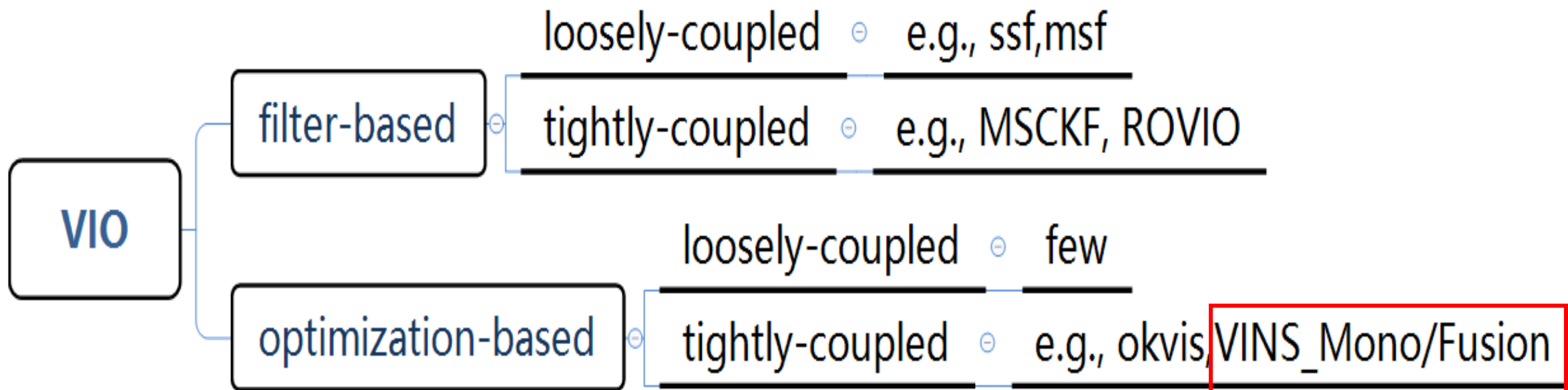# Exchange & Share about VINS

Tao Zhang

2019.6.6

# ◆VIO

- **VIO**: Vusual Inertial Odometer

- Classification

# ◆Background



HKUST Aerial Robotics Group

Home    Group ⌄    Publications    News
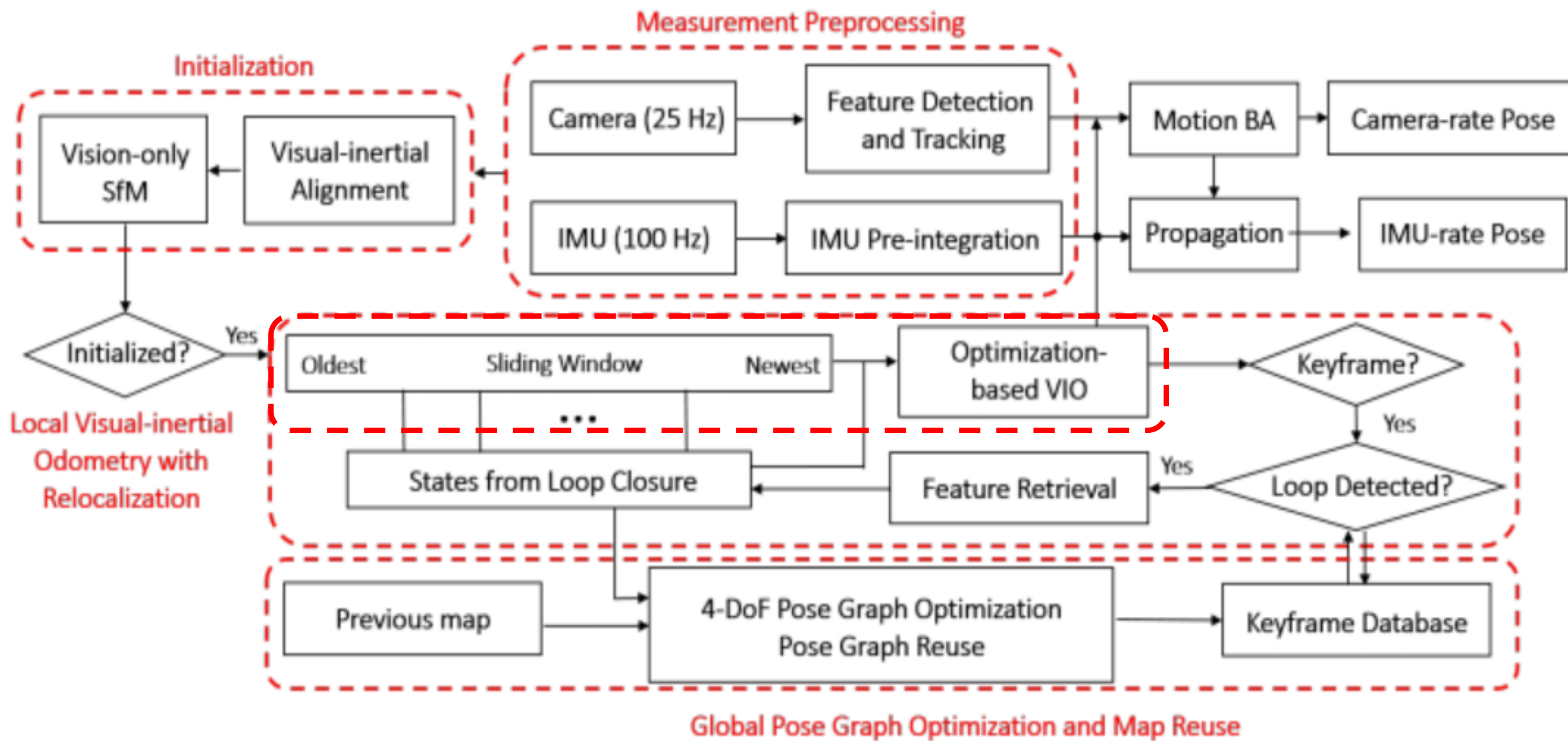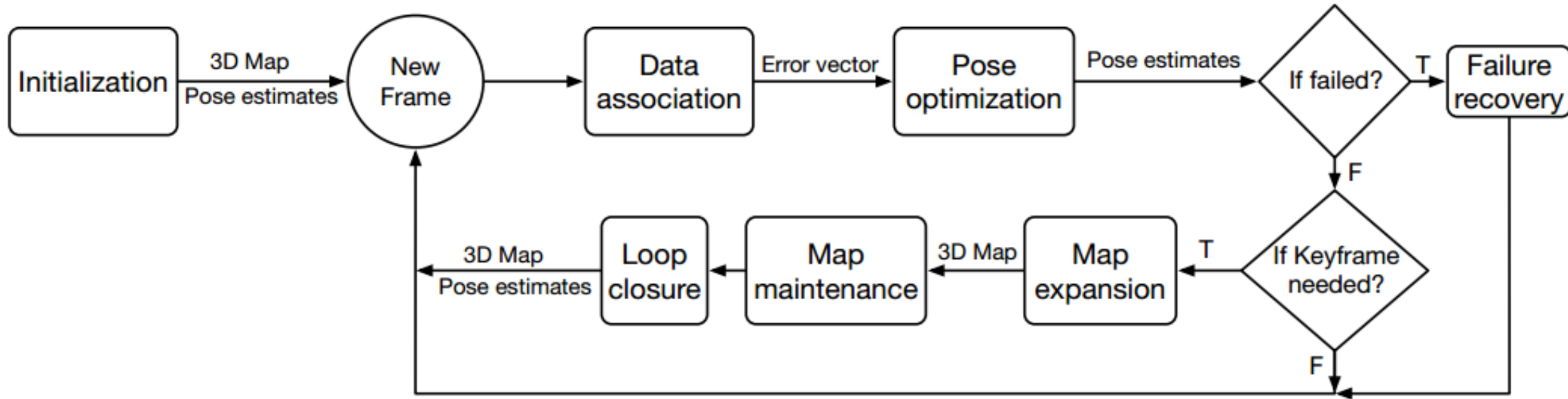
- VINS_Mono： VINS-Mono: A Robust and Versatile Monocular Visual-Inertial State Estimator. 2017.8
- VINS_Fusion： A General Optimization-based Framework for Local Odometry Estimation with Multiple Sensors
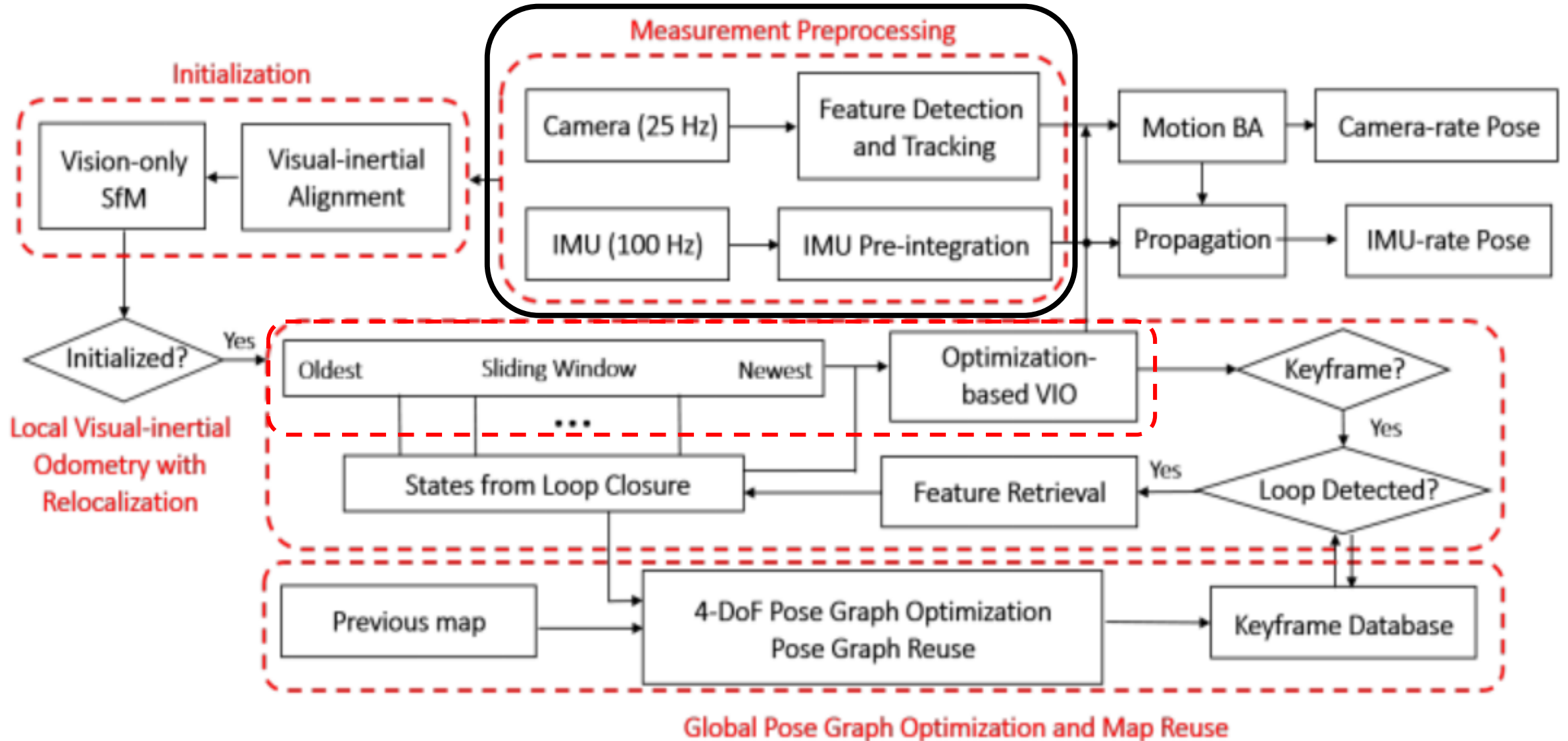  A General Optimization-based Framework for Global Pose Estimation with Multiple Sensors  2019.1

