

LVBA: LiDAR-Visual Bundle Adjustment for RGB Point Cloud Mapping (Supplementary)

Please note that equation numbers and section numbers from the main manuscript are labeled in this letter in red.

APPENDIX A COVARIANCE OF SIMPLIFIED PHOTOMETRIC CAMERA MODEL

Our method follows [1] and makes some simplification with it. When a ray is projected to a pixel on the camera, the photometric camera model can be represented as

$$\mathbf{I}(\rho) = \mathbf{f}(\tau V(\rho)\gamma) \quad (1)$$

where ρ is the pixel position, γ is the radiance of the ray, τ is the camera exposure time, $\mathbf{f}(\cdot)$ is the camera response function (CRF), and $V(\cdot)$ is the vignetting factor accounting for the lens vignetting effect. Since the CRF and vignetting factor of the camera are small but hard to obtain after the sequence is collected, we simply ignore them by setting them to one, leading to a camera photometric model as below

$$\mathbf{I} = \mathbf{f}(\tau V(\rho)\gamma) \stackrel{V(\cdot)=1}{\approx} f(\tau\gamma) \stackrel{\mathbf{f}(\mathbf{x})=\mathbf{x}}{=} \tau\gamma \quad (2)$$

Further, we model all sources of measurement noise (e.g. noise from AD transition on camera CMOS) as a Gaussian random noise $\delta_c \sim \mathcal{N}(0, \Sigma_c)$, which will lead to our simplified measurement model:

$$\mathbf{I} = \gamma\tau + \delta_c; \delta_c \sim \mathcal{N}(0, \Sigma_c); \Sigma_c = \begin{bmatrix} \sigma_c^2 & 0 & 0 \\ 0 & \sigma_c^2 & 0 \\ 0 & 0 & \sigma_c^2 \end{bmatrix} \quad (3)$$

where \mathbf{I} is the pixel color we actually measured. σ_c is the covariance of measurement noise for each channel.

In LVBA, we assume that all scene points lie on Lambertian surfaces, hence their radiance γ is the same in all directions. For two camera frames (i.e. the reference frame and target frames), when their poses are the ground truth and observing the same scene point, the observation radiance should be the same, which indicates:

$$\mathbf{0} = \gamma_r - \gamma_t = \tau_r^{-1}\mathbf{I}_r - \tau_t^{-1}\mathbf{I}_t + (\tau_t^{-1}\delta_{ct} - \tau_r^{-1}\delta_{cr}) \quad (4)$$

where γ_r and γ_t are the rays projected to the reference and target frames, respectively. τ_r and τ_t are the exposure times of the reference and target frames. \mathbf{I}_r and \mathbf{I}_t are the corresponding pixel colors when the rays are captured by the camera. δ_{cr} and δ_{ct} are the respective measurement noise. Thus we define the error function \mathbf{e} as:

$$\mathbf{e} := \tau_r^{-1}\mathbf{I}_r - \tau_t^{-1}\mathbf{I}_t \sim \mathcal{N}\left(0, \left(\frac{1}{\tau_t^2} + \frac{1}{\tau_r^2}\right)\Sigma_c\right) \quad (5)$$

The cost item L can be expressed as:

$$L = \|\mathbf{e}\|_{(\tau_t^{-2} + \tau_r^{-2})\Sigma_c}^2 \quad (6)$$

The cost function will be the summation of all generated cost items. It should be noted that if we use (6) to construct the cost function and optimize it, the degrees of freedom for the exposure time is not sufficiently constrained. Within a sequence, there is always a general solution: when all exposure time approaches infinity, all the cost items L will tend to zero. To address this problem, we define a ‘‘Relative Exposure Time (RET)’’ and optimize it instead of the exposure time. For the i -th camera frame with the exposure time τ_i , its relative exposure time is defined as:

$$\epsilon_i = \frac{\tau_i}{\tau_1} \quad (i = r, t) \quad (7)$$

where τ_1 is the exposure time of the first camera frame, and there is always $\epsilon_1 = 1$. Then, the error \mathbf{e} could be written as

$$\mathbf{e} = \tau_r^{-1}\mathbf{I}_r - \tau_t^{-1}\mathbf{I}_t = \tau_1^{-1}\left(\epsilon_r^{-1}\mathbf{I}_r - \epsilon_t^{-1}\mathbf{I}_t\right) \quad (8)$$

