

Incorporating Learnt Local and Global Embeddings into Monocular Visual SLAM

Huaiyang Huang · Haoyang Ye · Yuxiang Sun · Lujia Wang · Ming Liu*

the date of receipt and acceptance should be inserted later

Abstract Traditional approaches for Visual Simultaneous Localization and Mapping (VSLAM) rely on low-level vision information for state estimation, such as handcrafted local features or the image gradient. While significant progress has been made through this track, under more challenging configuration for monocular VSLAM, e.g., varying illumination, the performance of state-of-the-art systems generally degrades. As a consequence, robustness and accuracy for monocular VSLAM are still widely concerned. This paper presents a monocular VSLAM system that fully exploits learnt features for better state estimation. The proposed system leverages both learnt local features and global embeddings at different modules of the system: direct camera pose

estimation, inter-frame feature association, and loop closure detection. With a probabilistic explanation of keypoint prediction, we formulate the camera pose tracking in a direct manner and parameterize local features with uncertainty taken into account. To alleviate the quantization effect, we adapt the mapping module to generate 3D landmarks better to guarantee the system's robustness. Detecting temporal loop closure via deep global embeddings further improves the robustness and accuracy of the proposed system. The proposed system is extensively evaluated on public datasets (Tsukuba, EuRoC, and KITTI), and compared against the state-of-the-art methods. The competitive performance of camera pose estimation confirms the effectiveness of our method.

Keywords Visual simultaneous localization and mapping (SLAM), visual-based navigation, mapping

1 INTRODUCTION

Visual Simultaneous Localization and Mapping (VSLAM) is a fundamental building block for robotic state estimation [1], providing critical information for high-level applications. With the recent progress in this field, VSLAM has been well developed for various real-world applications, e.g., visual-based navigation for autonomous vehicles, ranging from path planning to decision making. Moreover, VSLAM can model the environment and build a corresponding map apart from ego-motion estimation. This capability can further benefit the robots in scene understanding and object recognition.

Despite the impressive progress under moderate scenarios, it is noticeable that diverse challenges [2], for instance, adverse illumination conditions [3], could cause failure to most of the current VLSAM systems. Therefore, the robustness and accuracy under challenging situations remain a problem demanding further investigation. Emerging advances in

This work was supported by the National Natural Science Foundation of China, under grant No. U1713211, Collaborative Research Fund by Research Grants Council Hong Kong, under Project No. C4063-18G, and HKUST-SJTU Joint Research Collaboration Fund, under project SJTU20EG03, awarded to Prof. Ming Liu. Ming Liu is the corresponding author.

· Huaiyang Huang
E-mail: hhuangat@connect.ust.hk

Haoyang Ye
E-mail: hy.ye@connect.ust.hk

Lujia Wang
E-mail: eewanglj@ust.hk

Ming Liu (*corresponding author*)
E-mail: eelium@ust.hk

Department of Electronic and Computer Engineering,
The Hong Kong University of Science and Technology,
Hong Kong, SAR, China

· Yuxiang Sun
E-mail: yx.sun@polyu.edu.hk, sun.yuxiang@outlook.com

Department of Mechanical Engineering,
The Hong Kong Polytechnic University, Hung Hom
Hong Kong, SAR, China

deep learning (DL) provide an alternative perspective in resolving the aforementioned challenges. For examples, leveraging learning-based depth prediction [4], object recognition [5], or high-level semantics [6] benefits the visual state estimation in several aspects, varying from local reconstruction [7], scale recovery [8, 9] to dynamic environment handling [10]. Especially in recent years, the exploration of learning local feature extraction and description is prevailing [11]. Following a *detect-and-describe* scheme to predict the repeatability and description within one forward pass, several works [12–14] demonstrate a superior efficiency against the traditional *detect-then-describe* pipelines, e.g., SIFT [15]. The larger receptive fields from convolutional neural networks (CNNs) and metric learning techniques further contribute to the competitive performance of learning-based methods in both detection and description. These evolutions pave the way for the introduction of deep features into VSLAM systems. While some existing works generally recover pose in an end-to-end manner [16, 17], limiting the generalization ability for different scenarios, we consider a more practical approach to utilize the deep feature.

Built on top of ORB-SLAM2 [18], we propose a monocular SLAM system powered by learnt local and global embeddings. Leveraging the repeatability probability predicted from a CNN, the initial camera pose can be tracked robustly in a direct manner. Pre-triangulated landmarks can then be associated with local observations efficiently, which establishes correspondences for the pose refinement. We discuss the parameterization of local features to facilitate different modules from pose estimation to mapping. To achieve global consistency, we further leverage deep global embeddings to detect temporal loop closure.

Experimental results on public datasets, especially accurate pose estimations on challenging scenarios where state-of-the-art methods fail, demonstrate the effectiveness of the proposed system. Fig. 1 provides an overview of our system. Our contributions are summarized as follows:

- A monocular visual SLAM system leveraging learned local and global embeddings.
- A hybrid tracking scheme along with uncertainty modeling based on the network predictions.
- A mapping module adapted from the traditional pipeline to fit the nature of the frontend.
- Extensive evaluations on public datasets to demonstrate the accuracy and robustness of the proposed method.

Some parts of this work has been presented in [19]. In this paper, we extend the conference paper [19] with the following contributions:

- We improve the camera tracking scheme and provide a probabilistic explanation.

- We extend the previous visual odometry system into a complete SLAM system with global features predicted by a neural network.
- We perform more experiments to verify the effectiveness of different modules and the overall system.

The rest of the article is organized as follows: In Sec. 2, we review recent progress and challenges in VSLAM along with methods integrating deep learning. In Sec. 3, we demonstrate commonly used notations for simplicity. Then, Sec. 4, Sec. 5 and Sec. 6 describe the three main modules in our system, namely feature extraction, camera tracking and association and mapping and loop closing, respectively. In Sec. 7, we report the experimental results for the ablation and comparative studies. Finally, Sec. 8 concludes this paper.

2 Related Works

2.1 Monocular Visual SLAM and its Challenges

The state-of-the-art approaches for monocular VO/VSLAM can be categorized into two dominant classes, namely *indirect* and *direct* methods. Indirect, or feature-based methods [20–22], generally extract points of interest along with their descriptors as a sparse representation of the input image. With multiview correspondence association, camera motion and sparse structure are recovered via minimizing the reprojection error. Among these works, ORB-SLAM [22] exploits the covisibility relationships to strengthen the map reuse and frame management, balancing the accuracy and the computational demand.

Direct methods [23–25], on the contrary, directly minimize the photometric error on input images to track camera poses. In particular, DSO [25] optimizes camera poses, sparse scene structure, and camera model parameters jointly. To combine the advantages of both methods, Froster *et al.* propose SVO [26], which tracks the camera poses via sparse image alignment and utilizes hierarchical bundle adjustment (BA) as the backend for optimizing the structure and camera motion. Inspired by these works, we propose a hybrid scheme for the camera motion estimation, which tracks the pose initially on the predicted repeatability map and then refines it in an indirect manner. In the context of challenging conditions, robust VO/VSLAM remains unsolved [1, 3, 27]. Potential solutions to these issues can be found in [28–30]. These works typically specialize in a particular problem (e.g., handling high dynamic range (HDR) environment), which would introduce certain overhead under common scenarios [28]. In contrast, our system is developed for a general VSLAM system, which does not require any preprocessing on the camera input for the image quality enhancement. At the same time, in the experiments, it works well under most conditions.

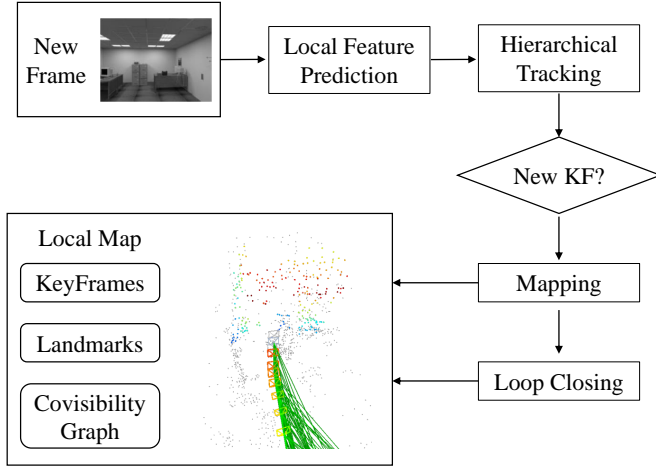


Fig. 1: System overview. Our system leverages deep networks to predict both local and global features. We follow the systematic design of ORB-SLAM2 [18], which can be divided into three modules: tracking, mapping and loop-closure. In the frontend, we propose a hierarchical scheme to recover the camera poses with network predictions. The mapping module generates new landmarks and maintains a covisibility graph. We further leverage learnt global embeddings to predict loop-closure and correct the global inconsistency.

