

DPLVO: Direct Point-Line Monocular Visual Odometry

Lipu Zhou¹, Shengze Wang² and Michael Kaess³

Abstract—In this paper, we present a direct visual odometry (VO) using points and lines. Direct methods generally choose pixels with sufficient gradients to minimize the photometric error for the status estimation. Pixels on lines are generally involved in this process. But the collinear constraint among these points are generally ignored, which may result in less accurate depth estimation. This paper introduces the collinear constraint into the state-of-the-art direct visual odometry DSO [1] to overcome this problem. The 3D lines, points and poses within a sliding window are jointly optimized. DSO implicitly establishes the data association for points among the keyframes within a sliding window by direct image alignment. This scheme is typically suitable for points, as points are generally only visible within a short time window. However, as lines are unbounded entities, they can be observed by a camera significantly longer than points. Thus, we seek to establish the long-term data association for lines among the keyframes. The 3D collinear points that are removed from the sliding window are served as collinear priors for the following windowed optimization. We prove that the prior collinear constraints of a 3D line can be compressed into six residuals in the optimization. This significantly reduces the computational complexity, and enables real-time performance for incorporating long 3D line segments into the windowed optimization. We present a new 3D line representation which reduces the four degrees of freedom (DoF) of a 3D line into two DoF by parameterizing the 3D line in the back-projection plane of its first 2D line observation. The experimental results show that our algorithm outperforms the state-of-the-art direct monocular VO algorithms.

Index Terms—SLAM, Mapping, Localization

I. INTRODUCTION

Simultaneous localization and mapping (SLAM) is a fundamental problem in robotics and computer vision community. Due to its importance, SLAM using various sensors has been extensively studied. Among them, monocular visual SLAM (VSLAM) and visual odometry (VO) have received tremendous attentions, as a monocular camera is widely available. Generally, VSLAM and VO algorithms can be classified into two categories: feature-based and direct methods. In the feature-based method, such as PTAM [2] and ORB-SLAM [3], the descriptors of salient features are extracted and matched along the video sequence. On the other hand, the direct approaches, such as LSD-SLAM [4] and DSO [1], leverage raw pixel intensities and minimize the photometric error. As the direct method can employ more information, it is generally more robust in poorly textured scenes [1]. Direct methods typically adopt pixels with sufficient

Manuscript received: February 24, 2021; Revised: May 31, 2021; Accepted: June 24, 2021.

This paper was recommended for publication by Editor Javier Civera upon evaluation of the Associate Editor and Reviewers' comments. This work was partially done when the first two authors were with the Robot Perception Lab at Carnegie Mellon University.

¹Lipu Zhou is with Magic Leap, 1376 Bordeaux Dr, Sunnyvale, CA 94089, USA lzhou@magic Leap.com

²Shengze Wang is with Department of Computer Science, University of North Carolina, 3175 Brooks, Chapel Hill, NC 27599, USA shengzew@cs.unc.edu

³Michael Kaess is with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA kaess@cmu.edu

Digital Object Identifier (DOI): see top of this page.



Fig. 1. An example of long-term line association. As lines are unbounded, they can generate long-term associations. This increases the correlation among poses, which benefits reducing drift, but may also significantly increase the runtime of the optimization if all the poses where a long 3D line segment is observed are jointly optimized. Thus it is important to design a scheme to balance the efficiency and accuracy. In addition, the line matching is also challenging. Since different parts of a line may be captured at different time, as demonstrated above, the feature-based matching scheme which is well-established for points is not suitable for lines.

gradients, which generally include pixels on lines. These pixels may have one dimensional freedom along the line during photometric tracking, which may result in less accurate depth estimation. We introduce collinear constraints into DSO [1] to improve the depth accuracy of the points on lines.

Lines widely exist in man-made environments. They are good complements to points. However, it is not trivial to exploit lines into VSLAM or VO. A line is an infinite object. Thus it can exist in the field of view (FoV) of a camera much longer than a point, which adds more correlation among poses, as demonstrated in Fig. 1. Although this benefits reducing the drift, it may increase the runtime of the optimization as well. In feature-based methods, the general strategy to employ lines is to detect and match lines [5]–[14]. As the line detection algorithm, such as LSD [15], generally cannot ensure to detect the same line segment in each frame, this may cause the line matching unstable. The deep learning based line detector [16] has got significant progress in recent years. But these algorithms require a powerful GPU, which is generally not available for embedded systems where the computation resource is limited. Additionally, once a part of a line segment moves out or emerge into the camera's FOV, the line matching may also fail. Furthermore, the line detection and matching significantly increase the computational load [5]. On the other hand, direct methods minimize the photometric error. This strategy is suitable for points, but it is difficult to be extended to lines. Thus, the number of studies on extending the direct method with lines is relatively small. In the literature, lines are generally not involved in the optimization in direct methods [17], [18]. Moreover, it is considered infeasible to add geometry prior on points in the direct method for real-time applications [1].

This paper seeks to improve the accuracy of the state-of-the-art direct method DSO [1] by establishing long-term line association and introducing collinear constraints into the windowed optimization. The main contributions of the paper are:

- We introduce the collinear constraint into the photometric error of DSO [19]. We establish long-term data association among different 2D line observations of a 3D line. We fix the poses of the keyframes that are removed from the sliding window and the depths of the collinear points in these removed keyframes to yield the prior collinear constraints for the following windowed optimization. We prove that the prior collinear constraints of a line are equivalent to six residuals in the optimization, which is independent of the number of prior collinear constraints. This makes incorporating long 3D line segments into the

optimization and keeping the real-time performance feasible.

- We present a new singular-free 3D line representation with two degrees of freedom (DoF). It is known that a 3D line has four DoF. The key point of our approach is to represent a line in the back-projection plane of the first 2D line observation.

II. RELATED WORK

Points are the most common features used in VSLAM and VO algorithms [1]–[3], [19], [20]. On the other hand, Lines are abundant in man-made scenes. Thus as an important complement to points, lines have received extensive attention in recent years.

Feature-based method with lines The feature-based method using points alone is prone to fail in low textured scenes, where abundant lines generally exist. As a counterpart of points, it seems natural to extend the well established feature-based framework of points to lines [5]–[14]. In these algorithms, lines are matched along the sequence as done for points. The drawback of this scheme is that there is no guarantee that the same segment of a 3D line can be repetitively detected along the sequence. Since the descriptors may represent different parts of a 3D line, the line matching may be unstable. Furthermore, this scheme cannot establish long-term data association. As a line is an infinite object, the segment of a 3D line captured by a moving camera generally changes. For instance, a new part of a 3D line may gradually emerge into the camera’s FoV, meanwhile an observed part may move out of the camera’s FoV, as demonstrated in Fig. 1. The feature-based method focuses on matching the first detected 2D line segment in the following frames, thus it probably fails once the captured line segment changes.

Direct method with lines The direct method minimizes the photometric error, which is difficult to be extended for lines. Thus the number of studies on the direct method with lines is relatively small. As direct methods generally employ pixels with sufficiently high gradients, pixels on lines are generally selected in the photometric cost. In previous works [17], [18], the 3D lines are used to force the collinear 3D points to meet the collinear constraint. Specifically, the depths of the collinear points are estimated as done for the normal points. The 3D lines are fitted from the estimated 3D points. Then the 3D points are projected into the corresponding 3D lines. Gomez-Ojeda *et al.* [21] introduce lines into the semi-direct method SVO [22]. In [21], a line is represented by two points, and the depth filter introduced in [22] is adopted to estimate the depth of the endpoints of lines. In these methods, the 3D lines themselves are not jointly optimized with the keyframe poses and points. This may result in suboptimal results.

Structural Line The indoor scenario is generally structured, known as the Manhattan world model that has three perpendicular line directions. The lines with these directions are called the structural lines [23], [24]. Vanishing points estimated from the parallel structural lines can be used to identify these directions, which in turn benefit the estimation of the camera orientation [25]–[27]. Some works [23], [24] introduce the structural lines into the extend Kalman filter (EKF) framework. Recently, Zou *et al.* [24] adopt the Atlanta world model to generalize the Manhattan world model. Non-structural lines, which are also important for pose estimation in the low texture areas, are ignored in these works. Besides, the Manhattan assumption limits its applicable scenarios.

