# Probabilistic Spatial Distribution Prior Based Attentional Keypoints Matching Network

Xiaoming Zhao, Jingmeng Liu, Xingming Wu, Weihai Chen*, *Member, IEEE*, Fanghong Guo, *Member, IEEE*, and Zhengguo Li*, *Senior Member, IEEE*

*Abstract*—**Keypoints matching is a pivotal component for many image-relevant applications such as image stitching, visual simultaneous localization and mapping (SLAM), and so on. Both handcrafted-based and recently emerged deep learning-based keypoints matching methods merely rely on keypoints and local features, while losing sight of other available sensors such as inertial measurement unit (IMU) in the above applications. In this paper, we demonstrate that the motion estimation from IMU integration can be used to exploit the spatial distribution prior of keypoints between images. To this end, a probabilistic perspective of attention formulation is proposed to integrate the spatial distribution prior into the attentional graph neural network naturally. With the assistance of spatial distribution prior, the effort of the network for modeling the hidden features can be reduced. Furthermore, we present a projection loss for the proposed keypoints matching network, which gives a smooth edge between matching and un-matching keypoints. Image matching experiments on visual SLAM datasets indicate the effectiveness and efficiency of the presented method.**

*Index Terms*—**Keypoints matching, Probabilistic, Motion prior, Attention, Graph neural network, Sensor fusion**

## I. Introduction

**K**EYPOINTS matching is an essential module in many image processing problems such as the visual SLAM, image stitching, and so on. It aims to establish 2D-2D matches (correspondences) of keypoints [1]–[4] between two images, so that the relative pose of cameras can be recovered with the multi-view geometry [5], [6] or a set of differently exposed images [7], [8]. Therefore, it becomes important to restore as many correct matches as possible.

The matching of two point sets is a permutation problem. Matching $N$ points to other $N$ points leads to $N!$ possible permutations [9]. The standard approach to relieve this problem in image keypoint matching is to obtain a discriminative feature for each keypoint which is invariant to viewpoint, scale, illumination, and so on. And then the keypoint matches are recovered based on the similarity of features. A widely used heuristic strategy is firstly restoring a set of putative matches with the mutual nearest neighbor (mNN), and then filtering out

Xiaoming Zhao, Jingmeng Liu, Xingming Wu, and Weihai Chen* are with the School of Automation Science and Electrical Engineering, Beihang University, Beijing, 100191 (e-mail: xmzhao@buaa.edu.cn, ljm-buaa110@163.com, wxmbuaa@163.com, and whchen@buaa.edu.cn).

Fanghong Guo is with the Department of Automation, Zhejiang University of Technology, Hangzhou 310014, China (email: fhguo@zjut.edu.cn).

Zhengguo Li* is with the SRO department, Institute for Infocomm Research, 1 Fusionopolis Way, Singapore (email: ezgli@i2r.a-star.edu.sg).

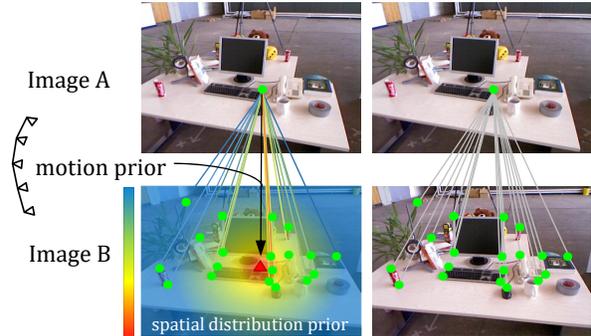*(Corresponding author: Weihai Chen and Zhengguo Li.)*

Fig. 1: We propose to utilize the motion prior to exploit the spatial distribution prior of keypoints in two successive images, and the prior is used to assist the keypoints matching process. In this figure, a keypoint in image A is going to match with the keypoints in image B. In the left column, the proposed spatial distribution prior of the keypoint in image A is displayed superimposed on image B, which is used to assist the matching attentional graph neural network. While the right column shows the matching process without any prior.

false positive matches [10]–[12] or regaining more consistent matches [13]–[15]. In addition to heuristic matching methods, several deep learning-based models [16]–[19] are proposed to exploit the local context based on the putative matches.

In this formulation, however, most local feature descriptors [10], [20]–[24] only encode an image patch with a limited region and ignore the global context. To track this problem, several end-to-end learned feature descriptors [25]–[27] are proposed to extract global keypoint features implicitly. Nonetheless, the keypoints matching involves the geometric distribution and feature similarity of keypoints in/between images. Humans usually obtain this information by checking the feature points back-and-forth when matching keypoints between two images. This behavior is imitated with the self- and cross- attentional graph neural network (GNN) in the recently proposed SuperGlue [28]. We extend this formulation with the observation that when humans have prior knowledge about the keypoints, they will compare and search for the matching keypoints based on the prior knowledge first, which will greatly reduce the matching effort. The search range is expanded only when no matching keypoints are found based on this prior knowledge.

As shown in the right column of Fig. 1, without prior

knowledge, all the keypoints to be matched are treated equally without discrimination. So the SuperGlue [28] has to refine the keypoint features with all contextual keypoints gradually through 18 attentional graph neural network (GNN) layers, which might be inefficient. On the other hand, the motion prior between two images could be obtained in SLAM systems using inertial measurement units (IMUs) or wheel odometry. For example, practical visual SLAM systems usually rely on a motion model to estimate an initial solution for the pose estimation between two consecutive images [29], [30]. It is thus desired to assist the exploitation of keypoints geometric distribution and feature similarity cross image by using the motion prior, as shown in the left column of Fig. 1.