2.2 Deep Learning in Visual SLAM

With the emergency of DL, recent researches demonstrate significant benefits by introducing relevant techniques into VSLAM systems. Compared with traditional methods in local feature detection [31–33] and description [15, 34, 35], learning-based approaches [12–14, 36, 37] exhibit a competitive performance in both matching accuracy and efficiency. Several works propose to use deep features in the camera pose estimation. Tang *et al.* [38] proposed BA-Net, where a feature-metric objective function is used for the direct image alignment, and the camera poses and predicted depth maps can be optimized in an end-to-end manner. In [39], Stumberg *et al.* proposed GN-Net, which introduced a Gauss-Newton loss for training the feature maps to be more invariant in the direct camera tracking. Works attempting to combine learnt local features with VO/VSLAM have been brought to our view [14, 40]. In [40], with labels from a VO backend, the original SuperPoint [12] is extended to predict the stability of local keypoints. In [14], to enhance the frontend efficiency for onboard RGB-D VO, Tang *et al.* proposed GCNv2. They supervise the detector with ground-truth keypoint locations labeled by Shi-Tomasi score. To further train the descriptor with triplet loss, positive and negative matches are retrieved via the projective geometry. The above works generally focus on leveraging traditional methods (e.g., multiview geometry) to train a more practical and

efficient frontend. On the contrary, here we consider bridging the gap of exploiting learned features to alleviate the challenges in monocular VSLAM.

3 Notations

Throughout the paper, we denote the image collected at the k -th time as I_k and the corresponding frame as \mathcal{F}_k . The world frame \mathcal{F}_w is set to be identical to the first camera frame \mathcal{F}_0 .

For I_k , the rigid transform $\mathbf{T}_k \in \mathbf{SE}(3)$ maps a 3D landmark $\mathbf{x}_i \in \mathbb{R}^3$ in \mathcal{F}_w to \mathcal{F}_k using:

$${}^{c_k}\mathbf{x}_i = \mathbf{R}_k \mathbf{x}_i + \mathbf{t}_k, \quad (1)$$

where $\mathbf{T}_k = [\mathbf{R}_k | \mathbf{t}_k]$. \mathbf{R}_k and \mathbf{t}_k are the rotational and translational components of \mathbf{T}_k , respectively. Accordingly, ${}^{c_k}\mathbf{x}_i$ denotes a 3D point in \mathcal{F}_k .

We denote a 2D pixel projected from a 3D landmark as $\tilde{\mathbf{u}}_{i,k}$. We use $\pi : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ to denote the projection function: $\tilde{\mathbf{u}}_{i,k} = \pi({}^{c_k}\mathbf{x}_i)$, where $\tilde{\mathbf{u}}_{i,k}$ is the coordinate in the pixel coordinate. π is defined as $\pi({}^{c_k}\mathbf{x}_i) = \mathbf{K} {}^{c_k}\mathbf{x}_i$, where \mathbf{K} is the intrinsic matrix for pinhole model.

The update of a camera pose is parameterized as an incremental twist $\boldsymbol{\xi} \in \mathfrak{se}(3)$. We use a left-multiplicative formulation $\oplus : \mathfrak{se}(3) \times \mathbf{SE}(3) \rightarrow \mathbf{SE}(3)$ for the update of \mathbf{T}_k , which is denoted as:

$$\boldsymbol{\xi} \oplus \mathbf{T}_k := \exp(\boldsymbol{\xi}^\wedge) \cdot \mathbf{T}_k, \quad (2)$$

where the mapping $\mathbf{SE}(3) \rightarrow \mathbb{R}^6$ is defined by $\mathbf{T} = \exp(\boldsymbol{\xi}^\wedge)$, where $(\cdot)^\wedge$ is the wedge operator.

4 Learning-based Feature Extraction

Here we adopt SuperPoint [12] as the feature extraction frontend. Recall the pipeline of SuperPoint, it first encodes the input image $I \in \mathbb{R}^{H \times W}$ with a single, shared encoder. Then two different heads decode a repeatability volume $\mathcal{H}' \in \mathbb{R}^{H_c \times W_c \times (C^2+1)}$ and a dense description $\mathcal{D}' \in \mathbb{R}^{H_c \times W_c \times 256}$, respectively. C is the size of the grid cell, which is 8 in our system, and $H_c = H/C$, $W_c = W/C$. \mathcal{H}' is then normalized by channel-wise softmax:

$$\mathcal{H}(h_c, w_c, y) = \frac{\exp(\mathcal{H}'(h_c, w_c, y))}{\sum_{k=1}^{65} \exp(\mathcal{H}'(h_c, w_c, k))}. \quad (3)$$

In this architecture, the keypoint detection problem is then formulated as a classification problem. Specifically, for $\mathcal{H}(h_c, w_c, :)$ $\in \mathbb{R}^{C^2+1}$, the first C^2 channels represent the probability of a pixel in a $C \times C$ grid to be a keypoint, while the last channel stands for the probability that no interest point lies in the current patch. For the details we refer the readers to [12].

We define the repeatability map, in patch-wise (\mathcal{R}_d) and in pixel-wise (\mathcal{R}) as:

$$\mathcal{R}_d = \mathcal{H}(\cdot, \cdot, C^2 + 1), \quad \mathcal{R} = -\log(s(\mathcal{H}(\cdot, \cdot, -1))). \quad (4)$$

where ‘.’ denotes the omitted channels, while ‘:’ denotes “from the first channel to the last channel”. $s()$ denotes the function that shuffles a feature vector into a 2D patch. In other words, it flattens the 3D volume to a 2D map. Here we further define the patch location where the pixel \mathbf{u} belongs to as \mathbf{u}_p . Then a probabilistic explanation is given by:

$$p(\mathbf{u}|I, \mathbf{w}) = \mathcal{R}(\mathbf{u}), \quad (5)$$

$$p(\mathbf{u}_p|I, \mathbf{w}) = 1 - \mathcal{R}_d(\mathbf{u}), \quad (6)$$

where $\mathcal{R}_d \in \mathbb{R}^{H_c \times W_c}$, the last channel of \mathcal{H} , stands for nonexistence of interest point in current patch. Accordingly, we reinterpret \mathcal{R}_d as the patch-wise repeatability prediction, which is used for the initial camera tracking, described in Sec. 5.1. $\mathcal{R} \in \mathbb{R}^{H \times W}$ is the repeatability map with full-resolution and $s: \mathbb{R}^{H_c \times W_c \times C^2} \rightarrow \mathbb{R}^{H \times W}$ maps the volume to a 2D heatmap. In the above formulation, we negate the repeatability response, so that a pixel \mathbf{u} (or patch \mathbf{u}_p) is more leaning to be a keypoint (or to contain a keypoint) if it has a smaller response $\mathcal{R}(\mathbf{u})$ (or $\mathcal{R}_d(\mathbf{u}_p)$).

In a non-maximum suppression (NMS) scheme, the locations of 2D features along with the final 2d grid $\mathcal{O} \in \mathbb{R}^{H_c \times W_c}$ are extracted from \mathcal{R} . A single cell in \mathcal{O} stores the index of the interest point in the current patch, zero if no salient point exists. For a local feature \mathbf{u}_i , the corresponding descriptor \mathbf{d}_i is sampled from $\mathcal{D} \in \mathbb{R}^{H \times W}$, which is interpolated from \mathcal{D}' . In addition, for the keyframes selected in our system, we predict the global embedding vector with NetVLAD [41], which is used to measure the image similarity across different frames for detecting the loop closure.

5 Camera Tracking and Association

5.1 Camera Pose Tracking

Assuming an initial sparse structure built from the backend, the frontend recovers the camera pose and associates 3D landmarks with 2D observations, which is illustrated in Fig. 2.

Initially, we track the current camera pose in a direct manner. The direct method of the camera pose estimation advances in their correspondence-free formulation, where there is no need for heuristic design of matching strategies. To take both the advantage of direct method and the network predictions, we propose to hierarchically track the camera pose on the repeatability map. Our intuition behind is that

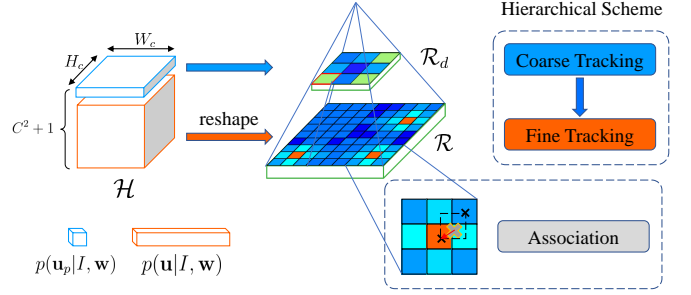


Fig. 2: Tracking scheme of the proposed system. We reinterpret the local feature detection as a two-level feature probability heatmap. Our method then tracks the repeatability map in a coarse-to-fine manner. With the tracking result, we associate the features with landmarks locally.