Line Representation Representing a 3D Line in the optimization is not trivial. It is known that a 3D line has four DoF [28]. A 3D line can be straightforwardly represented by two 3D points, which is adopted in some previous works [5], [7], [9]. This representation has six DoF. Solà *et al.* [29] present five 3D line representations. But none of them are minimal representation. The overparameterization may result in a rank-deficient Hessian matrix, which may lead to slow convergence and suboptimal results for the Levenberg-Marquardt (LM) algorithm [30]. In [31] and [32], two minimal representations for 3D lines depending on the representation of the rotation matrix are presented. As every minimal representation of the rotation matrix has singular cases, the two representations heritage these singularities. Furthermore, both representations are singular when a 3D line passes

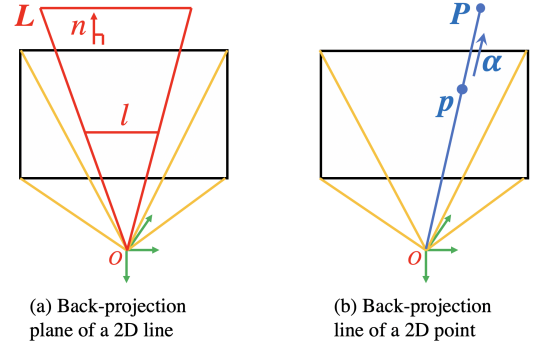


Fig. 2. The back-projections of a 2D line l and a point p on the image plane are a plane and a 3D line which pass through the origin of the camera coordinate system, respectively. l and p determine the normal n of the back-projection plane and the direction α of the back-projection line, respectively. The 3D line L is on the back-projection plane of l . The 3D point P is on the back-projection line of p .

through the origin of the coordinate system. Kottas *et al.* [33] adopt a quaternion and a distance scalar to represent a 3D line. Yang and Huang [34] multiply the quaternion and the distance scalar together to generate a four DoF representation, named the closed point (CP) representation. Like [31] and [32], the CP representation is singular when a 3D line passes through the origin. In [23] and [24], they customize the representation for the structural line in the Manhattan world. Li *et al.* [13] propose a parameterization for co-planar lines. This paper presents a 2D singularity-free representation for 3D lines using the back-projection plane of the first 2D line observation.

III. SYSTEM OVERVIEW

This work introduces lines into DSO [1]. Thus we extend the back end and front end of DSO for lines.

In the back end, we introduce a novel singular-free line parameterization which only has 2 DoF. We impose the collinear constraint on the depths of the 2D points sampled from a 2D line. To adapt the windowed optimization framework of DSO, we introduce the collinear prior for the collinear points that are removed from the sliding window. We prove that, no matter how many collinear priors there are, they can be compressed into six residuals, and can be calculated incrementally. This significantly reduces the computational complexity for long 3D line segments.

In the front end, we introduce the line management which consists of line detection, line initialization and line association. We adopt the LSD algorithm [15] to detect lines, and seek to combine lines with similar parameters. We only detect lines in the area that is not associated to an existing line. For each new 2D line segment, we sample some points from it according to the length of the line segment. The depth of each point is estimated as the ordinary points in DSO. These depths are then used to initialize the 3D line. To avoid the disadvantage of the descriptor-based line matching scheme mentioned above, we introduce a tracking-extension-redistribution method to establish the long-term line association among keyframes.

IV. NOTATIONS AND PRELIMINARIES

In this paper, we use boldfaced letters to represent vectors and matrices, and employ italic lowercase and italic uppercase to represent scalars and functions, respectively.

Camera Pose Let us denote the rotational and translational components of a camera pose transforming a point from the world coordinate system to the camera coordinate system as $\mathbf{R} \in SO(3)$ and $t \in \mathbb{R}^3$, respectively. The pose of a camera is then represented by $\mathbf{T} = \begin{bmatrix} \mathbf{R} & t \\ \mathbf{0} & 1 \end{bmatrix} \in SE(3)$.

Projection and Back-projection In this paper, the projection from a 3D point to the image plane is denoted as $\Pi_c : \mathbb{R}^3 \rightarrow \Omega$, where c represents the camera parameters. Accordingly, given a pixel

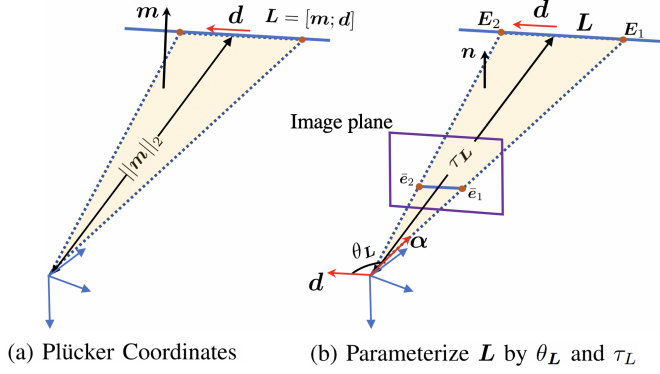


Fig. 3. A schematic of Plücker Coordinates (a) and our new line parameterization (b). Given a 2D line, the 3D line is on its back-projection plane, as demonstrated in Fig. 2. Thus we can represent a 3D line by 2 parameters (i.e., τ_L and θ_L). The parameterization is introduced in section V-B1.

and its depth, the back-projection is defined as $\Pi_c^{-1} : \Omega \times \mathbb{R} \rightarrow \mathbb{R}^3$, as demonstrated in Fig. 2 (b). Let us assume $\mathbf{p} \in \Omega_i$ is a point in a reference image I_i with the inverse depth d_p that is observed at \mathbf{p}' in a target image I_j . Suppose the poses of I_i and I_j are \mathbf{T}_i and \mathbf{T}_j , respectively. Then the relationship between \mathbf{p} and \mathbf{p}' has the form

$$\mathbf{p}' = \Pi_c(\mathbf{R}_{ij}\Pi_c^{-1}(\mathbf{p}, d_p) + \mathbf{t}_{ij}), \quad (1)$$

where \mathbf{R}_{ij} and \mathbf{t}_{ij} are the rotational and translational component of $\mathbf{T}_j\mathbf{T}_i^{-1}$.

Plücker Coordinates A 3D line L can be represented by the Plücker coordinates, as illustrated in Fig. 3 (a). The following formulas used in this paper are according to [28] and [35].

Let π denote the plane determined by L and the origin of the coordinate system, and τ_L the distance from the origin to L . Then the Plücker coordinates of L is a six-dimensional vector defined as

$$L = [m; d], \quad \text{s.t. } \|d\|_2 = 1, \quad d^T m = 0, \quad (2)$$

where d is the direction of L , and m is referenced to as the moment vector, which is perpendicular to π with $\|m\|_2 = \tau_L$. The Plücker coordinates are homogeneous coordinates. Here we impose $\|d\|_2 = 1$ on them to remove the scale ambiguity. Besides, this formulation can lead to simpler formulas in the following calculation.

Using the Plücker coordinates in (2), we can write the residual vector of the distance from a point $\mathbf{X} \in \mathbb{R}^3$ to L as

$$e(L, \mathbf{X}) = \mathbf{m} - \mathbf{X} \times \mathbf{d}. \quad (3)$$

Given \mathbf{R} and \mathbf{t} , the transformation matrix of the Plücker coordinates has the form

$$H(\mathbf{R}, \mathbf{t}) = \begin{bmatrix} \mathbf{R} & [\mathbf{t}]_{\times} \mathbf{R} \\ \mathbf{0} & \mathbf{R} \end{bmatrix}. \quad (4)$$

Using (4), we can write the transformation of the Plücker coordinates from one coordinate system to another as

$$L' = H(\mathbf{R}, \mathbf{t})L. \quad (5)$$

The Plücker coordinates are convenient for computation. However, due to the constraints in (2), this representation is not convenient in optimization. In this paper, we represent the Plücker coordinates by two parameters, which leads to unconstrained optimization.