# ◆VINS_Mono Frame

# ◆Keyframe-SLAM flowchart

# ◆Measurement Preprocessing

## ◆vison front-end

●Feature processing

-Harris corners (goodFeaturesToTrack() in OpenCV)

-KLT sparse optical flow tracker(calcOpticalFlowPyrLK() in OpenCV)

-RANSAC outlier rejection

● Keyframe selection

-average parallax method(平均视差法)

　Rotation-compensated average feature parallax is larger than a threshold

-tracking quality method(跟踪质量法)

　Number of tracked features in the current frame is less than a threshold

# ◆IMU Pre-integration

$$\hat{\mathbf{a}}_t = \mathbf{a}_t + \mathbf{b}_{a_t} + \mathbf{R}_w^t \mathbf{g}^w + \mathbf{n}_a$$
$$\hat{\boldsymbol{\omega}}_t = \boldsymbol{\omega}_t + \mathbf{b}_{w_t} + \mathbf{n}_w. \tag{1}$$

raw measurements

$$\mathbf{p}_{b_{k+1}}^w = \mathbf{p}_{b_k}^w + \mathbf{v}_{b_k}^w \Delta t_k$$
$$+ \iint_{t \in [t_k, t_{k+1}]} (\boxed{\mathbf{R}_t^w}(\hat{\mathbf{a}}_t - \mathbf{b}_{a_t} - \mathbf{n}_a) - \mathbf{g}^w) \, dt^2$$

$$\mathbf{v}_{b_{k+1}}^w = \mathbf{v}_{b_k}^w + \int_{t \in [t_k, t_{k+1}]} (\boxed{\mathbf{R}_t^w}(\hat{\mathbf{a}}_t - \mathbf{b}_{a_t} - \mathbf{n}_a) - \mathbf{g}^w) \, dt$$

$$\mathbf{q}_{b_{k+1}}^w = \boxed{\mathbf{q}_{b_k}^w} \otimes \int_{t \in [t_k, t_{k+1}]} \frac{1}{2} \boldsymbol{\Omega}(\hat{\boldsymbol{\omega}}_t - \mathbf{b}_{w_t} - \mathbf{n}_w) \mathbf{q}_t^{b_k} \, dt, \tag{3}$$

integration in word frame**(traditional method**)

# ◆IMU Pre-integration

$$\mathbf{R}_w^{b_k}\mathbf{p}_{b_{k+1}}^w = \mathbf{R}_w^{b_k}(\mathbf{p}_{b_k}^w + \mathbf{v}_{b_k}^w\Delta t_k - \frac{1}{2}\mathbf{g}^w\Delta t_k^2) + \boldsymbol{\alpha}_{b_{k+1}}^{b_k}$$

$$\mathbf{R}_w^{b_k}\mathbf{v}_{b_{k+1}}^w = \mathbf{R}_w^{b_k}(\mathbf{v}_{b_k}^w - \mathbf{g}^w\Delta t_k) + \boldsymbol{\beta}_{b_{k+1}}^{b_k}$$

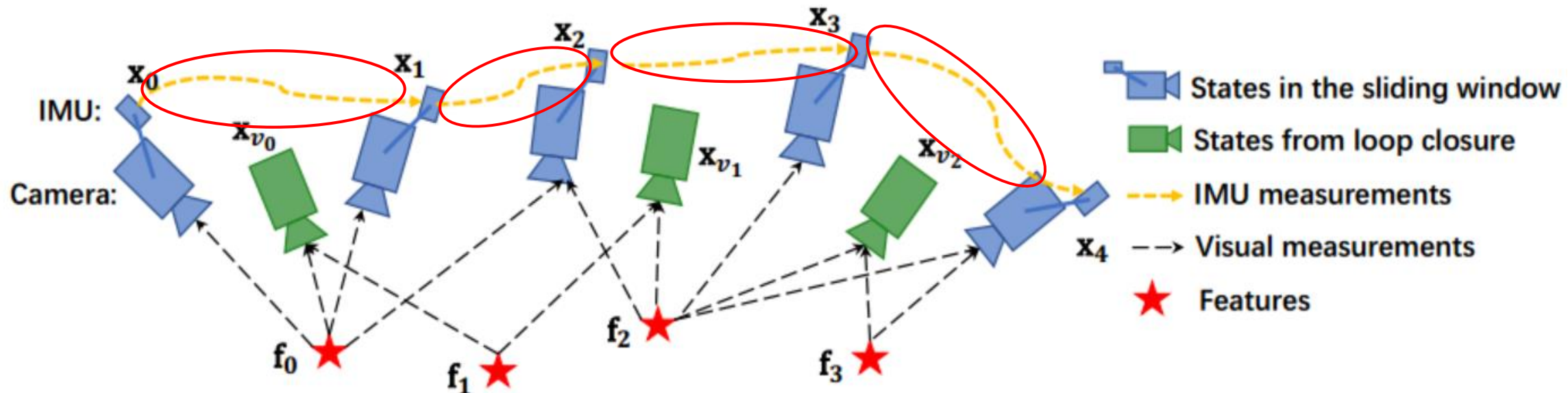$$\mathbf{q}_w^{b_k}\otimes\mathbf{q}_{b_{k+1}}^w = \boldsymbol{\gamma}_{b_{k+1}}^{b_k},$$

(5)

Process the integration from one imu frame to another imu frame

$$\boldsymbol{\alpha}_{b_{k+1}}^{b_k} = \iint_{t\in[t_k,t_{k+1}]}\mathbf{R}_t^{b_k}(\hat{\mathbf{a}}_t - \mathbf{b}_{a_t} - \mathbf{n}_a)dt^2$$

$$\boldsymbol{\beta}_{b_{k+1}}^{b_k} = \int_{t\in[t_k,t_{k+1}]}\mathbf{R}_t^{b_k}(\hat{\mathbf{a}}_t - \mathbf{b}_{a_t} - \mathbf{n}_a)dt$$

(6)

$$\boldsymbol{\gamma}_{b_{k+1}}^{b_k} = \int_{t\in[t_k,t_{k+1}]}\frac{1}{2}\boldsymbol{\Omega}(\hat{\boldsymbol{\omega}}_t - \mathbf{b}_{w_t} - \mathbf{n}_w)\boldsymbol{\gamma}_t^{b_k}dt.$$

Pre-integration(s.t., the imu measurement)