3D landmarks triangulated from 2D observations are consistently repeatable in different views. In other words, the reprojection location of a 3D landmark in $\mathcal{R}_{(\cdot)}$ should be the local extremum.

For the simplicity, here we denote the set of visible landmarks as $\mathcal{X} = \{\mathbf{x}_0, \mathbf{x}_1 \dots \mathbf{x}_n\}$. The current camera pose is denoted as \mathbf{T}_k and is parameterized as ξ_k . Given a landmark position \mathbf{x}_i and the current camera motion parameter ξ_k , the likelihood of local observation can be derived as:

$$p(\tilde{\mathbf{u}}_{ik}|\mathbf{w}, \xi_k, \mathbf{x}_i) = \mathcal{R}(\tilde{\mathbf{u}}_{ik}) = \mathcal{R}(\pi(\exp(\xi_k^\wedge) \cdot \mathbf{x}_i)). \quad (7)$$

To estimate ξ_k , our method seeks to minimize the total likelihood, which is given by:

$$\begin{aligned} \hat{\xi}_k &= \arg \max_{\xi_k} \prod_{\mathbf{x}_i \in \mathcal{X}} p(\cdot) = \arg \min_{\xi_k} - \sum_{\mathbf{x}_i \in \mathcal{X}} \log(p(\cdot)) \\ &= \arg \min_{\xi_k} - \sum_{\mathbf{x}_i \in \mathcal{X}} \log(\mathcal{R}(\pi(\exp(\xi_k^\wedge) \cdot \mathbf{x}_i))), \end{aligned} \quad (8)$$

where we abbreviate $p(\tilde{\mathbf{u}}_{ik}|\mathbf{w}, \xi_k, \mathbf{x}_i)$ to $p(\cdot)$. With a robust kernel function denoted as $\|\cdot\|_\gamma$, the total energy function is given by:

$$E = \sum_{\mathbf{x}_i \in \mathcal{X}} \|\log(\mathcal{R}(\pi(\exp(\xi_k^\wedge) \cdot \mathbf{x}_i)))\|_\gamma. \quad (9)$$

In the implementation, we choose Huber norm as the loss function. Based on the above derivation, this problem is converted into a standard least-squares problem, which can be solved using Gauss-Newton method.

Direct method is considered very sensitive to the initial guess. A common strategy to alleviate this issue is to leverage a pyramidal implementation. In our case, generating pyramidal prediction results require multiple forward passes, prohibiting its real-time performance. However, we disassemble the prediction into two levels, which represents the repeatability from coarse to fine, as described in Sec. 4 and shown in the Fig. 2. Therefore, similar to the pyramidal implementation in the optical flow estimation or direct

Algorithm 1: Tracking and Feature Association

Data: $I_k, \tilde{\xi}_k, \mathcal{O}_k, \mathcal{R}, \mathcal{R}_d, \mathcal{M}$.

```

1  $\hat{\xi}_k = \text{applyConstantVelocity}(\tilde{\xi}_k)$ 
2  $\mathcal{X} = \text{collectVisibleLandmarks}(\mathcal{M})$ 
3 //  $\mathcal{M}$  denotes the local map.
4  $\hat{\xi}_k = \text{coarseTracking}(\tilde{\xi}_k, \mathcal{X})$ 
5 foreach  $\tilde{\mathbf{u}}_{i,k}$  do
6    $\mathcal{V}_{i,k} = \text{checkAdjacentPatches}(\tilde{\mathbf{u}}_{i,k}, \mathcal{O}_k)$ 
7   if  $|\mathcal{V}_{i,k}| \equiv 1$  then
8      $\text{associate}(\mathbf{x}_i, \mathcal{V}_{i,k}[0])$ 
9   else
10     $j = \text{findBestAssociation}(\mathbf{x}_i, \mathcal{V}_{i,k})$ 
11     $\text{associate}(\mathbf{x}_i, \mathcal{V}_{i,k}[j])$ 
12  end
13 end
14  $\hat{\xi}_k = \text{poseRefinement}(\hat{\xi}_k, \mathcal{X})$ 
15  $\text{outlierFiltering}(\hat{\xi}_k, \mathcal{X})$ 

```

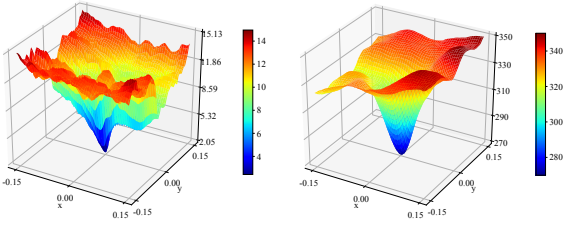


Fig. 3: Cost surface for camera pose tracking of *photometry*-based residual (left) and *repeatability*-based residual (right). For each plot, the x-y plane stands for different translational offsets to the ground-truth pose and the value on z-axis is the total cost with the corresponding transform.

image alignment, we hierarchically track the current camera in a coarse-to-fine manner. For estimating the camera pose in the coarsest level, oppositely we minimize the probability that a patch does not contain keypoint, causing a little subtlety in the formulation. This gives:

$$\begin{aligned} \hat{\xi}_k &= \arg \min_{\xi_k} p(\mathbf{u}_p | \xi_k, \mathbf{x}_i, \mathbf{w}), \\ &= \arg \min_{\xi_k} \mathcal{R}_d(\pi_d(\exp(\xi_k^\wedge) \cdot \mathbf{x}_i)). \end{aligned} \quad (10)$$

For each residual term, the jacobian can be derived as:

$$\begin{aligned} \mathbf{J}_i &= \mathbf{J}_{\text{repeat}} \cdot \mathbf{J}_{\text{proj}} \cdot \mathbf{J}_{\text{pose}} \\ &= \frac{\partial \mathcal{R}_{(\cdot)}(\tilde{\mathbf{u}}_{i,k})}{\partial \tilde{\mathbf{u}}_{i,k}} \cdot \frac{\partial \tilde{\mathbf{u}}_{i,k}}{\partial \mathbf{x}_i^c} \cdot \frac{\partial ((\xi_k \oplus \mathbf{T}_k) \mathbf{x}_i)}{\partial \xi_k}, \end{aligned} \quad (11)$$

where similar to direct tracking methods, $\mathbf{J}_{\text{repeat}}$ is the gradient of the corresponding repeatability map, which is evaluated at the projected pixel $\tilde{\mathbf{u}}_{i,k}$ and is calculated by bilinear interpolation at each iteration. \mathbf{J}_{proj} and \mathbf{J}_{pose} are the Jacobian of projection function w.r.t transformed point and the left-compositional derivative of the transformed point w.r.t the twist update ξ_k , respectively.

In practice, we find that tracking on the coarse level already produces a relatively good pose estimation result. When this estimation is fed to the fine level, our tracking module is able to fine-tune it for better accuracy. Analogous to the traditional optical flow pipeline, the pyramidal implementation is a significant design to alleviate the convergence basin issue and make the tracking module work better with predicted repeatability. Fig. 3 compares the cost surfaces of our method and the classical direct method. The cost surface of photometry-based tracking suffers from non-convexity responsible for the sensitivity to initial guess, the narrowness of convergence basin, and potential numerical issues in the optimization. In contrast, tracking on repeatability exhibits a smoother cost surface, which in practice yields better convergence property.

5.2 Feature Association

With the initial tracking, Here we propose a local feature association scheme, which minimizes the computational overhead by avoiding the descriptor matching as much as possible. Our observation is that, with the initial tracking, the reprojection location is considered to have an accuracy at sub-pixel level. That is to say, if the optimization converges, the reprojection location $\tilde{\mathbf{u}}_{i,k}$ of an inlier \mathbf{x}_i is supposed to have sub-patch accuracy. In addition, as we described before, we perform NMS or gridding to store one feature in a 8×8 patch. Under these observations, the 2D association of \mathbf{x}_i should belonging to the four adjacent patches or grid cells of $\tilde{\mathbf{u}}_{i,k}$. Therefore, instead of exhaustively matching descriptors or adopting some heuristics for tracking, we actively attempt to associate features in four adjacent patches. Based on the actual number of adjacent patches which contain a keypoint, we associate landmarks and local features as follows:

- If there is only one keypoint candidate, we directly associate it with the landmark.
- If there are multiple keypoint candidates, we then associate it according to descriptor distance, which is equivalent to finding the best local association.

It is noticeable that this association step is achieved with minimal computational overhead, which might be caused by matching high-dimensional descriptors.

