

# S-PTAM stereo parallel tracking and mapping

Sunday, December 27, 2020 8:07 AM

Notation:

$SE(3)$  transformation:

$$T = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}$$

$E^{cw}$ : transformation representing a camera pose that transforms a point in world frame

$$x^w = [x^w \ y^w \ z^w \ 1]^T$$

to a point in camera frame

$$x^c = [x^c \ y^c \ z^c \ 1]^T$$

that is

$$x^c = E^{cw} x^w$$

Motion matrix  $M$ : a  $4 \times 4$  matrix  $SE(3)$  representing the changes in camera pose by left-multiplication:

$$E^{cw} = M^c E^{cw}_{\text{prev}}$$

In Lie Groups, the motion matrix  $M$  can be represented by a six-vector

$$\mu = [tx, ty, tz, \theta_{roll}, \theta_{pitch}, \theta_{yaw}]$$

$\mu, M$  are related by

$$M = \exp(\mu) = e^{\sum_{j=1}^6 \mu_j G_j}$$

$G_j, j=1 \dots 6$  is the group generator matrix.

Measurement  $z = [u \ v]$ , is the true 2D position that matches with the projected 3D point on the image plane.

Map point:  $p_i$  is ordered pair  $(x^*, d)$ , contains

- 3D point  $x^w$

- its descriptor  $d$

Stereo keyframe  $k$  is a stereo pair of images with the associated Stereo Camera pose.

Map is defined as the set of map points and the set of Stereo keyframes.

A point in camera frame  $x^c$  projects into the image as

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = p(x^c)$$

$$p(x^c) = \begin{bmatrix} f_u & 0 & u_0 \\ 0 & f_v & v_0 \end{bmatrix} \begin{bmatrix} \frac{x^c}{z^c} \\ \frac{y^c}{z^c} \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{f_u x^c}{z^c} + u_0 \\ \frac{f_v y^c}{z^c} + v_0 \end{bmatrix}$$

The global reference frame is the camera pose at the first frame.

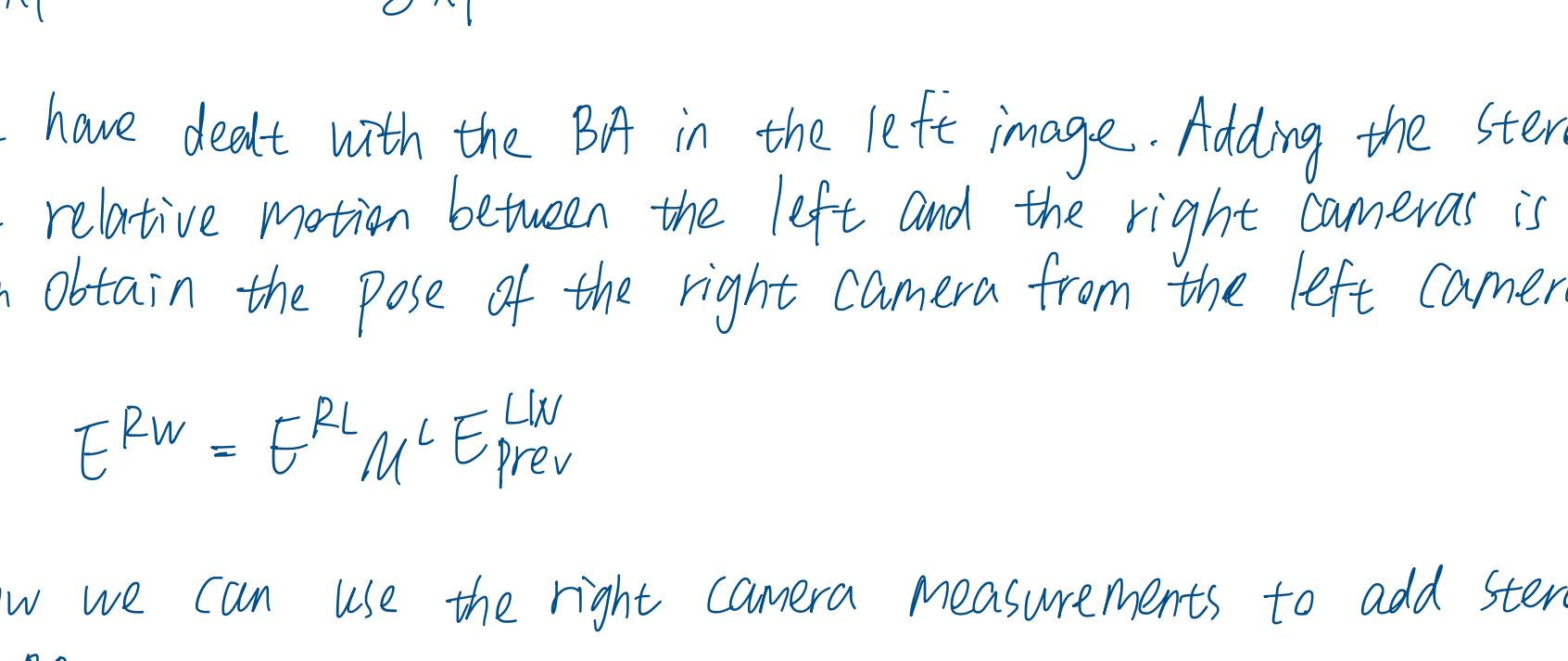
An initial map is estimated by matching and triangulating SIFT point features in the first stereo pair.