V. DIRECT POINT-LINE MODEL

A. Photometric Error

Direct method minimizes the photometric error. We adopt the photometric error model introduced in [1]. The photometric error for \mathbf{p} with inverse depth d_p in I_i observed by I_j is defined as:

$$E_{p_j} = \sum_{\mathbf{p} \in \mathbb{N}_p} w_p \left\| (I_j[\mathbf{p}'] - b_j) - \frac{t_j e^{a_j}}{t_i e^{a_i}} (I_i[\mathbf{p}] - b_i) \right\|_{\gamma}, \quad (6)$$

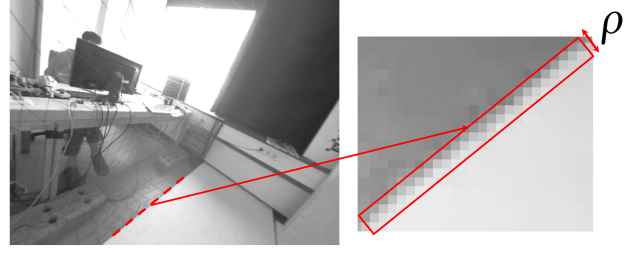


Fig. 4. Line support region. As introduced in [15], a real line segment has a support region due to the quantization or motion blur, which can represent the uncertainty of a line. We use this information to weigh the collinear constraints. ρ is the width of the line support region.

where a_i, a_j, b_i and b_j are the affine brightness transform parameters, t_i and t_j are the exposure time of I_i and I_j , w_p is a gradient-dependent weighting factor, \mathbb{N}_p is a set of neighborhoods around \mathbf{p} , and $\|\cdot\|_{\gamma}$ is the Huber norm. The relationship between \mathbf{p} and \mathbf{p}' is given in (1). Consequently, the full photometric error over all frames \mathbb{F} has the form as

$$E_{photo} = \sum_{i \in \mathbb{F}} \sum_{\mathbf{p} \in \mathbb{P}_i} \sum_{j \in obs(\mathbf{p})} E_{p_j}, \quad (7)$$

where \mathbb{P}_i is the set of points in frame i , and $obs(\mathbf{p})$ is the set of frames where \mathbf{p} is visible.

B. Collinear Constraint

We impose the collinear constraints on the points on lines. We first present a novel 3D line representation. Next, we introduce the cost function of the collinear constraint for the keyframes in the sliding window. Finally, we discuss the collinear prior yielded by the removed collinear points.

1) 3D Line Representation: We call the image where a 3D line L is first observed as the **anchor image** of L . On the other hand, we call the subsequent images observing L as the **associate images**. Here we show that L in (2) can be parameterized by two unknowns in the coordinate system of the anchor image of L , as demonstrated in Fig. 3 (b).

We assume the image has been undistorted and the camera matrix is \mathbf{K} , and suppose l_{\perp} is the captured image of L in the anchor image. The back-projection of l_{\perp} is a plane demonstrated in Fig. 2. Assume \bar{e}_1 and \bar{e}_2 are the homogeneous coordinates of the two endpoints of l_{\perp} , and E_1 and E_2 are the corresponding 3D points on L . Then we have $l_{\perp} = \bar{e}_1 \times \bar{e}_2$, and the direction of L is $d = \frac{E_2 - E_1}{\|E_2 - E_1\|_2}$. The norm of the back-projection plane $\pi_{l_{\perp}}$ of l_{\perp} in the camera coordinate system is $n = \frac{\mathbf{K}^T l_{\perp}}{\|\mathbf{K}^T l_{\perp}\|_2}$. The direction of the back-projection ray of \bar{e}_1 is $\alpha = \frac{\mathbf{K}^{-1} \bar{e}_1}{\|\mathbf{K}^{-1} \bar{e}_1\|_2}$. Assume θ_L is the angle from α to d , and τ_L is the distance from the camera origin to the 3D line. Then the Plücker coordinates L defined in (2) can be parameterized by θ_L and τ_L as below

$$\begin{aligned} m &= \tau_L n, \\ d &= \cos(\theta_L) \alpha + \sin(\theta_L) \beta, \end{aligned} \quad (8)$$

where $\beta = n \times \alpha$.

2) Collinear Constraint from Anchor Image of L : Suppose p_{\perp}^L with depth $d_{p_{\perp}^L}$ is a point on l_{\perp} in the anchor image I_i of L . Then the corresponding 3D point of p_{\perp}^L in the coordinate system of I_i has the form $\mathbf{X}_{\perp}^L = \Pi_c^{-1}(p_{\perp}^L, d_{p_{\perp}^L})$. Here \mathbf{X}_{\perp}^L and L are in the same coordinate system. Directly using (3), we get the collinear constraint on \mathbf{X}_{\perp}^L as below

$$E_{p_{\perp}^L}^L = \left\| e(L, \mathbf{X}_{\perp}^L) \right\|_2^2. \quad (9)$$

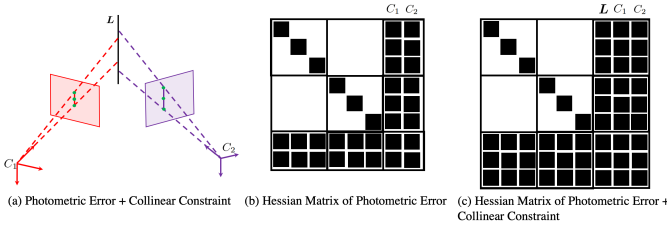


Fig. 5. A schematic of the Hessian matrix of the photometric error (7) and the photometric error with the collinear constraint (15). 3 points are sampled from each 2D line. By properly aligning the variables, the Hessian matrix of (15) also has a sparse pattern. Thus the resulting linear system of the iterative algorithm can be efficiently solved [36].

Let us denote the set of sampled points on \mathbf{l}_\perp as \mathbb{A}^L . The collinear constraint on the 3D points of \mathbb{A}^L is given by

$$E_\perp^L = \sum_{\mathbf{p}_\perp^L \in \mathbb{A}^L} E_{\mathbf{p}_\perp^L}^L. \quad (10)$$

3) *Collinear Constraint from Associate Image of \mathbf{L}* : Let us use \mathbf{l}_j^L to denote the 2D line of \mathbf{L} in the associate image I_j . Assume \mathbf{p}_j^L is a pixel on \mathbf{l}_j^L . Given its depth $d_{\mathbf{p}_j^L}$, we can get the corresponding 3D point coordinates $\mathbf{X}_j^L = \Pi_c^{-1}(\mathbf{p}_j^L, d_{\mathbf{p}_j^L})$ in the camera coordinate system of I_j . Here \mathbf{X}_j^L and \mathbf{L} are in different coordinate systems. Using (3) and (5), we get the collinear constraint on \mathbf{X}_j^L

$$E_{\mathbf{p}_j^L}^L = \left\| e(\mathbf{L}_j, \mathbf{X}_j^L) \right\|_2^2, \quad \mathbf{L}_j = H(\mathbf{R}_{i_j}, \mathbf{t}_{i_j}) \mathbf{L}, \quad (11)$$

where \mathbf{R}_{i_j} and \mathbf{t}_{i_j} are defined in (1). Assume \mathbb{B}_j^L denotes the set of sampled points on \mathbf{l}_j . Then the constraint on the 3D points of \mathbb{B}_j^L has the form as

$$E_j^L = \sum_{\mathbf{p}_j^L \in \mathbb{B}_j^L} E_{\mathbf{p}_j^L}^L. \quad (12)$$

4) *Collinear Constraints*: As shown in Fig. 4, a line segment has a support region [15] due to the quantization or motion blur, which represents the uncertainty of the line segment. Let us denote the width of the support region as ρ . We use $\frac{1}{\rho}$ to weigh the collinear constraints. Unlike [18] that samples the same number of points for each line, we split the support region into several segments with a fixed length ϵ . In each segment, we select the point with the largest gradient. As a longer line generally has a lower uncertainty on its parameters, this sampling method can increase the weight of a longer line by introducing more points.