5.3 Feature Parameterization and Pose Refinement

For the indirect methods, the *reprojection error* is generally used for pose estimation or Bundle Adjustment. A generic formulation is given by:

$$\mathbf{u}_{i,k} = \mathbf{h}_{i,k}(\mathbf{x}_i, \xi_k) + \mathbf{n}_{i,k}, \quad (12)$$

where $\mathbf{h}_{i,k}$ is the measurement function between the i -th landmark position \mathbf{x}_i and k -th keyframe pose ξ_k ; and $\mathbf{n}_{i,k}$ is zero-mean white Gaussian noise that perturbs the measurement, i.e., $\mathbf{n}_{i,k} \sim \mathcal{N}(\mathbf{0}, \Sigma_{i,k})$. In other words, this formulation models local observations as 2D Gaussian components, i.e., $\tilde{\mathbf{u}}_{i,k} \sim \mathcal{N}(\bar{\mathbf{u}}_{i,k}, \bar{\Sigma}_{i,k})$. As long as the majority of methods follow this generic formulation, the problem is how to define the parameters for each local feature, i.e., $\bar{\mathbf{u}}_{i,k}$ and $\bar{\Sigma}_{i,k}$. Previous works resolve this issue in two approaches. In the early stage, methods model local features as 2D Gaussian distributions with anisotropic covariance. Some interesting discussions can be found in [42, 43], where the authors generally conclude modelling such anisotropic covariance can improve the visual state estimation. On the contrary, modern SLAM systems [18, 22, 26] generally assign an isotropic covariance for each feature. This gives $\bar{\Sigma}_{i,k} = \sigma^2 \mathbf{I}_{2 \times 2}$, with σ proportional to the pyramidal level of the corresponding feature. However, as described above, generating a pyramidal prediction is prohibitive, making this parameterization method not practical in our implementation.

Following previous works, we seek to weigh the contribution of different local features to the overall estimation. Compared to traditional feature extraction techniques, CNNs have a larger receptive field and are capable of generating more consistent local features. Therefore, we model their 2D distribution from keypoint prediction heatmap with a probabilistic explanation. Considering the prediction as a mixture of multivariate Gaussian distributions, to approximate the single variate parameters, we have:

$$\bar{\mathbf{u}} = \mathbb{E}[\mathbf{u}] = \sum_{\mathbf{u} \in \mathcal{P}} \mathbf{u} \cdot p(\mathbf{u}|\mathbf{w}) \quad (13)$$

$$\begin{aligned} \bar{\Sigma} &= \mathbb{E}[(\mathbf{u} - \bar{\mathbf{u}})(\mathbf{u} - \bar{\mathbf{u}})^T] \\ &= \sum_{\mathbf{u} \in \mathcal{P}} (p_i^2 \Sigma_{\mathbf{u}} + p_i \mathbf{u} \mathbf{u}^T) - \bar{\mathbf{u}} \bar{\mathbf{u}}^T \end{aligned} \quad (14)$$

For the simplicity of implementation, here we represent \mathcal{P} as a 3×3 patch. This parameterization scheme is further verified by Sec. 7.2.

After associating landmarks with the local observations, the current camera pose \mathbf{T}_k is refined by minimizing the reprojection error. For \mathbf{x}_i , the error function is defined as:

$$\mathbf{e}_{i,k}^{\text{repro}} = \pi(\mathbf{R}_k \mathbf{x}_i + \mathbf{t}_k) - \tilde{\mathbf{u}}_{i,k}. \quad (15)$$

Similar to the initial tracking, the refined pose is solved by iterative least-squares. The overall energy function to minimize is defined as:

$$E^{\text{repro}} = \sum_{i \in \mathcal{P}_k} \left\| \left(\mathbf{e}_{i,k}^{\text{repro}} \right)^T \Sigma_{i,k}^{-1} \mathbf{e}_{i,k}^{\text{repro}} \right\|_{\gamma}, \quad (16)$$

where $\Sigma_{i,k}$ is the covariance of 2D feature location $\tilde{\mathbf{u}}_{i,k}$ for weighting different features' contribution to the optimization. The total pipeline of our tracking and association module is also demonstrated in algorithm 1.

6 Mapping and Loop Closing

6.1 Mapping

6.1.1 Landmark Creation

When a new keyframe is constructed, the mapping module first creates new landmarks with previous observations. For correspondence search, the top- n keyframes sharing the most covisibility score with the current keyframe are chosen to be the candidate frames. ORB-SLAM2 [18] exploits bag-of-words (BoW) [44] to accelerate feature matching between keyframes for the landmark generation. Inevitably, this strategy introduces quantization effect that fewer correspondences would be found compared to brute-force search. For original ORB-SLAM2, we observe that such effect is affordable as it extracts quantities of features (e.g., 1000 for a image with VGA resolution). However, in our case, two factors amplify this quantization effect: first, we only extract the most representative features, limiting the number of total local features; second, float descriptors lie in a much larger space than binary descriptors, making quantization effect more significant as more clusters are required for a representative quantization. This degrades both the robustness and the accuracy in practice. Therefore, two methods are adopted for feature association between keyframes. The first one is the approximate nearest neighbor (ANN) search. At the creation of each keyframe, a database of untracked local features is established. In the association step, we query each feature from the databases of candidate keyframes. Moreover, after the association step, the database is reindexed to guarantee the search efficiency with future keyframes. The problem of ANN-based association is that for the repetitive pattern (e.g. checkboard), ambiguity exists and increases the number of outliers in the triangulation step.

To resolve this problem, we further search the epipolar line of the 2d grid \mathcal{O}_j of the target keyframe \mathcal{F}_j to associate the features. For an untracked local feature $\mathbf{u}_{i,k}$ in the current keyframe \mathcal{F}_k , the epipolar line ${}^i\mathbf{l}_{k,j} = [l_0, l_1, l_2]^T$ on the image plane of \mathcal{F}_j is derived as:

$$({}^i\mathbf{l}_{k,j})^T = \mathbf{u}_{i,k}^T \mathbf{K}^{-T} \mathbf{t}_{k,j}^{\wedge} \mathbf{R}_{k,j} \mathbf{K}^{-1},$$

$$\text{with } \mathbf{R}_{k,j} = \mathbf{R}_k \mathbf{R}_j^{-1} \quad \mathbf{t}_{k,j} = -\mathbf{R}_k \mathbf{R}_j^T \mathbf{t}_j + \mathbf{t}_k,$$

where \mathbf{K} is the projection matrix. We search \mathcal{O}_j along the entire epipolar line ${}^i\mathbf{l}_{k,j}$. If current grid has an unmatched feature, we further check the epipolar distance $d_{i,k,l,j}$ as the geometry constraint:

$$d_{i,k,l,j} = \frac{1}{\det(\Sigma_{i,k})} \frac{\mathbf{u}_{l,j}^T \cdot {}^i\mathbf{l}_{k,j}}{\sqrt{l_0^2 + l_1^2}},$$

which is weighed by $1/\det(\Sigma_{i,k})$ for the consideration of uncertainty. The inlier with the best descriptor distance is

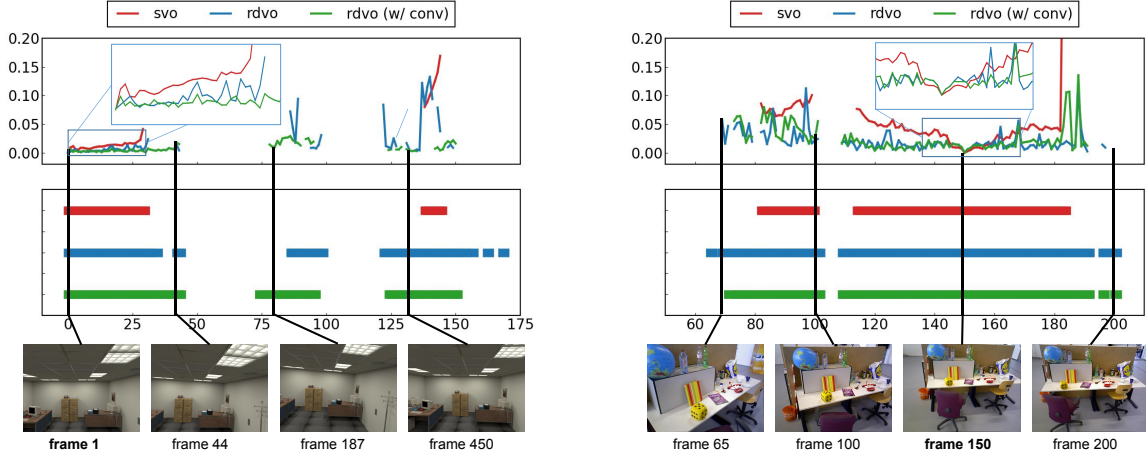


Fig. 4: Convergence radius and accuracy of different direct image alignment methods (indicated by color). All frames of the respective sequence are tracked on frame 1 (left) and frame 150 (right), using the identity as initialization. The bottom plots show for which frames tracking succeeds; the top plots show the final translational error. We observe that apart from the traditional photometry-based direct tracking (svo), the proposed repeatability-based tracking advances in larger convergence basin and better tracking precision.

considered as a successful match for the candidate feature. Finally, the 3D positions are recovered via mid-point triangulation. For rejecting outliers in the triangulation, we follow [22] to further verify the sign of depth and the reprojection error in both keyframes.

