

Learned Monocular Depth Priors in Visual-Inertial Initialization Supplemental Material

Yunwen Zhou, Abhishek Kar, Eric Turner, Adarsh Kowdle, Chao X. Guo,
Ryan C. DuToit, and Konstantine Tsotsos

Google AR

{verse,abhiskar,elturner,adarshkowdle,chaoguo,rdutoit,ksotsos}@google.com

1 Depth Residual function

Reiterating our paper, the proposed depth residual in VI-BA for keyframe k and feature point i takes the form of the \log of the ratio between the measured depth scaled/shifted by \mathbf{S}_k and the feature point’s estimated depth:

$$r_{\mathcal{L}_{ik}} = \log \left((a_k d_{ik} + b_k) \cdot \Omega(\mathcal{C}^j \mathbf{f}_i, \mathbf{q}_j, \mathbf{p}_j, \mathbf{q}_k, \mathbf{p}_k) \right) \quad (1)$$

Cross-Keyframe Depth Transformation $\Omega(\cdot)$ in Eq. (1) is the function transforming the feature point depth w_{ij}^{-1} from its first observed j^{th} camera keyframe to the depth measurement camera k^{th} keyframe $\{C_k\}$. Since feature point $\mathcal{C}^j \mathbf{f}_i$ is parameterized as inverse depth, to recover it into euclidean space $\mathcal{C}^j \mathbf{l}_i$, we have

$$\mathcal{C}^j \mathbf{l}_i = w_{ij}^{-1} [u_{ij} \ v_{ij} \ 1]^T$$

where inverse parameterization is known to be unstable in linearizations when used directly. Instead, we define $\mathcal{C}^j \mathbf{h}_i = [u_{ij}, v_{ij}, 1]^T$, so

$$\mathcal{C}^j \mathbf{h}_i = w_{ij} \mathcal{C}^j \mathbf{l}_i$$

where we are able to deal with w_{ij} separately from the geometry transform without explicitly unfolding $\mathcal{C}^j \mathbf{h}_i$ using w_{ij} . We define

$${}^B \mathbf{y}_i = w_{ij} {}^B \mathbf{l}_i \quad (2)$$

where ${}^B \mathbf{l}_i$ is the i^{th} feature position in $\{B\}$ frame other than global $\{G\}$. If it is global $\{G\}$, Eq. (2) becomes $\mathbf{y}_i = w_{ij} \mathbf{l}_i$.

We can infer $\mathcal{C}^k \mathbf{y}_i$ through a series of transforms only on top of $\mathcal{C}^j \mathbf{h}_i$ using the IMU-camera extrinsics $(\mathbf{q}_C, \mathbf{p}_C)$, j^{th} and k^{th} keyframe poses. Firstly, we show how to get $w_{ij} {}^j \mathbf{l}_i$ that is ${}^j \mathbf{y}_i$ using $\mathcal{C}^j \mathbf{h}_i$, where ${}^j \mathbf{l}_i$ is in j^{th} IMU keyframe coordinate system

$${}^j \mathbf{y}_i = R(\mathbf{q}_C) \mathcal{C}^j \mathbf{h}_i + w_{ij} \mathbf{p}_C$$

Similarly, we can keep transforming all the way to the k^{th} camera frame and get ${}^{C_k}\mathbf{y}_i$

$$\begin{aligned}\mathbf{y}_i &= R(\mathbf{q}_j)^j \mathbf{y}_i + w_{ij} \mathbf{p}_j \\ {}^k \mathbf{y}_i &= R(\mathbf{q}_k^{-1})(\mathbf{y}_i - w_{ij} \mathbf{p}_k) \\ {}^{C_k} \mathbf{y}_i &= R(\mathbf{q}_C^{-1})({}^k \mathbf{y}_i - w_{ij} \mathbf{p}_C)\end{aligned}$$

Reiterating Eq. (2), along with $[\cdot]_3$ that extracts the 3^{rd} element from the vector, $\Omega(\cdot)$ then becomes

$$\Omega({}^{C_j} \mathbf{f}_i, \mathbf{q}_j, \mathbf{p}_j, \mathbf{q}_k, \mathbf{p}_k) = \frac{[{}^{C_k} \mathbf{y}_i]_3}{w_{ij}} \quad (3)$$

As shown below, Eq. (1) can be simply written as followings by substitution using Eq. (3)

$$r_{\mathcal{L}_{ik}} = \log(a_k d_{ik} + b_k) + \log([{}^{C_k} \mathbf{y}_i]_3) - \log(w_{ij}) \quad (4)$$

Eq. (4) fits our problem space since our ML model yields inverse depth and our feature point is parameterized as inverse depth, which results in a more stable first-order approximation that is helpful for optimization purposes.