The 3D line \mathbf{L} is observed in its anchor image and the subsequent associate images. Let \mathbb{H}^L denote the set of associate images of \mathbf{L} . Using (10) and (12), we get the collinear constraint from \mathbf{L} as

$$E^L = \frac{1}{\rho_\perp^L} E_\perp^L + \sum_{j \in \mathbb{H}^L} \frac{1}{\rho_\perp^L + \rho_j^L} E_j^L, \quad (13)$$

where ρ_\perp^L and ρ_j^L are the widths of the support region of the line segments in the anchor and associate images. Let \mathbb{L} denote the set of all the 3D lines. Then the full cost function over \mathbb{L} has the form

$$E_{line} = \sum_{\mathbf{L} \in \mathbb{L}} E^L. \quad (14)$$

C. Full Cost Function

In our algorithm, the camera poses, 3D points and lines are jointly adjusted to minimize the photometric error as well as to meet the collinear constraint. Specifically, the full cost function over all frames, points and lines is the combination of (7) and (14)

$$E = E_{photo} + E_{line}, \quad (15)$$

Fig. 5 demonstrates the Hessian matrix of (15). It is of the general sparse pattern, thus the resulting linear system of the iterative algorithm can be efficiently solved [36].

D. Windowed Optimization

Minimizing the full cost (15) is computationally demanding. Thus we adopt the sliding window approach in [1] to balance the computational complexity and accuracy. For the removed keyframes, we adopt different strategies for E_{photo} and E_{line} .

1) *Marginalization*: For E_{photo} , we adopt the marginalization strategy used in [1]. Briefly, inactive keyframes and points are marginalized using the Schur complement, and all the remaining observations related to the inactive keyframes are removed to preserve the sparsity of the Hessian matrix. This leads to a quadratic function on the remaining unknowns that can be used in the subsequent optimization and marginalization. For 3D lines, we want to use the old observations to constrain the subsequent optimization, and update the line parameters when new observations emerge. Thus, we do not marginalize them. Instead, we introduce prior collinear constraints for a 3D line when its anchor image or associate images are removed from the sliding window.

2) *Prior Collinear Constraint*: We fix the poses of the keyframes removed from the sliding window and the depths of the collinear pixels in these keyframes. Given a fixed pose, we can transform the collinear points in this local camera coordinate system to the global coordinate system. These fixed points provide collinear priors to the corresponding lines in the subsequent windowed optimization. We denote the set of fixed points on \mathbf{L} as \mathbb{X}_f^L . Suppose χ_k^L is the k th point in \mathbb{X}_f^L . Based on (3) and (5), the prior collinear constraint from χ_k^L on \mathbf{L} has the form

$$E_{\chi_k^L}^L = \left\| e(\chi_k^L, \mathbf{L}') \right\|_2^2, \quad \mathbf{L}' = H(\mathbf{R}_i^T, -\mathbf{R}_i^T \mathbf{t}_i) \mathbf{L}, \quad (16)$$

where \mathbf{R}_i and \mathbf{t}_i are the rotation matrix and translation vector from the global coordinate system to the coordinate system of the anchor image I_i of \mathbf{L} , respectively. Stacking the constraints in (16) for all $\chi_k^L \in \mathbb{X}_f^L$, we get the collinear prior on \mathbf{L}

$$E_{prior}^L = \left\| \mathbf{f}^L \right\|_2^2, \quad \mathbf{f}^L = \left[\dots, e(\chi_k^L, \mathbf{L}')^T, \dots \right]^T. \quad (17)$$

Suppose the size of \mathbb{X}_f^L is N^L . Then the size of \mathbf{f}^L is $3N^L$.

3) *Compressed Prior Collinear Constraint*: As keyframes are removed, the size of \mathbf{f}^L in (17) enlarges, which seemingly inevitably increases the runtime of the optimization. Here we show that the computational complexity of the collinear prior in the optimization can be independent of N^L . We first consider the form of \mathbf{f}^L and its Jacobian matrix in terms of the introduced line parameterization (8). We have the following lemma:

Lemma 1: For the parameterization (8) of \mathbf{L} , \mathbf{f}^L and its Jacobian matrix $\mathbf{J}_{\mathbf{f}^L}$ can be respectively factorized as follows

$$\mathbf{f}^L = \mathbf{A}^L \mathbf{L} \text{ and } \mathbf{J}_{\mathbf{f}^L} = \mathbf{A}^L \mathbf{B}^L, \quad (18)$$

where \mathbf{L} is the Plücker coordinates given in (8), and \mathbf{A}^L and \mathbf{B}^L are defined in (25) and (23) with size $3N^L \times 6$ and 6×2 , respectively.

Let us define $\mathbf{M}^L = (\mathbf{A}^L)^T \mathbf{A}^L$. We have the following lemma for \mathbf{M}^L :

Lemma 2: \mathbf{M}^L can be factorized as $\mathbf{M}^L = (\mathbf{C}^L)^T \mathbf{C}^L$, where \mathbf{C}^L is a 6×6 matrix.

We define the compressed residual vector for \mathbf{f}^L as

$$\mathbf{g}^L = \mathbf{C}^L \mathbf{L}. \quad (19)$$

The size of \mathbf{g}^L is independent of the number of points in \mathbb{X}_f^L . In the following theorem, we show that \mathbf{g}^L is sufficient for the Gauss-Newton or LM algorithm.

Theorem 1: \mathbf{g}^L can replace \mathbf{f}^L in the Gauss-Newton or LM algorithm.

The computational complexity for calculating the compressed residual vector \mathbf{g}^L , the Jacobian matrix $\mathbf{J}_{\mathbf{g}^L}$, $\mathbf{J}_{\mathbf{g}^L}^T \mathbf{J}_{\mathbf{g}^L}$ and $\mathbf{J}_{\mathbf{g}^L}^T \mathbf{g}^L$ are $\frac{6}{N^L}$ of the corresponding ones for \mathbf{f}^L . This can significantly reduce the computational cost when N^L is large as demonstrated in Fig. 1. We prove Lemma 1, 2 and Theorem 1 in the appendix.

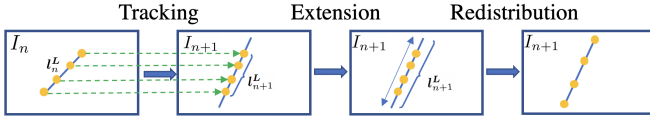


Fig. 6. Our line association between the keyframe I_n and I_{n+1} consists of three steps: tracking, extension and redistribution. In the tracking step, each point on line l_n^L is individually tracked by minimizing the photometric error. Here we only require that the matched point in I_{n+1} is on the corresponding line, and we do not require an exact point-to-point correspondence. Then we seek to extend the line l_{n+1}^L obtained from the tracking to reflect the new observation of the line in I_{n+1} . Finally, we resample some points from l_{n+1}^L .

4) *Incremental Calculation of M^L* : The most important step to get g^L in (19) is to calculate M^L . The computational complexity will increase when N^L enlarges. Here we show that it is not necessary to compute M^L from scratch every time. Instead, M^L can be calculated incrementally. Assume I_k is removed from the sliding window and N_k^L 3D points on L are observed in I_k . We can calculate A_k for the new coming N_k^L 3D points by (22). Next we append A_k^L to the current A^L to generate the new one $A_{new}^L = [A^L; A_k^L]$. Then M_{new}^L has the form

$$\begin{aligned} M_{new}^L &= (A_{new}^L)^T A_{new}^L = (A^L)^T A^L + A_k^T A_k \\ &= M^L + A_k^T A_k \end{aligned} \quad (20)$$

VI. FRONT END

The front end determines and initializes the points, lines and keyframe poses used in the back end. We introduce the line management into the front end of DSO, and slightly change the point and frame management of DSO.

A. Line Management

The line management includes line detection, initialization and association.

1) *Line Detection*: We adopt the LSD algorithm [15] to detect lines. Except for the first keyframe, the line detection is conducted after the line association described in the section VI-A3. Thus, starting from the second keyframe, we only detect lines in the regions which are not associated with the existing lines. As a line segment may be detected as several shorter ones, we seek to merge the line segments with similar parameters. Specifically, assume the two line segments l_1 and l_2 have the directions v_1 and v_2 , and the distances from the origin to l_1 and l_2 are c_1 and c_2 . If $\text{angle}(v_1, v_2) < \vartheta$ and $|c_1 - c_2| < \varrho$, we calculate a new line l_3 using all the points of l_1 and l_2 . Here $\text{angle}(\cdot, \cdot)$ represents the angle between two vectors. If more than $q\%$ of the points are within a distance ζ to the new line, we merge l_1 and l_2 . In our experiments, we set $\vartheta = 10^\circ$, $\varrho = 10$ pixels, $q = 95$ and $\zeta = 2$ pixels. The new 3D lines are initialized as described below.