6.1.2 Backend Optimization

Similar to ORB-SLAM2, batched BA is utilized for backend optimization and map management. The variables in the local map for optimization consisting of updates to keyframe poses $\mathcal{T} = [\xi_{k_0}, \xi_{k_1}, \dots, \xi_{k_n}]$ and positions of landmarks $\mathcal{M} = [\mathbf{x}_{i_0}, \mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_m}]$. We denote the full state vector as $\mathcal{X} = [\mathcal{T}, \mathcal{M}]$, which is solved by:

$$\mathcal{X} = \arg \min_{\mathcal{X}} \sum_{\mathbf{T}_k} \sum_{\mathbf{x}_i} \text{obs}(i, k) \|(\mathbf{e}_{i,k}^{\text{repro}})^T \Sigma_{i,k}^{-1} \mathbf{e}_{i,k}^{\text{repro}}\|_{\gamma}, \quad (17)$$

where $\text{obs}(i, k) = 1$ if \mathbf{x}_i is observed by \mathcal{F}_k , $\text{obs}(i, k) = 0$ otherwise. Note that, as in ORB-SLAM2 [18], the poses of keyframes that share no observations of visible landmarks with the current keyframe are fixed to maintain a consistent scale estimation. The map management operations are inherited from [18]. Briefly, it culls redundant keyframes, removes outliers in the local BA and fuses the landmarks.

6.2 Loop Closing

We extend the built-in loop-closure module to work with learnt global embeddings. When a keyframe is created, we generate the global embedding from NetVLAD [41]. The advantage of using traditional BoW is that we can leverage

the inverted file formulation for efficient keyframe retrieval. Yet in our case, the network directly outputs a global embedding vector instead of preserving the visual words information. Therefore, in our implementation, we directly compare the l_2 -distance between the global embeddings between the current keyframe and the past keyframes that does not share any common observations with the current keyframe. Then loop-closure candidates are considered to have embeddings close to that of the current keyframe. With the loop closure candidates, we further adopt the built-in modules of ORB-SLAM2 [18] to perform candidate verification and loop closure correction. Briefly, if the loop closure is detected successfully, it first corrects the **Sim(3)** pose graph generated by the minimal spanning tree from the global covisibility graph. Then, global BA is performed with a similar pipeline described in Sec. 6.

7 Experimental Results

In this section, we first perform ablation studies on different modules in our system, including the proposed direct tracking frontend and the feature parameterization. Then, we report the general state estimation accuracy on several public datasets and compare our system against other popular methods. Finally, we provide several qualitative results along with runtime analysis to further demonstrate the effectiveness of the proposed method.

7.1 Evaluation on the Camera Tracking

In this section, we validate the proposed direct tracking scheme. For the baselines to compare with, we first select the front-end of SVO [26]. The tracking front-end of SVO detects FAST corners in a grid, and tracks recovered 3D landmarks in a sparse image registration manner, which shares a similar scheme to our system and is considered to be a good competitor. In addition, inspired by the recently proposed GN-Net in [39], we propose another baseline as competitor, which tracks the current frame with deep features as input in a direct manner. The objective function is given by:

$$E_{\text{feat}} = \sum \|\mathbf{F}_{k'}(\mathbf{u}_{i,k'}) - \mathbf{F}_k(\pi(\mathbf{R}_k \mathbf{x}_i + \mathbf{t}_k))\|, \quad (18)$$

where \mathbf{F}_k is the feature map of the k -th frame. This formulation is based on the insight that CNNs can produce invariant feature for the correspondence across different frames. Similarly, we add this *feature-metric constraint* into our repeatability formulation. In the implementation, we leverage the feature volume predicted by SuperPoint and reduce its dimension using PCA to achieve an applicable real-time implementation on CPU. We denote this baseline as **rdvo (w/ conv)** and our method as **rdvo**.

Here we mainly evaluate the convergence radius and pose estimation accuracy of the tracking module. For the evaluation, given a test sequence, we first generate a sparse point-cloud with a pair of RGB image and depth map. Then, we evaluate different methods by registering the sparse visual structure to each frame in the sequence. The experiments are conducted on the ICL-NUIM [45] and TUM Dataset [46]. The results are shown in Fig. 4. We observe that compared to SVO, the proposed method generally enlarges the convergence radius for direct tracking. In other words, starting from the same initialization point, more frames with larger translation and rotation variances are successfully registered by our method. This confirms the effectiveness of the proposed direct tracking scheme with predicted repeatability.

It is also interesting that, introducing high-level descriptions from the convolutional feature, **rdvo (w/ conv)** provides a close localization accuracy and robustness against **rdvo**. While both of them outperform the tracking scheme in SVO, **rdvo (w/ conv)** enlarges the convergence radius a little on the ICL-NUIM dataset, which indicates that using the *feature-metric* constraints can increase the robustness of direct tracking to some extent. Despite that the convergence radius is increased as expected, we notice that, **rdvo (w/ conv)** does not consistently outperform **rdvo** and **svo**. We think the reason could be that the image alignment generally requires a subpixel accuracy for an accurate alignment result. Unlike repeatability predictions that implicitly provide pixel-level keypoint location, convolutional feature slices of correspondences in different frames and have some subtleties with viewpoint variances. These subtleties yield

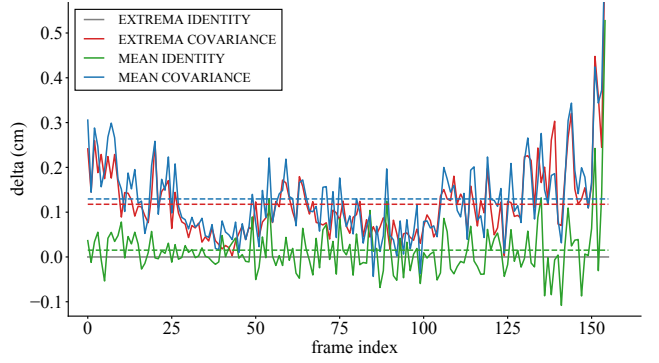


Fig. 5: Relative improvement on different methods compared to the baseline method. The dashed line represents mean error of the corresponding method.

Table 1: Translational RMSE (cm) and tracking success rate (%) on the New Tsukuba dataset. The tracking success is defined as (#tracked frame/#frame in the sequence).

Sequence Name	Ours	DSO	ORB-SLAM w/o loop
fluorescent	10.9 ± 5.7	59.0 ± 39.0	(18.2 ± 15.3)
	87.6 ± 0.3	98.2 ± 0.0	88.4 ± 11.6
daylight	9.4 ± 4.9	50.0 ± 29.7	14.7 ± 9.7
	96.2 ± 0.2	98.1 ± 0.0	82.4 ± 16.9
lamps	9.3 ± 4.6	×	×
	94.5 ± 7.1	44.5 ± 37.4	0.0 ± 0.0
flashlight	18.3 ± 10.3	×	×
	94.4 ± 0.1	58.4 ± 10.9	0.0 ± 0.0

some degradation to the registration results. Besides, although we interpolate between feature slices, the convolutional feature map is considered to provide patch-level information, which is downsampled 8 times in our case. As the residual dimension is very large (256 in our case), implementing it on CPU leads to significant computational overhead, especially for interpolating the feature map.

There exist some image alignment methods based on raw gradient input [47–50]. However, these methods differs from baselines in the experiments in that they densely select pixels with high gradient and generally transform the gradient map into another space, e.g., Distance Transform (DT), for smoother gradient information [48, 50]. Some also assume structural constraints between landmarks for better structural representation [24]. As a consequence, we only focus on *sparse* image alignment methods here.