For every frame after the tracking thread estimates the 6DOF pose for each stereo frame by minimizing the re-projection error between the projected Map points and their correspondences.

The system selects a subset of keyframes that will be used in a second thread to estimate the map at a lower rate.

The map points are triangulated from the stereo matches of each keyframe, and added to the map.

The mapping thread is constantly minimizing the local reprojection error by refining all the map points, and the stereo poses using Bundle Adjustment. We use a pose graph to maintain the global consistency of the map.



Feature extraction and description.

- Keypoint detector: GFTT

- Keypoint descriptor: BRISK.

Pose tracking:

- Matching.

- Project each map point into the viewing frustum of the predicted stereo pose, and search for a match in the neighborhood of the point.

- a decaying velocity model is used for pose prediction.

- binary descriptors are used, and the Hamming distance is calculated.

• Pose refinement

- $E^{cw}$ : current camera pose in global frame  $w$ .

- Previous camera pose  $E^{cw}_{\text{prev}}$  with the relative motion  $M^c$  in the local camera frame.

$$E^{cw} = M^c E^{cw}_{\text{prev}}$$

- To find the  $M^c$ ,

$$J \mu = \Delta z (\mu_{\text{prev}}) \quad (6)$$

$$\mu = (tx, ty, tz, \theta_{roll}, \theta_{pitch}, \theta_{yaw})^T$$

$\Delta z$ : re-projection error that only depends on the camera motion  $\mu$ , as we consider the map fixed.

J: Jacobian of re-projection error w.r.t camera motion.

$$J_{ij} = \frac{\partial \Delta z_i(\mu)}{\partial \mu_j} = \frac{\partial \left( \begin{bmatrix} u \\ v \end{bmatrix}_i - P(\exp(\mu_j) E^{cw}_{\text{prev}} x^w_i) \right)}{\partial \mu_j}$$

$$= - \frac{\partial P(x^c_i)}{\partial x^c_i} \frac{\partial x^c_i}{\partial \mu_j}$$

where

$$\frac{\partial P(x^c_i)}{\partial x^c_i} = \begin{bmatrix} \frac{f_u}{z^c} & 0 & -\frac{f_u x^c}{z^{c2}} \\ 0 & \frac{f_v}{z^c} & -\frac{f_v y^c}{z^{c2}} \end{bmatrix}$$

$$\frac{\partial x^c_i}{\partial \mu_j} = G_j E^{cw}_{\text{prev}} x^w_i$$

$\mu$  is found by solving (6). In order to do this, given a set

$$S = \{z_1, \dots, z_n\}$$

of matched measurements, the new value for  $\mu$  is obtained by minimizing an objective function:

$$\mu' = \arg \min_{\mu} \sum_{i \in S} l(J_i \mu - \Delta z_i(\mu_{\text{prev}}))$$

$l(\cdot)$  is the Huber function, the optimization is done via LM.

• Keyframes selection and map points creation.

- A frame is selected to be a keyframe if the number of tracked points is less than 90% of the points tracked in last keyframe.

- The remaining unmatched features from the stereo pair are triangulated to create new map points.

- The keyframe is queued into the map refinement thread.

• Local mapping.

- The refinement of the keyframe poses, and the 3D map points, is done via BA, that minimizes the reprojection error of every point in every image.

- The problem can be stated as follows:

Given an initial set of  $N$  keyframe poses  $\{E_1, \dots, E_N\}$ , an initial set of  $M$  3D points  $\{x^w_1, \dots, x^w_M\}$ , and measurements  $\{z_1, \dots, z_N\}$ , each set  $S_j$  contains the measurements  $z_{ij}$  of the  $i$ -th point in the  $j$ -th keyframe, the simultaneous estimation of the multiple cameras and the point cloud is achieved by solving

$$J \mu = \Delta z(\mu_{\text{prev}}, x^w)$$

$$\Delta z(\mu, x^w) = z - P(\exp(\mu) E^{cw}_{\text{prev}} x^w)$$

We must minimize

$$J \mu = \Delta z(\mu_{\text{prev}}, x^w) \quad \text{as } M_i \text{ is fixed in BA, as it defines the world frame.}$$

$$\{ \mu'_j \}_{j=1}^N, \{ x^w_i \}_{i=1}^M = \arg \min_{\mu} \sum_{j=1}^N \sum_{i \in S_j} l_b(\mu_j, z_{ji})$$

$$\mu'_j = J_{ji} \left[ \begin{array}{c} \mu_j \\ x^w_i \end{array} \right] - \Delta z_i(\mu_{\text{prev}}, j, x^w_i)$$

Given that the vector of parameters is divided in two groups, (cameras and points), the Jacobian can be decomposed as:

• Keypoint detector: GFTT

• Keypoint descriptor: BRISK.