In this paper, we utilize the motion prior to compute the keypoint distribution prior, and present a keypoint distribution prior integration strategy of attentional GNN based on the probabilistic perspective of attention. More specifically, the initial pose of accurate IMU integration [31], [32] is regraded as the motion prior of two successive images. Then, the spatial distribution prior of keypoints is obtained by warping keypoints across images with the motion prior. To integrate the spatial distribution prior into the attentional GNN, we propose to regard the attention as a Gaussian distribution that models the probability of the feature correlation. As such, the spatial distribution prior can be naturally and efficiently integrated into the attention module based on the conditional independence hypothesis. With the assistance of the prior, the attentional GNN can pay less effort to recover the correct matches. Thus, we streamline the network by decreasing some of the attentional layers. Furthermore, different from the matching loss, in which the ground-truth matches are obtained with a hard threshold of the keypoint projection errors, we propose to relax the hard threshold by directly utilizing the projection error in a margin.

The main contributions of this paper are in three folds:

- With the probabilistic interpretation of attention, we exploit the motion prior from IMU measurements to obtain the spatial distribution prior of keypoints. Thus the keypoints matching network can be streamlined with fewer attentional GNN layers, resulting in less computational cost.
- We proposed to use the projection errors instead of the hard-threshold ground-truth matches to supervise the training of attentional GNN, so that the network can achieve better overall matching performance.
- Our approach can estimate keypoint matches efficiently and accurately on SLAM datasets such as InteriorNet [33], TUM-RGBD [34], and ETH3D [35].

The rest of this paper is organized as follows. Section II reviews the handcrafted and deep learned keypoints matching methods for image pairs. Section III begins by formulating the attentional GNN, introduces the direct and probabilistic spatial distribution prior integration to attention, and proposes the projection loss. In Section IV, we evaluate and discuss the probabilistic prior integration and projection loss on different datasets. Finally, a conclusion is given in Section V.

## II. RELATED WORKS

Keypoints matching for an image pair usually follows the following steps: a) detecting keypoints and extracting local features, b) finding the putative matches with brutally nearest neighbor matching, and then c) filtering out the false positive matches. For the first step, vast handcrafted [1], [3], [10], [20], [36] and learned [4], [21]–[23], [25]–[27], [37]–[41] methods have been proposed, devoting efforts to extract repeatable keypoints and discriminative local features. Whereas the latter two steps concentrate on retrieving accurate and exhaustive keypoint matches in an image pair based on the extracted keypoints and local features in the first step. They can also roughly be divided into two categories: the handcrafted and the learned.

### A. Handcrafted keypoints matching

Researchers have designed massive approaches for keypoints matching based on heuristic experiences. The most successful and widely used one is the ratio test [10], which is proposed along with the SIFT descriptor by Lowe. It removes the false positive matches by testing the similarity between the nearest and next-nearest neighbor. Subsequent research on ratio test improves the similarity measurement by using the Earth Mover's Distance [11].

To recover more robust matches, various methods explores local keypoints distribution such as triangle constraint [42], vector field consensus [14], local neighborhood structures [9], [43], coherence-based separability constraint [15], and motion smoothness [44]. In addition to local keypoints distribution exploration, other strategies such as suitable local feature selection [45] and optical flow guided matching [46] have also been explored to recover more matches. However, the above methods could ignore the global consistency since they only focus on the local distribution.

On the other hand, the global properties of keypoints have also been explored. In such approaches, the matching problem is formulated as correspondence function [12], bounded distortion transformation [13], rigid transformation [47], or graph matching [48]–[51], and they can be solved iteratively. Our proposed method is based on graph neural network and is most similar to graph matching methods [48]–[51] which are designed based on human experiences. For example, [50] focuses on the second- and high-order graph matching and [51] introduces a dual calibration strategy to model the correspondence relationship in points and edges respectively. Moreover, the RANSAC [52] procedure in subsequent tasks could also be regarded as a false positive removing process, in which the potential matches are iteratively sampled to fit a model and the putative matches are classified as inliers and outliers based on the fitness to the model. However, the iterative solution of global formulations is less computational efficient, and the result could degrade when the underlying matching model differs from the predefined model.
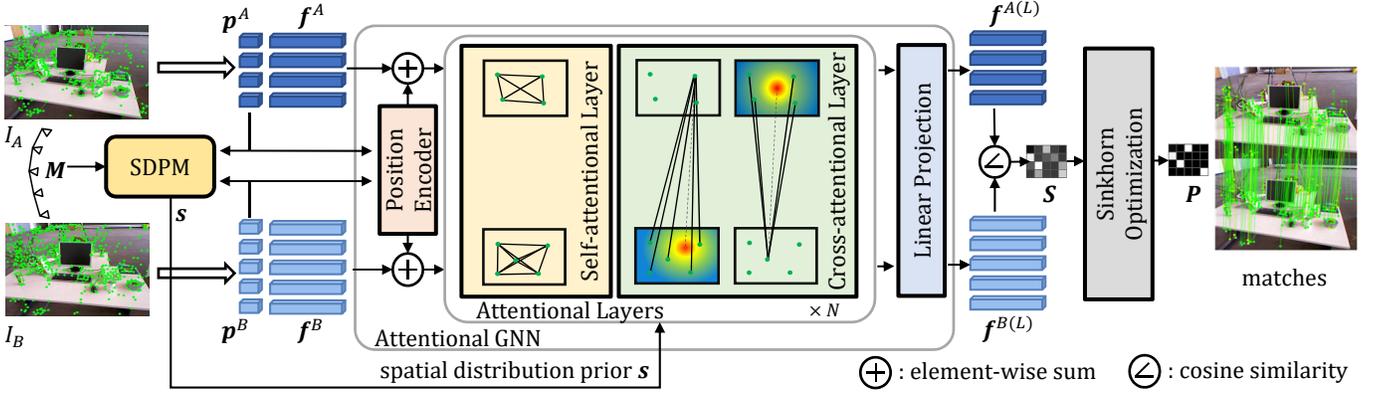
Fig. 2: The proposed spatial distribution prior assisted matching pipeline. The $(\boldsymbol{p}^A, \boldsymbol{f}^A)$ and $(\boldsymbol{p}^B, \boldsymbol{f}^B)$ are local features extracted from two successive image frames $I_A$ and $I_B$ respectively, and $\boldsymbol{M}$ is the IMU measurements between $I_A$ and $I_B$. The **S**patial **D**istribution **P**rior **M**odule (SDPM) takes the $\boldsymbol{p}^A$, $\boldsymbol{p}^B$, and $\boldsymbol{M}$ as inputs and computes the spatial distribution prior $\boldsymbol{s}$. The attentional GNN then can utilize the spatial distribution prior $\boldsymbol{s}$ to assist the optimization of the keypoint features. At last, the Sinkhorn algorithm is adopted to optimize the cosine similarity matrix $\boldsymbol{S}$ of the features, and it outputs the matching result $\boldsymbol{P}$.