7.2 Evaluation on the Local Feature Parameterization

In this section, we evaluate the effectiveness of different feature parameterization methods in pose estimation. As the local feature is parameterized as a gaussian distribution, with

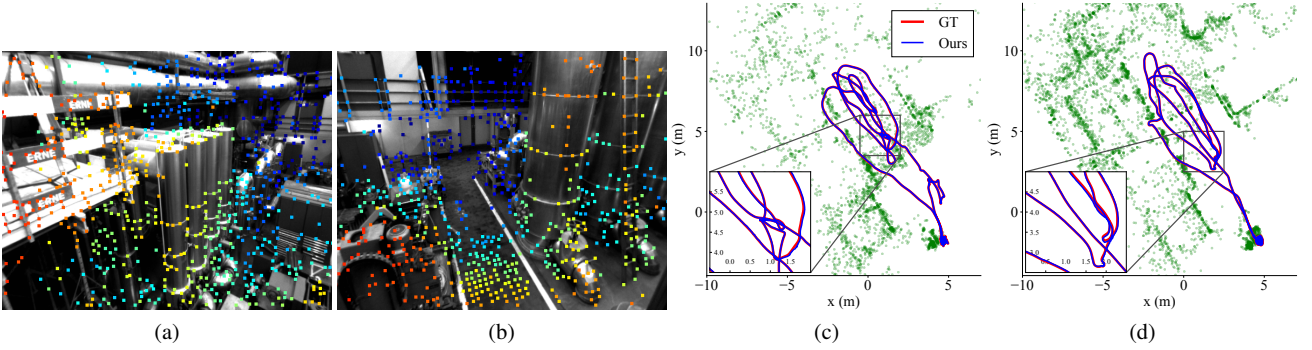


Fig. 6: (a) and (b): excerpts of sequence Machine Hall 01 and 02 the EuRoC dataset, keypoints are colored by the inverse depth. (c) and (d): trajectories of ours against the ground truth on Machine Hall 01 and 02.

Table 2: Translational RMSE (cm) on the EuRoC dataset. Left column: performances of pure visual odometry. Right column: performance of SLAM system, where results on sequences without explicit loop-closure are excluded.

	Sequence	Ours w/o loop	DSO	ORB-SLAM2 w/o loop	Ours	ORB-SLAM2
easy	Machine Hall 01	1.63 ± 0.86	5.73 ± 2.28	1.77 ± 0.91	-	-
	Machine Hall 02	1.46 ± 0.77	4.53 ± 2.71	1.72 ± 0.81	-	-
	Machine Hall 03	3.20 ± 1.53	20.89 ± 10.23	3.19 ± 1.52	-	-
	Vicon Room 1 01	3.34 ± 4.34	13.52 ± 8.72	3.28 ± 1.20	-	-
	Vicon Room 2 01	1.85 ± 0.83	4.70 ± 2.41	2.59 ± 1.34	1.60 ± 0.15	4.44 ± 1.31
	Vicon Room 2 02	4.73 ± 2.49	12.31 ± 6.37	(5.08 ± 2.95)	2.59 ± 0.37	1.69 ± 0.10
hard	Machine Hall 04	8.59 ± 3.37	20.03 ± 8.96	9.78 ± 4.05	-	-
	Machine Hall 05	4.34 ± 2.09	10.47 ± 3.78	(12.67 ± 5.86)	-	-
	Vicon Room 1 02	23.82 ± 9.38	32.83 ± 27.41	(14.06 ± 5.27)	1.45 ± 0.11	1.59 ± 0.35
	Vicon Room 1 03	(17.17 ± 10.71)	94.15 ± 39.78	×	4.44 ± 1.31	7.26 ± 5.78

different combination of mean and variance, we mainly compare four different parameterization methods, given by:

- EXTREMA IDENTITY (EI): $\mathbf{u} = \mathbf{u}_{\text{peak}}, \Sigma = \mathbf{I}_{2 \times 2}$
- EXTREMA COVARIANCE (EC): $\mathbf{u} = \mathbf{u}_{\text{peak}}, \Sigma = \bar{\Sigma}$
- MEAN IDENTITY (MI): $\mathbf{u} = \bar{\mathbf{u}}, \Sigma = \mathbf{I}_{2 \times 2}$
- MEAN COVARIANCE (MC): $\mathbf{u} = \bar{\mathbf{u}}, \Sigma = \bar{\Sigma}$

where the $\bar{\mathbf{u}}$ and $\bar{\Sigma}$ are described in Sec. 5.1.

For the dataset selection, we notice that the accuracy of SLAM system (or BA) is mainly bounded by the data association, and thus the local feature parameterization is not the most important factor. Therefore, using a real-world dataset can hardly compare the effectiveness of different parameterization methods. This is also because of the systematic error in the real-world dataset, for example, the uncertainty of calibration between cameras and sensors providing absolute pose measurements. As a consequence, here we perform the evaluations on the ICL-NUIM dataset.

To examine how different parameterization methods work in the SLAM system, for each frame in a sequence, we back-project the detected keypoints to the 3D space as landmarks.

Then these landmarks are associated with N (5 in the experiment) adjacent frames. Combining the poses along with the landmark positions, we jointly optimize these variables following a standard pipeline of BA. We maintain an optimization problem with identical optimizable parameters and constraints, and introduce different local feature parameterization methods to examine their effectiveness.

The results are represented in Fig. 5, from which we have the following conclusions: First, the results show that considering the prediction with an anisotropic covariance assigning to different keypoints could generally increase the bundle adjustment accuracy. For methods that do not consider the uncertainty (EI and MI), the covariance matrices are all set to be identity. This means that in the optimization, different residual factors contribute to the final results equally. On the contrary, for methods taking uncertainty into account (EC and MC), the anisotropic covariance serve as a reflection of keypoint prediction uncertainty. Through modelling the prediction uncertainty, different factors are weighed in the joint optimization according to their observation confidence. In other words, factors with higher uncertainty should

Table 3: Translational RMSE (m) on the KITTI Odometry Dataset.

Sequence	Environment	Ours w/o loop	DSO	ORB-SLAM2 w/o loop	Ours	ORB-SLAM2
00	Urban	96.23 \pm 44.76	118.58 \pm 0.51	73.39 \pm 32.66	5.31 \pm 0.57	6.11 \pm 0.32
01	Highway	\times	105.14 \pm 204.86	\times	-	-
02	Urban+Country	76.94 \pm 50.34	122.67 \pm 6.93	27.14 \pm 14.04	22.54 \pm 1.46	28.92 \pm 4.26
03	Country	1.34 \pm 0.29	2.08 \pm 0.08	0.96 \pm 0.53	-	-
04	Country	0.21 \pm 0.08	0.99 \pm 0.01	1.14 \pm 0.58	0.51 \pm 0.18	0.87 \pm 0.31
05	Urban	48.74 \pm 26.33	50.93 \pm 0.69	39.42 \pm 21.05	4.75 \pm 0.61	5.65 \pm 1.04
06	Urban	54.37 \pm 25.51	65.81 \pm 1.29	54.45 \pm 25.21	10.16 \pm 1.74	16.13 \pm 0.83
07	Urban	16.96 \pm 9.37	18.37 \pm 0.54	16.71 \pm 9.47	4.16 \pm 0.63	2.45 \pm 0.31
08	Urban+Country	64.82 \pm 39.95	119.12 \pm 12.74	51.55 \pm 31.78	-	-
09	Urban+Country	60.82 \pm 31.60	61.21 \pm 7.72	47.31 \pm 25.96	9.52 \pm 0.47	51.63 \pm 12.70
10	Urban+Country	8.27 \pm 0.86	34.84 \pm 38.41	9.56 \pm 1.02	-	-

have less contribution to the optimization framework. As a consequence, comparing EC and MC against EI and MI, the accuracy of BA is increased with modelling the uncertainty. Second, comparing methods that consider keypoint location as extrema (EI and EC) or weighted mean (MI and MC), we find that the pose accuracy is improved slightly. This indicates that most of the keypoint predictions are already of good quality, so that in a BA framework, considering the refinement of keypoint location with prediction uncertainty can only provide trivial improvement in the pose estimation. Accordingly, these experimental results guide us to finally choose the local feature parameterization in our system described in Sec. 5.1.

7.3 Evaluation on Trajectory Estimation

In this section, we evaluate our system on several public datasets, the New Tsukuba [51], the EuRoC Mav dataset [52] and KITTI Odometry dataset [53]. We compare our method to both the state-of-the-art indirect (ORB-SLAM2 [18]) and direct (DSO [25]) VO/VSLAM algorithms. GCN-SLAM [14], the recently proposed method using a learned binary descriptor, is expected to be one of the competitors. However, the monocular version adapted from the open-source implementation¹ fails to produce competitive results. Thus we exclude it for further evaluation. A major reason for the failure of GCN-SLAM is that it is designed for RGB-D inputs, while for the monocular VO, maintaining the scale consistency and estimating the depth of local feature raise more challenges.

