



# StructVIO: Visual-Inertial Odometry with Structural Regularity of Man-Made Environments

Danping Zou

VALSE online seminar  
2019年7月10日



上海交通大學  
SHANGHAI JIAO TONG UNIVERSITY

# Visual SLAM



- Visual SLAM techniques have been widely applied to unmanned vehicles.



Stereo camera



Fisheye camera

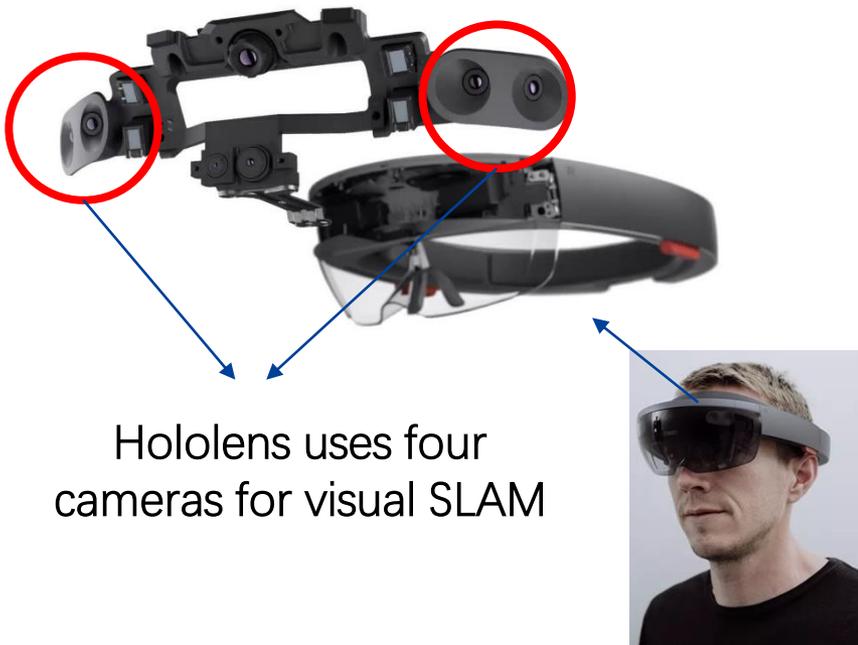


Stereo camera

# Visual SLAM



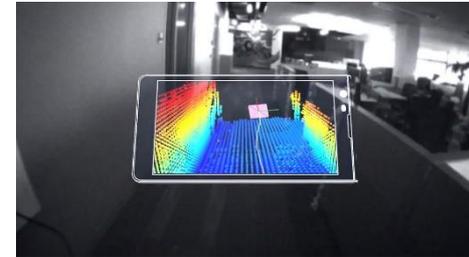
- Augmented reality (AR) (Hololens Glass, Project Tango Tablet)



Hololens uses four cameras for visual SLAM



Tango use one fisheye camera for visual SLAM



# Visual SLAM



- Operation system on cellphones
  - Google and Apple integrate visual SLAM into their OS (iOS, Android).



# Visual SLAM

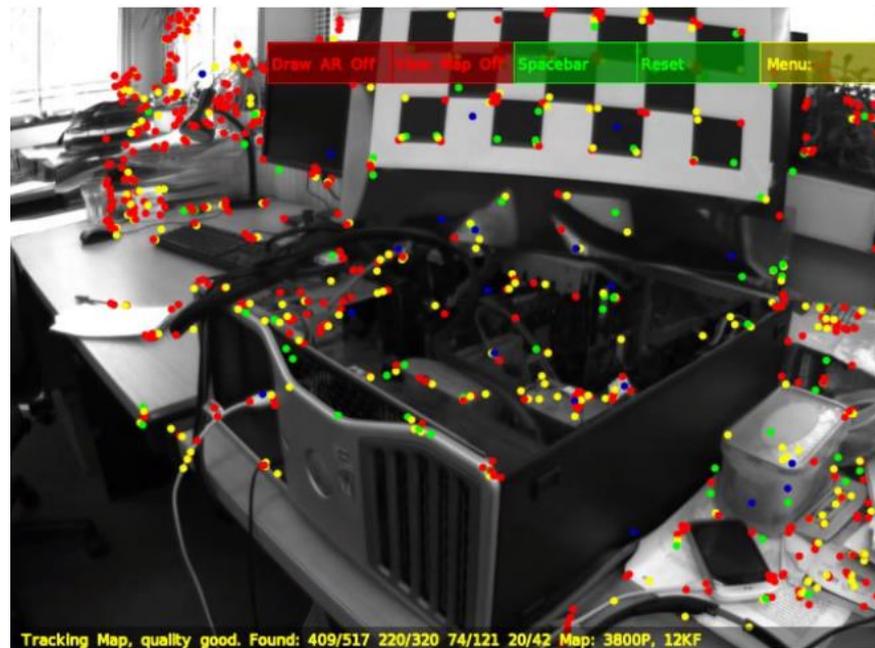
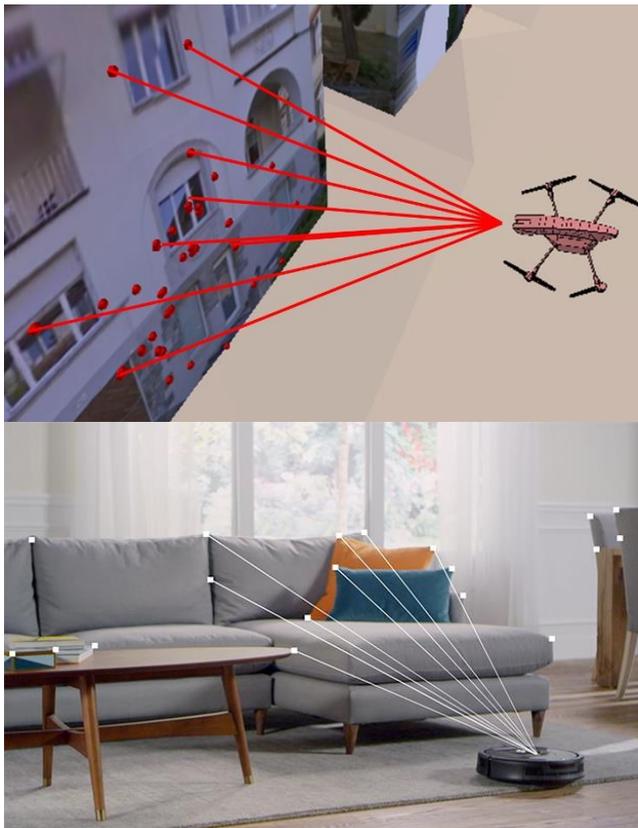


- A lot of algorithms have been proposed for visual SLAM in the past 15 years.
  - MonoSLAM (2003), StructSLAM(2014)
  - PTAM(2007), ORB-SLAM(2015)
  - SVO(2014), LSD-SLAM(2014), DSO(2016)
- Pure visual SLAM system is not robust in practical applications.
- Visual-inertial systems become predominant for real applications.
  - MSCKF (2007), ROVIO (2009)
  - OKVIS (2015), VINS(2017), ICE-BA(2018)

# Features in v/vi-SLAM systems



- Most visual-SLAM or visual-inertial systems choose points as the landmarks.



# Features in v/vi-SLAM systems



- Man made environments exhibit **strong regularity on geometry**.



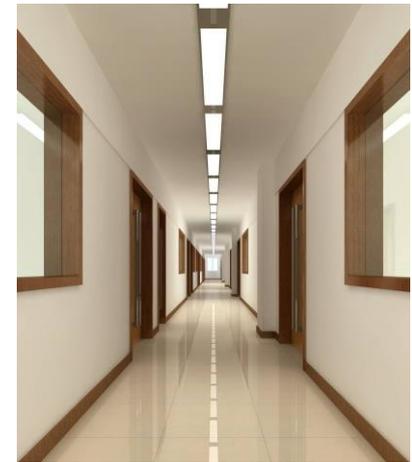
Natural scenes



Street

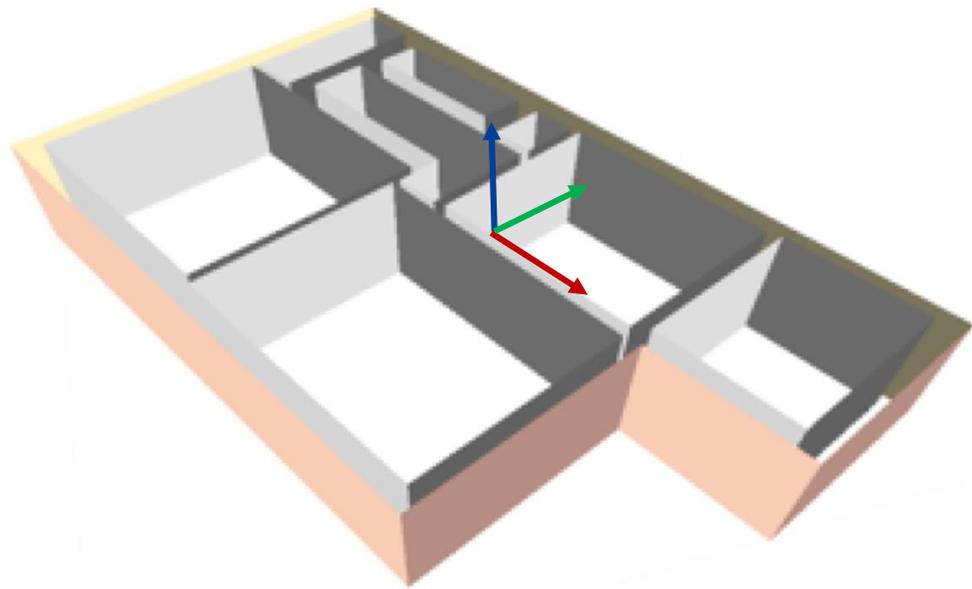


Underground parking



Indoor

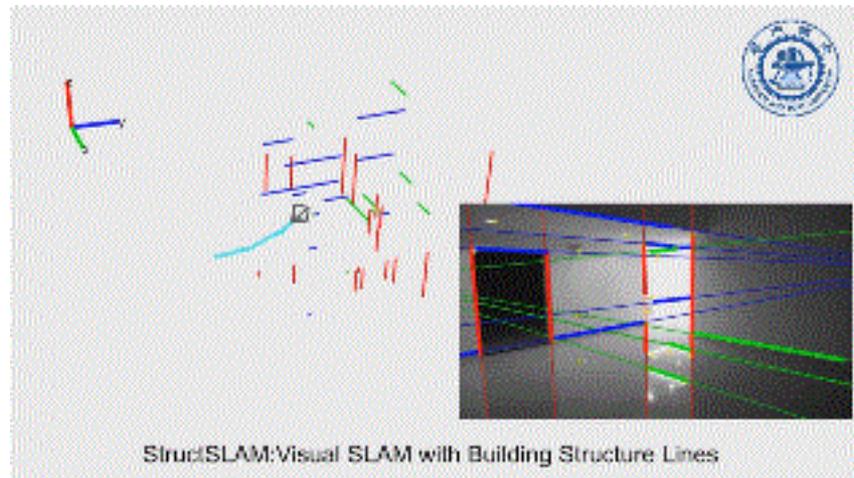
# Structural regularity - Manhattan word