## B. Learned keypoints matching

In recent years, deep networks reveal their superior performance in various computer vision tasks including keypoints matching. The deep learned image keypoints matching can be roughly divided into two categories, i.e. semantic and geometric matching. The semantic keypoints matching approach matches the semantic keypoints in different instances of the same category of objects [53]–[55], while the geometric matching matches the keypoints on different images of the same scene and does not consider high-level semantic information [16]–[19], [28]. Both of them are robust to lighting conditions [56], [57]. The keypoints matching existing in SLAM systems usually refers to the latter, thus we focus on the geometric keypoints matching in this paper.

The seminal work of the deep network for the point set is the PointNet [58], in which the multi-layer perception (MLP) extracts the point feature and the max-pooling layer aggregates the global feature. Although the PointNet is originally invented for the classification and segmentation of unordered 3D points, it also encourages the design of putative matches filtering networks [16]–[19]. They firstly pair the putative keypoint matches as 4D quads. Then the networks take the 4D quads as input and output a score for each putative match. The false matches are filtered out with a threshold on the scores. Specifically, PointCN [16] aggregates the context of all matching features with Context-Normalization (CN) modules and estimates the essential matrix of an image pair with differentiable weighted eight points model. Following PointCN, NM-Net [17] mines compatibility-specific locality of keypoints to discover reliable local neighbors. And OA-Net [18] introduces differentiable pooling and unpooling layers to exploit the global context. More recently, ACNe [19] extends the CN by introducing the local and global attentive weights. However, the putative matches filtering networks merely learn the keypoints geometric features, as they only take the keypoint positions as input and ignore the local features.

Recently, SuperGlue [28] takes advantage of both keypoint positions and local features. It utilizes the attentional GNN to propagate and aggregate the contextual features both in intra-image and inter-image. In this way, each local feature in an image is refined with features in both images. Then the Sinkhorn [59] algorithm optimizes the matching score matrix and produces an assignment matrix representing the matching result. Our method is built on the architecture of attentional GNN introduced in [28]. Besides the keypoint positions and local features, our method also utilizes the initial pose from IMU integration. The spatial distribution prior of keypoints is computed based on the initial pose to reduce the matching efforts.

## III. METHOD

In this section, we first present a brief overview of the proposed pipeline. To incorporate with the prior, a spatial distribution prior module is then introduced. After formulating the attentional GNN, the direct and probabilistic spatial prior integration method of attentional GNN is proposed. Moreover, the matching loss is investigated and a more reasonable projection loss is presented.

## A. Pipeline overview

*1) Problem formulation:* Formally, we consider two images $\{I_A, I_B\}$, their depth maps $\{\boldsymbol{D}^A, \boldsymbol{D}^B\}$ and the IMU measurements $\{\boldsymbol{M}_i = (\boldsymbol{\omega}_B, \boldsymbol{a}_B)_i | i \in [1, M]\}$ between $I_A$ and $I_B$, where $\boldsymbol{\omega}_B = (\omega_x, \omega_y, \omega_z)$ and $\boldsymbol{a}_B = (a_x, a_y, a_z)$ denote the angular velocity and linear acceleration in the IMU body frame respectively. The keypoint features $\{\boldsymbol{F}_i = (\boldsymbol{p}_i, \boldsymbol{f}_i) | i \in [1, N]\}$ of each image are first extracted by the feature extractor, where $\boldsymbol{p}_i = (u, v)$ is the position of the i-th keypoint, $\boldsymbol{f}_i$ denotes the corresponding feature, and $N$ is the number of extracted keypoints in the image. Assuming there are $N_A$ keypoints in $I_A$ and $N_B$ keypoints in $I_B$, thus the keypoint features of $\{I_A, I_B\}$ are represented as $\{\boldsymbol{F}_i^A\}_{i \in [1, N_A]}$ and
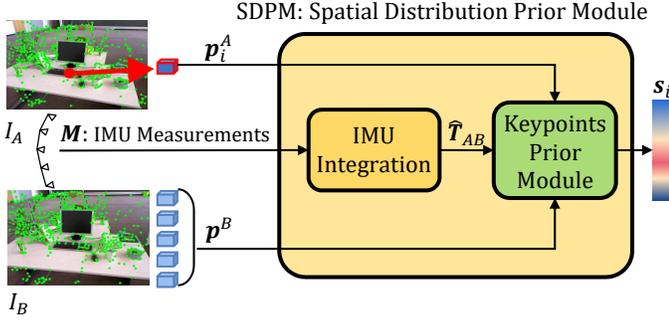
Fig. 3: The Spatial Distribution Prior Module (SDPM). In this module, the IMU measurements $M$ are first integrated into an initial pose $\hat{T}_{AB}$. Then the Keypoints Prior Module computes the spatial distribution prior $s$. For clarity, we take the $i$-th keypoint $p_i^A$ in image $I_A$ as an example, and its spatial distribution prior $s_i$ related to $p^B$ is a normalized vector.