All the experiments are done using the same desktop with i7-8700K and NVIDIA 1080Ti. For quantitative evaluation, translational RMSE of absolute trajectory error (ATE) [46] is used. We run different algorithms on each sequence 5 times and average the evaluation metrics. The tracking failure is either reported by the system itself or determined af-

terward if the error is larger than 1 meter (typically caused by scale inconsistency). If tracking failure exists for one or more runs, the corresponding result is shown in parentheses (\cdot), while failure for all runs is marked as \times . The results on both datasets are reported in Table. 1, Table. 2 and Table. 3, respectively.

7.3.1 The New Tsukuba Dataset

The New Tsukuba dataset provides synthetic images rendered by computer graphics techniques. It is challenging for monocular visual odometry as 1) the illumination condition of sequences *lamps* and *flashlight* is extreme, and 2) the camera rotation is relatively aggressive. As ORB-SLAM2 does not provide the setting for the Tsukuba dataset, we use the setting of another indoor dataset with the same image resolution and adjust the camera parameters accordingly. Note that the New Tsukuba Dataset does not contain any loopy motion, therefore all methods can be considered as pure visual odometry rather than full SLAM systems.

Table. 1 reports the translational RMSE and success rate of tracking for each sequence. As expected, challenging illumination conditions of the *lamps* and *flashlight* lead both ORB-SLAM and DSO to tracking failure. Besides, to accelerate feature association for camera pose tracking, ORB-SLAM searches correspondences in a local window with a motion prior, making it sensitive to aggressive rotation change. As a consequence, it occasionally fails even under the *fluorescent* with a moderate illumination configuration.

Besides the robustness and accuracy under these test scenarios, our system is capable of maintaining a relatively consistent trajectory estimation accuracy, regardless of the illumination variances. Especially on *lamps*, the RMSE is even slightly smaller than on *daylight*. Larger error under *flashlight* over other sequences indicates that learning-based descriptors share a similar characteristic with hand-crafted ones in degradation under photometric noise [25].

¹ https://github.com/jiexiong2016/GCNv2_SLAM

7.3.2 EuRoC Mav Dataset

The EuRoC dataset contains several sequences that are challenging for monocular ego-motion estimation. Accordingly, for the comparison, we divide the dataset into *moderate* (easy) and *challenging* (hard) sequences. For the sequences MH01, MH02, MH03, V101, V201 and V202, the camera motion is slow, and the illumination conditions are moderate, making them relatively easy for monocular VO/VSLAM. On the contrary, in MH04, MH05, and V103, the illumination varies in a wide range. In V102, V103, and V203, the camera motion is aggressive, including nearly pure rotation or fast movement. Additionally, for that currently a pinhole camera model is assumed by our system, we pre-rectify all the images in the evaluation, which limits the field-of-view (FOV) of the camera and brings more challenges for monocular VO/VSLAM, especially indirect methods.

Here we report performances for both VO/VSLAM systems. Due to all the methods fail on V203, the corresponding results are not included. For the pure VO, as shown in Table. 2, for the moderate sequences, the proposed system has comparable accuracy with the state-of-the-art methods. For the challenging sequences, our method increases the robustness and accuracy compared to the baselines. Especially for V103, ORB-SLAM is unable to track the camera motion with severe exposure change and aggressive motion continuously, while our method only fails in 1/10 runs, indicating a more robust performance. Considering the full SLAM system, the proposed system successfully increase the robustness and accuracy compared to the VO-only solution. It is noticeable that for the difficult sequence V103, we achieve an RMSE of 4.44 cm, which significantly improves the result from the pure VO system. Compared our method against ORB-SLAM2, we observe our system generally provides a comparable performance with improvement on V103.

7.3.3 KITTI Odometry Dataset

The KITTI Odometry Dataset [53] is commonly used for evaluating visual state estimation systems in large-scale outdoor scenario. The result comparisons are reported in Table. 3 for both pure odometry and SLAM. As in Table. 3, the pure visual odometry results do not outperform ORB-SLAM2 while the full SLAM results is more competitive, which validate the improvement on introducing loop closure with NetVLAD. Our system takes more time in data association due to the usage of floating descriptors, which limits the computation time for batched optimization. However, as more representative local features are selected, our system manages a more lightweight map and advances in global BA with loop closure. We notice that for the sequences with loops (e.g., 00, 02, 05-07, 09, 10), generally the proposed system performs better than ORB-SLAM2. Especially for

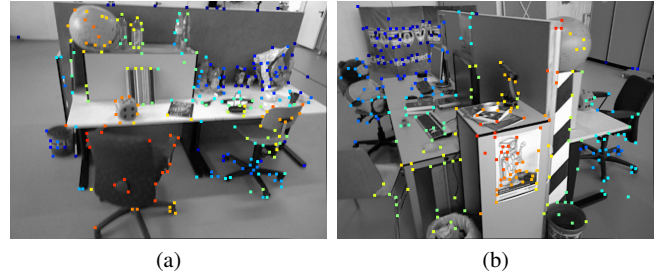


Fig. 7: Excerpts on the TUM Dataset, keypoints are colored by the inverse depth.

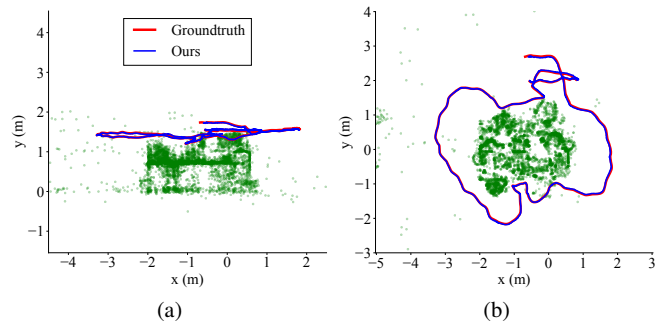


Fig. 8: Qualitative results on TUM Dataset.

sequence 09, we notice that ORB-SLAM2 failed to detect loop-closure with BoW, leading the trajectory to be inconsistent globally. On the contrary, our system successfully detects the loop with learnt global embedding. Along with global BA to optimize the trajectory, not surprisingly we achieve a better estimation result. Exception occurs on 08, where our system exhibits more severe scale drift in trajectory estimation, resulting in larger trajectory error. For the sequence 01, both methods fail to maintain a correct scale estimation, while only DSO succeeds.

7.4 Qualitative Results

Here we also provide some qualitative results in Fig. 6, Fig. 7 and Fig. 8. They show the reconstructed sparse map and estimated trajectories on several sequences. In Fig. 6, even some repetitive features can be well associated and triangulated with relatively accurate depth, for example, features on the floor. As in Fig. 7, in this indoor environment, we notice that our estimated trajectory has only trivial drift without using loop-closure. In addition, the sparse map well recovers the geometry in the test desk scene. Fig. 8 shows that the trajectory along with recovered visual structure on TUM Dataset, where we notice the estimated trajectory is well aligned with the ground truth.

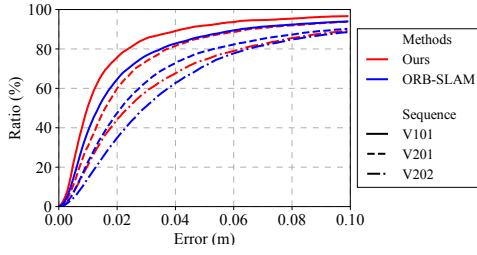


Fig. 9: Evaluation on the local reconstructions. The higher the percentage of points towards a zero distance the better.

7.5 Evaluations on Local Reconstruction

To evaluate the reconstruction accuracy, the sparse point-clouds generated by different methods are scaled and transformed via the alignment results. The error metric is defined as the RMSE of distances to the nearest neighborhood. The similar evaluation process can be found in [54]. Three sequences V101, V201, and V202 of EuRoC that provide dense ground-truth structure are selected for the evaluation.

As shown in Fig. 9, our method recovers a more accurate local structure, which in turn guarantees the accuracy of local trajectory estimation. Additionally, in V101, although the ATE of our method does not outperform ORB-SLAM2, the local structure is more accurate than ORB-SLAM2, indicating fewer outliers in our system. One advantage of indirect methods over direct ones is the covisibility connections established by powerful descriptors, with which local map is fully reused and the VO drift can be reduced [25]. ORB-SLAM2 does a great job of associating features cross views and fusing redundant points. However, under the cases of wider baseline, our system better associates landmarks with local features. Especially, comparing the sparse reconstruction of certain objects in the scene (e.g. the check-board), drastically fewer outliers and redundancies are generated from our system.

8 Conclusions and Future Work

In this paper, we presented a monocular SLAM system leveraging learnt description, repeatability and global embedding. Different from previous work, we focused on tightly coupling the learning-based frontend with indirect multiview geometry to fully exploit the network predictions. By interpreting the repeatability prediction in a probabilistic manner, we proposed a two-step tracking scheme to estimate the camera pose: direct tracking on the repeatability maps and refining the pose with the patch-wise association. We further discuss the local feature parameterization to consider the uncertainty in the network prediction. With loop closure detection from deep global embeddings, we maintain the consistency via global bundle adjustment. The experimen-

tal results demonstrated that the proposed system is capable to handle challenging situations, where both the state-of-the-art indirect and direct methods suffer from strong degradation.