2) *Line Initialization*: Assume l_m^L is the m th new 2D line segment detected in the latest keyframe I_n . We sample N 2D points from it. In subsequent frames, the N points are then individually tracked along the epipolar line, and their depths are estimated as done for the ordinary point in DSO. After their depths are estimated, we check the quality of these depths. Specifically, given the depths, we calculate the 3D points using the back-projection function Π_c^{-1} . Then we conduct the principal component analysis (PCA) on these 3D points. Assume λ_1 , λ_2 and λ_3 are the 3 resulting principal component variances in descending order. For an idea 3D line, we have $\lambda_2 = \lambda_3 = 0$. Thus, we consider the initialization succeeds if $\frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3} > \kappa$. If the initialization succeeds, we compute the parameters of the 3D line as introduced in section V-B1, and all the 3D points are projected onto the line and their depths are revised accordingly. We set $\kappa = 0.7$ in our experiments.

3) *Line Association*: We conduct line association for each initialized lines in a new keyframe I_{n+1} . Instead of adopting the feature-based line matching strategy [37], we introduce a **tracking-extension-redistribution** method to establish the line association among keyframes, as shown in Fig. 6.

Tracking: Let l_n^L denote the image of the 3D line L in the n th keyframe I_n . Assume $p_{n,k}^L$ is the k th point on l_n^L . We compute the corresponding point of $p_{n,k}^L$ in I_{n+1} by (1), and denote it as $\hat{p}_{n+1,k}^L$. If $\hat{p}_{n+1,k}^L$ is out of I_{n+1} , it is ignored. Inspired by [22], we compute a correct $\delta u \in \mathbb{R}^2$ for $\hat{p}_{n+1,k}^L$ by minimizing the photometric error (6). Unlike [22], here the corrected point $\tilde{p}_{n+1,k}^L = \hat{p}_{n+1,k}^L + \delta u$ does not have to match $p_{n,k}^L$. Instead, we only require $\tilde{p}_{n+1,k}^L$ can get closer to the line. Then we fit a line l_{n+1}^L in I_{n+1} using these corrected points.

Extension: As a new part of L may be captured in I_{n+1} , we seek to extend l_{n+1}^L to include the new observation. Motivated by [15], we introduce a directional region growing method to extend l_{n+1}^L . Assume the length of l_{n+1}^L is d_{n+1}^L . Let \bar{g} , δ and g_{min} denote the mean, standard deviation and minimum of the norms of the gradients at $\lceil d_{n+1}^L \rceil$ points uniformly distributed on l_{n+1}^L , respectively. Here $\lceil \cdot \rceil$ represents the ceiling of a real number. Let e_{n+1}^L denote one of the endpoints of l_{n+1}^L , and η_{n+1}^L represent the direction of l_{n+1}^L . For each step, we first seek to extend e_{n+1}^L with ν pixels. Specifically, we check the candidate endpoint $e_{n+1}^L = e_{n+1}^L + \nu \eta_{n+1}^L$. If the norm of the gradient at e_{n+1}^L is larger than $\max(\bar{g} - 2\delta, g_{min})$ and the angle between its gradient and the normal of l_{n+1}^L is less than 22.5° , e_{n+1}^L is considered as the new endpoint of l_{n+1}^L . The threshold 22.5° is recommended by [15]. We continue this process until the test fails. We set $\nu = 5$ pixels in our experiments.

Redistribution: We sample N points from l_{n+1}^L . Let $p_{n+1,k}^L$ denote the k th point on l_{n+1}^L . The depth of $p_{n+1,k}^L$ is initialized as the depth of the 3D point on the back-projection ray of $p_{n+1,k}^L$ that is closest to L .

B. Point and Frame Management

We slightly change the point and frame management in DSO [1]. For the point management, when selecting the candidate points, we avoid choosing the point whose distance to any of the lines is less than f pixels. For the frame management, we introduce a new criterion to determine whether a new keyframe is needed for lines. In detail, we add a new keyframe if there are more than k new lines have been initialized or the total length of the initialized lines exceeds γ pixels. We sample some points from each 2D line segment according to the length of the 2D line segment. For frame pose estimation, these collinear points are treated as ordinary points and combined with the ordinary points to get the pose of a new frame by direct image alignment. So lines do not impact on the runtime of the motion estimation. In our experiments, we set $f = 5$, $k = 3$ and $\gamma = 100$.

VII. EXPERIMENTAL RESULTS

In this section, we compare the performance of our algorithm with the state-of-the-art VO or SLAM algorithms using the TUM monoVO dataset [38], ICL_NUM [39], EuRoC MAV Dataset [40] and TUM RGB-D dataset [41].

A. The TUM monoVO Dataset

The TUM monoVO dataset contains 50 photometrically calibrated and loop-closed sequences without the ground-truth pose. We run each sequence forwards and backwards 10 times to account for the nondeterministic behavior. As this dataset does not provide the ground-truth trajectory but has the loop-closure-ground-truth, we adopt the alignment error e_{align} , rotational drift e_r and scale drift e_s introduced in [38] to evaluate the performance. As in [1], we adopt the cumulative error plots to present the results, which reveal both accuracy and robustness of an algorithm by counting how many runs its errors are below a certain threshold.

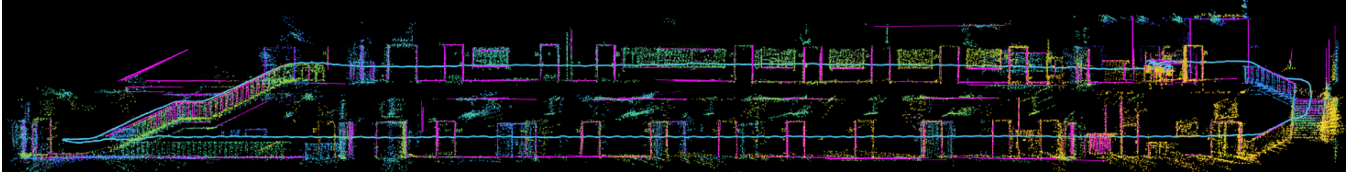


Fig. 7. The trajectory, lines and points generated by our algorithm on the TUM monoVO dataset [38].

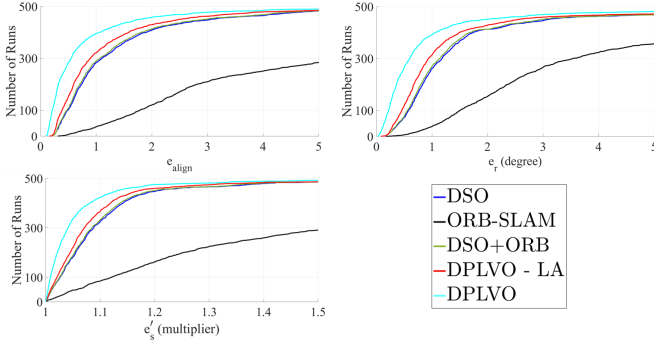


Fig. 8. The cumulative error of alignment error e_{align} , rotational drift e_r and the revised scale drift $e'_s = \max(e_s, e_s^{-1})$ of different algorithms on the TUM monoVO dataset [38].

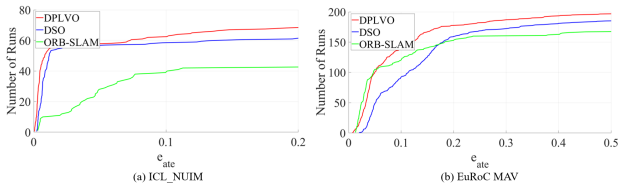


Fig. 9. The cumulative error of the absolute trajectory error e_{ate} of our algorithm, DSO [1] and ORB-SLAM [3] on the ICL_NUM [39] dataset and the EuRoC MAV dataset [40].