Fig. 4: The keypoint prior module. The red and green dots represent the detected keypoints in image $I_A$ and $I_B$ respectively, and the blue pentacles denote the landmarks in the 3D world. For the $i$-th keypoints in $I_A$, this module warps it with the initial pose $\hat{T}_{AB}$ into $I_B$ (the red triangle in $I_B$), and the spatial distribution prior $s_i$ is formulated as a Gaussian distribution in $I_B$.

$\{F_i^B\}_{i\in[1,N_B]}$. Our goal is to match the pixel coordinate $\{p_i^A\}_{i\in[1,N_A]}$ and $\{p_i^B\}_{i\in[1,N_A]}$, with the initial pose $T_{AB}$ integrated from the IMU measurement $\{M_i\}_{i\in[1,M]}$. Note that under this formulation, the keypoint features can be extracted with any existing feature extractor such as the handcrafted [10], [20], [36] or the learned [21], [23]–[27] feature.

*2) Prior assisted attentional GNN matching pipeline:* The overview of the proposed matching pipeline is shown in Fig. 2. First, the Spatial Distribution Prior Module (SDPM), which will be presented in detail in Section III-B, computes the spatial distribution prior $s$ by using the keypoints position $p^A$, $p^B$, and the IMU measurements $M$. Before the attentional GNN layers, the Keypoint Encoder encodes the keypoints position $p$ into feature space $f$. Then, the spatial distribution prior $s$ is fed to the attentional GNN (Section III-C and Section III-D) to assist the keypoint features optimization. Each attentional GNN layer has two types of attention: self-attention and cross-attention. The spatial distribution prior $s$ is integrated into these attention mechanisms in a probabilistic way. Next, the cosine similarity matrix $S$, which is computed by the linear transformed features of the attentional GNN, is optimized by the Sinkhorn algorithm [59] (Section III-E) to obtain the assignment matrix $P$ representing the matching result.

### B. Spatial distribution prior module

As shown in Fig. 3, the spatial distribution prior module (SDPM) contains two sub-modules, the IMU integration and the keypoints prior module. The spatial distribution prior of the $i$-th keypoint in $I_A$ and all the keypoints in $I_B$ is illustrated for clarity.

*1) IMU Integration:* A basic motion assumption is to consider the velocity between two images is constant [29], [30]. To obtain an accurate prior motion, IMU measurements have been widely used in recent SLAM systems to constrain the pose optimization problem [31], [32], [60]–[62]. In this paper, we adopt the switched linear system based IMU integration model [31], [32]. It integrates all IMU measurements $\{M_i =$
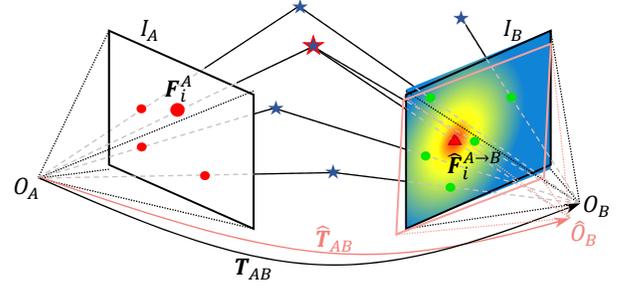
$(\boldsymbol{\omega}_B, \boldsymbol{a}_B)_i | i \in [1, M]\}$ between $I_A$ and $I_B$ and outputs motion prior, that is, the initial relative pose $\hat{T}_{AB}$ of $I_A$ and $I_B$.

*2) Keypoints Prior Module:* As illustrated in Fig. 4, the keypoints prior module first warps the keypoint $F_i^A = (p_i, f_i)^A$ in $I_A$ to $I_B$ with the motion prior $\hat{T}_{AB}$:

$$\hat{p}_i^{A\to B} = \Pi(\hat{T}_{AB}\Pi^{-1}(p_i^A, D^A)), \tag{1}$$

where $\Pi(p)$ is the camera projection function that projects a 3D landmark $l \in \mathbb{R}^3$ into the camera image plane.