1. Rich of line features
2. Three known directions (x, y, z)

# Visual SLAM with Manhattan world model

- StructSLAM (Presented VALSE online seminar, 2016, 30<sup>th</sup>, Mar)
  - Point + structural lines (lines aligned with x, y, z directions)
  - The direction of lines improves the observability of camera orientation



Zhou, Huizhong, Danping, Zou, et al. "StructSLAM: Visual SLAM with building structure lines." *Vehicular Technology, IEEE Transactions on* 64.4 (2015): 1364-1375. - Special session for indoor localization

# Real world is full of diversity



- A lot of man made environments can not be well described by Manhattan world model.
- Oblique/curvy structures.



# StructVIO



- A novel visual-inertial odometry method is presented
  - Use **Atlanta world** model to better describe irregular scenes.
  - Made **several improvements** to existing VIO approach.
  - A VIO dataset that can be used evaluate different methods.



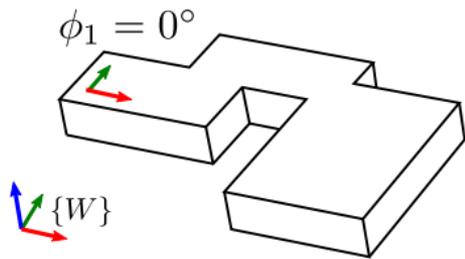
Zou, Danping, et al. “StructVIO: Visual-inertial Odometry with Structural Regularity of Man-made Environments.” IEEE Trans. on Robotics, 2019

Executable, tools & dataset : <http://drone.sjtu.edu.cn/dpzou/project/structvio.html>

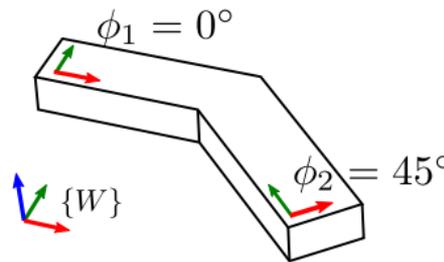
# Key idea – Atlanta world model



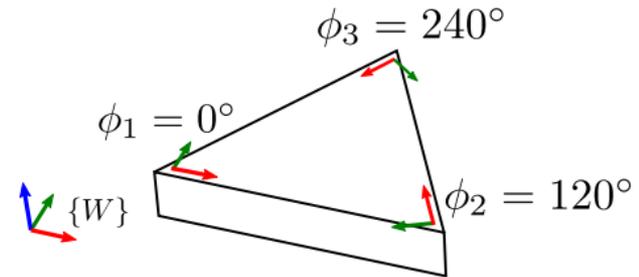
- We can approximate an irregular world by a group of local Manhattan worlds.
- Each one of them can be represented by a heading direction :  $\phi$ .



One Manhattan world



Two Manhattan worlds



Three Manhattan worlds

# Key idea – Atlanta world model



- Locally, the world is a Manhattan world. We can still use
  - Three directions
  - Structural line features
- to improve the performance of the VIO system.

Three directions

X,Y directions – Render the Yaw angle observable (locally)

Z direction – Render the gravity direction observable

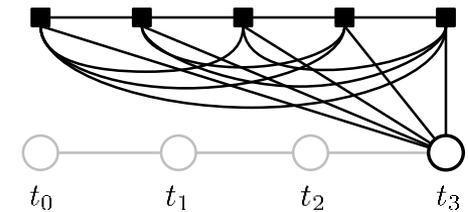
Line features

A good complementary to point features in texture-less scenes.

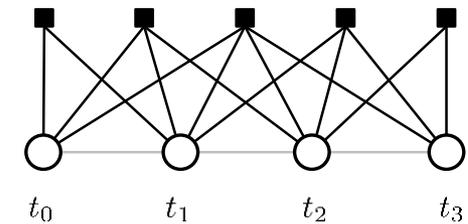
# The framework of StructVIO



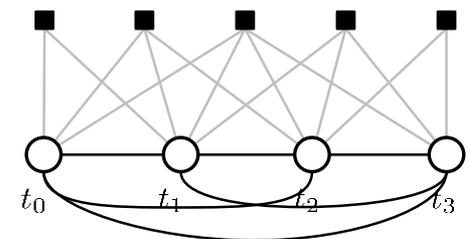
- We adopt the multi-state EKF filter based framework.
- Comparing with classic EKF filter
  - Much faster since the features are not included in the state vector.
- Comparing with key-frame optimization
  - Short feature trajectories are fully explored.
  - State update using a single feature trajectory.
  - Efficient but without losing much accuracy.



Classic EKF filter



Key-frame optimization

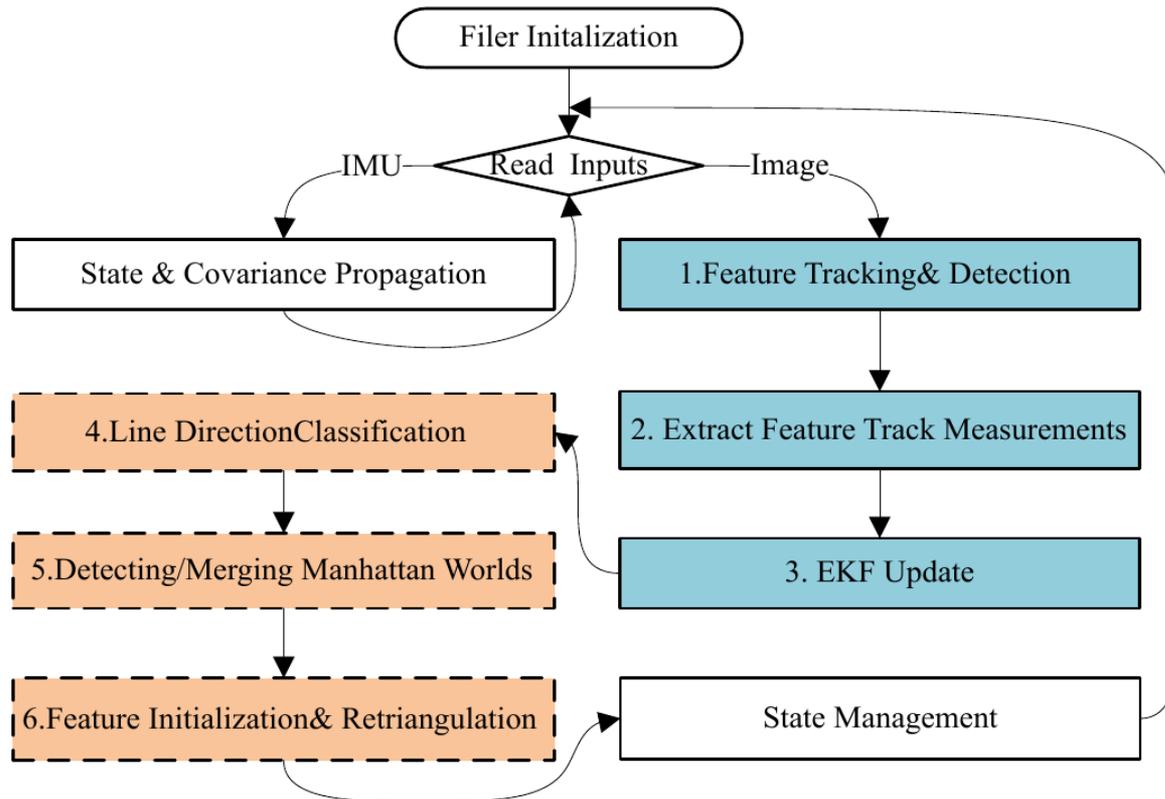


Multi-state EKF filter

# The framework of StructVIO



- The pipeline of StructVIO is as the following:



# State definition of StructVIO



- The state vector consists of the *current IMU state*, *historical IMU poses*, *calibration parameters*, and the *heading directions* of local Manhattan worlds

$$x_k = [x_{I_k}, {}^I C \tau, \phi_1, \dots, \phi_N, \begin{matrix} W \\ I_1 \tau, \dots, I_M \tau \end{matrix}]$$

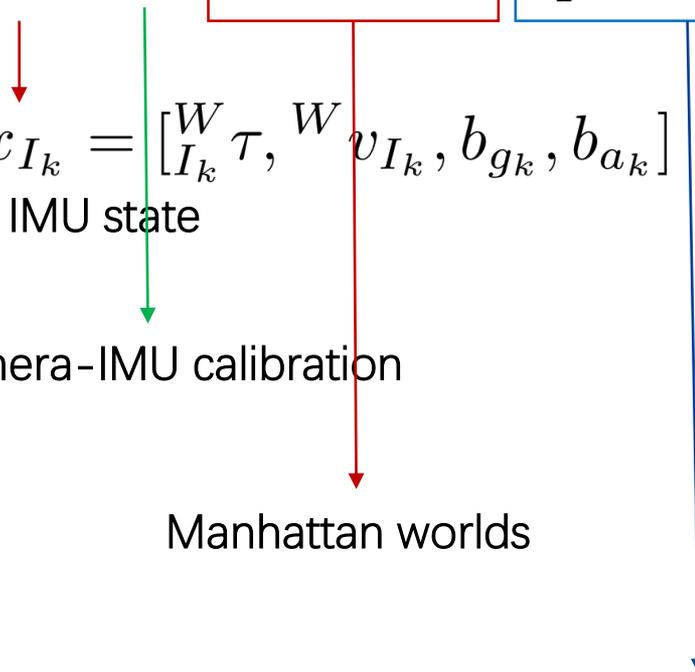
$$x_{I_k} = [{}^W I_k \tau, {}^W v_{I_k}, b_{g_k}, b_{a_k}]$$