discretization、convariance matrix、jacobian matrix
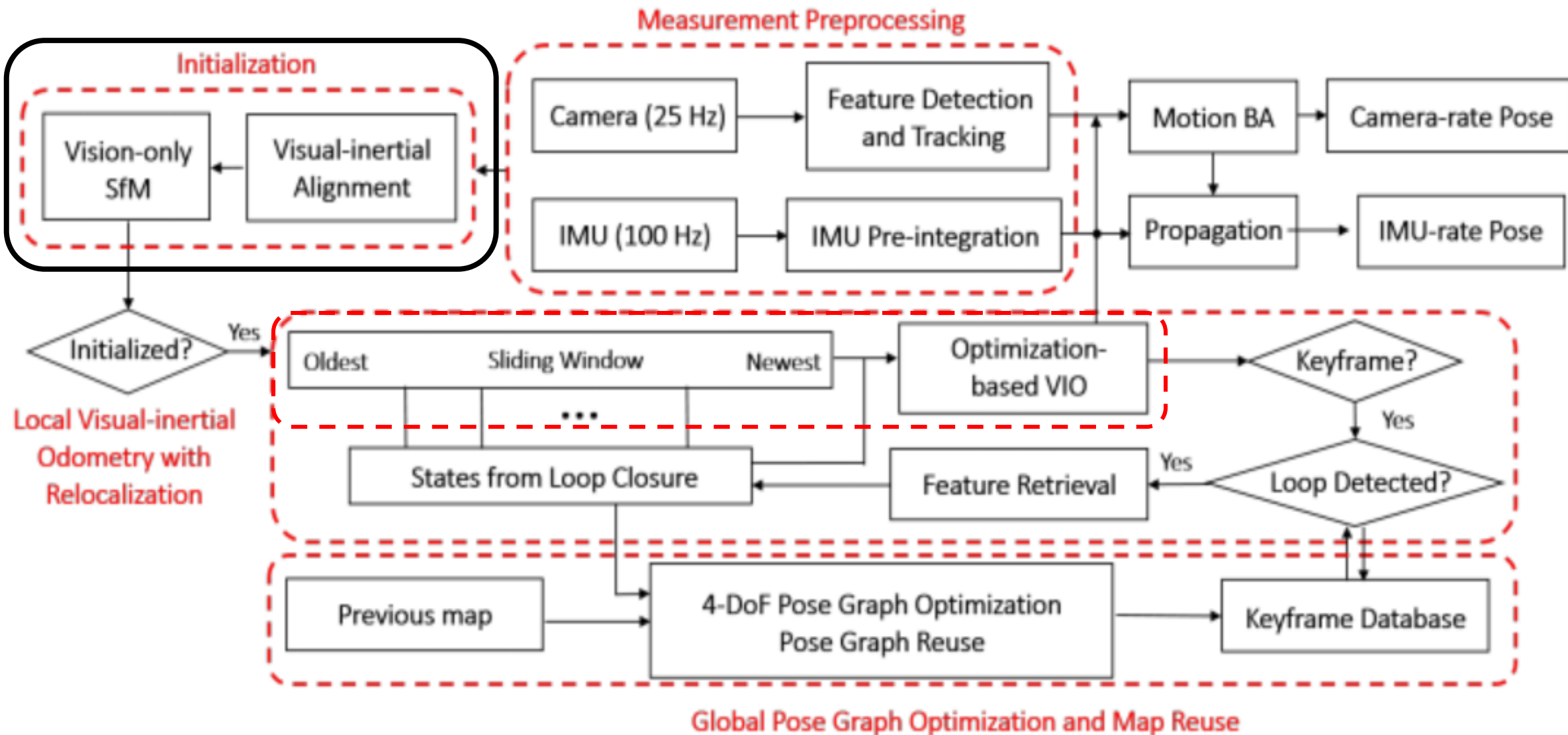………and so much work

# ◆IMU Pre-integration



Legend:
- **States in the sliding window** (blue camera)
- **States from loop closure** (green camera)
- **IMU measurements** (yellow dashed arrow)
- **Visual measurements** (black dashed arrow)
- **Features** (red star)

$$\begin{bmatrix} \hat{\boldsymbol{\alpha}}_{b_{k+1}}^{b_k} \\ \hat{\boldsymbol{\beta}}_{b_{k+1}}^{b_k} \\ \hat{\boldsymbol{\gamma}}_{b_{k+1}}^{b_k} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{R}_w^{b_k}(\mathbf{p}_{b_{k+1}}^w - \mathbf{p}_{b_k}^w + \frac{1}{2}\mathbf{g}^w \Delta t_k^2 - \mathbf{v}_{b_k}^w \Delta t_k) \\ \mathbf{R}_w^{b_k}(\mathbf{v}_{b_{k+1}}^w + \mathbf{g}^w \Delta t_k - \mathbf{v}_{b_k}^w) \\ \mathbf{q}_{b_k}^{w^{-1}} \otimes \mathbf{q}_{b_{k+1}}^w \\ \mathbf{b}_{ab_{k+1}} - \mathbf{b}_{ab_k} \\ \mathbf{b}_{wb_{k+1}} - \mathbf{b}_{wb_k} \end{bmatrix}. \quad (13)$$

the IMU measurement model based on pre-integration

measurement part          estimated part

# ◆Initialization



Measurement Preprocessing

Initialization

Vision-only SfM ← Visual-inertial Alignment

Camera (25 Hz) → Feature Detection and Tracking

IMU (100 Hz) → IMU Pre-integration

Motion BA → Camera-rate Pose

Propagation → IMU-rate Pose

Local Visual-inertial Odometry with Relocalization

Initialized? — Yes → Oldest | Sliding Window | Newest ... → Optimization-based VIO → Keyframe?

States from Loop Closure ← Feature Retrieval ← Yes — Loop Detected?

Global Pose Graph Optimization and Map Reuse

Previous map → 4-DoF Pose Graph Optimization Pose Graph Reuse → Keyframe Database
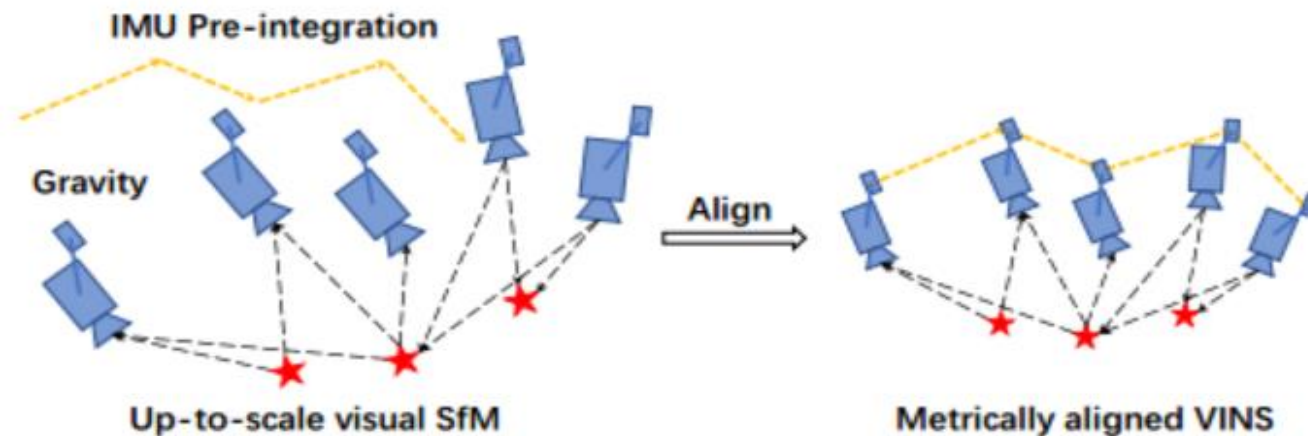
# ◆Initialization

- Why? Tightly-coupled visual-inertial odometry is a highly nonlinear system;
- Assumption:
  1. known camera intrinsic and camera-IMU extrinsic before initialization
  2. known accelerometer and gyroscope biases initial value
- **Initialization: a loosely-coupled sensor fusion process**(vision-only SLAM + visual-inertial alignment)
-  ignore accelerometer bias terms in the initial step

## ◆vision-only SLAM(SfM)

● Sliding Window Vision-Only SLAM
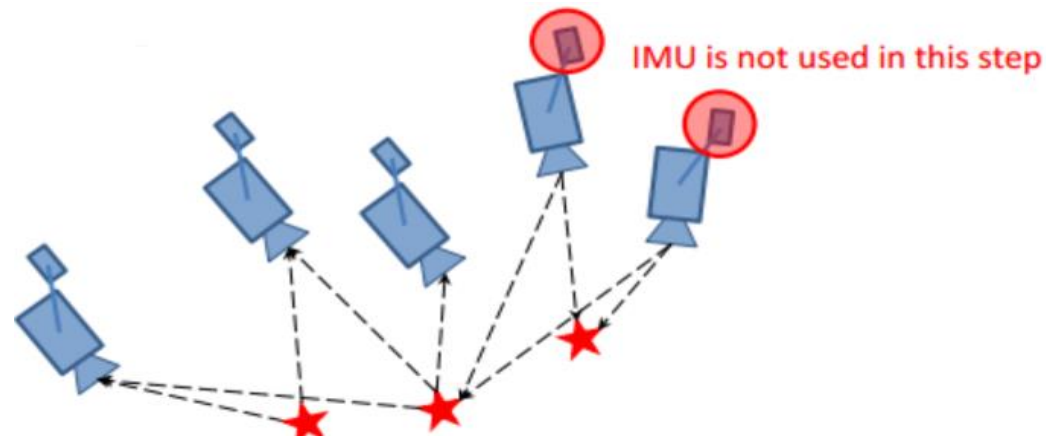  -in small window(10 fps)
  -stable feature tracking (more than 30 tracked features) and sufficient parallax (more than 20 rotation-compensated pixels)
  -5-point method 、PnP、**global full BA**
            →up-to-scale camera poes ＋ feature positions
  -the **1$^{st}$ camera frame** as the reference frame, i.e., not aligned with gravity
  -**IMU not used in this step**



IMU is not used in this step

# ◆VI Alignment

- Metric scale, world frame
- Gyroscope Bias Calibration

$$\min_{\delta b_w} \sum_{k \in \mathcal{B}} \left\| \mathbf{q}_{b_{k+1}}^{c_0}{}^{-1} \otimes \mathbf{q}_{b_k}^{c_0} \otimes \gamma_{b_{k+1}}^{b_k} \right\|^2 \qquad (15)$$

$$\gamma_{b_{k+1}}^{b_k} \approx \hat{\gamma}_{b_{k+1}}^{b_k} \otimes \begin{bmatrix} 1 \\ \frac{1}{2} \mathbf{J}_{b_w}^{\gamma} \delta \mathbf{b}_w \end{bmatrix},$$