Ablation Study We first use the TUM monoVO dataset [38] to evaluate the different design options of our algorithm. We compare the following two variants of our algorithm:

- **DPLVO - LA:** The line association is removed from our algorithm. That is to say this version does not include E_j^L defined in (12) in the cost function.
- **DSO + ORB:** We add the feature-based point matching for the corners and the resulting reprojection errors in DSO.

The results in Fig. 8 show that the collinear constraint alone can improve the performance of DSO. The proposed line association scheme can further increase the accuracy. The difference between the results of DSO and DSO + ORB is marginal. The optical flow tracking for the corner feature is well defined. Additional feature-based matching cannot significantly improve the performance. On the other hand, the number of edge points is generally much larger than the number of corner features in the man-made scenarios, as demonstrated in Fig. 1, and the photometric tracking for the points on the line is problematic. The collinear constraint and line association can alleviate this problem, thus they improve the results.

Comparison with the State-of-the-art We compare the performance of our algorithm with DSO [1] and ORB-SLAM [3]. Fig. 8 presents the results. It is clear that our algorithm outperforms DSO and ORB-SLAM. Our algorithm establishes long-term line association among keyframes which reduces the drift.

B. The EuRoC MAV and ICL_NUM Dataset

The ICL_NUM dataset [39] and EuRoC MAV dataset [40] consist of 11 and 8 sequences, respectively. In these sequences, the

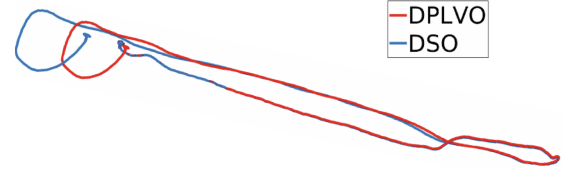


Fig. 10. The trajectories estimated by our algorithm and DSO [1] on sequence 40 of the TUM monoVO dataset [38]. The start point and the end point of this sequence are at the small place. The gap between the start and the end points of DSO is clearly larger than the one of our algorithm.

TABLE I
ABSOLUTE TRAJECTORY ERROR (CM) ON THE TUM RGBD BENCHMARK [41]. THE RESULTS ARE THE MEDIAN OVER 5 EXECUTIONS FOR EACH SEQUENCE.

Sequence	Ours	DLGO [18]	DSO [1]	PL-SLAM [5]	ORB-SLAM [3]
fr1_xyz	3.85	5.38	6.37	1.21	0.9
fr3_walk_xyz	9.7	27.5	15.5	1.54	1.24
fr3_walk_half	23.4	37.4	32.6	1.60	1.74

photometric calibration and the exposure time are not available. As the two datasets provide the ground-truth trajectories, we adopt the absolute trajectory error e_{ate} (ATE) to evaluate the performance. Again, we run each sequence 10 times forwards and backwards, and adopt the cumulative error plot to present the results in Fig. 9. The results show that our algorithm outperforms ORB-SLAM [3] and DSO [1] on the ICL_NUM dataset. ORB-SLAM achieves a better accuracy but less robust than DSO [1] and our algorithm. But our algorithm significantly reduces the gap between DSO and ORB-SLAM. One drawback of DSO is that it does not establish the accurate data correspondence. Our line association scheme overcomes this drawback, thus our algorithm improves the accuracy.

C. The TUM RGB-D Dataset

We also use the TUM RGB-D dataset [41] to evaluate the performance of our algorithm. The results are listed in table I. It is clear that our algorithm outperforms the state-of-the-art direct methods [1] and [18]. The TUM RGB-D dataset is challenging for the direct method. This dataset does not provide the photometric calibration and the exposure time. In addition, since the camera in the TUM RGB-D dataset is a rolling shutter camera, simply employing the direct method does not work well [42]. Thus the feature-based methods PL-SLAM [5] and ORB-SLAM surpass the direct methods in this dataset.

D. Runtime

We ran the runtime experiments on a laptop with an i7 3.4GHZ CPU and 16G memory. On the TUM monoVO dataset, the average runtime for the back end of our algorithm and DSO is 166.5 ms and 140.5 ms, respectively. Lines introduce more parameters, thus the time for optimization increases. The average runtime for the direct frame alignment, and point and frame management in the front end is 18.2 ms and 15.9 ms for our algorithm and DSO, respectively. Our algorithm treats the collinear points as ordinary points in the front

end, thus the runtime is similar. As our algorithm may introduce more points into these steps, thus the runtime is a bit longer. Our algorithm requires an additional 11.3 *m.s* for the line management on average, but these computation only requires when a new keyframe is added. In addition, as the line management is independent of the point and frame management, we can run them in parallel. Therefore, the front end of our algorithm can achieve the real-time performance.

E. Qualitative Results

Fig. 7 demonstrates the 3D lines and point could generated by our algorithm. Fig. 10 shows the trajectories of our algorithm and DSO on sequence 40 of the TUM monoVO dataset. This sequence contains a loop. The camera returns back to the original location. The gap between the start and the final points of DSO is larger than the one of our algorithm.

VIII. CONCLUSION

In this paper, we introduce the collinear constraint into DSO [1]. We establish long-term data association for lines. The collinear constraints in the removed keyframes are treated as priors in the following windowed optimization. We prove that no matter how many prior collinear constraints a line has, they can be compressed into six constraints. So the computation will not increase as new observations emerge. Furthermore, the resulting Hessian Matrix is still sparse, thus the optimization is efficient. Experimental results show that our algorithm significantly improves DSO and outperforms the state-of-the-art direct VO algorithms.

APPENDIX

A. Proof of Lemma 1

Proof. Let us first expand $e(\chi_k^L, L')$ in (16). Substituting the definition of H in (4) into (16), we have

$$\begin{aligned} e(\chi_k^L, L') &= \mathbf{R}_i^T \mathbf{m} - \left[\mathbf{R}_i^T \mathbf{t}_i \right]_{\times} \mathbf{R}_i^T \mathbf{d} - \left[\chi_k^L \right]_{\times} \mathbf{R}_i^T \mathbf{d} \\ &= \mathbf{R}_i^T \mathbf{m} - \underbrace{\left(\left[\mathbf{R}_i^T \mathbf{t}_i \right]_{\times} \mathbf{R}_i^T + \left[\chi_k^L \right]_{\times} \mathbf{R}_i^T \right)}_{\mathbf{U}_k} \mathbf{d} \\ &= \mathbf{R}_i^T \mathbf{m} - \mathbf{U}_k \mathbf{d} \end{aligned} \quad (21)$$

Substituting the line parameterization (8) into (21), we get

$$\begin{aligned} e_k(\tau_L, \theta_L) &= \tau_L \mathbf{R}_i^T \mathbf{n} - \cos(\theta_L) \mathbf{U}_k \boldsymbol{\alpha} - \sin(\theta_L) \mathbf{U}_k \boldsymbol{\beta} \\ &= \underbrace{\left[\mathbf{R}_i^T, -\mathbf{U}_k \right]}_{\mathbf{A}_k^L} \underbrace{\begin{bmatrix} \tau_L \mathbf{n} \\ \cos(\theta_L) \boldsymbol{\alpha} + \sin(\theta_L) \boldsymbol{\beta} \end{bmatrix}}_{\mathbf{L}} = \mathbf{A}_k^L \mathbf{L}. \end{aligned} \quad (22)$$

Then we can calculate the Jacobian matrix \mathbf{J}_k of $e_k(\tau_L, \theta_L)$. It has the form

$$\mathbf{J}_k^L = \underbrace{\left[\mathbf{R}_i^T, -\mathbf{U}_k \right]}_{\mathbf{A}_k^L} \underbrace{\begin{bmatrix} \mathbf{n}, & \mathbf{0} \\ \mathbf{0}, & -\sin(\theta_L) \boldsymbol{\alpha} + \cos(\theta_L) \boldsymbol{\beta} \end{bmatrix}}_{\mathbf{B}^L} = \mathbf{A}_k^L \mathbf{B}^L \quad (23)$$

Stacking all the residual vector from \mathbb{X}_f^L , we can rewrite \mathbf{f}^L defined in (17) as

$$\mathbf{f}^L = \underbrace{\left[\dots, (\mathbf{A}_k^L)^T, \dots \right]^T}_{\mathbf{A}^L} \mathbf{L} = \mathbf{A}^L \mathbf{L}. \quad (24)$$