We regard the warped position $\hat{p}_i^{A\to B}$ in $I_B$ as keypoint prior knowledge across images. We name it as spatial distribution prior $s_{i,j}$ of keypoints and encode it with Gaussian distribution

$$s_{i,j} = \exp \frac{-(d_{ij}^{A,B})^2}{\sigma}, \tag{2}$$

where $d_{ij}^{A,B} = \|\hat{p}_i^{A\to B} - p_j^B\|$ denotes the re-projection distances of keypoint $\hat{p}_i^{A\to B}$ and $p_j^B$. Similarly, the spatial distribution prior of keypoints in the same image can be obtained with $d_{ij}^{k,k} = \|p_i^k - p_j^k\|$ $(k \in A, B)$. In Section III-D, we present two methods to integrate the spatial distribution prior $s_i$ into attentional GNN.

### C. Attentional GNN for keypoints matching

The attentional GNN contains three modules, namely, Position Encoder, Attentional Layers, and the Linear Projection layer.

*1) Position Encoder:* As shown in Fig. 2, given the features $F^A = \{p^A, f^A\}$ and $F^B = \{p^B, f^B\}$ of an image pair, we first embed the keypoint position $p$ into its feature space $f$ with the position encoder. This enables the attentional GNN to utilize the position in an implicit way. Following [28], an MLP module is adopted as the position encoder for each feature $F_i \in [1, N]$:

$$f_i^{(1)} = f_i + \text{MLP}(p_i) \quad i \in [1, N]. \tag{3}$$

*2) Attentional Layers:* The attentional GNN layers for keypoints matching [28] formulates the vertexes as keypoint feature $f$, edges as attention $\alpha$. It optimizes each keypoint
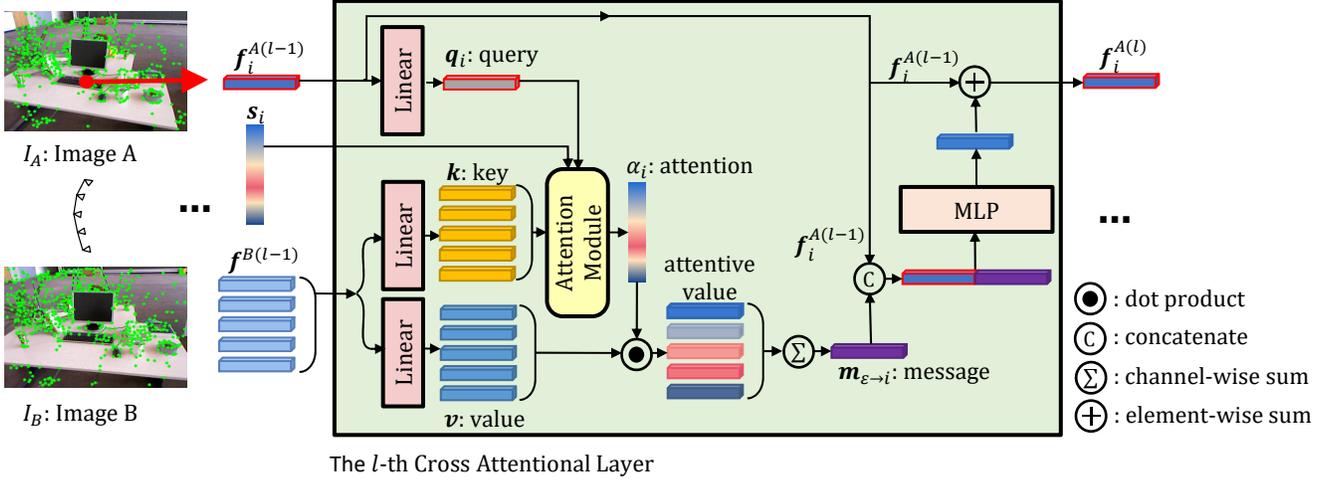
Fig. 5: The cross attentional layer. The $i$-th feature $\boldsymbol{f}_i^{A(l-1)}$ in $I_A$ is optimized by all the features $\boldsymbol{f}^{B(l-1)}$ in $I_B$ and the spatial distribution prior in the $l$-th cross attentional GNN layer.

feature $\boldsymbol{f}$ iteratively to obtain the contextual information in-and-cross images. Fig. 5 illustrates the process that the $i$-th feature in $I_A$ is optimized by all the features in $I_B$. In the following, to facilitate the notation, the superscripts $A$ and $B$ are omitted to denote arbitrary keypoint pairs of arbitrary images.

As in transformer [63], the input features are first linear-transformed to the query $\boldsymbol{q}_i$, key $\boldsymbol{k}$, and value $\boldsymbol{v}$:

$$
\begin{aligned}
\boldsymbol{q}_i &= \boldsymbol{W}_1 f_i^{(l-1)} + \boldsymbol{b}_1 \\
\boldsymbol{k}_j &= \boldsymbol{W}_2 f_j^{(l-1)} + \boldsymbol{b}_2 \\
\boldsymbol{v}_j &= \boldsymbol{W}_3 f_j^{(l-1)} + \boldsymbol{b}_3
\end{aligned}
\tag{4}
$$

where $j$ denotes the $j$-th feature in images. Then, the attention

$$
\alpha_{ij} = \text{atten}(\boldsymbol{q}_i, \boldsymbol{k}_j, s_{i,j})
\tag{5}
$$

is obtained to focus on different values $\boldsymbol{v}$. So that the contextual messages $\text{m}_{\varepsilon \to i}$ for the feature $\boldsymbol{f}_i$ are propagated through all the edges $\varepsilon \to i$ connected with vertex $\boldsymbol{f}_i$:

$$
\text{m}_{\varepsilon \to i} = \sum_{j,(ij) \in \varepsilon} \alpha_{ij} \boldsymbol{v}_j,
\tag{6}
$$