Current IMU state

Camera-IMU calibration

Manhattan worlds

Historical IMU poses



# StructVIO – Technical details



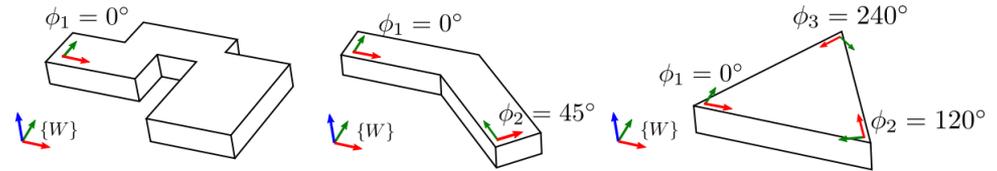
- Inside of the filter
  - Parameterization
  - Measurement equation
- Outside of the filter
  - Structural line related tasks:
    - Line detection & tracking
    - Classification of structural lines
    - initialization & triangulation
    - Handling long feature tracks
  - Manhattan world :
    - Detection
    - Merging
- Other details
  - Outlier rejection

# Coordinate frames



- World frame :  $\{W\}$

- Z axis aligned with gravity
- Starting point as the origin



- Local Manhattan frame:  $\phi_i \in [0, \pi/2), i = 1, \dots, N$

- Camera frame :  $\{C\}$

- Z axis aligned with the optical axis toward the viewing direction.
- X, Y axes aligned with x,y axes of the image

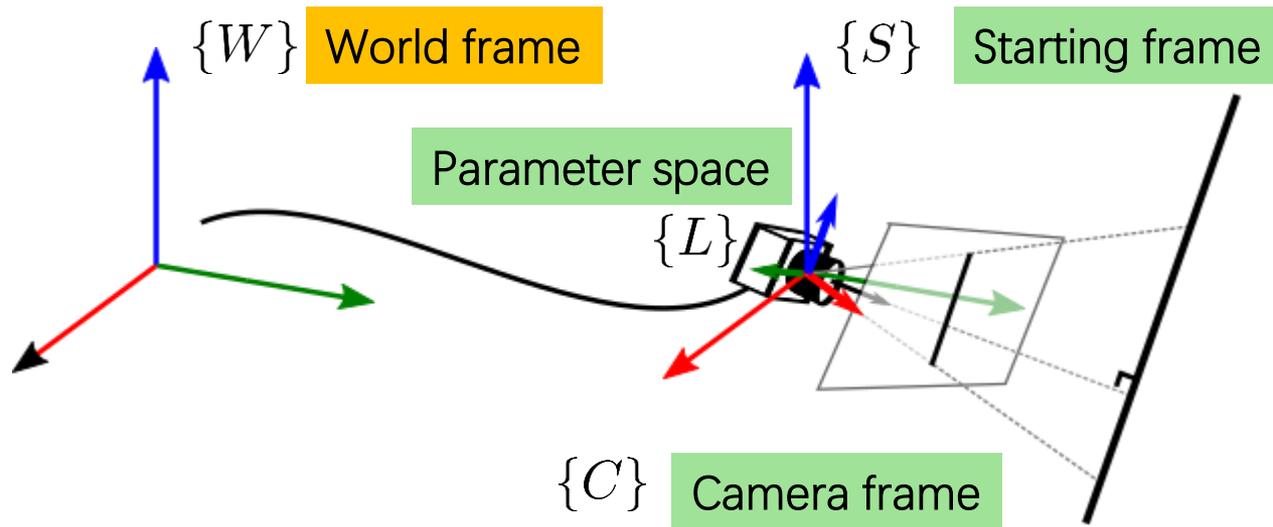
- Starting frame:  $\{S\}$  - Moving Manhattan frame

- The origin is located at the camera center.
- Three axes aligned with those of local Manhattan frame

# Representation of a structural line



- We use a camera-centric representation.



- Parameter space :  $\{L\}$  - use for line representation

# Structural line parameter space



- In parameter space  $\{L\}$ , a structural line can be represented by a point and a vertical direction.

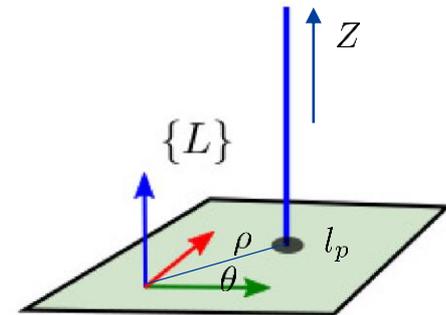
$$l_p \sim (a, b, 0)^T \quad Z = (0, 0, 1)^T$$

- To achieve better linearization, the intersection point can be represented using inverse-depth approach. We have

$$l_p \sim (\theta, \rho, 0)^T$$

$$\theta = \text{atan2}(b, a)$$

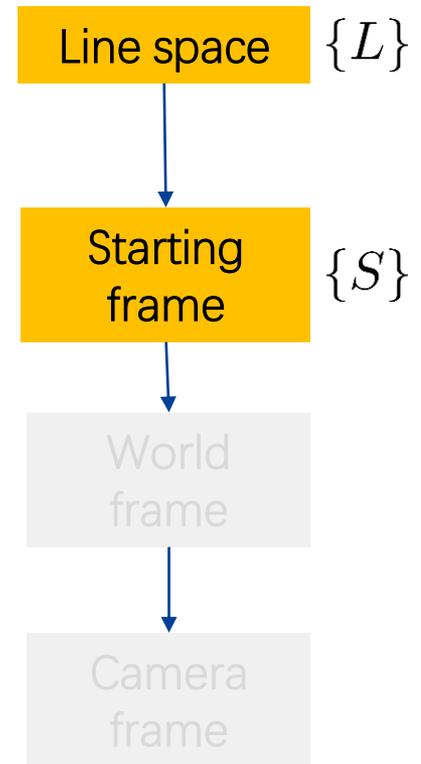
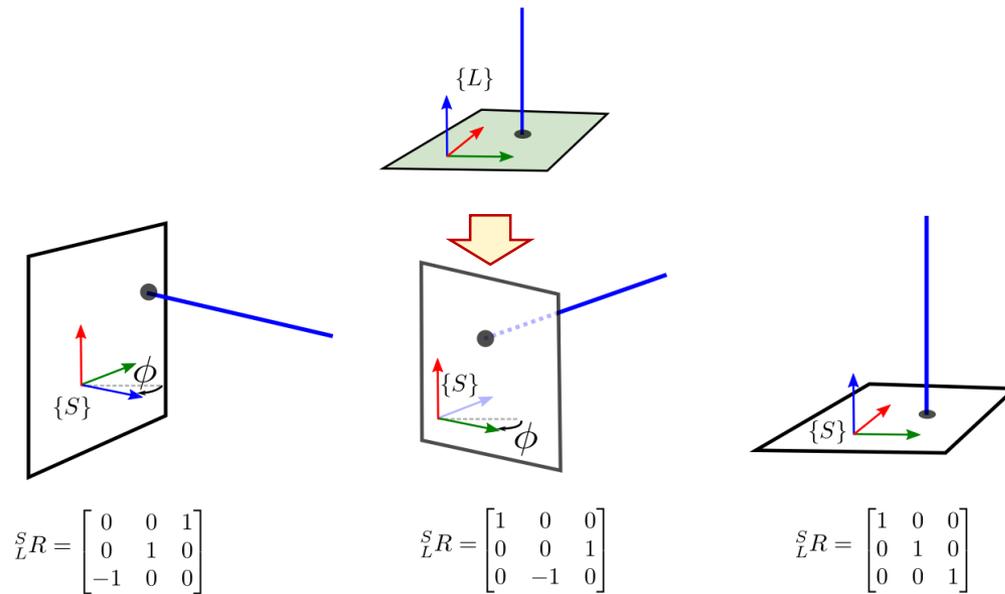
$$\rho = 1/\sqrt{a^2 + b^2}$$



# Line space -> Starting frame



- The structural line can be transformed into three axes of the starting frame  $\{S\}$  by the rotation  ${}^S_L R$ .



# Starting frame -> World frame

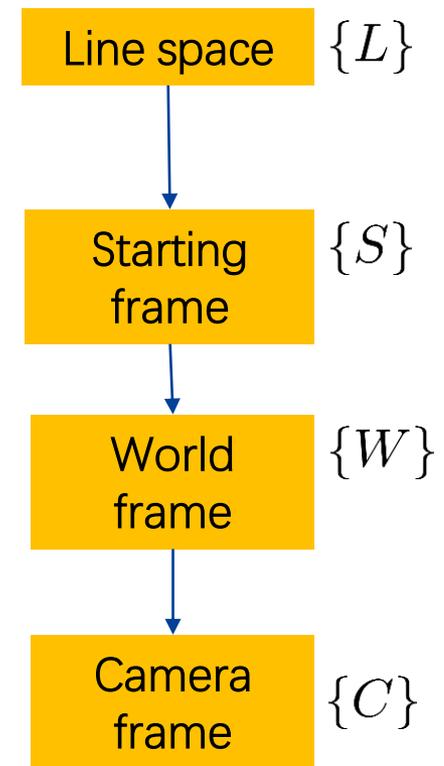


- The structural line can be further transformed into the world frame by using the heading direction ( $\phi_i$ ) of the local Manhattan world.

$${}^W_S R(\phi_i) = \begin{bmatrix} \cos(\phi_i) & \sin(\phi_i) & 0 \\ -\sin(\phi_i) & \cos(\phi_i) & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

- The structural line is then transformed to the current camera frame by.