$$\hat{\mathbf{z}}_{b_{k+1}}^{b_k} = \begin{bmatrix} \hat{\boldsymbol{\alpha}}_{b_{k+1}}^{b_k} - \mathbf{p}_c^b + \mathbf{R}_{c_0}^{b_k} \mathbf{R}_{b_{k+1}}^{c_0} \mathbf{p}_c^b \\ \hat{\boldsymbol{\beta}}_{b_{k+1}}^{b_k} \end{bmatrix} = \mathbf{H}_{b_{k+1}}^{b_k} \mathcal{X}_I + \mathbf{n}_{b_{k+1}}^{b_k}$$
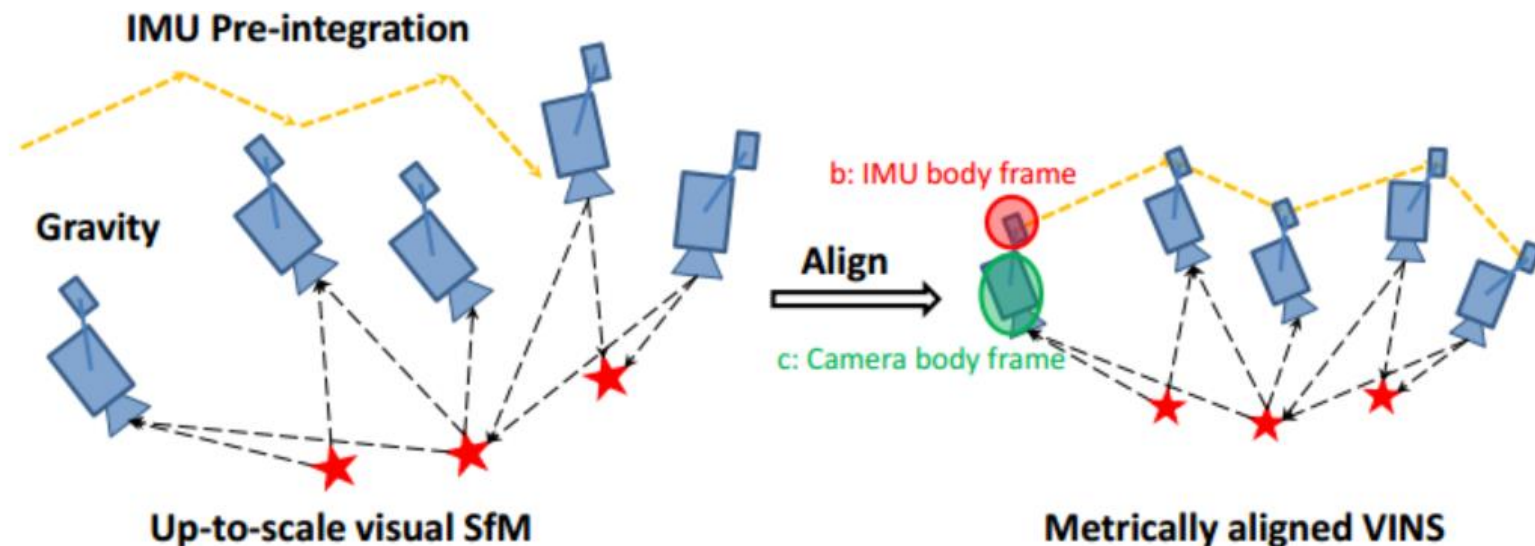
$$(18)$$

$$\mathbf{H}_{b_{k+1}}^{b_k} = \begin{bmatrix} -\mathbf{I}\Delta t_k & \mathbf{0} & \frac{1}{2} \mathbf{R}_{c_0}^{b_k} \Delta t_k^2 & \mathbf{R}_{c_0}^{b_k} (\bar{\mathbf{p}}_{c_{k+1}}^{c_0} - \bar{\mathbf{p}}_{c_k}^{c_0}) \\ -\mathbf{I} & \mathbf{R}_{c_0}^{b_k} \mathbf{R}_{b_{k+1}}^{c_0} & \mathbf{R}_{c_0}^{b_k} \Delta t_k & \mathbf{0} \end{bmatrix}$$

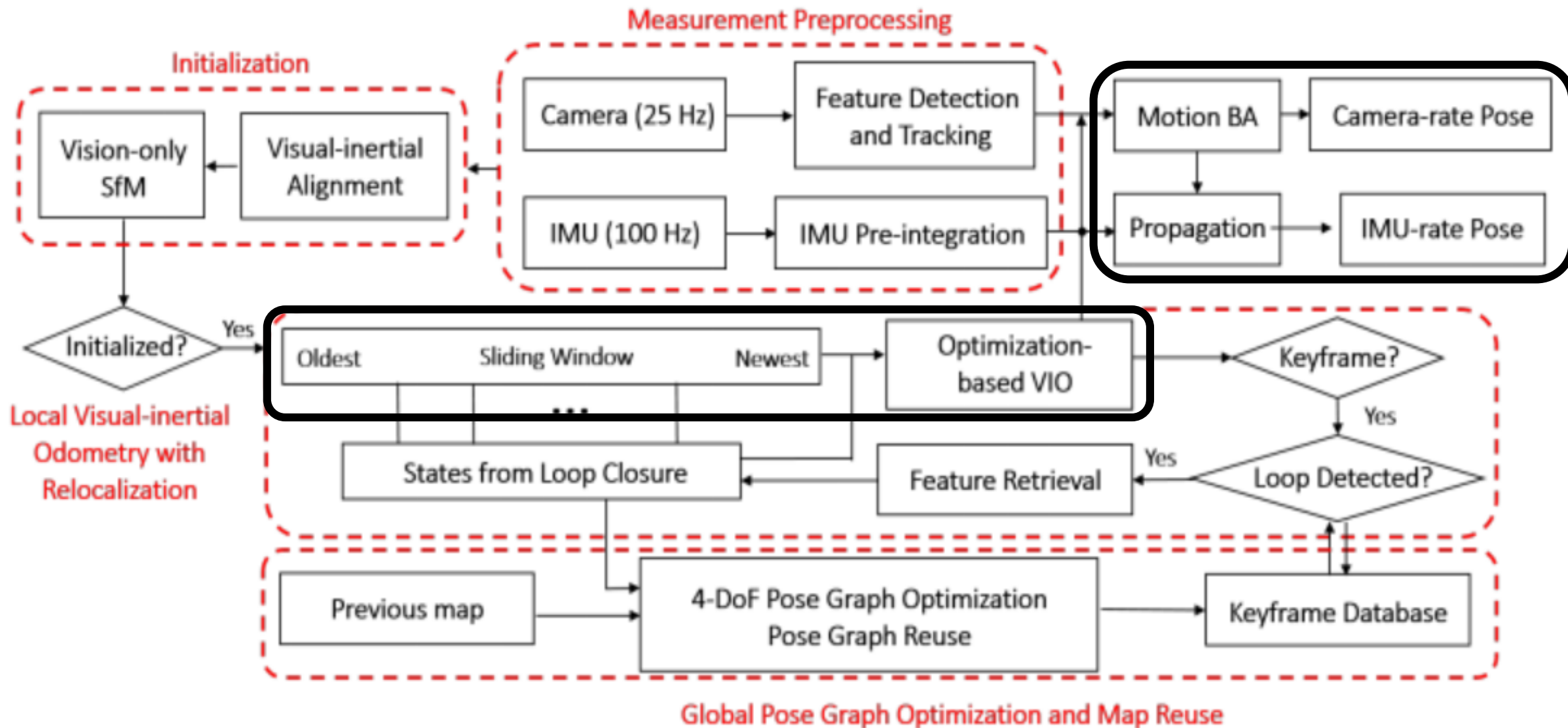- Velocity, Gravity Vector and Metric Scale Initialization

$$\mathcal{X}_I = \begin{bmatrix} \mathbf{v}_{b_0}^{b_0}, \mathbf{v}_{b_1}^{b_1}, \cdots \mathbf{v}_{b_n}^{b_n}, \mathbf{g}^{c_0}, s \end{bmatrix}$$

$$\min_{\mathcal{X}_I} \sum_{k \in \mathcal{B}} \left\| \hat{\mathbf{z}}_{b_{k+1}}^{b_k} - \mathbf{H}_{b_{k+1}}^{b_k} \mathcal{X}_I \right\|^2$$

Linear least squre problem
because R and T is known



**IMU Pre-integration**

**Gravity**

**Up-to-scale visual SfM**

b: IMU body frame

**Align**

c: Camera body frame

**Metrically aligned VINS**

# ◆ VIO



**Initialization**

**Measurement Preprocessing**

- Vision-only SfM
- Visual-inertial Alignment
- Camera (25 Hz) → Feature Detection and Tracking
- IMU (100 Hz) → IMU Pre-integration
- Motion BA → Camera-rate Pose
- Propagation → IMU-rate Pose

**Local Visual-inertial Odometry with Relocalization**

- Initialized? — Yes
- Oldest | Sliding Window | Newest → Optimization-based VIO → Keyframe?
- States from Loop Closure ← Feature Retrieval ← Loop Detected?