The message $\boldsymbol{m}_{\varepsilon \to i}$ is aggregated into the $i$-th feature $\boldsymbol{f}_i^{(l-1)}$ in a residual style. The residual is the MLP dimensionality reduction of the concatenated contextual messages $[\boldsymbol{f}_i^{(l-1)} \| \text{m}_{\varepsilon \to i}]$:

$$
\boldsymbol{f}_i^{(l)} = \boldsymbol{f}_i^{(l-1)} + \text{MLP}([\boldsymbol{f}_i^{(l-1)} \| \text{m}_{\varepsilon \to i}]).
\tag{7}
$$

*3) Linear Projection:* In the above processes, the self- and cross- attentional GNN layers are used iteratively to imitate the back-and-forth behaviors when asking a human to match keypoints. For a specific feature $\boldsymbol{f}_i$ in one image, the self-attention GNN layer propagates contextual messages from the same image, while the cross-attention GNN layer obtains contextual messages from another image. After several self-and cross- attentional GNN layers, the final matching feature

$\boldsymbol{f}_i$ are linear projections of the last attentional GNN layer outputs $\boldsymbol{f}_i^{(L)}$:

$$
\boldsymbol{f}_i = \boldsymbol{W} \boldsymbol{f}_i^{(L)} + \boldsymbol{b}.
\tag{8}
$$

*D. Prior Assisted Attentional GNN for keypoints matching*

For the attention module in Fig. 5, no prior knowledge is considered in the vanilla formulation [28]

$$
\alpha_{ij} = \text{atten}(\boldsymbol{q}_i, \boldsymbol{k}_j, s_{i,j}) = \text{softmax}_j(\boldsymbol{q}_i^T \boldsymbol{k}_j) = \frac{e^{\boldsymbol{q}_i^T \boldsymbol{k}_j}}{\sum_{j=1}^n e^{\boldsymbol{q}_i^T \boldsymbol{k}_j}}, \tag{9}
$$

it has to learn to extract appropriate contextual cues from scratch. Thus multiple attentional GNN layers are necessary to enable the network to obtain correct context message gradually. This leads to more computational cost, which is not suitable for time-sensitive visual SLAM systems. To solve this problem, we propose to integrate the spatial distribution prior into the attention module.

*1) Direct spatial prior integration:* It is rational that a keypoint should have strong correlations with the ones close to it. Based on this intuition, the correlation strength of each keypoint should decrease with the distance to others. A straightforward approach to integrating the prior $s_{ij}$ into attentional GNN is using it to weigh the attention in Equation (6). But directly weighting on $\alpha_{ij}$ will break the normalization property, thus we can weigh on the variable of softmax $\alpha_{ij} = \text{softmax}_j(s_{ij}\boldsymbol{q}_i^T \boldsymbol{k}_j)$. However, when the distance $d_{ij}$ of feature $i$ and $j$ is greater than $3\sigma$, the prior $s_{ij}$ becomes very close to zero. Thus the attention on other hidden features will be weakened unexpectedly. So the direct spatial prior integration is formulated as

$$
\alpha_{ij} = \text{atten}(\boldsymbol{q}_i, \boldsymbol{k}_j, s_{i,j}) = \text{softmax}_j((1 + s_{ij})\boldsymbol{q}_i^T \boldsymbol{k}_j). \tag{10}
$$

It can be seen that the weight $(1 + s_{ij})$ ranges from one to two, so it can preserve the attention on other hidden features even if $s_{ij} = 0$. Under this formulation, the spatial prior could hinder the original attention extraction as they are tightly coupled together. So a probabilistic spatial prior integration is proposed to tackle this problem.
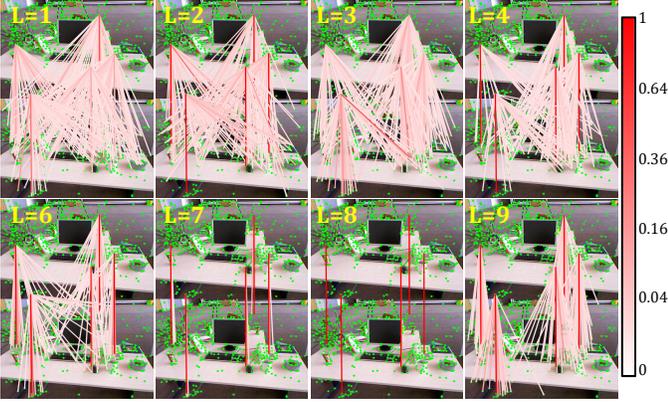
Fig. 6: The visualization of the cross-attention of the vanilla attentional GNN layers in [28]. The cross-attention values of five keypoints in the first image are visualized, and the attentions lower than 0.01 are omitted for clarity. The first several vanilla attention layers focus on a large region, and it is not until the seventh layer that the network focus on the corresponding keypoints. Moreover, the ninth attention layer slightly expands its focus region.