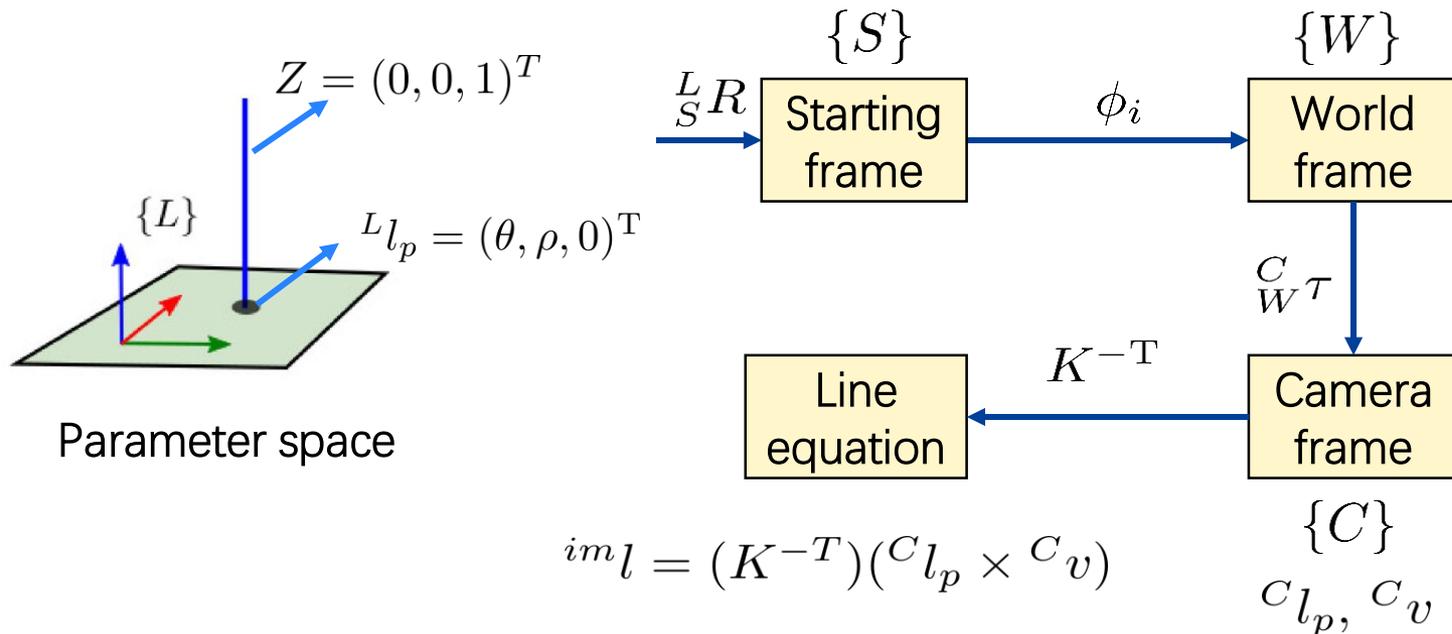
$${}^C_W \tau = ({}^C_W R, {}^C p_W)$$



# Line projection on the image



- Apply the transformations to both the point  $l_p$  and the vertical direction  $Z$



# Line projection on the image



- Line projection can be written as the following functions

$${}^{im}l = \Pi(l, \phi_i, {}^S_L R, {}^W_C \tau)$$

$${}^{im}l = \Pi(l, \phi_i, {}^S_L R, {}^I_C \tau, {}^W_I \tau) \quad (\text{unknown camera-IMU calibration } {}^I_C \tau)$$

, where  ${}^S_L R$  are known constants after line direction classification.

- Hence we further write

$${}^{im}l = \Pi(l, \phi_i, {}^W_C \tau)$$

$${}^{im}l = \Pi(l, \phi_i, {}^I_C \tau, {}^W_C \tau)$$

- We can use the above functions to derive the measurement equations.

# Measurement equations



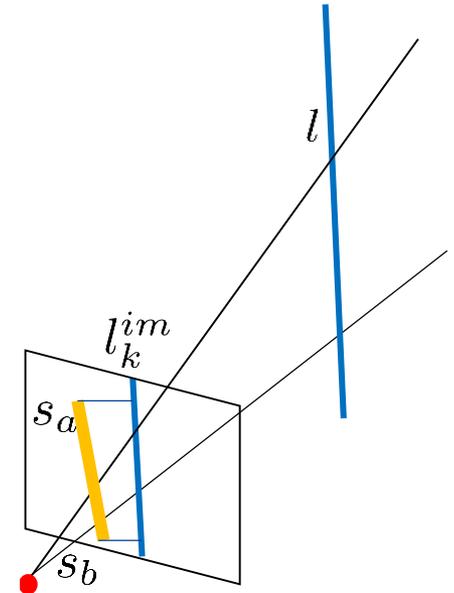
- Measurement equation by re-projection errors

- The line projection at time  $k$  is given by:

$${}^{im}l_k = \Pi(l, \phi_i, {}^I_C \tau, {}^W_C \tau)$$

- The line segment detected on the image is denoted by :  $s_a \leftrightarrow s_b$
- Hence the re-projection error can be computed as the signed distances between the line projection and the two end points:

$$r_k = D(s_a, s_b, {}^{im}l) = D(s_a, s_b, \Pi(l, \phi_i, {}^I_C \tau, {}^{I_k} \tau))$$



# Measurement equations



- After local linearization, we have

$$r_k = h_0 + J_l \delta l + J_\phi \delta \phi + J_{IC} \delta_C^I \tau + J_{WI_k} \delta_{I_k}^W \tau$$

- By stacking all observations from time 1 to time M

$$r_1 = h_0 + J_l \delta l + J_\phi \delta \phi + J_{IC} \delta_C^I \tau + J_{WI_1} \delta_{I_1}^W \tau$$

... ..

$$r_k = h_0 + J_l \delta l + J_\phi \delta \phi + J_{IC} \delta_C^I \tau + J_{WI_k} \delta_{I_k}^W \tau$$

... ..

$$r_M = h_0 + J_l \delta l + J_\phi \delta \phi + J_{IC} \delta_C^I \tau + J_{WI_M} \delta_{I_M}^W \tau$$



$$z = H_l \boxed{\delta l} + H_\phi \boxed{\delta \phi} + H_{CI} \boxed{\delta_C^I \tau} + H_{WI} \boxed{[\delta_{I_1}^W \tau \cdots \delta_{I_M}^W \tau]}$$

Line  
parameters

Heading of  
Manhattan  
world

Camera-  
IMU  
calibration

Historical IMU  
poses

# Measurement equations



- Project the residual to the left null space of  $H_l$ , we can get rid of the line parameters:

$$z = H_l \delta l + H_\phi \delta \phi + H_{CI} \delta_C^I \tau + H_{WI} [\delta_{I_1}^W \tau \cdots \delta_{I_M}^W \tau]$$



$$z^{(0)} = H_\phi^{(0)} \delta \phi + H_{CI}^{(0)} \delta^C x_I + H_{WI}^{(0)} [\delta_{I_1}^W \tau \cdots \delta_{I_M}^W \tau],$$

- The measurement equation involves
  - 1. Heading direction of the local Manhattan world
  - 2. IMU-camera relative pose
  - 3. Historical IMU poses

# Outside of the filter



- Structural line related tasks:
  - Line detection & classification of structural lines
  - initialization & triangulation
  - Line tracking
  - Handling long feature tracks

# Outside of the filter

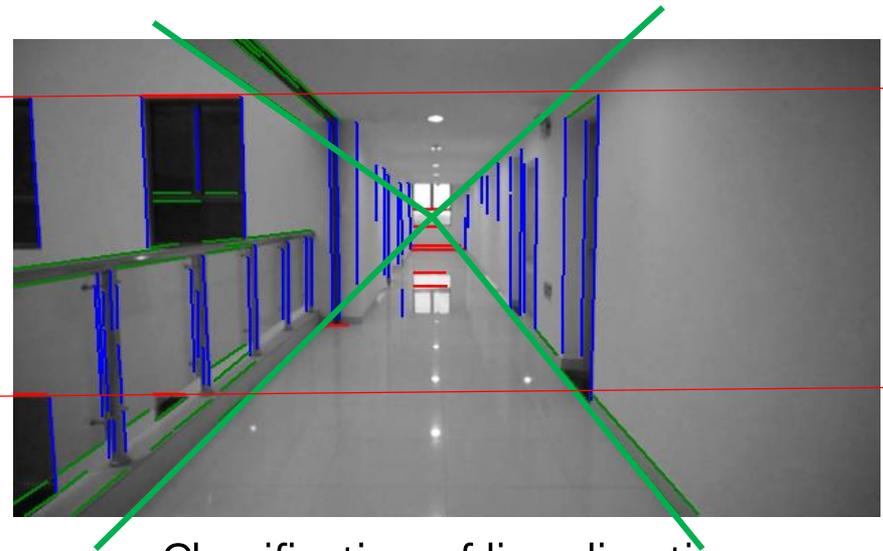


- Structural line related tasks:
  - Line detection & classification of structural lines

For a line segment  $s_a \leftrightarrow s_b$  find its Manhattan world  $\phi_i$  and its direction (X,Y,orZ)



Detection of line segments



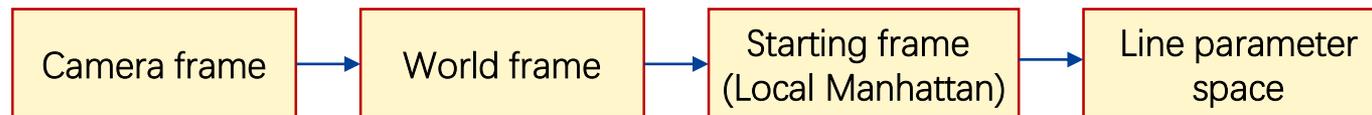
Classification of line directions (X,Y,or Z) and identify the Manhattan world  $\phi_i$

# Outside of the filter



- Structural line related tasks:
  - Initialization
    1. Longer line segment first
    2. Establish the starting frame (in which Manhattan world it lies)
    3. Use the middle point  $m$  of the line segment for initialization