Pose tracking:

- Matching.

- Project each map point into the viewing frustum of the predicted stereo pose, and search for a match in the neighborhood of the point.

- a decaying velocity model is used for pose prediction.

- binary descriptors are used, and the Hamming distance is calculated.

• Pose refinement

- $E^{cw}$ : current camera pose in global frame  $w$ .

- Previous camera pose  $E^{cw}_{\text{prev}}$  with the relative motion  $M^c$  in the local camera frame.

$$E^{cw} = M^c E^{cw}_{\text{prev}}$$

- To find the  $M^c$ ,

$$J \mu = \Delta z (\mu_{\text{prev}}) \quad (6)$$

$$\mu = (tx, ty, tz, \theta_{roll}, \theta_{pitch}, \theta_{yaw})^T$$

$\Delta z$ : re-projection error that only depends on the camera motion  $\mu$ , as we consider the map fixed.

J: Jacobian of re-projection error w.r.t camera motion.

$$J_{ij} = \frac{\partial \Delta z_i(\mu)}{\partial \mu_j} = \frac{\partial \left( \begin{bmatrix} u \\ v \end{bmatrix}_i - P(\exp(\mu_j) E^{cw}_{\text{prev}} x^w_i) \right)}{\partial \mu_j}$$

$$= - \frac{\partial P(x^c_i)}{\partial x^c_i} \frac{\partial x^c_i}{\partial \mu_j}$$

where

$$\frac{\partial P(x^c_i)}{\partial x^c_i} = \begin{bmatrix} \frac{f_u}{z^c} & 0 & -\frac{f_u x^c}{z^{c2}} \\ 0 & \frac{f_v}{z^c} & -\frac{f_v y^c}{z^{c2}} \end{bmatrix}$$

$$\frac{\partial x^c_i}{\partial \mu_j} = G_j E^{cw}_{\text{prev}} x^w_i$$

$\mu$  is found by solving (6). In order to do this, given a set

$$S = \{z_1, \dots, z_n\}$$

of matched measurements, the new value for  $\mu$  is obtained by minimizing an objective function:

$$\mu' = \arg \min_{\mu} \sum_{i \in S} l(J_i \mu - \Delta z_i(\mu_{\text{prev}}))$$

$l(\cdot)$  is the Huber function, the optimization is done via LM.

• Keyframes selection and map points creation.

- A frame is selected to be a keyframe if the number of tracked points is less than 90% of the points tracked in last keyframe.

- The remaining unmatched features from the stereo pair are triangulated to create new map points.

- The keyframe is queued into the map refinement thread.

• Local mapping.

- The refinement of the keyframe poses, and the 3D map points, is done via BA, that minimizes the reprojection error of every point in every image.

- The problem can be stated as follows:

Given an initial set of  $N$  keyframe poses  $\{E_1, \dots, E_N\}$ , an initial set of  $M$  3D points  $\{x^w_1, \dots, x^w_M\}$ , and measurements  $\{z_1, \dots, z_N\}$ , each set  $S_j$  contains the measurements  $z_{ij}$  of the  $i$ -th point in the  $j$ -th keyframe, the simultaneous estimation of the multiple cameras and the point cloud is achieved by solving

$$J \mu = \Delta z(\mu_{\text{prev}}, x^w)$$

$$\Delta z(\mu, x^w) = z - P(\exp(\mu) E^{cw}_{\text{prev}} x^w)$$

We must minimize

$$J \mu = \Delta z(\mu_{\text{prev}}, x^w) \quad \text{as } M_i \text{ is fixed in BA, as it defines the world frame.}$$

$$\{ \mu'_j \}_{j=1}^N, \{ x^w_i \}_{i=1}^M = \arg \min_{\mu} \sum_{j=1}^N \sum_{i \in S_j} l_b(\mu_j, z_{ji})$$

$$\mu'_j = J_{ji} \left[ \begin{array}{c} \mu_j \\ x^w_i \end{array} \right] - \Delta z_i(\mu_{\text{prev}}, j, x^w_i)$$

Given that the vector of parameters is divided in two groups, (cameras and points), the Jacobian can be decomposed as:

• Keypoint detector: GFTT

• Keypoint descriptor: BRISK.

Pose tracking:

- Matching.

- Project each map point into the viewing frustum of the predicted stereo pose, and search for a match in the neighborhood of the point.

- a decaying velocity model is used for pose prediction.

- binary descriptors are used, and the Hamming distance is calculated.

• Pose refinement

- $E^{cw}$ : current camera pose in global frame  $w$ .

- Previous camera pose  $E^{cw}_{$