**Global Pose Graph Optimization and Map Reuse**

- Previous map → 4-DoF Pose Graph Optimization Pose Graph Reuse → Keyframe Database
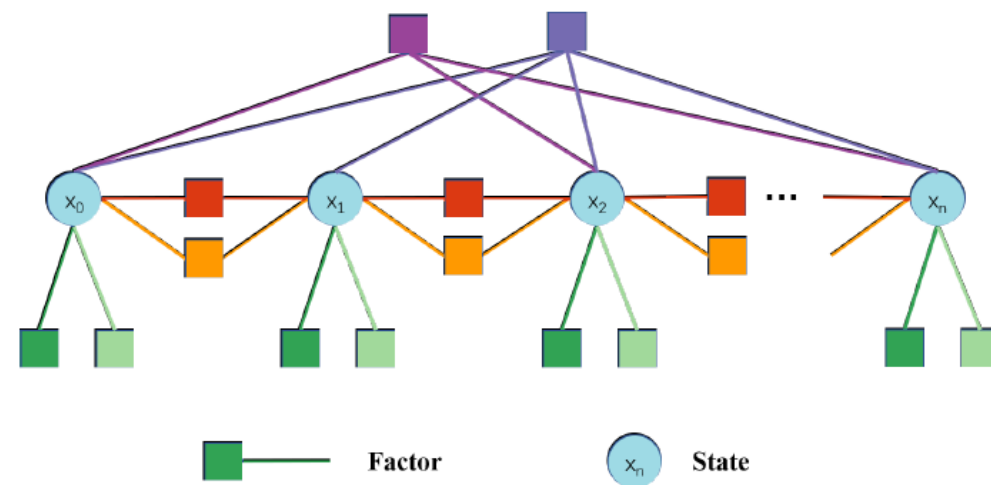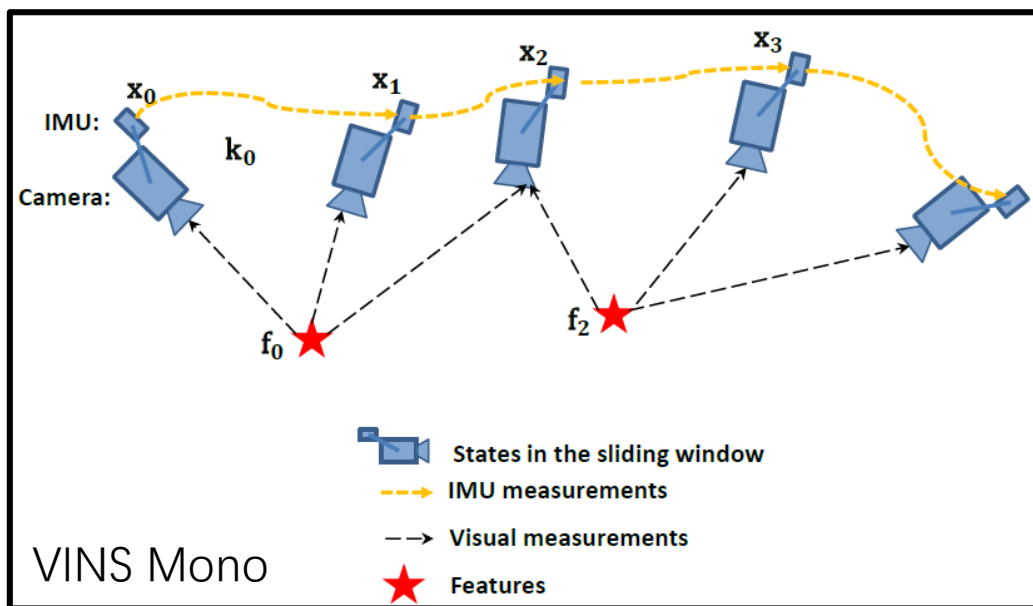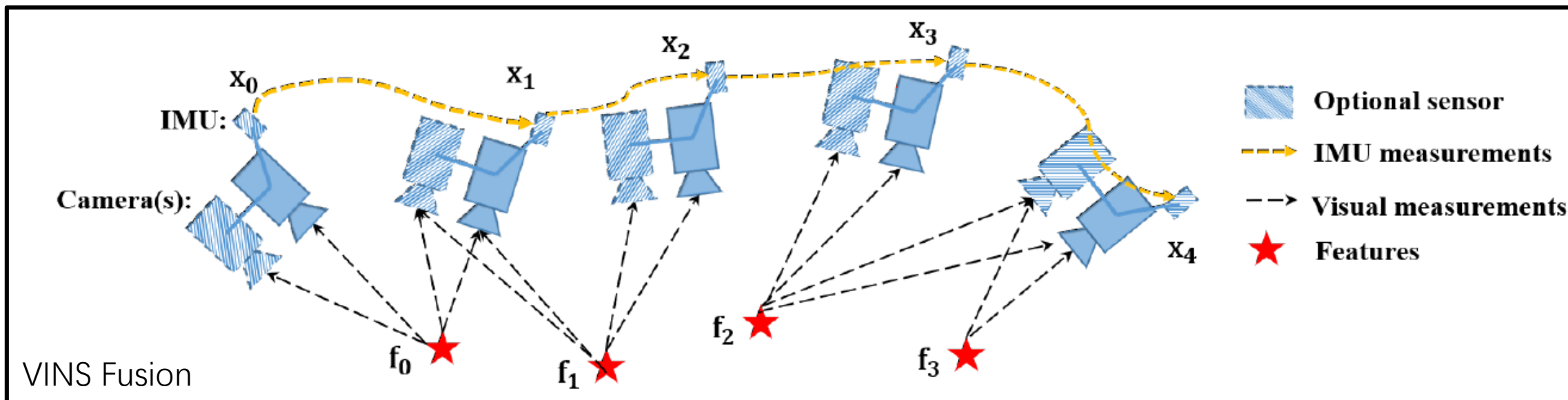
# ◆VIO



VINS Mono



Fig. 2. A graphic illustration of the pose graph. Each node represents states (position, orientation, velocity and so on) at one moment. Each edge represents a factor, which is derived by one measurement. Edges constrain one state, two states or multiple states.



VINS Fusion

# ◆VIO

- Nonlinear graph optimization-based, visual-inertial bundle adjustment
- **tightly-coupled**,Sliding window
- The full state vector

$$\mathcal{X} = \left[\mathbf{x}_0, \mathbf{x}_1, \cdots \mathbf{x}_n, \mathbf{x}_c^b, \lambda_0, \lambda_1, \cdots \lambda_m\right]$$

$$\mathbf{x}_k = \left[\mathbf{p}_{b_k}^w, \mathbf{v}_{b_k}^w, \mathbf{q}_{b_k}^w, \mathbf{b}_a, \mathbf{b}_g\right], k \in [0, n]$$

$$\mathbf{x}_c^b = \left[\mathbf{p}_c^b, \mathbf{q}_c^b\right], \qquad \text{VINS Mono}$$

$$\mathcal{X} = [\mathbf{p}_0, \mathbf{R}_0, \mathbf{p}_1, \mathbf{R}_1, ..., \mathbf{p}_n, \mathbf{R}_n, \mathbf{x}_{cam}, \mathbf{x}_{imu}]$$

$$\mathbf{x}_{cam} = [\lambda_0, \lambda_1, ..., \lambda_l] \qquad \text{VINS Fusion}$$

$$\mathbf{x}_{imu} = [\mathbf{v}_0, \mathbf{b}_{a_0}, \mathbf{b}_{g_0}, \mathbf{v}_1, \mathbf{b}_{a_1}, \mathbf{b}_{g_1}, ..., \mathbf{v}_n, \mathbf{b}_{a_n}, \mathbf{b}_{g_n}]$$

n——number of keyframes in sliding window
m——number of features in sliding window

lamda——inverse depth
Xk——IMU state when the kth image is captured

- Minimize residuals from all sensors

$$\min_{\mathcal{X}} \left\{ \|\mathbf{r}_p - \mathbf{H}_p\mathcal{X}\|^2 + \sum_{k \in \mathcal{B}} \left\|\mathbf{r}_{\mathcal{B}}(\hat{\mathbf{z}}_{b_{k+1}}^{b_k}, \mathcal{X})\right\|_{\mathbf{P}_{b_{k+1}}^{b_k}}^2 + \right.$$

$$\left. \sum_{(l,j) \in \mathcal{C}} \rho(\|\mathbf{r}_{\mathcal{C}}(\hat{\mathbf{z}}_l^{c_j}, \mathcal{X})\|_{\mathbf{P}_l^{c_j}}^2) \right\},$$

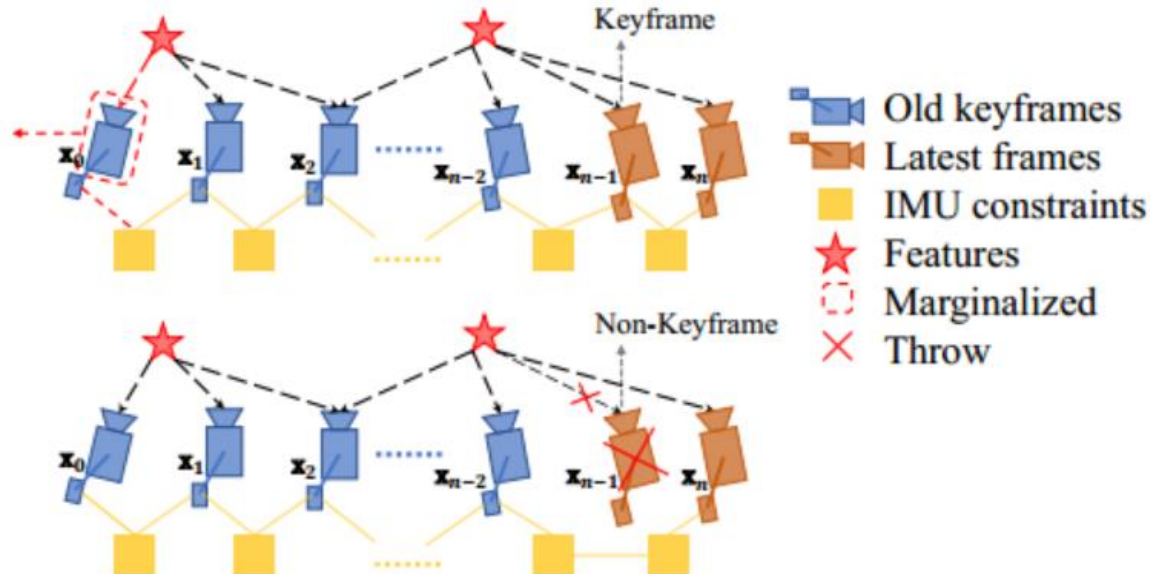Prior from marginalization
IMU measurement residual
Vision measurement residual

$$\|\mathbf{r}\|_{\Omega}^2 = \mathbf{r}^T \Omega^{-1} \mathbf{r}.$$

Mahalanobis norm

# ◆Marginalization

- Purpose: bound the computational complexity of optimization-based VIO
- Principles:
  - Add all frames into the sliding window, and remove non-keyframes after the nonlinear optimization
  - keep as many keyframes with sufficient parallax as possible;
  - Maintain matrix sparsity by throwing away visual measurements from nonkeyframes.
- Method: using the Schur complement

$$\begin{bmatrix} \mathbf{H}_{mm} & \mathbf{H}_{mr} \\ \mathbf{H}_{rm} & \mathbf{H}_{rr} \end{bmatrix} \begin{bmatrix} \delta\mathcal{X}_m \\ \delta\mathcal{X}_r \end{bmatrix} = \begin{bmatrix} \mathbf{b}_m \\ \mathbf{b}_r \end{bmatrix}$$