*2) Probabilistic spatial prior integration:* To integrate the spatial prior into the attention module in a natural way, we exploit the probabilistic perspective of the attention. In the formulation of the vanilla attention module, i.e. Equation (9), the attention $\alpha_{ij}$ is a normalized Gaussian distribution [64] of the cosine similarity of query $q_i$ and key $k_j$. With this perspective, the vanilla attention $\alpha_{ij}$ can be regarded as a joint distribution for each embedded features $q_i$ and $k_j$ as

$$\Pr(f_i, f_j|\text{appearance, spatial, ...}) = \alpha_{ij} \propto \exp(q_i^T k_j). \quad (11)$$

In fact, it is a distribution conditioned by appearance, spatial characteristics, and other hidden contexts. Our intuition is that the vanilla attention module has to model all the hidden contexts from scratch, which could be difficult. This intuition can be further confirmed with the closer investigations on the attention $\alpha_{ij}$ in each layer. As shown in Fig. 6, the network struggles to find a proper attentional relationship of features across images. With the assistance of spatial distribution prior, the attentional GNN can focus on the modeling of other contexts such as the appearance, and the attentional GNN layers can be streamlined to only two attentional GNN layers, as shown in Fig. 8.

Then, we consider the spatial distribution prior $s_{ij}$ for keypoints $p_i$ and $p_j$ as a Gaussian distribution

$$\Pr(f_i, f_j|\text{spatial}) \propto s_{ij} = \exp(\frac{-d_{ij}^2}{\sigma}), \quad (12)$$

and it is independent of other hidden contexts. As such, with the spatial distribution prior $\Pr(f_i, f_j|\text{spatial})$, the network can be relaxed to solely model the appearance and other hidden contexts $\Pr(f_i, f_j|\text{appearance, ...})$. Thus the attentional

distribution can be formulated as

$$\Pr(f_i, f_j|\text{appearance, spatial, ...})$$
$$\propto \Pr(f_i, f_j|\text{spatial}) \Pr(f_i, f_j|\text{appearance, ...})$$
$$\propto \exp(\frac{-d_{ij}^2}{\sigma}) \exp(q_i^T k_j)$$
$$\propto \exp(\frac{-d_{ij}^2}{\sigma} + q_i^T k_j). \quad (13)$$

With the normalization of the above distribution, this process can be efficiently implemented as

$$\Pr(f_i, f_j|\text{appearance,spatial, ...})$$
$$= \text{softmax}(\frac{-d_{ij}^2}{\sigma} + q_i^T k_j). \quad (14)$$

*E. Sinkhorn optimization*

As shown in Fig. 2, after the attentional GNN, a pairwise score matrix $S = \{S_{ij} = (f_i^A)^T f_j^B \mid i \in [1, N_A], j \in [1, N_B]\}$ of features is first computed in the optimization layer. Following [28], two dustbins are added to the last column and row of $S$ to form $\bar{S}$ to cope with the unmatched keypoints. As the keypoint matching in image pair is a bipartite matching formulation in optimal transport [65], its assignment matrix $\bar{P}$ is given by

$$L(a, b) = \min_{P \in U(a,b)} \bar{P} \odot \bar{S} \quad (15)$$

where $\odot$ denotes the Hadamard product operation, $a = [\mathbf{1}_{N_A}^T \quad N_B]^T$ and $b = [\mathbf{1}_{N_B}^T \quad N_A]^T$ are the mass of two discrete measures, $U(a, b)$ denotes the assemble couplings under the constraints $\bar{P}\mathbf{1}_{N_B+1} = a$ and $\bar{P}^T\mathbf{1}_{N_A+1} = b$. The classical solution of bipartite matching is the Hungarian algorithm [66]. Its differentiable version Sinkhorn algorithm [59] entropy-regularizes the soft assignment and can be solved on GPU efficiently.

The assignment matrix $\bar{P}$ indicates the matching confidence of a keypoints pair. Then the matches are recovered by finding the row and column minimum and filtering out the false positives with a confidence threshold.

*F. Loss functions*

To train the proposed attentional GNN, we explore the matching loss proposed in [28] and propose a projection loss in this section. The former utilizes hard-threshold matches while the latter uses projection errors to supervise the attentional GNN.
*1) Matching loss:* The elements in the assignment matrix $\bar{P}_{1:N_A,1:N_B}$ indicate the matching confidence for each possible keypoint pairs, and the elements in dustbins $\bar{P}_{N_A+1,1:N_B}$ and $\bar{P}_{1:N_A,N_B+1}$ suggest the unmatched confidence. The direct way to constrain $\bar{P}$ is by using the ground-truth keypoint matching $\mathcal{M} \subseteq \{(i,j)|i \in [1, N_A], j \in [1, N_B]\}$, unmatching $\mathcal{I} \subseteq [1, N_A]$, and $\mathcal{J} \subseteq [1, N_B]$.

To obtain the ground-truth matches $\mathcal{M}$, keypoints $p_i^A$ in the image $I_A$ are warped to $I_B$ as $p_i^{A \to B}$ with its depth map and ground-truth pose $T_{AB}$. The distance matrix $\mathcal{D} \in \mathbb{R}^{N_A \times N_B}$
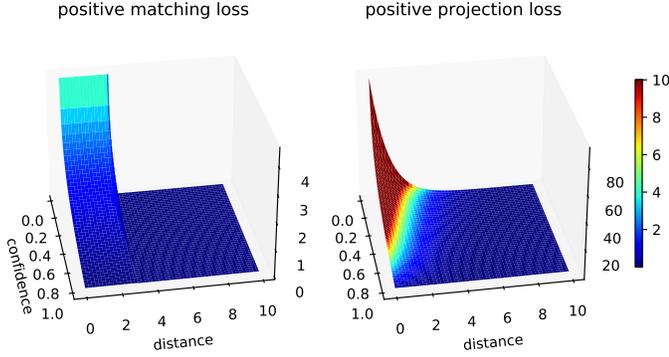
Fig. 7: Positive losses with respect to estimating matching confidence and ground-truth matching distance (with $th = 3$ and $mg = 10$). The matching loss truncates to zero when ground truth matching distance $d > th$, while the surface of the projection loss is much softer at the $th$.

between the warped keypoints $\boldsymbol{p}_i^{A \to B}$ and $\boldsymbol{p}_j^B$ is computed. Then, the correct matching $\mathcal{M}$ is defined as

$$\mathcal{M} = \{(i,j) | d_{ij} < th, d_{ij} = \min d_{:j}, d_{ij} = \min d_{i:}, d_{ij} \in \mathcal{D}\}, \quad (16)$$

where $th$ is a pixel threshold.