For horizontal lines (aligned with X, Y axes of a certain Manhattan frame)



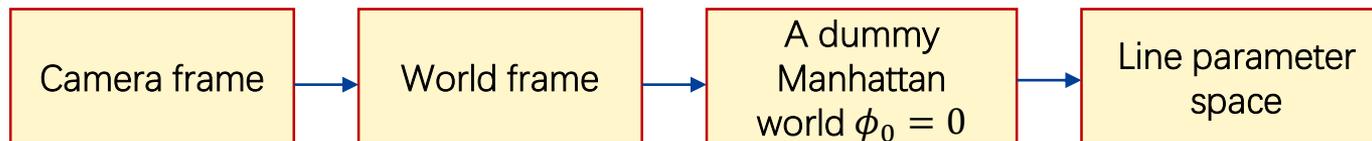
$${}^L m = {}^L_S R {}^S_W R(\phi_i) {}^W_C R K^{-1} m$$

# Outside of the filter



- Structural line related tasks:
  - Initialization
    1. Longer line segment first
    2. Establish the starting frame (in which Manhattan world it lies)
    3. Use the middle point  $m$  of the line segment for initialization

For vertical lines (aligned with Z axis of any Manhattan frames)



$$L_m = {}^W_C R K^{-1} m$$

# Outside of the filter



- We can write the initialization process as a function:

$$l_0 = \begin{bmatrix} \theta_0 \\ \rho_0 \end{bmatrix} = \begin{bmatrix} \Pi^{-1}(s, \phi_i, {}^W_C R) \\ \rho_0 \end{bmatrix}$$

Initial parameters

$s$ : line segment

$\phi_i$ : Local Manhattan frame

${}^W_C R$ : Current camera orientation

$\rho_0$ : Initial inverse depth

$$\Sigma_0 = \begin{bmatrix} \sigma_{\theta_0}^2 & 0 \\ 0 & \sigma_{\rho_0}^2 \end{bmatrix} \quad \begin{array}{l} \sigma_{\theta_0}^2: \text{small value to account line detection error (2-4 pixels)} \\ \sigma_{\rho_0}^2: \text{uncertainty of inverse depth (5 by default)} \end{array}$$

Initial covariance

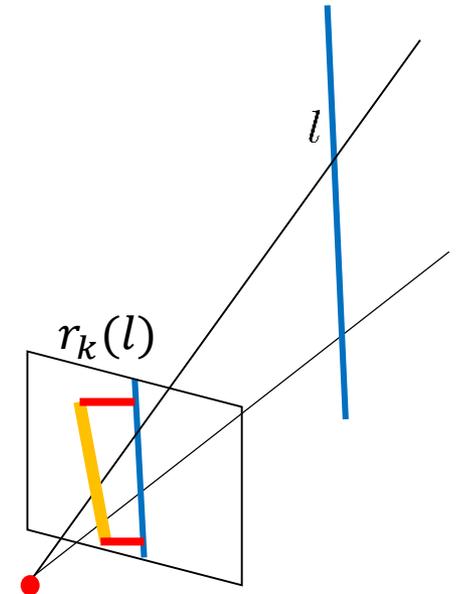
# Outside of the filter



- Line triangulation with prior Knowledge

	Line projection error	Prior knowledge
$\arg \min_{l=(\theta, \rho)^T}$	$\sum_{k \in \mathcal{V}} r_k^2(l) / \sigma_{im}^2$	$+ (l - l_0)^T \Sigma_0^{-1} (l - l_0)$

$l_0$  : Prior line parameters  
 $r_k(l)$  : line projection error  
 $\mathcal{V}$  : set of visible views

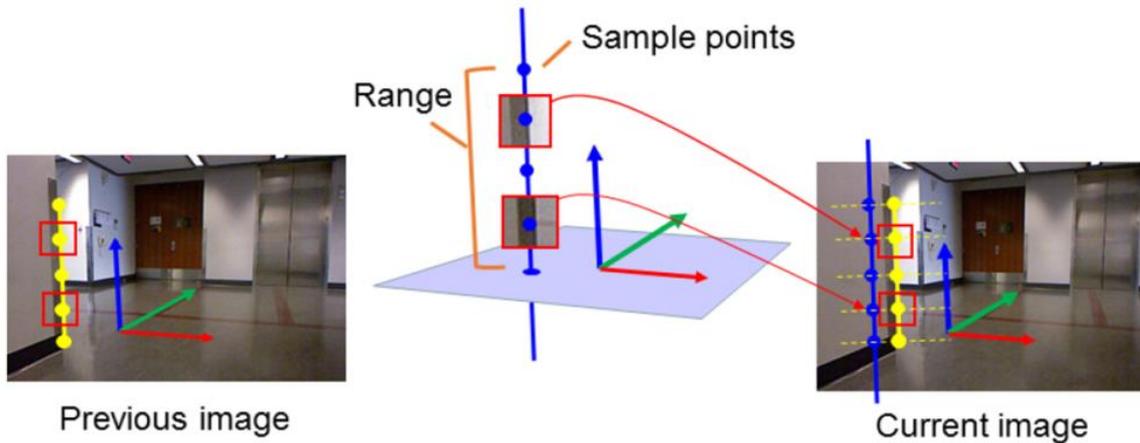


# Outside of the filter



## Line tracking

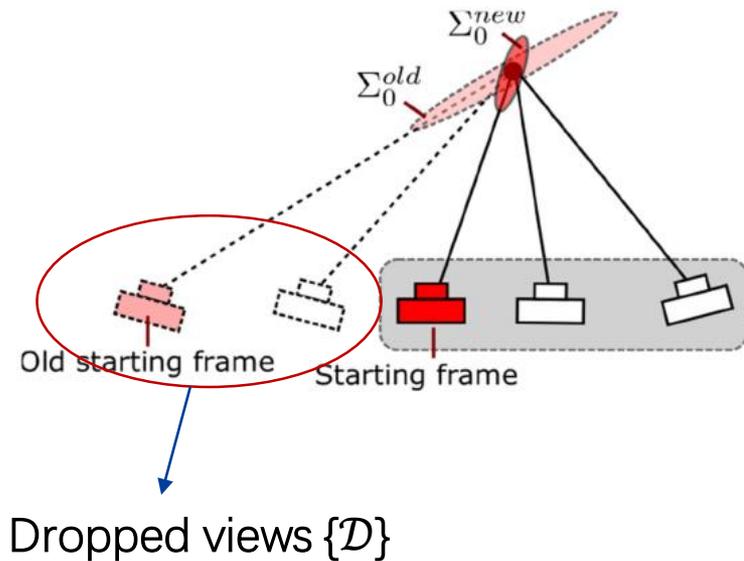
- 1. Sample **several points** on the line
- 2. Project those points onto the image, searching corresponding points perpendicular to the line projection.
- 3. Use the small patches around those points as the descriptor



# Outside of the filter



- Handling long feature tracks



Step1 – Absorb dropped measurements into priori information :

$$\arg \min_{l=(\theta, \rho)^T} \sum_{k \in D} r_k^2(l) / \sigma_{im}^2 + (l - l_0^{\text{old}})^T (\Sigma_0^{\text{old}})^{-1} (l - l_0^{\text{old}})$$



$$\Lambda \delta l = Y \quad \text{Normal equation in the last Gauss-Newton iteration}$$



$$\Sigma_0^{\text{new}} \leftarrow \Lambda^{-1}$$

Step2 – Change the starting frame  $S \rightarrow S'$

$l' \leftarrow \frac{S'}{S} \mathcal{T}(l), l'_0 \leftarrow \frac{S'}{S} \mathcal{T}(l_0)$	$\Sigma' \leftarrow J \Sigma J^T, \Sigma'_0 \leftarrow J \Sigma_0 J^T$
--	--

Current estimate

Prior information

# Outside of the filter



- Manhattan world detection :
  - 1. starts once vertical lines are identified
  - 2. compute the horizontal line  $l_\infty = K^{-T} C_W R [0,0,1]^T$
  - 3. run 1-line RANSAC to detect one of the two horizontal directions (X or Y)
    - Randomly select one line, extended it to intersect  $l_\infty$  to get a vanishing point  $v_x$
    - Compute the other vanishing point  $v_y$
    - Check the consistent line segments aligned with  $v_x$  or  $v_y$
    - Repeat the aforementioned steps
  - It is a **possible Manhattan world** if the maximum consensus set contains sufficient inliers.

# Outside of the filter



- Manhattan world merging :

- The heading direction of two Manhattan worlds could be very close.

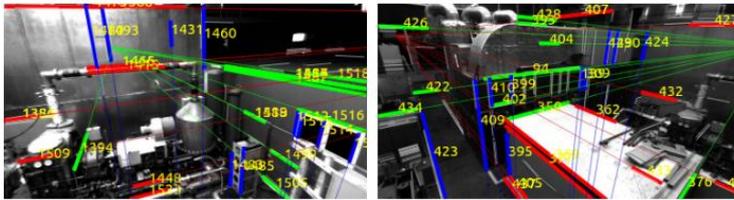
$$|\phi_i - \phi_j| < \Delta\phi$$

- We merge them by removing the newly detected one and update the information of related structural lines