2 Formulation and Derivation of Closed-Form Solver For Initializing VI-BA

It is well-known that VI-SFM is a highly non-linear problem, therefore it is important to find an accurate initial linearization point. Our solution is based on [4]. We employ visual reprojection error to approximate the 0^{th} keyframe's pose and velocity in gravity-aligned global coordinate frame $\{G\}$. Then, the remaining keyframe poses and velocities are inferred from IMU integration.

Given N keyframes, we define our estimated states \mathbf{X} for closed-form solver as following,

$$\mathbf{X} = [\mathbf{v}; \mathbf{g}; \Delta \mathbf{p}_0; \Delta \mathbf{p}_1; \dots; \Delta \mathbf{p}_{N-1}] \quad (5)$$

where \mathbf{v} , \mathbf{g} are the initial velocity and gravity vector, $\Delta \mathbf{p}_k$ is the k^{th} keyframe position estimation difference. All parameters are expressed with respect to the 0^{th} keyframe, and we re-express them with respect to $\{G\}$ through \mathbf{g} after the solve. The rest of this section shows the derivation of the linear equation to solve the problem.

As in [2, 5], we marginalize feature points, yielding constraints among keyframes. Then, for feature points initial values in VI-BA, they are triangulated after the closed-form solve.

IMU Measurement Model. First, we recall the IMU measurement model:

$$\boldsymbol{\omega}^m(t) = \boldsymbol{\omega}(t) + \mathbf{b}^w(t) + \boldsymbol{\eta}^w(t) \quad (6)$$

$$\mathbf{a}^m(t) = R(t)^T (\mathbf{G} \mathbf{a}(t) - \mathbf{G} \mathbf{g}) + \mathbf{b}^a(t) + \boldsymbol{\eta}^a(t) \quad (7)$$

where $\boldsymbol{\omega}^m(t)$ and $\mathbf{a}^m(t)$ are the gyro and acceleration measurement at timestamp t , $\mathbf{G} \mathbf{g} = [0; 0; -G]$ is the gravity vector in gravity-aligned global $\{G\}$. $\mathbf{b}^a(t)$

and $\mathbf{b}^w(t)$ are IMU accelerometer bias and gyro bias respectively, and their corresponding noises are $\boldsymbol{\eta}^a(t)$, $\boldsymbol{\eta}^g(t)$. $R(t)$ is the IMU body frame rotation w.r.t gravity-aligned global $\{G\}$.

Integrated Keyframe Positions. According to [4] and IMU measurement model Eq. (6), Eq. (7), we can write k^{th} keyframe position as

$$\mathbf{p}_k = \mathbf{v}\Delta t_k + \mathbf{g}\frac{\Delta t_k^2}{2} + \boldsymbol{\xi}_k \quad (8)$$

where $\Delta t_k = t_k - t_0$ is the elapse timestamp up to t_k since t_0 . $\boldsymbol{\xi}_k$ is the k^{th} keyframe integrated position in the 0^{th} keyframe at timestamp t_k without the impact of gravity. The integrator is described in [5]

$$\boldsymbol{\xi}_k = \int_{t_0}^{t_k} \int_{t_0}^{\tau} R_0(\eta)(\mathbf{a}^m(\eta) - \mathbf{b}^a(\eta))d\eta d\tau$$

where $R_0(\eta)$ is integrated using gyro measurement from t_0 to timestamp η in 0^{th} keyframe coordinate system. To remove the non-linearity of the estimation, $\mathbf{b}^w(t)$ is assumed to be 0. $\mathbf{b}^a(t)$ is also assumed to be 0, since the small baseline closed-form solution is resilient to accelerometer bias, which is studied in [3]. The bias random walk noises $\boldsymbol{\eta}^a(t)$, $\boldsymbol{\eta}^w(t)$ are treated as zero mean so it doesn't appear in the equation.