In the future, we would like to investigate more lightweight network architectures or binary descriptor learning to accelerate our system. In the current system implementation, descriptor matching causes certain overhead for the map management, which we believe can be further optimized. In addition, supervising the learning frontend in an end-to-end manner might help our system better exploit from data-driven approaches. For example, using learnt relative pose as a prior for camera pose tracking could boost the general robustness under challenging conditions, in which cases direct formulations could not work well. Last but not least, extending the current system into a stereo or visual-inertial system is considered to be of much more practical value. As they can provide absolute scale information, such extension could better help robotic navigation in the long-term operations.

References

1. C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
2. Y. Sun, M. Liu, and M. Q.-H. Meng, "Motion removal for reliable rgb-d slam in dynamic environments," *Robotics and Autonomous Systems*, 2018.
3. S. Park, T. Schöps, and M. Pollefeys, "Illumination change robustness in direct visual slam," in *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2017, pp. 4523–4530.
4. R. Mahjourian, M. Wicke, and A. Angelova, "Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5667–5675.
5. W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
6. J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
7. S. Y. Loo, A. J. Amiri, S. Mashohor, S. H. Tang, and H. Zhang, "CNN-SVO: Improving the mapping in semi-direct visual odometry using single-image depth prediction," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 5218–5223.
8. K. Tateno, F. Tombari, I. Laina, and N. Navab, "Cnn-slam: Real-time dense monocular slam with learned depth prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6243–6252.
9. S. Yang and S. Scherer, "Cubeslam: Monocular 3-d object slam," *IEEE Transactions on Robotics*, 2019.
10. M. Runz, M. Buffier, and L. Agapito, "Maskfusion: Real-time recognition, tracking and reconstruction of multiple moving objects," in *2018 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2018, pp. 10–20.

11. V. Balntas, K. Lenc, A. Vedaldi, and K. Mikolajczyk, "Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5173–5182.
12. D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 224–236.
13. M. Dusmanu, I. Rocco, T. Pajdla, M. Pollefeys, J. Sivic, A. Torii, and T. Sattler, "D2-net: A trainable cnn for joint detection and description of local features," *arXiv preprint arXiv:1905.03561*, 2019.
14. J. Tang, L. Ericson, J. Folkesson, and P. Jensfelt, "GCNv2: Efficient correspondence prediction for real-time slam," *arXiv preprint arXiv:1902.11046*, 2019.
15. D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
16. S. Wang, R. Clark, H. Wen, and N. Trigoni, "Deepvo: Towards end-to-end visual odometry with deep recurrent convolutional neural networks," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 2043–2050.
17. R. Li, S. Wang, Z. Long, and D. Gu, "Undeepvo: Monocular visual odometry through unsupervised deep learning," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 7286–7291.
18. R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: An open-source slam system for monocular, stereo, and rgb-d cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
19. Y. S. Huaiyang Huang, Haoyang Ye and M. Liu, "Monocular visual odometry using learned repeatability and description," in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020.
20. G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *Proceedings of the 2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. IEEE Computer Society, 2007, pp. 1–10.
21. H. Strasdat, A. J. Davison, J. M. Montiel, and K. Konolige, "Double window optimisation for constant time visual slam," in *2011 International Conference on Computer Vision*. IEEE, 2011, pp. 2352–2359.
22. 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.
23. R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "Dtam: Dense tracking and mapping in real-time," in *2011 international conference on computer vision*. IEEE, 2011, pp. 2320–2327.
24. J. Engel, T. Schöps, and D. Cremers, "Lsd-slam: Large-scale direct monocular slam," in *European conference on computer vision*. Springer, 2014, pp. 834–849.
25. 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, 2018.
26. C. Forster, M. Pizzoli, and D. Scaramuzza, "Svo: Fast semi-direct monocular visual odometry," in *2014 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2014, pp. 15–22.
27. N. Yang, R. Wang, X. Gao, and D. Cremers, "Challenges in monocular visual odometry: Photometric calibration, motion bias, and rolling shutter effect," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 2878–2885, 2018.
28. H. Alismail, M. Kaess, B. Browning, and S. Lucey, "Direct visual odometry in low light using binary descriptors," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 444–451, 2016.
29. Z. Zhang, C. Forster, and D. Scaramuzza, "Active exposure control for robust visual odometry in hdr environments," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3894–3901.
30. P. Kim, B. Coltin, O. Alexandrov, and H. J. Kim, "Robust visual localization in changing lighting conditions," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 5447–5452.
31. C. G. Harris, M. Stephens *et al.*, "A combined corner and edge detector," in *Alvey vision conference*, vol. 15, no. 50. Citeseer, 1988, pp. 10–5244.
32. J. Shi and C. Tomasi, "Good features to track," Cornell University, Tech. Rep., 1993.
33. E. Rosten and T. Drummond, "Machine learning for high-speed corner detection," in *European conference on computer vision*. Springer, 2006, pp. 430–443.
34. H. Bay, T. Tuytelaars, and L. Van Gool, "Surf: Speeded up robust features," in *European conference on computer vision*. Springer, 2006, pp. 404–417.
35. E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *Proc. Int. Conf. Computer Vision*, Nov. 2011, pp. 2564–2571.
36. J. Revaud, P. Weinzaepfel, C. De Souza, N. Pion, G. Csurka, Y. Cabon, and M. Humenberger, "R2d2: Repeatable and reliable detector and descriptor," *arXiv preprint arXiv:1906.06195*, 2019.
37. P.-E. Sarlin, C. Cadena, R. Siegwart, and M. Dymczyk, "From coarse to fine: Robust hierarchical localization at large scale," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 716–12 725.
38. C. Tang and P. Tan, "Ba-net: Dense bundle adjustment networks," in *2019 International Conference on Learning Representations*.
39. L. von Stumberg, P. Wenzel, Q. Khan, and D. Cremers, "Gn-net: The gauss-newton loss for multi-weather relocalization," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 2, pp. 890–897, 2020.
40. D. DeTone, T. Malisiewicz, and A. Rabinovich, "Self-improving visual odometry," *arXiv preprint arXiv:1812.03245*, 2018.
41. R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307.
42. Y. Kanazawa and K. Kanatani, "Do we really have to consider covariance matrices for image feature points?" *Electronics and Communications in Japan (Part III: Fundamental Electronic Science)*, vol. 86, no. 1, pp. 1–10, 2003.
43. M. J. Brooks, W. Chojnacki, D. Gawley, and A. Van Den Hengel, "What value covariance information in estimating vision parameters?" in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, vol. 1. IEEE, 2001, pp. 302–308.
44. D. Gálvez-López and J. D. Tardós, "Bags of binary words for fast place recognition in image sequences," *IEEE Transactions on Robotics*, vol. 28, no. 5, pp. 1188–1197, October 2012.
45. A. Handa, T. Whelan, J. McDonald, and A. Davison, "A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM," in *IEEE Intl. Conf. on Robotics and Automation, ICRA*, Hong Kong, China, May 2014.
46. J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012.
47. Y. Ling, M. Kuse, and S. Shen, "Edge alignment-based visual-inertial fusion for tracking of aggressive motions," *Autonomous Robots*, vol. 42, no. 3, pp. 513–528, 2018.
48. Y. Zhou, H. Li, and L. Kneip, "Canny-vo: Visual odometry with rgb-d cameras based on geometric 3-d–2-d edge alignment," *IEEE Transactions on Robotics*, vol. 35, no. 1, pp. 184–199, 2018.
49. I. Nurutdinova and A. Fitzgibbon, "Towards pointless structure from motion: 3d reconstruction and camera parameters from general 3d curves," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2363–2371.

50. M. Kuse and S. Shen, "Robust camera motion estimation using direct edge alignment and sub-gradient method," in *2016 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2016, pp. 573–579.
51. S. Martull, M. Peris, and K. Fukui, "Realistic cg stereo image dataset with ground truth disparity maps," in *ICPR workshop Trak-Mark2012*, vol. 111, no. 430, 2012, pp. 117–118.
52. 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*, 2016. [Online]. Available: <http://ijr.sagepub.com/content/early/2016/01/21/0278364915620033.abstract>
53. A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 3354–3361.
54. A. Millane, Z. Taylor, H. Oleynikova, J. Nieto, R. Siegwart, and C. Cadena, "C-blox: A scalable and consistent tsdf-based dense mapping approach," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 995–1002.