$$\underbrace{\left(\mathbf{H}_{rr} - \mathbf{H}_{rm}\mathbf{H}_{mm}^{-1}\mathbf{H}_{mr}\right)}_{\mathbf{H}_p}\delta\mathcal{X}_r = \underbrace{\mathbf{b}_r - \mathbf{H}_{rm}\mathbf{H}_{mm}^{-1}\mathbf{b}_m}_{\mathbf{b}_p}$$
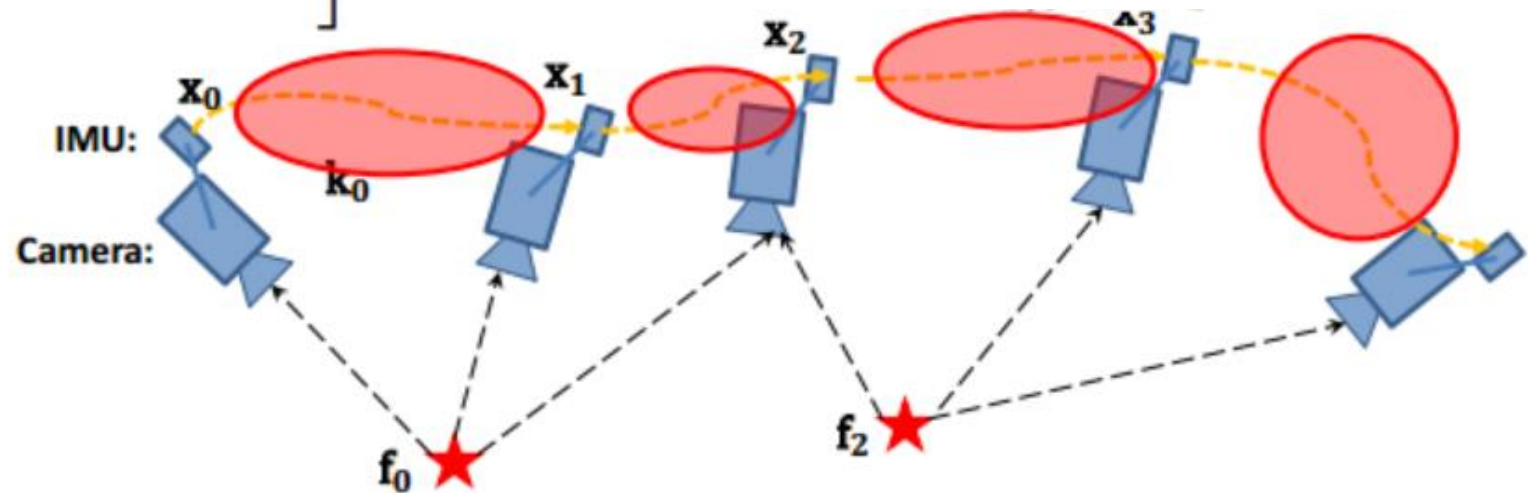


Old keyframes
Latest frames
IMU constraints
Features
Marginalized
Throw

# ◆IMU Measurement Residual

$$
\mathbf{r}_{\mathcal{B}}(\hat{\mathbf{z}}_{b_{k+1}}^{b_k}, \mathcal{X}) = \begin{bmatrix} \delta\boldsymbol{\alpha}_{b_{k+1}}^{b_k} \\ \delta\boldsymbol{\beta}_{b_{k+1}}^{b_k} \\ \delta\boldsymbol{\theta}_{b_{k+1}}^{b_k} \\ \delta\mathbf{b}_a \\ \delta\mathbf{b}_g \end{bmatrix}
$$
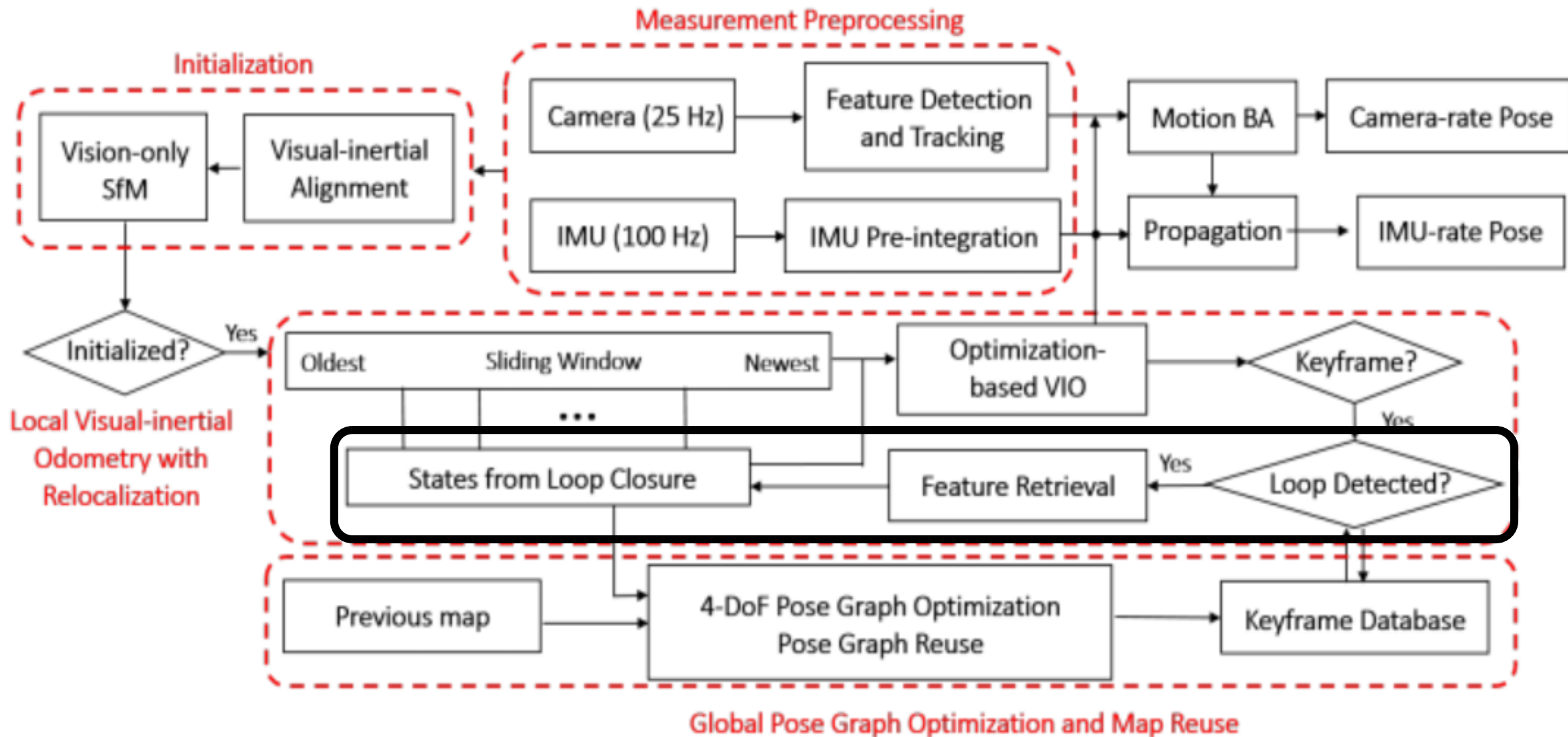
$$
= \begin{bmatrix} \mathbf{R}_w^{b_k}(\mathbf{p}_{b_{k+1}}^w - \mathbf{p}_{b_k}^w + \frac{1}{2}\mathbf{g}^w\Delta t_k^2 - \mathbf{v}_{b_k}^w\Delta t_k) - \hat{\boldsymbol{\alpha}}_{b_{k+1}}^{b_k} \\ \mathbf{R}_w^{b_k}(\mathbf{v}_{b_{k+1}}^w + \mathbf{g}^w\Delta t_k - \mathbf{v}_{b_k}^w) - \hat{\boldsymbol{\beta}}_{b_{k+1}}^{b_k} \\ 2\begin{bmatrix} \mathbf{q}_{b_k}^{w^{-1}} \otimes \mathbf{q}_{b_{k+1}}^w \otimes (\hat{\boldsymbol{\gamma}}_{b_{k+1}}^{b_k})^{-1} \end{bmatrix}_{xyz} \\ \mathbf{b}_{ab_{k+1}} - \mathbf{b}_{ab_k} \\ \mathbf{b}_{wb_{k+1}} - \mathbf{b}_{wb_k} \end{bmatrix}
$$

IMU pre-integration "blocks"(observation part)
IMU estimation "blocks"(estimation part)

# ◆Visual Measurement Residual

$$\mathbf{r}_C(\hat{\mathbf{z}}_l^{c_j}, \mathcal{X}) = \begin{bmatrix} \mathbf{b}_1 & \mathbf{b}_2 \end{bmatrix}^T \cdot (\hat{\mathcal{P}}_l^{c_j} - \frac{\mathcal{P}_l^{c_j}}{\|\mathcal{P}_l^{c_j}\|})$$

- ●Reprojection error
- ●Inverse depth model

$\hat{\mathcal{P}}_l^{c_j} = \pi^{-1}\left(\begin{bmatrix} \hat{u}_l^{c_j} \\ \end{bmatrix}\right)$

$\mathcal{P}_l^{c_j} = $

This factor is universal for both left camera and right camera. We can project a feature from the left image to the left image in temporal space, also <mark>we can project a feature from the left image to the right image in spatial space.</mark>
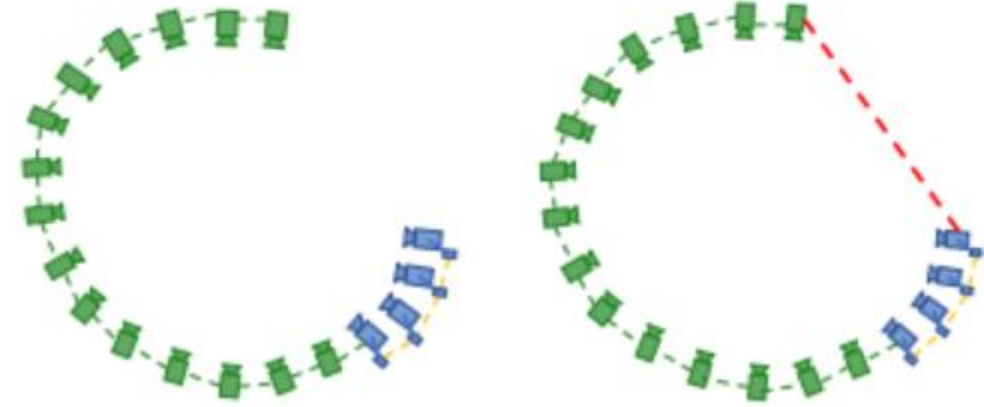
VINS Fusion

# ◆RELOCALIZATION

## ◆Loop Detection

- Features described by the BRIEF descriptor
  - Harris corners used in VIO
  - Additional 500 more FAST corners
- DBoW2 returns loop closure candidates

consistency check. We keep all BRIEF descriptors for feature retrieving, but discard the raw image to reduce memory consumption.