Defining integrated keyframe positions in 0^{th} keyframe coordinate system as $\mathbf{P} = [\mathbf{p}_0; \mathbf{p}_1; \dots; \mathbf{p}_{N-1}]$, we form a linear equation with states in \mathbf{X} and $\boldsymbol{\Xi} = [\mathbf{0}_{3 \times 1}; \boldsymbol{\xi}_0; \dots; \boldsymbol{\xi}_{N-1}]$

$$\mathbf{P} = F\mathbf{g} + W \begin{bmatrix} \mathbf{v} \\ \Delta \mathbf{p}_0 \\ \vdots \\ \Delta \mathbf{p}_{N-1} \end{bmatrix} + \boldsymbol{\Xi} \quad (9)$$

where

$$F = \begin{bmatrix} \mathbf{0}_{3 \times 3} \\ \vdots \\ \frac{\Delta t_{N-1}^2}{2} \mathbf{I}_3 \end{bmatrix}, \quad W = \begin{bmatrix} \mathbf{0}_{3 \times 3} & \mathbf{I}_3 & & \\ & \vdots & \ddots & \\ \Delta t_{N-1} \mathbf{I}_3 & & & \mathbf{I}_3 \end{bmatrix}$$

Visual Constraint. We define the estimated i^{th} feature point position at k^{th} camera frame as ${}^{C_k}\mathbf{l}_i = [x, y, z]^T$. With undistorted 2D perspective projection measurement $[u_{ik}, v_{ik}]^T$, we can formulate the visual constraint for a single feature point in a linear equation:

$$\begin{bmatrix} 1 & 0 & -u_{ik} \\ 0 & 1 & -v_{ik} \end{bmatrix} {}^{C_k}\mathbf{l}_i = K_{ik} {}^{C_k}\mathbf{l}_i = 0 \quad (10)$$

Then we transform ${}^{C_k}\mathbf{l}_i$ to \mathbf{l}_i in 0^{th} keyframe coordinate system though pre-calibrated IMU-camera extrinsics $[R_C, \mathbf{p}_C]$ and k^{th} IMU keyframe pose $[R_k, \mathbf{p}_k]$, in which R_k is computed by zero bias gyro measurements integration.

$${}^{C_k}\mathbf{l}_i = R_C^T R_k^T \mathbf{l}_i - R_C^T R_k^T \mathbf{p}_k - R_C^T \mathbf{p}_C \quad (11)$$

Substitute $C_k \mathbf{l}_i$ in Eq. (10) with Eq. (11), we can write Eq. (10) as

$$A_{ik} \mathbf{p}_k + H_{ik} \mathbf{l}_i = \mathbf{b}_{ik} \quad (12)$$

where

$$\begin{aligned} A_{ik} &= -K_{ik} R_C^T R_k^T, \quad H_{ik} = K_{ik} R_C^T R_k^T, \\ \mathbf{b}_{ik} &= K_{ik} R_C^T \mathbf{p}_C \end{aligned}$$

The i^{th} feature point should be observed by at least by 2 keyframes among N keyframes, and we can stack Eq. (12) together for all visual constraints w.r.t one feature point as

$$A_i \mathbf{P} + H_i \mathbf{l}_i = \mathbf{b}_i \quad (13)$$

where

$$\begin{aligned} A_i &= \text{Diag}(A_{ik}), \quad H_i = [H_{i0}; \dots; H_{iN-1}] \\ \mathbf{b}_i &= [\mathbf{b}_{i0}; \dots; \mathbf{b}_{iN-1}] \end{aligned}$$