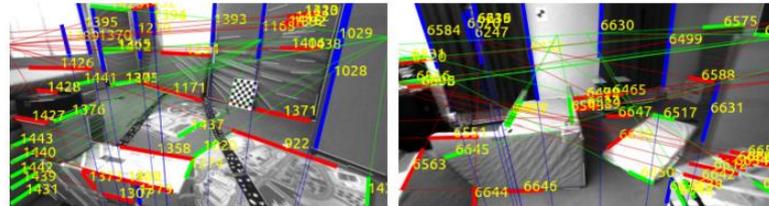
# Results



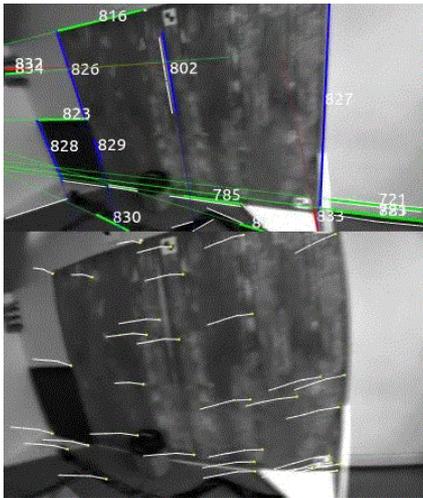
- Benchmark tests on Euroc dataset



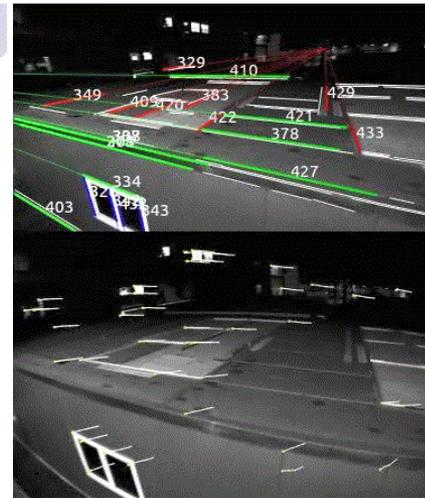
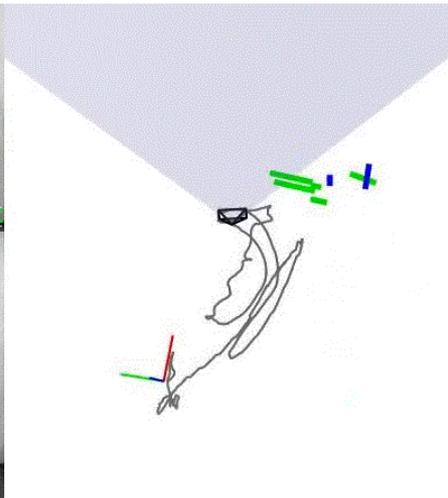
(a) Machine hall



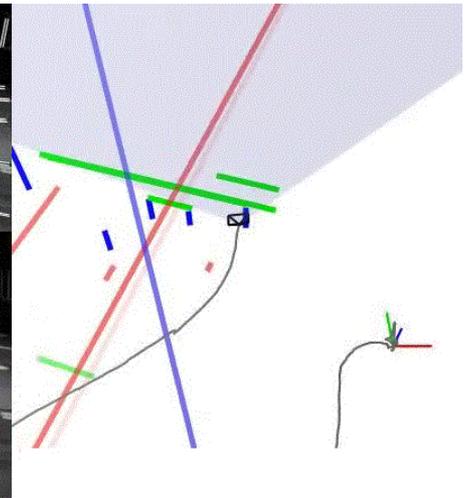
(b) Vicon room



V2\_03\_difficult



MH\_05\_difficult



# Results



- Euroc dataset

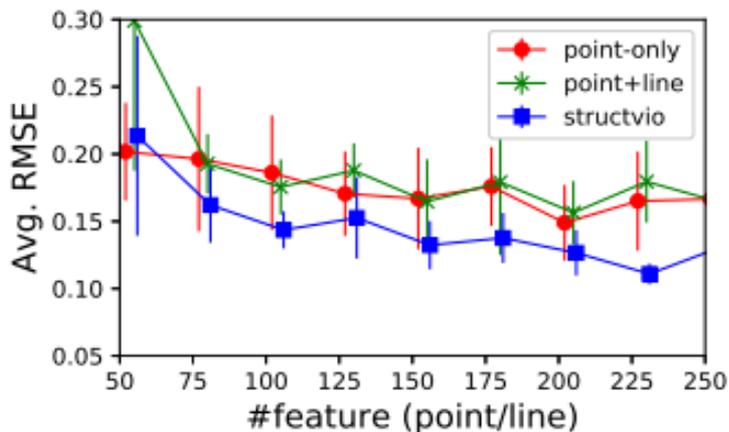
Dataset	OKVIS[5]		VINS[6](w/o loop)		StructVIO	
	RMSE	Max.	RMSE	Max.	RMSE	Max.
MH_01_easy	0.308	0.597	<b>0.157</b> <sup>2</sup>	0.349	<b>0.079</b> <sup>1</sup>	0.251
MH_02_easy	0.407	0.811	<b>0.181</b> <sup>2</sup>	0.533	<b>0.145</b> <sup>1</sup>	0.267
MH_03_medium	0.241	0.411	<b>0.196</b> <sup>2</sup>	0.450	<b>0.103</b> <sup>1</sup>	0.271
MH_04_difficult	0.363	0.641	<b>0.345</b> <sup>2</sup>	0.475	<b>0.130</b> <sup>1</sup>	0.286
MH_05_difficult	0.439	0.751	<b>0.303</b> <sup>2</sup>	0.434	<b>0.182</b> <sup>1</sup>	0.358
V1_01_easy	<b>0.076</b> <sup>2</sup>	0.224	0.090	0.201	<b>0.060</b> <sup>1</sup>	0.180
V1_02_medium	0.141	0.254	<b>0.098</b> <sup>1</sup>	0.334	<b>0.130</b> <sup>2</sup>	0.260
V1_03_difficult	0.240	0.492	<b>0.183</b> <sup>2</sup>	0.376	<b>0.090</b> <sup>1</sup>	0.263
V2_01_easy	0.134	0.308	<b>0.080</b> <sup>2</sup>	0.232	<b>0.045</b> <sup>1</sup>	0.140
V2_02_medium	0.187	0.407	<b>0.149</b> <sup>2</sup>	0.379	<b>0.066</b> <sup>1</sup>	0.157
V2_03_difficult	<b>0.255</b> <sup>2</sup>	0.606	0.268	0.627	<b>0.110</b> <sup>1</sup>	0.231

RMSE-Rooted Mean Squared Error

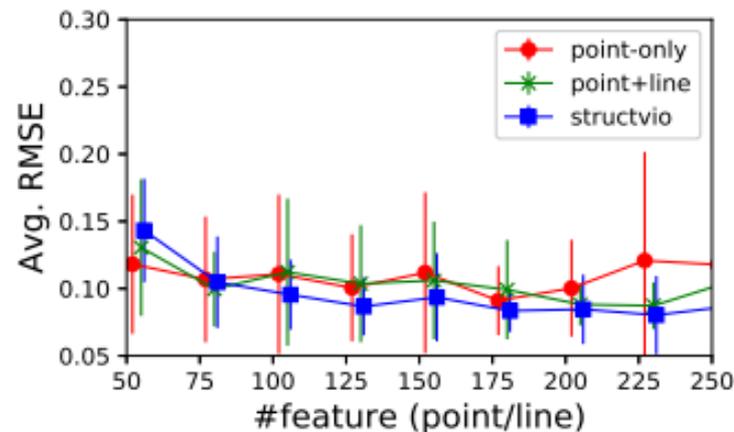
# Results



- Euroc datasets
  - StructVIO performs better in Machine hall, since it exhibit stronger structural regularity.



(a) Machine hall

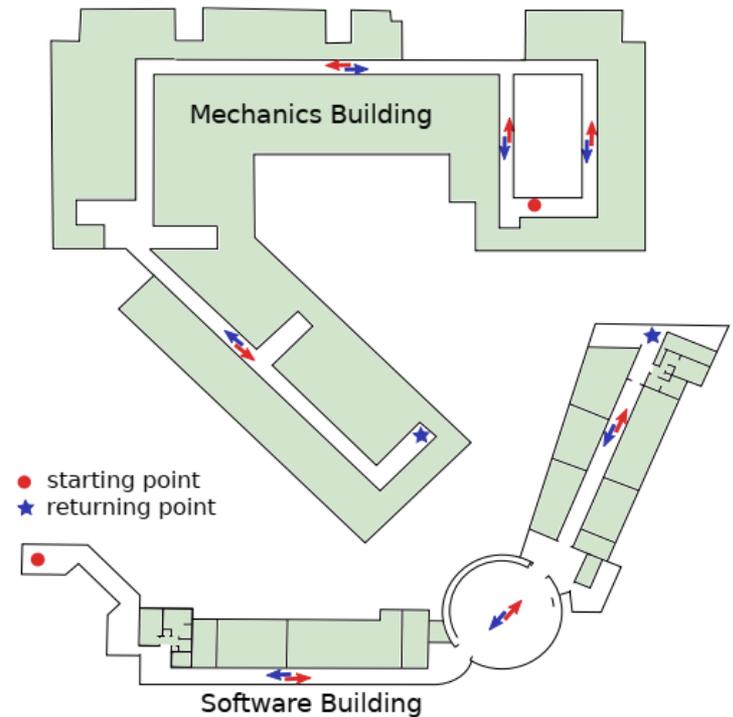


(b) Vicon room

# Results



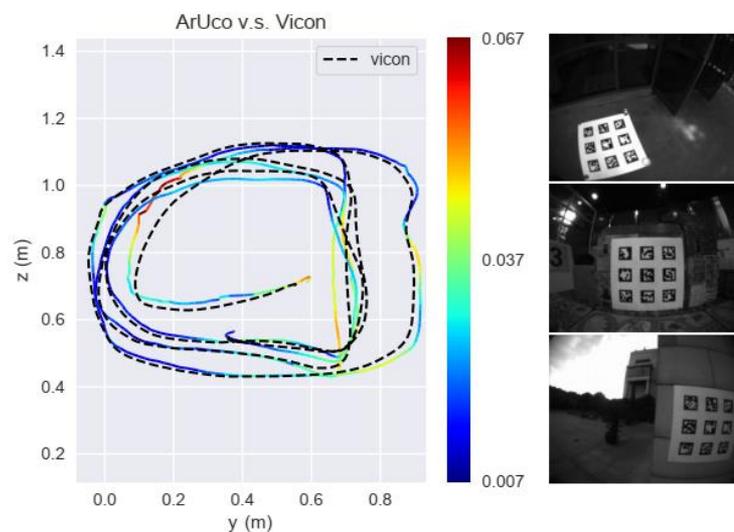
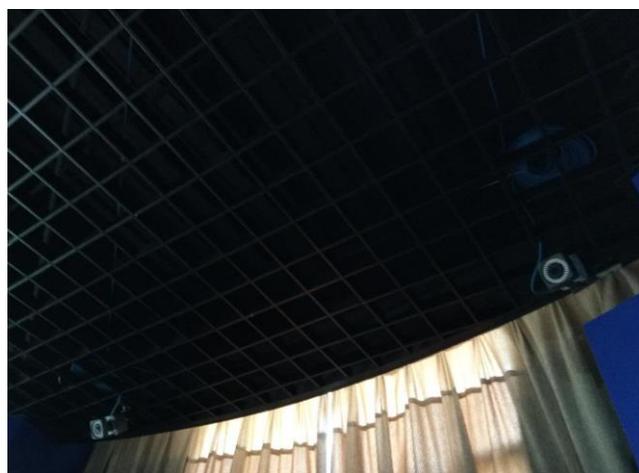
- Visual-inertial data collected by Google Tango Tablet (16 test sequences)
- Different buildings in SJTU campus.
- Indoor/Outdoor, Large illumination changes, 5~10 minutes walking



# Results



- Ground truth data were collected by either Vicon or ArUco code.