Similarly, stacking the Jacobian matrix \mathbf{J}_k^L for $\chi_k^L \in \mathbb{X}_f^L$, we get the Jacobian matrix of \mathbf{f}^L

$$\mathbf{J}_{f^L} = \left[\dots, (\mathbf{A}_k^L)^T, \dots \right]^T \mathbf{B}^L = \mathbf{A}^L \mathbf{B}^L. \quad (25)$$

B. Proof of Lemma 2

Proof. For any non-zero $\mathbf{x} \in \mathbb{R}^6$, we have $\mathbf{x}^T \mathbf{M}^L \mathbf{x} = \mathbf{x}^T (\mathbf{A}^L)^T \mathbf{A}^L \mathbf{x} \geq 0$. Thus \mathbf{M}^L is a symmetric positive semi-definite matrix. Let us denote the eigendecomposition of \mathbf{M}^L as $\mathbf{Q}^T \boldsymbol{\Lambda} \mathbf{Q}$. Here $\boldsymbol{\Lambda}$ is a diagonal matrix. Denote the i th diagonal element of $\boldsymbol{\Lambda}$ as λ_i . As \mathbf{M}^L is positive semi-definite, we have $\lambda_i \geq 0$. Let us define $\mathbf{C} = \boldsymbol{\Lambda}^{\frac{1}{2}} \mathbf{Q}$, where $\boldsymbol{\Lambda}^{\frac{1}{2}}$ is a diagonal matrix whose i th diagonal element is $\sqrt{\lambda_i}$. It is easy to verify that $\mathbf{M}^L = (\mathbf{C}^L)^T \mathbf{C}^L$. Here \mathbf{C} is a 6×6 matrix, whose size is independent of the number of points. \square

C. Proof of Theorem 1

Proof. Let us divide the residual vector $\boldsymbol{\delta}$ into two parts. One is \mathbf{f}^L defined in (17), and the other one is \mathbf{h} containing the remaining residuals. We use \mathbf{g}^L to replace \mathbf{f}^L to form a compressed residual vector $\boldsymbol{\delta}^r$. Denote as \mathbf{J}_{f^L} , \mathbf{J}_{g^L} and \mathbf{J}_h the Jacobian matrix of \mathbf{f}^L , \mathbf{g}^L and \mathbf{h} , respectively. Then $\boldsymbol{\delta}$ and $\boldsymbol{\delta}^r$, and their corresponding Jacobian matrices \mathbf{J} and \mathbf{J}^r have the form:

$$\boldsymbol{\delta} = \begin{bmatrix} \mathbf{h} \\ \mathbf{f}^L \end{bmatrix}, \quad \boldsymbol{\delta}^r = \begin{bmatrix} \mathbf{h} \\ \mathbf{g}^L \end{bmatrix}, \quad \mathbf{J} = \begin{bmatrix} \mathbf{J}_h \\ \mathbf{J}_{f^L} \end{bmatrix}, \quad \mathbf{J}^r = \begin{bmatrix} \mathbf{J}_h \\ \mathbf{J}_{g^L} \end{bmatrix}. \quad (26)$$

The key step of the LM algorithm to minimize the least-squares problems $\|\boldsymbol{\delta}\|_2^2$ and $\|\boldsymbol{\delta}^r\|_2^2$ is to calculate the following linear equations

$$\mathbf{J}^T \mathbf{J} + \lambda \mathbf{I} = \mathbf{J}^T \boldsymbol{\delta} \quad \text{and} \quad \mathbf{J}^{rT} \mathbf{J}^r + \lambda \mathbf{I} = \mathbf{J}^{rT} \boldsymbol{\delta}^r. \quad (27)$$

For the Gauss-Newton algorithm, we set $\lambda = 0$. According to (26), $\mathbf{J}^T \mathbf{J}$ and $\mathbf{J}^{rT} \mathbf{J}^r$ can be written as

$$\mathbf{J}^T \mathbf{J} = \mathbf{J}_h^T \mathbf{J}_h + \mathbf{J}_{f^L}^T \mathbf{J}_{f^L}, \quad \mathbf{J}^{rT} \mathbf{J}^r = \mathbf{J}_h^T \mathbf{J}_h + \mathbf{J}_{g^L}^T \mathbf{J}_{g^L}. \quad (28)$$

Based on Lemma 1, the Jacobian matrices of \mathbf{f}^L and \mathbf{g}^L in the residual vectors $\boldsymbol{\delta}$ and $\boldsymbol{\delta}^r$ respectively have the form:

$$\mathbf{J}_{f^L} = \begin{bmatrix} \mathbf{0} & \dots & \mathbf{A}^L \mathbf{B}^L & \dots & \mathbf{0} \end{bmatrix}, \quad \mathbf{J}_{g^L} = \begin{bmatrix} \mathbf{0} & \dots & \mathbf{C}^L \mathbf{B}^L & \dots & \mathbf{0} \end{bmatrix}. \quad (29)$$

Let us first focus on $\mathbf{J}_{f^L}^T \mathbf{J}_{f^L}$. Using (29), we have

$$\mathbf{J}_{f^L}^T \mathbf{J}_{f^L} = \begin{bmatrix} \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & (\mathbf{B}^L)^T (\mathbf{A}^L)^T \mathbf{A}^L \mathbf{B}^L & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \mathbf{0} & \dots & \mathbf{0} \end{bmatrix}. \quad (30)$$

$\mathbf{J}_{g^L}^T \mathbf{J}_{g^L}$ has the same form as $\mathbf{J}_{f^L}^T \mathbf{J}_{f^L}$. It only has one non-zero block $(\mathbf{B}^L)^T (\mathbf{C}^L)^T \mathbf{C}^L \mathbf{B}^L$. As $(\mathbf{C}^L)^T \mathbf{C}^L = (\mathbf{A}^L)^T \mathbf{A}^L$, we have $(\mathbf{J}_f^L)^T \mathbf{J}_f^L = (\mathbf{J}_g^L)^T \mathbf{J}_g^L$. So we have

$$\mathbf{J}^T \mathbf{J} = \mathbf{J}^{rT} \mathbf{J}^r. \quad (31)$$

We next consider $\mathbf{J}^T \boldsymbol{\delta}$ and $\mathbf{J}^{rT} \boldsymbol{\delta}^r$. Substituting the definitions in (26) into $\mathbf{J}^T \boldsymbol{\delta}$ and $\mathbf{J}^{rT} \boldsymbol{\delta}^r$, we have

$$\mathbf{J}^T \boldsymbol{\delta} = \mathbf{J}_h^T \mathbf{h} + \mathbf{J}_{f^L}^T \mathbf{f}^L, \quad \mathbf{J}^{rT} \boldsymbol{\delta}^r = \mathbf{J}_h^T \mathbf{h} + \mathbf{J}_{g^L}^T \mathbf{g}^L. \quad (32)$$

Substituting (29) and (24) into $\mathbf{J}_{f^L}^T \mathbf{f}^L$, we get

$$\mathbf{J}_{f^L}^T \mathbf{f}^L = \begin{bmatrix} \mathbf{0} & \dots & (\mathbf{B}^L)^T (\mathbf{A}^L)^T \mathbf{A}^L \mathbf{L} & \dots & \mathbf{0} \end{bmatrix}. \quad (33)$$

Using (29) and the definition of \mathbf{g}^L in (19), we know that $\mathbf{J}_{g^L}^T \mathbf{g}^L$ has the same form as $\mathbf{J}_{f^L}^T \mathbf{f}^L$. Specifically, there only exists one non-zero block $(\mathbf{B}^L)^T (\mathbf{C}^L)^T \mathbf{C}^L \mathbf{L}$. As $(\mathbf{C}^L)^T \mathbf{C}^L = (\mathbf{A}^L)^T \mathbf{A}^L$, we have $\mathbf{J}_{f^L}^T \mathbf{f}^L = \mathbf{J}_{g^L}^T \mathbf{g}^L$. Using (32), we get

$$\mathbf{J}^T \boldsymbol{\delta} = \mathbf{J}^{rT} \boldsymbol{\delta}^r. \quad (34)$$