$A_{ik} = 0_{3 \times 3}$, $H_{ik} = 0_{3 \times 2}$, $\mathbf{b}_{ik} = \mathbf{0}_{2 \times 1}$ if k^{th} keyframe doesn't observe the feature. As described in [2, 5], H_i in Eq. (13) can be projected on its left-nullspace to marginalize out \mathbf{l}_i . Then we have

$$A_i^* \mathbf{P} = \mathbf{b}_i^* \quad (14)$$

Reiterating Eq. (9), for i^{th} feature point observed by N keyframes, we can form the linear equation for the visual constraint

$$F_i^* \mathbf{g} + W_i^* \begin{bmatrix} \mathbf{v} \\ \Delta \mathbf{p}_0 \\ \vdots \\ \Delta \mathbf{p}_{N-1} \end{bmatrix} = \mathbf{r}_i \quad (15)$$

where

$$\begin{aligned} F_i^* &= A_i^* F, \quad W_i^* = A_i^* W \\ \mathbf{r}_i &= \mathbf{b}_i^* - A_i^* \boldsymbol{\Xi} \end{aligned}$$

Closed-form VI-SFM Least Square Problem. Suppose we have \mathcal{F} feature points in total, the closed-form VI-SFM solver is essentially to solve a least square problem with a number of visual constraints as Eq. (15)

$$\begin{aligned} &\text{minimize } \sum_{i \in \mathcal{F}} \left\| F_i^* \mathbf{g} + W_i^* \begin{bmatrix} \mathbf{v} \\ \Delta \mathbf{p}_0 \\ \vdots \\ \Delta \mathbf{p}_{N-1} \end{bmatrix} - \mathbf{r}_i \right\|^2, \\ &\text{subject to } \|\mathbf{g}\|^2 = G \end{aligned} \quad (16)$$

3 Experimental Results for 10KFs under 10/4Hz settings

In Tab. 1, we present results of 10KFs under $10Hz/4Hz$ settings as specified in [1]. We partition the datasets into 1.6 second trajectories for $10Hz$ and 3 seconds for $4Hz$ to run the same exhaustive initialization benchmark. Ours performs best in the $10Hz$ setting, while results are mixed at $4Hz$ where ours performs similarly to the baseline. This is expected, as lower framerates with the same number of keyframes results in overall larger baselines (i.e., more motion). Note that by construction, the $4Hz$ sequences result in slower initialization time (at least 2.25s vs. 0.897s and 0.399s for $10KFs/10Hz$ and $5KFs/10Hz$). For practical applications, faster initialization is preferred.

Table 1: Aggregated initialization benchmark for Inertial-Only, Baseline and our proposed method using various framerates on all EuRoC datasets. For each metric, lower is better.

Metrics	10KFs 10Hz			10KFs 4Hz		
	Inertial-Only	Baseline	Ours	Inertial-Only	Baseline	Ours
Scale Error (%) $\ \bar{a}\ > 0.005G$	35.67	16.36	13.78	14.65	8.76	8.41
Position RMSE (meters)	0.123	0.038	0.034	0.066	0.061	0.064
Gravity RMSE (degrees)	1.89	1.42	1.38	1.36	1.32	1.36
$\log(\mathbf{Condition\ Num}) \ \bar{a}\ < 0.005G$	n/a	12.12	11.23	n/a	10.25	10.54

References

1. Campos, C., Montiel, J.M., Tardós, J.D.: Inertial-only optimization for visual-inertial initialization. In: 2020 IEEE International Conference on Robotics and Automation (ICRA). pp. 51–57. IEEE (2020)
2. Forster, C., Carlone, L., Dellaert, F., Scaramuzza, D.: On-manifold preintegration theory for fast and accurate visual-inertial navigation. CoRR **abs/1512.02363** (2015), <http://arxiv.org/abs/1512.02363>
3. Kaiser, J., Martinelli, A., Fontana, F., Scaramuzza, D.: Simultaneous state initialization and gyroscope bias calibration in visual inertial aided navigation. IEEE Robotics and Automation Letters **2**(1), 18–25 (2017). <https://doi.org/10.1109/LRA.2016.2521413>
4. Li, M., Mourikis, A.I.: A convex formulation for motion estimation using visual and inertial sensors. In: Proceedings of the Workshop on Multi-View Geometry, held in conjunction with RSS. Berkeley, CA (July 2014)
5. Mourikis, A.I., Roumeliotis, S.I.: A multi-state constraint kalman filter for vision-aided inertial navigation. In: Proceedings 2007 IEEE International Conference on Robotics and Automation. pp. 3565–3572. IEEE (2007)