Starting segment:  $s$   
Ending segment:  $e$

Align the starting segments :  $T^* = \arg \min_T \sum_{t \in s} (\|T(p^t) - g_s^t\|^2)$

Compute the ending segment's RMSE and Max errors:

$$RMSE = \sqrt{\frac{1}{|e|} \sum_{t \in e} \|T(p^t) - g_e^t\|^2} \quad Max. = \max |T(p^t) - g_e^t|$$

# Results



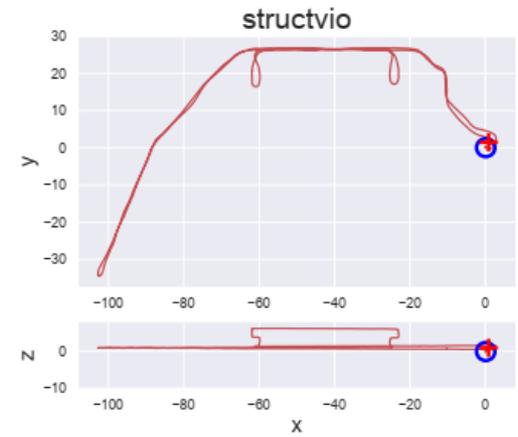
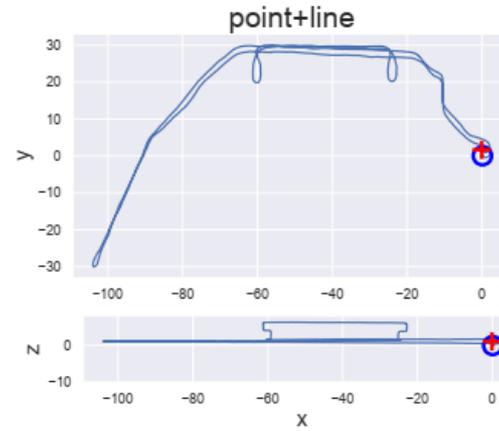
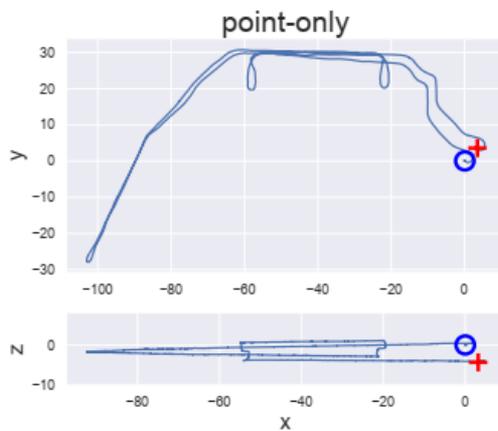
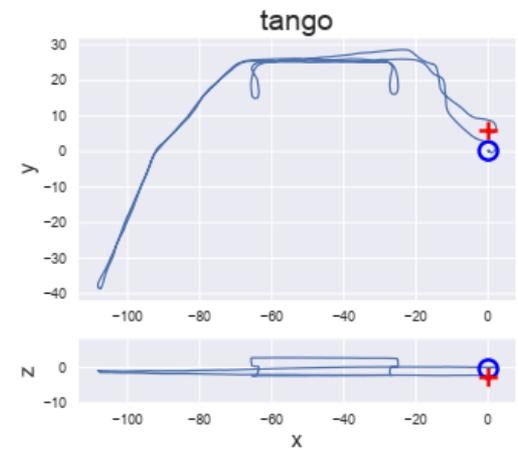
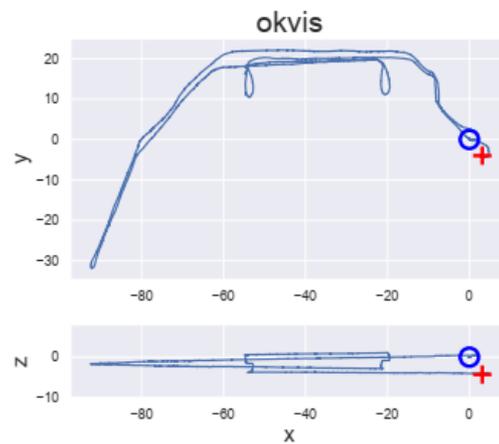
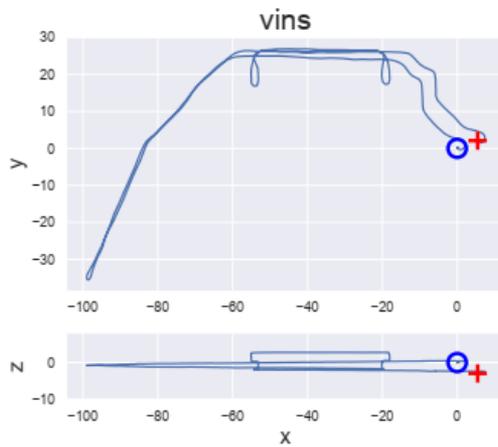
- Methods: OKVIS, VINS, Project Tango, Point-only, Point-line, StructVIO

Seq. Name	Traveling Dist. [m]	OKVIS <sup>[5]</sup>		VINS <sup>[6]</sup> (w/o loop)		Project Tango		Point-only		Point+Line		StructVIO	
		RMSE	Max.	RMSE	Max.	RMSE	Max.	RMSE	Max.	RMSE	Max.	RMSE	Max
Soft-01	315.167	6.702	9.619	4.861	6.688	5.715	8.181	<b>2.153</b> <sup>2</sup>	2.728	2.262	2.842	<b>1.931</b> <sup>1</sup>	2.437
Soft-02	438.198	4.623	6.713	2.713	4.086	4.238	6.226	3.905	5.243	<b>1.468</b> <sup>2</sup>	2.026	<b>1.429</b> <sup>1</sup>	1.984
Soft-03	347.966	<b>4.505</b> <sup>2</sup>	6.223	7.270	9.832	167.825	228.630	6.515	8.119	8.618	10.790	<b>0.325</b> <sup>1</sup>	1.020
Soft-04	400.356	3.993	5.784	28.667	75.479	2.453	3.544	<b>1.550</b> <sup>1</sup>	2.028	4.051	5.262	<b>1.722</b> <sup>2</sup>	2.241
Mech-01	340.578	3.627	4.745	2.452	3.260	<b>1.948</b> <sup>2</sup>	2.726	3.298	3.961	4.323	5.181	<b>0.909</b> <sup>1</sup>	1.165
Mech-02	388.548	3.079	4.195	3.570	4.754	<b>1.596</b> <sup>2</sup>	2.217	1.663	2.108	2.317	2.927	<b>0.779</b> <sup>1</sup>	1.022
Mech-03	317.974	3.875	5.324	4.682	9.113	4.220	5.781	<b>2.384</b> <sup>2</sup>	3.020	4.193	5.272	<b>1.161</b> <sup>1</sup>	1.532
Mech-04	650.430	-	-	3.002	8.592	1.915	5.808	1.785	4.663	<b>1.425</b> <sup>2</sup>	3.729	<b>0.742</b> <sup>1</sup>	1.940
MicroA-01	257.586	2.485	3.382	<b>0.654</b> <sup>2</sup>	1.148	45.599	61.058	2.849	3.505	2.189	2.721	<b>0.642</b> <sup>1</sup>	1.225
MicroA-02	190.203	3.428	5.186	14.222	57.172	<b>1.145</b> <sup>1</sup>	1.692	1.964	2.514	<b>1.723</b> <sup>2</sup>	2.207	2.089	2.661
MicroA-03	388.730	0.078	0.779	<b>1.800</b> <sup>1</sup>	2.578	4.400	6.253	3.824	5.169	3.072	4.232	<b>1.884</b> <sup>2</sup>	2.892
MicroA-04	237.856	6.136	8.532	<b>0.994</b> <sup>2</sup>	1.765	55.200	75.318	2.056	2.897	2.406	2.879	<b>0.350</b> <sup>1</sup>	0.448
MicroB-01	338.962	2.898	4.025	<b>1.856</b> <sup>2</sup>	2.944	38.197	50.572	7.084	8.576	7.337	8.913	<b>1.477</b> <sup>1</sup>	1.902
MicroB-02	306.316	2.240	3.490	<b>1.030</b> <sup>2</sup>	2.431	5.660	8.652	2.521	3.714	3.197	4.610	<b>0.470</b> <sup>1</sup>	0.799
MicroB-03	485.291	-	-	2.132	3.368	<b>2.009</b> <sup>2</sup>	2.960	6.490	8.978	4.507	6.301	<b>0.445</b> <sup>1</sup>	0.675
MicroB-04	357.251	4.064	6.481	<b>1.332</b> <sup>2</sup>	2.068	13.962	22.028	5.078	7.713	1.977	3.074	<b>0.473</b> <sup>1</sup>	0.777
Mean Drift Err.(%)		1.078%		1.410%		6.180%		0.957%		<b>0.956%</b> <sup>2</sup>		<b>0.292%</b> <sup>1</sup>	
Median Drift Err.(%)		0.781%		<b>0.538%</b> <sup>2</sup>		0.900%		0.559%		0.570%		<b>0.176%</b> <sup>1</sup>	