We note that our monocular VIO is able to render roll and pitch angles observable. As such, we do not need to rely on rotation-invariant features, such as the ORB feature used in ORB SLAM [4].

# ◆Feature Retrieval

- BRIEF descriptor matching

- 2D-2D: fundamental matrix test with RANSAC

- 3D-2D: PnP test with RANSAC

- the number of inliers beyond a certain threshold, then loop closure OK
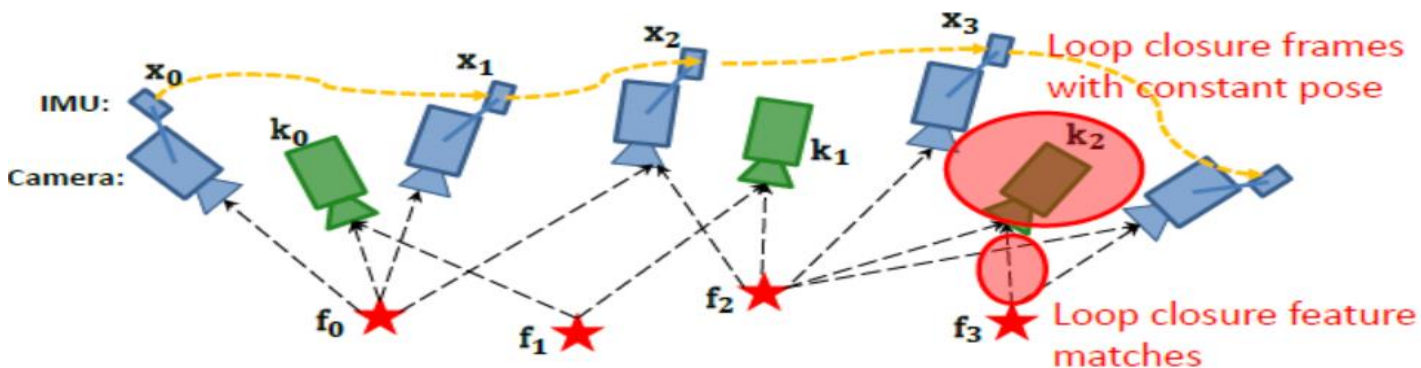
# ◆Tightly-Coupled Relocalization

$$\min_{\mathcal{X}} \left\{ \left\| \mathbf{r}_p - \mathbf{H}_p\mathcal{X} \right\|^2 + \sum_{k\in\mathcal{B}} \left\| \mathbf{r}_{\mathcal{B}}(\hat{\mathbf{z}}_{b_{k+1}}^{b_k}, \mathcal{X}) \right\|_{\mathbf{P}_{b_{k+1}}^{b_k}}^2 \right. $$
$$\left. + \sum_{(l,j)\in\mathcal{C}} \rho(\left\| \mathbf{r}_{\mathcal{C}}(\hat{\mathbf{z}}_l^{c_j}, \mathcal{X}) \right\|_{\mathbf{P}_l^{c_j}}^2) \right.$$

$$\left. + \sum_{(l,v)\in\mathcal{L}} \rho(\left\| \mathbf{r}_{\mathcal{C}}(\hat{\mathbf{z}}_l^v, \mathcal{X}, \hat{\mathbf{q}}_v^w, \hat{\mathbf{p}}_v^w) \right\|_{\mathbf{P}_l^{c_v}}^2) \right\},$$
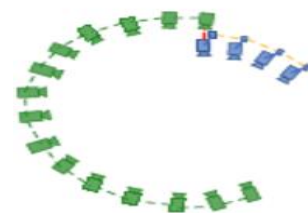
- ● VIO Residuals
- ● Loop closure vision measurement residual
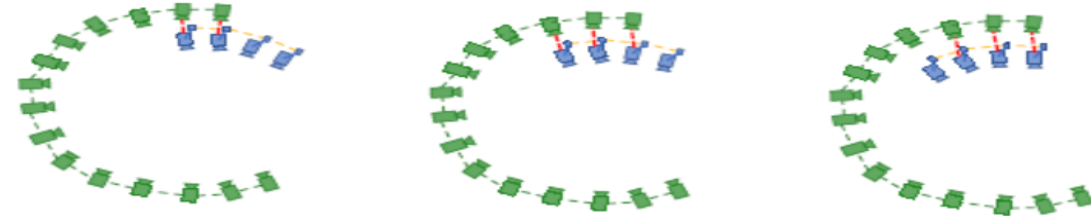- ● Poses of loop closure frames are constant
- ● multi-view constraints for relocalization



IMU:

Camera:

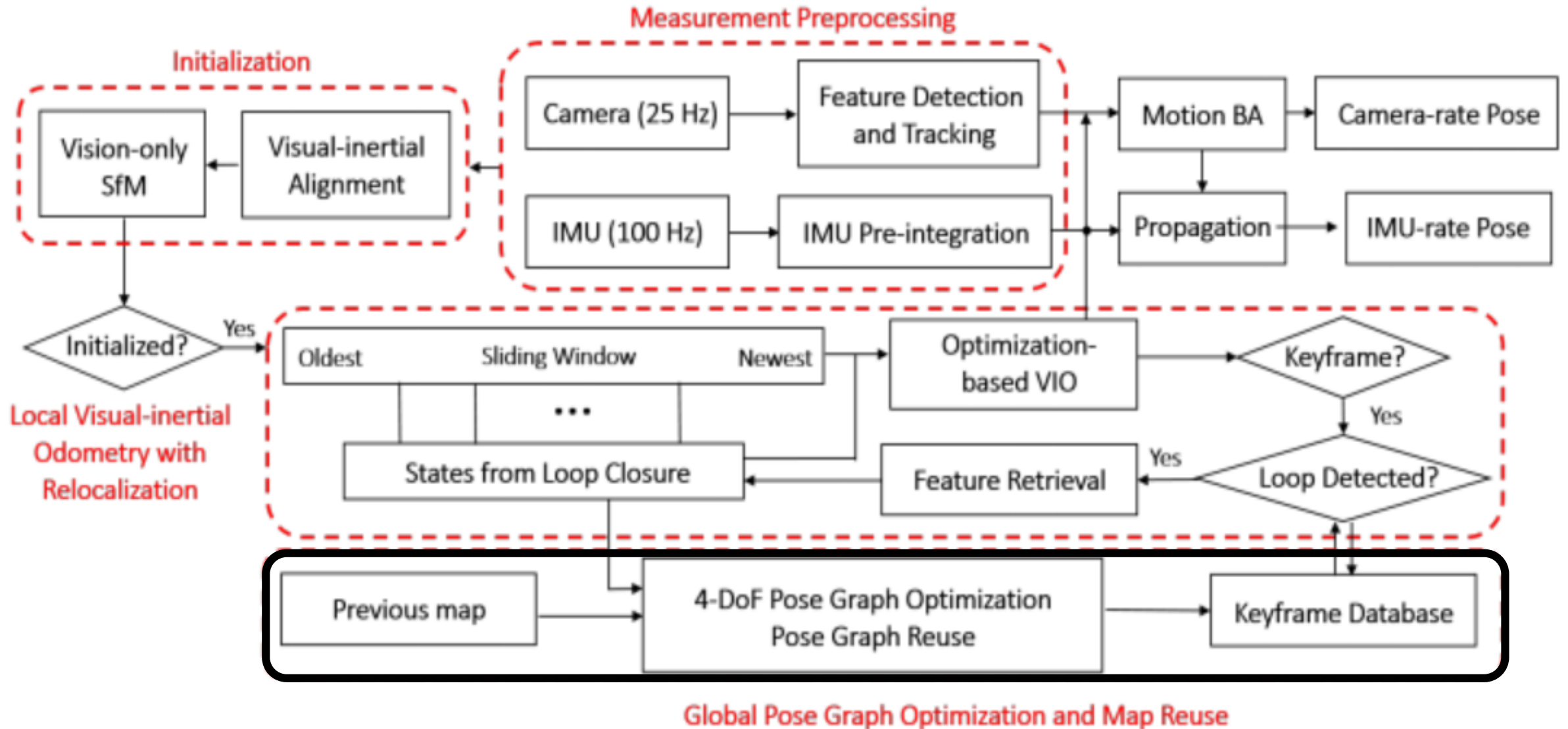Loop closure frames with constant pose

Loop closure feature matches

- 🎥 States in the sliding window
- 🟩 States from loop closure
- - - → IMU measurements
- - - → Visual measurements
- ★ Features

3. Relocalization

4. Relocalization with Multiple Constraints

# ◆GLOBAL POSE GRAPH OPTIMIZATION



Measurement Preprocessing

Initialization

Vision-only SfM ← Visual-inertial Alignment

Camera (25 Hz) → Feature Detection and Tracking → Motion BA → Camera-rate Pose

IMU (100 Hz) → IMU Pre-integration → Propagation → IMU-rate Pose

Local Visual-inertial Odometry with Relocalization

Initialized? — Yes →

Oldest — Sliding Window — Newest → Optimization-based VIO → Keyframe?

...