Since the constant scale factor τ_1^{-1} only scaled the whole cost function and does not affect the optimal camera state during optimization. We simply ignore it, leading to the error \mathbf{e}_ϵ :

$$\begin{aligned} \mathbf{e}_\epsilon &:= \tau_1\mathbf{e} = \epsilon_r^{-1}\mathbf{I}_r - \epsilon_t^{-1}\mathbf{I}_t \\ &\sim \mathcal{N}\left(0, \tau_1^2\left(\frac{1}{\epsilon_t^2} + \frac{1}{\epsilon_r^2}\right)\Sigma_c\right) \end{aligned} \quad (9)$$

Also, since the value of τ_1 and σ_c in our covariance only scaled the whole cost function, which does not affect the optimal camera states during optimization, we simply set $\sigma_c\tau_1 = 1$. And the final covariance matrix Σ will be

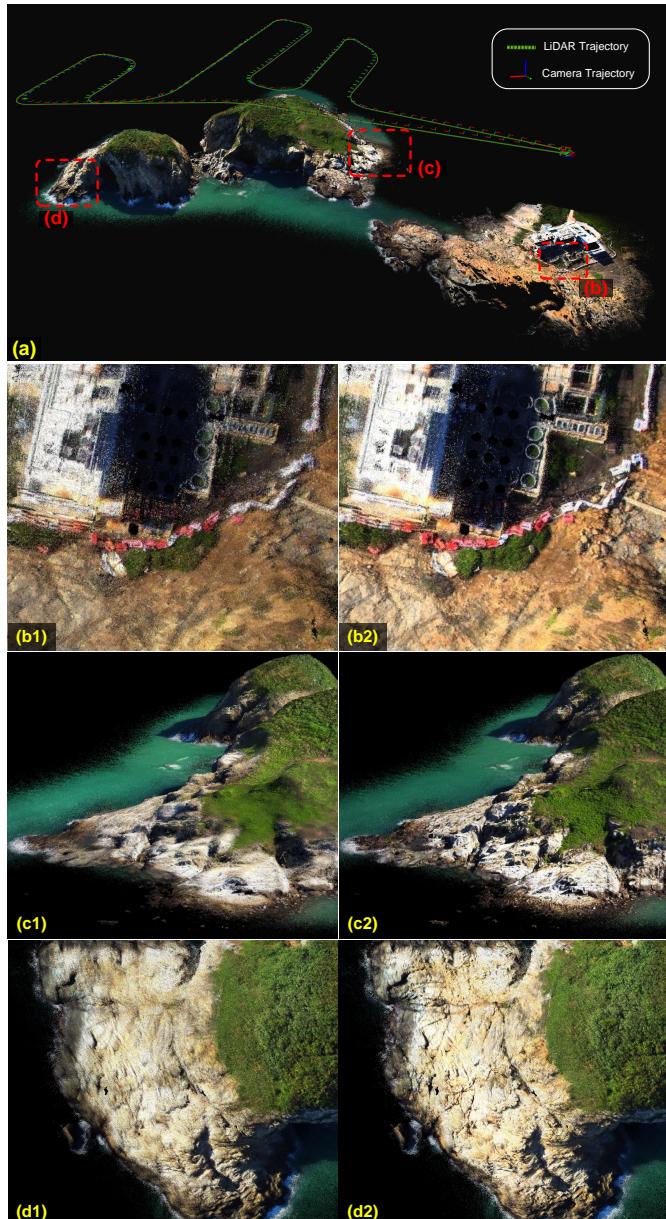
$$\Sigma = \left(\frac{1}{\epsilon_t^2} + \frac{1}{\epsilon_r^2}\right)\mathbf{E} \quad (10)$$

where \mathbf{E} is the 3×3 identity matrix. Combining (9) and (10), the final cost item is $\|\mathbf{e}_\epsilon\|_\Sigma^2$ as illustrated in (8).

It is worthy mention that implementing the covariance (9) in the cost function $\|\epsilon_t^{-1}\mathbf{I}_t - \epsilon_r^{-1}\mathbf{I}_r\|^2$ is necessary. Otherwise, an apparently wrong optimal solution to the cost function would exist, which is $\epsilon_i = \infty, \forall i$. Incorporating the covariance in the cost function would rectify this issue by excluding the wrong optimum.

APPENDIX B MAPPING RESULT FOR ASYNCHRONOUS LiDAR AND CAMERA SEQUENCES

The results of our UAV mapping experiment. Panel (a) offers an overview of the outcome following our optimization process. Panels (b1), (c1), and (d1) show detailed views of the initial values before optimization. In contrast, panels (b2), (c2), and (d2) display the corresponding areas after our optimization.



REFERENCES

- [1] J. Engel, V. Koltun, and D. Cremers, “Direct sparse odometry,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 3, pp. 611–625, 2017.