# Results



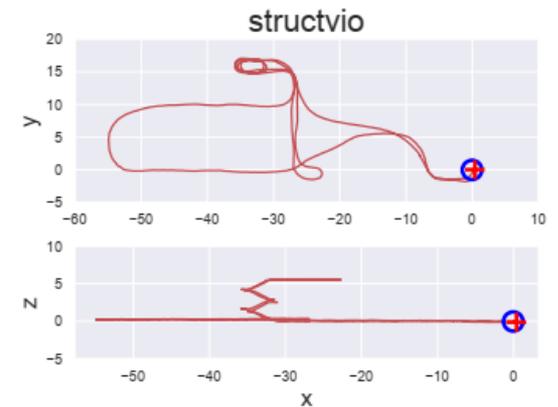
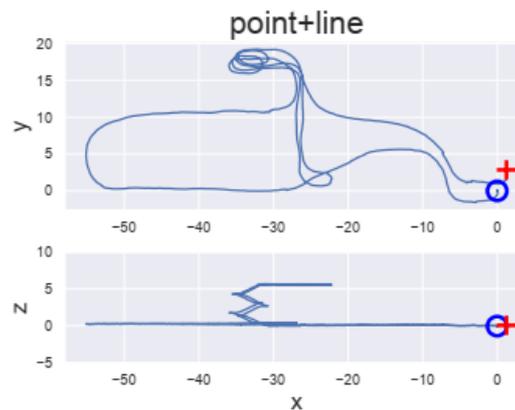
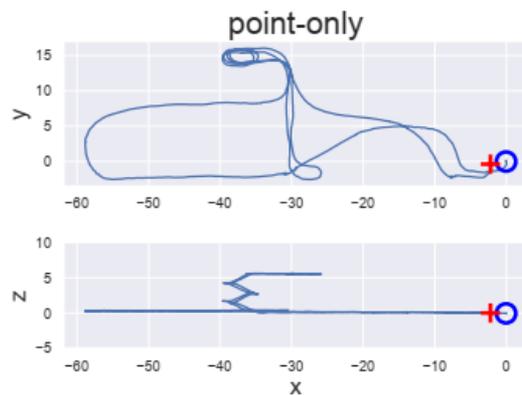
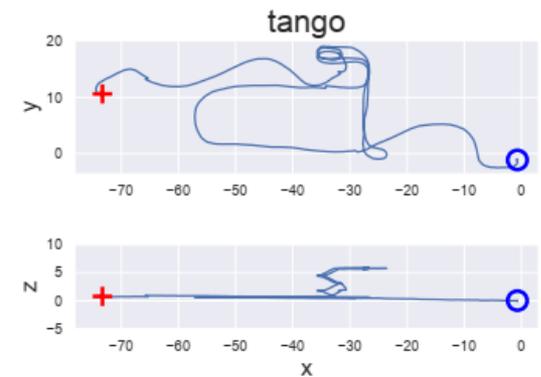
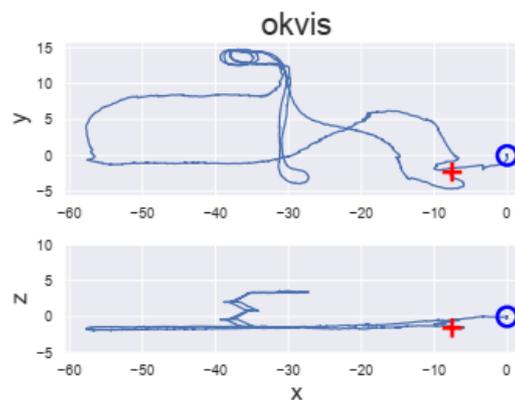
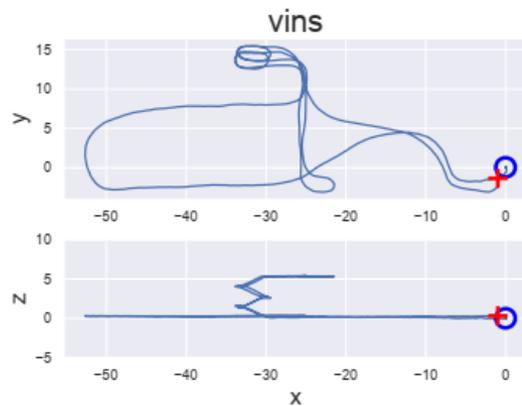
- Software building (Soft-02)



# Results



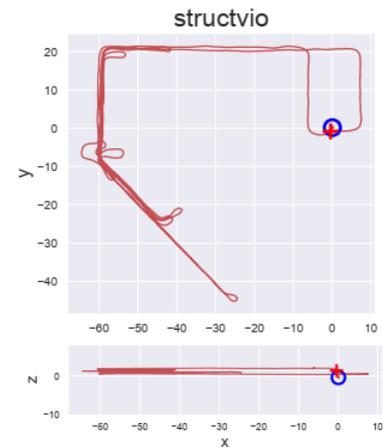
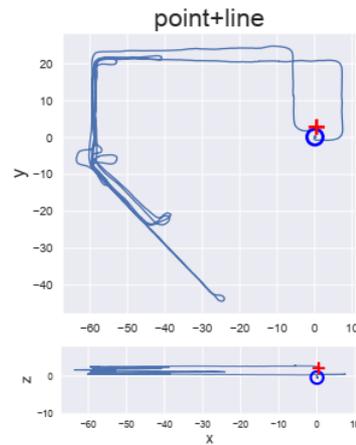
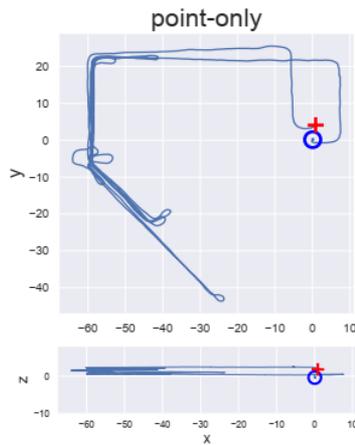
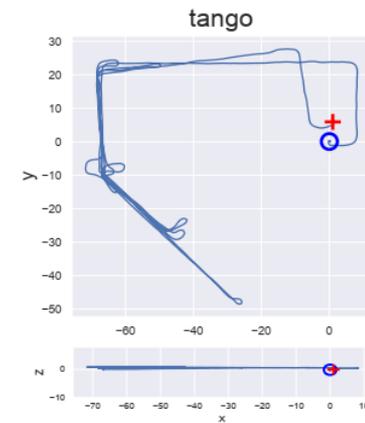
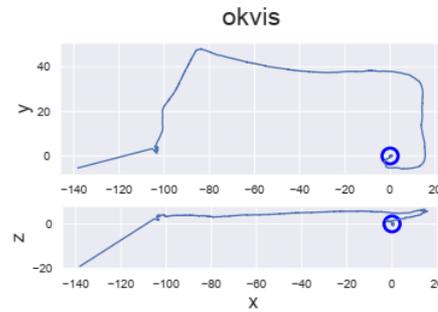
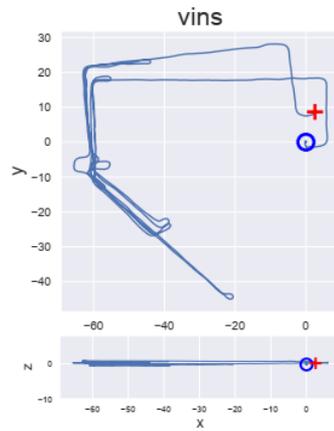
- Micro Electronic Engineering Building (MicroA-04)



# Results



- Mechanical Engineering Building (Mech-04)



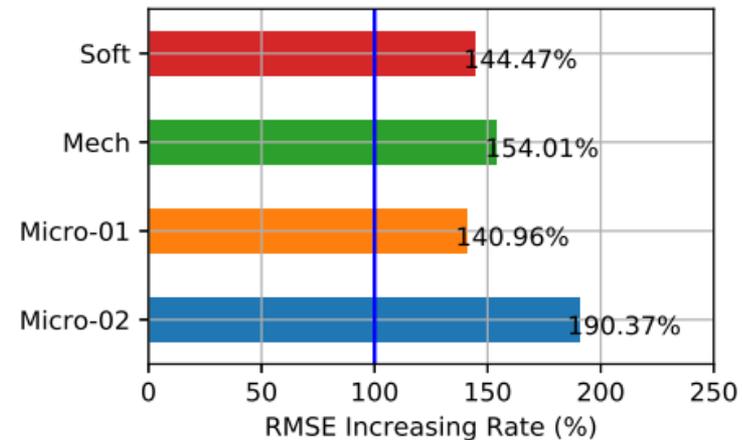
# Results



- Other tests

Seq. Name	Atlata world		Manhattan world	
	RMSE	Max.	RMSE	Max.
Mech-01	0.909	1.165	1.144	1.524
Mech-02	0.779	1.022	1.286	1.061
Mech-03	1.161	1.532	2.029	1.211
Mech-04	0.742	1.940	1.822	2.193
Soft-01	1.931	2.437	2.896	2.397
Soft-02	1.429	1.984	3.092	4.149
Soft-03	0.325	1.020	3.352	4.236
Soft-04	1.722	2.241	3.178	4.120

Atlanta world vs Manhattan world



Without dealing with  
dropped measurements for  
long feature tracks

- Software, tools, and dataset:

<http://drone.sjtu.edu.cn/dpzou/project/structvio.html>

# Software usage



- Supports the following modes
  - Point/Line/Structural line-only
  - Point+line
  - Point+Structural line

```
-l <0|1|2>, --line_type <0|1|2>  
  Type of lines used: 0-structlines, 1-general lines, 2-both  
  
-f <0|1|2>, --feature_type <0|1|2>  
  Type of features used: 0 - point, 1 - line, 2 - both
```

```
./structvio -i ./Soft-01 -n Soft-01 -r Soft-01-res -c structvio_data.yaml
```

- Script for evaluation (revised from evo)

```
vio_eva.py -r <path to the result file - 'state.txt'> -d <path to the folder of data - e.g. 'Mech-01'>
```



謝 謝