States from Loop Closure ← Feature Retrieval ← Loop Detected? — Yes

Previous map → 4-DoF Pose Graph Optimization Pose Graph Reuse → Keyframe Database
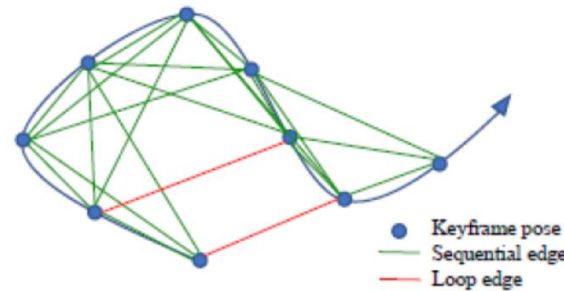
Global Pose Graph Optimization and Map Reuse

# ◆GLOBAL POSE GRAPH OPTIMIZATION

- 4-DOF pose graph optimization(x, y, z and yaw angle)

    -Roll and pitch are observable from VIO

- Relocalization and pose graph optimization run in different threads and different rate

- Adding Keyframes into the Pose Graph

    -When a keyframe is marginalized out from the sliding window

    -two types of edges:

    1) Sequential Edge

    2) Loop Closure Edge

- 4-DOF relative pose residual

    -obtained directly from VIO



- ● Keyframe pose
- — Sequential edge
- — Loop edge

$$\mathbf{r}_{i,j}(\mathbf{p}_i^w, \psi_i, \mathbf{p}_j^w, \psi_j) = \begin{bmatrix} \mathbf{R}(\hat{\phi}_i, \hat{\theta}_i, \psi_i)^{-1}(\mathbf{p}_j^w - \mathbf{p}_i^w) - \hat{\mathbf{p}}_{ij}^i \\ \psi_j - \psi_i - \hat{\psi}_{ij} \end{bmatrix},$$

$$(28)$$

# ◆GLOBAL POSE GRAPH OPTIMIZATION

● minimizing the following cost function

$$\min_{\mathbf{P},\psi} \left\{ \sum_{(i,j)\in\mathcal{S}} \|\mathbf{r}_{i,j}\|^2 + \sum_{(i,j)\in\mathcal{L}} \rho\left(\|\mathbf{r}_{i,j}\|^2\right) \right\}$$
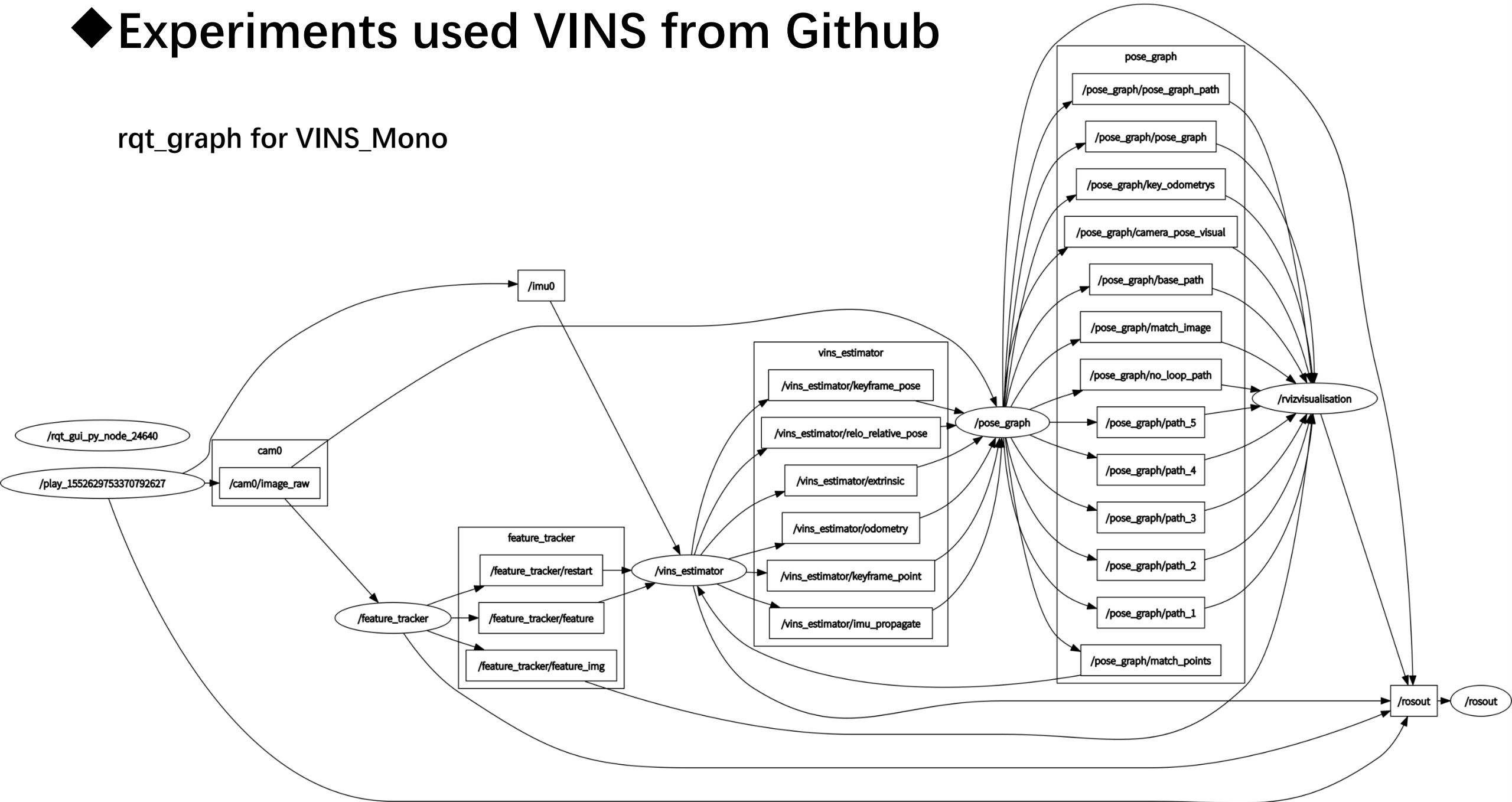
Sequential edges
Loop closure edges
Huber norm, used to further reduce the impact
of any possible wrong loops

● Pose Graph Management

of the system in the long run. To this end, we implement a downsample process to maintain the pose graph database to a limited size. All keyframes with loop closure constraints will be kept, while other keyframes that are either too close or have very similar orientations to its neighbors may be removed. The probability of a keyframe being removed is proportional to its spatial density to its neighbors.
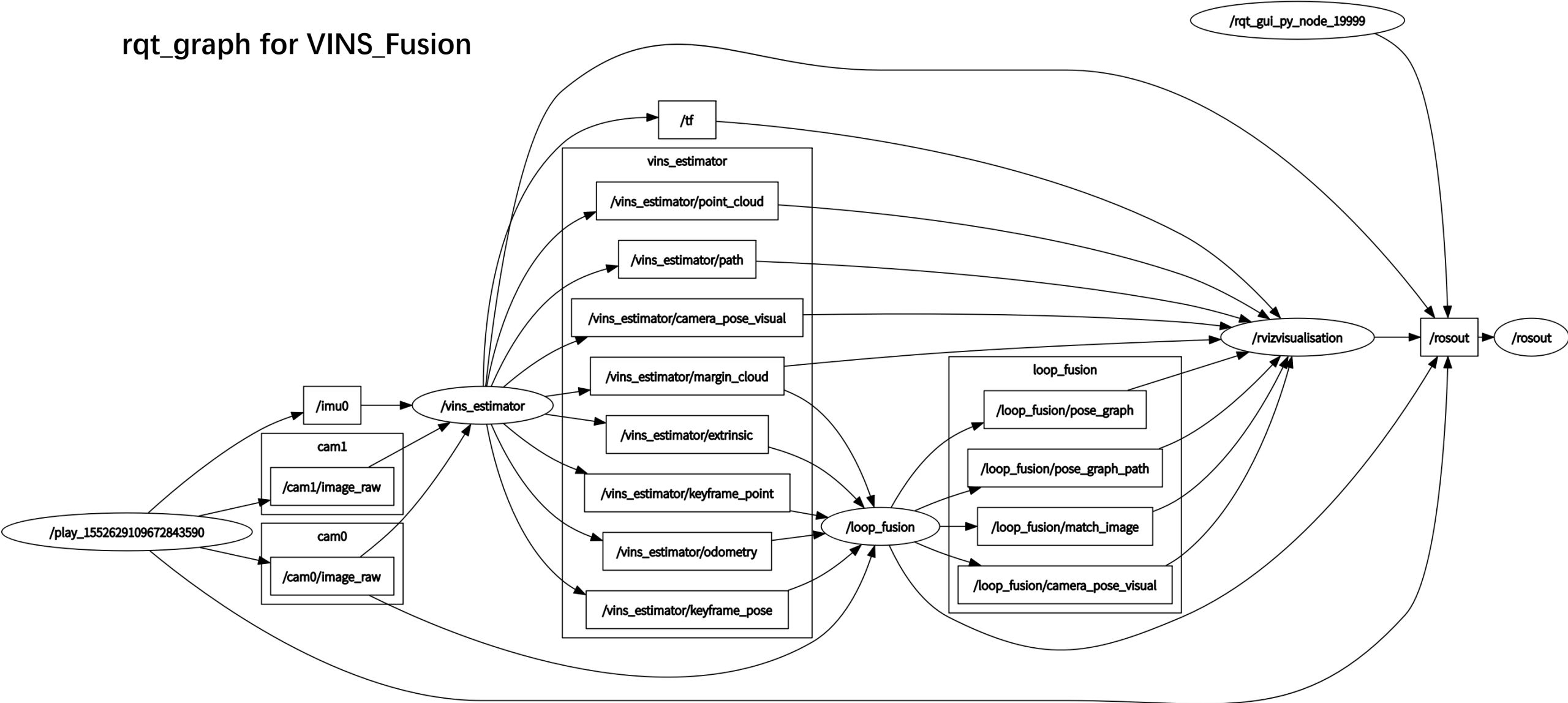
# ◆Experiments used VINS from Github

rqt_graph for VINS_Mono

# ◆Experiments used VINS from Github

rqt_graph for VINS_Fusion

# Thanks!