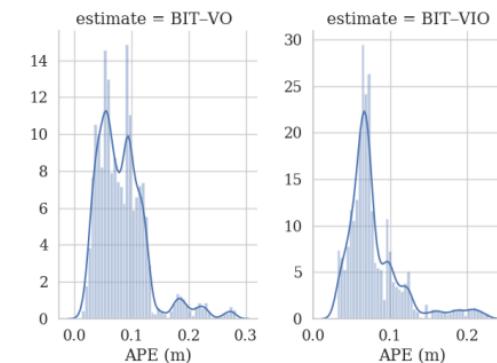


TABLE I  
IMU AT 400 Hz AND BIT-VO AT 300 FPS. BIT-VIO HAS LOWER APE.

APE with IMU at 400 Hz, BIT-VO at 300 FPS					
Traj.	Type	BIT-VO APE (m)	IMU APE (m)	BIT-VIO APE (m)	Length (m)
A	RMSE:	0.156247	0.156138	<b>0.156138</b>	3.6
	median:	0.138958	0.138315	<b>0.133202</b>	
B	RMSE:	<b>0.144674</b>	0.151913	0.150545	3.42
	median:	<b>0.096717</b>	0.106351	0.09959	
C	RMSE:	0.134617	0.121777	<b>0.12071</b>	2.6
	median:	0.119079	0.112582	<b>0.111856</b>	
D	RMSE:	0.05454	0.06597	<b>0.063498</b>	1.1
	median:	0.050929	0.053554	<b>0.049605</b>	
E	RMSE:	0.094479	0.08533	<b>0.086911</b>	1.7
	median:	0.07561	0.067952	<b>0.068756</b>	
F	RMSE:	0.175323	0.154923	<b>0.153335</b>	2.12
	median:	0.174444	0.142438	<b>0.140952</b>	
G	RMSE:	0.206866	0.196163	<b>0.195263</b>	2.4
	median:	0.15714	0.149788	<b>0.149103</b>	
H	RMSE:	0.088185	0.090189	<b>0.090376</b>	3.3
	median:	0.070431	0.07014	<b>0.070212</b>	
I	RMSE:	<b>0.134361</b>	0.13838	0.134328	2.01
	median:	<b>0.116587</b>	0.12805	0.124618	
J	RMSE:	0.057402	0.058275	<b>0.057461</b>	1.6
	median:	0.044689	0.044147	<b>0.042597</b>	
K	RMSE:	0.075214	0.092716	<b>0.092355</b>	2.48
	median:	0.067571	0.062949	<b>0.062812</b>	
L	RMSE:	0.10864	0.104469	<b>0.104366</b>	2.55
	median:	0.089689	0.088151	<b>0.08788</b>	
M	RMSE:	<b>0.094062</b>	0.108688	0.109156	2.9
	median:	<b>0.075766</b>	0.093151	0.09367	
N	RMSE:	0.101385	0.107546	<b>0.105227</b>	2.71
	median:	0.09016	0.089743	<b>0.087536</b>	

The Absolute Trajectory Error (ATE or APE) of our BITVIO <<< of BIT-VO.



# V. Conclusion and Next Steps

We have presented BIT-VIO, the first-ever 6-Degrees of Freedom (6-DOF) Visual Inertial Odometry (VIO) algorithm, which utilizes the advantages of the FPSP for vision-IMU-fused state estimation.

Our BIT-VIO algorithm operates and corrects by loosely-coupled sensor-fusion iterated Extended Kalman Filter (iEKF) at 300 FPS with an IMU at 400 Hz.

We evaluate BIT-VIO qualitatively against BIT-VO and demonstrate improvements in ATE across many trajectories. Moreover, the high-frequency noise evident in BIT-VO is effectively filtered out, resulting in a smoother estimated trajectory.

We plan to take the next steps toward a tightly-coupled VIO approach with the FPSP. Feature extraction would be integrated directly with IMU instead of two separate pose estimates, making the trajectories more robust and loop-closure better achievable.

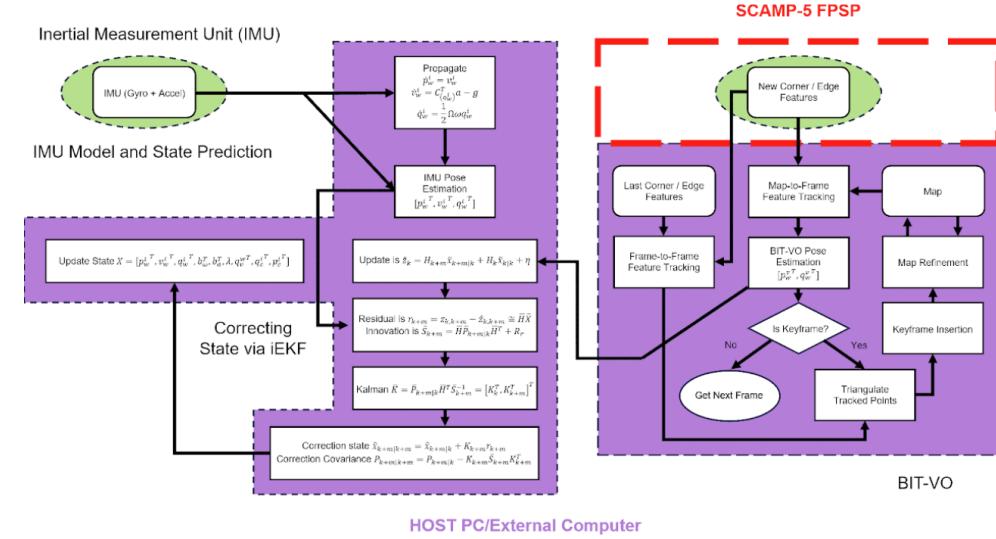


FIGURE 1: Pipeline of BIT-VIO. Multi-Sensor Fusion (iEKF-MSF) is to the left. BIT-VO is to the right, the vision sensor utilizes the FPSP, highlighted in red. New Corner/Edge Features are detected via the FPSP, off-loading computational load by allowing some image and signal processing to be done on the chip before transferring to a PC host or other external device to be further processed, i.e., some of the framework is alleviated on FPSP.