Thus, using (27), (31) and (34), we know that \mathbf{g}^L can replace \mathbf{f}^L in the Gauss-Newton or LM algorithm. \square

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.
- [2] G. Klein and D. Murray, "Parallel tracking and mapping for small AR workspaces," in *2007 6th IEEE and ACM international symposium on mixed and augmented reality*. IEEE, 2007, pp. 225–234.
- [3] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "ORB-SLAM: a versatile and accurate monocular slam system," *IEEE transactions on robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [4] J. Engel, T. Schöps, and D. Cremers, "LSD-SLAM: Large-scale monocular SLAM," in *European conference on computer vision*. Springer, 2014, pp. 834–849.
- [5] A. Pumarola, A. Vakhitov, A. Agudo, A. Sanfeliu, and F. Moreno-Noguer, "PL-SLAM: Real-time monocular visual SLAM with points and lines," in *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2017, pp. 4503–4508.
- [6] X. Zuo, X. Xie, Y. Liu, and G. Huang, "Robust Visual SLAM with Point and Line Features," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 1775–1782.
- [7] F. Zheng, G. Tsai, Z. Zhang, S. Liu, C.-C. Chu, and H. Hu, "Trifovio: Robust and efficient stereo visual inertial odometry using points and lines," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 3686–3693.
- [8] S. J. Lee and S. S. Hwang, "Elaborate monocular point and line slam with robust initialization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1121–1129.
- [9] R. Gomez-Ojeda, F.-A. Moreno, D. Zuñiga-Noël, D. Scaramuzza, and J. Gonzalez-Jimenez, "PL-SLAM: A stereo SLAM system through the combination of points and line segments," *IEEE Transactions on Robotics*, vol. 35, no. 3, pp. 734–746, 2019.
- [10] Y. Yang, P. Geneva, K. Eickenhoff, and G. Huang, "Visual-Inertial Odometry with Point and Line Features," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 2447–2454.
- [11] S.-S. Huang, Z.-Y. Ma, T.-J. Mu, H. Fu, and S.-M. Hu, "Lidar-monocular visual odometry using point and line features," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 1091–1097.
- [12] B. Fang and Z. Zhan, "A visual slam method based on point-line fusion in weak-matching scene," *International Journal of Advanced Robotic Systems*, vol. 17, no. 2, p. 1729881420904193, 2020.
- [13] X. Li, Y. Li, E. P. Örnek, J. Lin, and F. Tombari, "Co-Planar Parametrization for Stereo-SLAM and Visual-Inertial Odometry," *IEEE Robotics and Automation Letters*, 2020.
- [14] X. Li, Y. He, J. Lin, and X. Liu, "Leveraging planar regularities for point line visual-inertial odometry," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 1775–1782.
- [15] R. G. Von Gioi, J. Jakubowicz, J.-M. Morel, and G. Randall, "LSD: A fast line segment detector with a false detection control," *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 4, pp. 722–732, 2008.
- [16] K. Zhao, Q. Han, C. Zhang, J. Xu, and M. Cheng, "Deep hough transform for semantic line detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 01, pp. 1–1, may 5555.
- [17] S. Yang and S. Scherer, "Direct Monocular Odometry using Points and Lines," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3871–3877.
- [18] S.-J. Li, B. Ren, Y. Liu, M.-M. Cheng, D. Frost, and V. A. Prisacariu, "Direct Line Guidance Odometry," in *2018 IEEE international conference on Robotics and automation (ICRA)*. IEEE, 2018, pp. 1–7.
- [19] X. Gao, R. Wang, N. Demmel, and D. Cremers, "LDSO: Direct Sparse Odometry with Loop Closure," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 2198–2204.
- [20] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM3: An Accurate Open-Source Library for Visual, Visual-Inertial, and Multimap SLAM," *IEEE Transactions on Robotics*, pp. 1–17, 2021.
- [21] R. Gomez-Ojeda, J. Briales, and J. Gonzalez-Jimenez, "PL-SVO: Semidirect monocular visual odometry by combining points and line segments," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 4211–4216.
- [22] C. Forster, Z. Zhang, M. Gassner, M. Werlberger, and D. Scaramuzza, "SVO: Semidirect visual odometry for monocular and multicamera systems," *IEEE Transactions on Robotics*, vol. 33, no. 2, pp. 249–265, 2016.
- [23] H. Zhou, D. Zou, L. Pei, R. Ying, P. Liu, and W. Yu, "Structslam: Visual slam with building structure lines," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 4, pp. 1364–1375, 2015.
- [24] D. Zou, Y. Wu, L. Pei, H. Ling, and W. Yu, "StructVIO: visual-inertial odometry with structural regularity of man-made environments," *IEEE Transactions on Robotics*, vol. 35, no. 4, pp. 999–1013, 2019.
- [25] Y. H. Lee, C. Nam, K. Y. Lee, Y. S. Li, S. Y. Yeon, and N. L. Doh, "VPass: Algorithmic compass using vanishing points in indoor environments," in *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2009, pp. 936–941.
- [26] G. Zhang, D. H. Kang, and I. H. Suh, "Loop closure through vanishing points in a line-based monocular SLAM," in *2012 IEEE International Conference on Robotics and Automation*. IEEE, 2012, pp. 4565–4570.
- [27] F. Camposeco and M. Pollefeys, "Using vanishing points to improve visual-inertial odometry," in *2015 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2015, pp. 5219–5225.
- [28] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [29] J. Sola, T. Vidal-Calleja, J. Civera, and J. M. M. Montiel, "Impact of landmark parametrization on monocular ekf-slam with points and lines," *International journal of computer vision*, vol. 97, no. 3, pp. 339–368, 2012.
- [30] J. J. Moré, "The levenberg-marquardt algorithm: implementation and theory," in *Numerical analysis*. Springer, 1978, pp. 105–116.
- [31] L. Zhang and R. Koch, "Structure and Motion from Line Correspondences: Representation, Projection, Initialization and Sparse Bundle Adjustment," *Journal of Visual Communication and Image Representation*, vol. 25, no. 5, pp. 904–915, 2014.
- [32] A. Bartoli and P. Sturm, "Structure-from-motion using lines: Representation, triangulation, and bundle adjustment," *Computer vision and image understanding*, vol. 100, no. 3, pp. 416–441, 2005.
- [33] D. G. Kottas and S. I. Roumeliotis, "Efficient and consistent vision-aided inertial navigation using line observations," in *2013 IEEE International Conference on Robotics and Automation*. IEEE, 2013, pp. 1540–1547.
- [34] Y. Yang and G. Huang, "Aided inertial navigation: Unified feature representations and observability analysis," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 3528–3534.
- [35] A. Bartoli and P. Sturm, "The 3d line motion matrix and alignment of line reconstructions," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, vol. 1. IEEE, 2001, pp. I–I.
- [36] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, "Bundle adjustment—a modern synthesis," in *International workshop on vision algorithms*. Springer, 1999, pp. 298–372.
- [37] L. Zhang and R. Koch, "An efficient and robust line segment matching approach based on LBD descriptor and pairwise geometric consistency," *Journal of Visual Communication and Image Representation*, vol. 24, no. 7, pp. 794–805, 2013.
- [38] J. Engel, V. Usenko, and D. Cremers, "A photometrically calibrated benchmark for monocular visual odometry," in *arXiv:1607.02555*, July 2016.
- [39] A. Handa, T. Whelan, J. McDonald, and A. J. Davison, "A benchmark for rgb-d visual odometry, 3d reconstruction and slam," in *2014 IEEE international conference on Robotics and automation (ICRA)*. IEEE, 2014, pp. 1524–1531.
- [40] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, and R. Siegwart, "The euroc micro aerial vehicle datasets," *The International Journal of Robotics Research*, vol. 35, no. 10, pp. 1157–1163, 2016.
- [41] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of RGB-D SLAM systems," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 573–580.
- [42] D. Schubert, N. Demmel, V. Usenko, J. Stueckler, and D. Cremers, "Direct sparse odometry with rolling shutter," in *European Conference on Computer Vision (ECCV)*, September 2018.