As such, the matching loss could be formulated as the mean negative log-likelihood of the assignment $\bar{\boldsymbol{P}}$:

$$\begin{cases} \mathcal{L}_{positive} = -\dfrac{\sum_{(ij) \in \mathcal{M}} \log \bar{\boldsymbol{P}}_{ij}}{\#\mathcal{M}} \\ \mathcal{L}_{negative} = -\dfrac{\sum_{i \in \mathcal{I}} \log \bar{\boldsymbol{P}}_{i,N_B+1} + \sum_{j \in \mathcal{J}} \log \bar{\boldsymbol{P}}_{N_A+1,j}}{\#\mathcal{I} + \#\mathcal{J}} \end{cases}, \quad (17)$$

where $\#\mathcal{M}$ denotes the total number of matched keypoint pairs, $\#\mathcal{I}$ and $\#\mathcal{J}$ are the number of unmatched keypoints in $I_A$ and $I_B$ respectively. Thus the total matching loss is given as

$$\mathcal{L}_{total} = 2\mathcal{L}_{positive} + \mathcal{L}_{negative}. \quad (18)$$

*2) Projection loss:* However, in our experiments, we find that the model could not always converge to its best weights with the matching loss. We hypothesize that the hard-threshold ground-truth definition of Equation (16) with a fix $th$ could be sub-optimal. As shown in Fig. 7, the positive matching loss truncates the loss at $th$, regarding all the loss as the same when $d_{ij} < th$, while neglecting the loss for those $d_{ij} > th$.

To address this problem, we relax the threshold $th$ in Equation (16) to a much larger margin $mg$ ($mg > th$), and the loss should be defined with the distance $d_{ij}$. To this end, a subset of $\mathcal{D}$ is defined as positive matches

$$\tilde{\mathcal{D}} = \{d_{ij} | d_{ij} < mg, d_{ij} = \min d_{:j}, d_{ij} = \min d_{i:}, d_{ij} \in \mathcal{D}\}. \quad (19)$$

And the unmatched keypoints $\tilde{\mathcal{I}}$ of $I_a$ and $\tilde{\mathcal{J}}$ of $I_B$ are those not in the positive matches $\tilde{\mathcal{D}}$.

Under this formulation, each element in $\tilde{\mathcal{D}}$ is a possible match and its value indicates the ground truth matching distance.

Then the projection loss of all possible matches is given as

$$\mathcal{L}_{positive} = -\frac{\sum_{(ij) \in \tilde{\mathcal{D}}} \mathcal{L}_{projection}}{\#\tilde{\mathcal{D}}}, \quad (20)$$

where

$$\mathcal{L}_{projection} = \exp(th - d_{ij}) \log \bar{\boldsymbol{P}}_{ij}. \quad (21)$$

The projection loss can constrain the loss of positive samples softer and more reasonably, as shown in Fig. 7. On the one hand, when the matching distance of keypoints is close to zero, they are more likely to be matched. In this case, the loss will push the output confidence to one. On the other hand, a larger ground truth matching distance indicates a weak tie, thus the loss constraint softly decreases with the distance. So that the loss can recall more potential matches.

## IV. EXPERIMENTS

In this section, we evaluated the proposed method on the one simulated and two real SLAM datasets. We first introduce experimental configurations including the datasets and implementation details. Then the proposed model is compared with state-of-the-art matching networks. To further inspect the proposed model, we discuss two types of prior integration and loss functions in the ablation studies.

### A. Datasets

We train and evaluate the proposed model on three indoor SLAM datasets: InteriorNet [33], TUM-RGBD [34], and ETH3D [35]. To check the generalization performance of the model, evaluations are also conducted on two homography datasets: HPatches [67] and Oxford-Paris [68]. Detailed introductions of these datasets are as follows:

**InteriorNet dataset** [33] provides 20 million synthetic RGBD images (RGB color images and corresponding depth maps), IMU data, and ground-truth camera poses. This dataset was created by professional designers based on real-world decorations. As it provides accurate depth maps and IMU measurements, there is no need to filter the data. We selected 21 trajectories for training, and 4 trajectories for testing.

**TUM-RGBD dataset** [34] contains the data captured by Kinect-V1 and motion capture system in three different indoor environments. The data includes well-calibrated RGBD images, ground truth camera poses, and accelerator readings for some sequences. As almost all the images were collected in a good light condition and rich textures, it is suitable for training the proposed matching model. To do so, we removed the sequences with dynamic objects and metallic spheres and selected 26 sequences for training, 7 sequences for testing.

**ETH3D dataset** [35] also provides the well-calibrated RGBD images and ground truth camera poses. Moreover, it includes time-aligned IMU measurements with linear acceleration and liner angular velocity for each sequence. Thus it is a proper training dataset for the proposed model. We selected the training and testing sequences from its training split. Those sequences with illumination changes, black color images, and
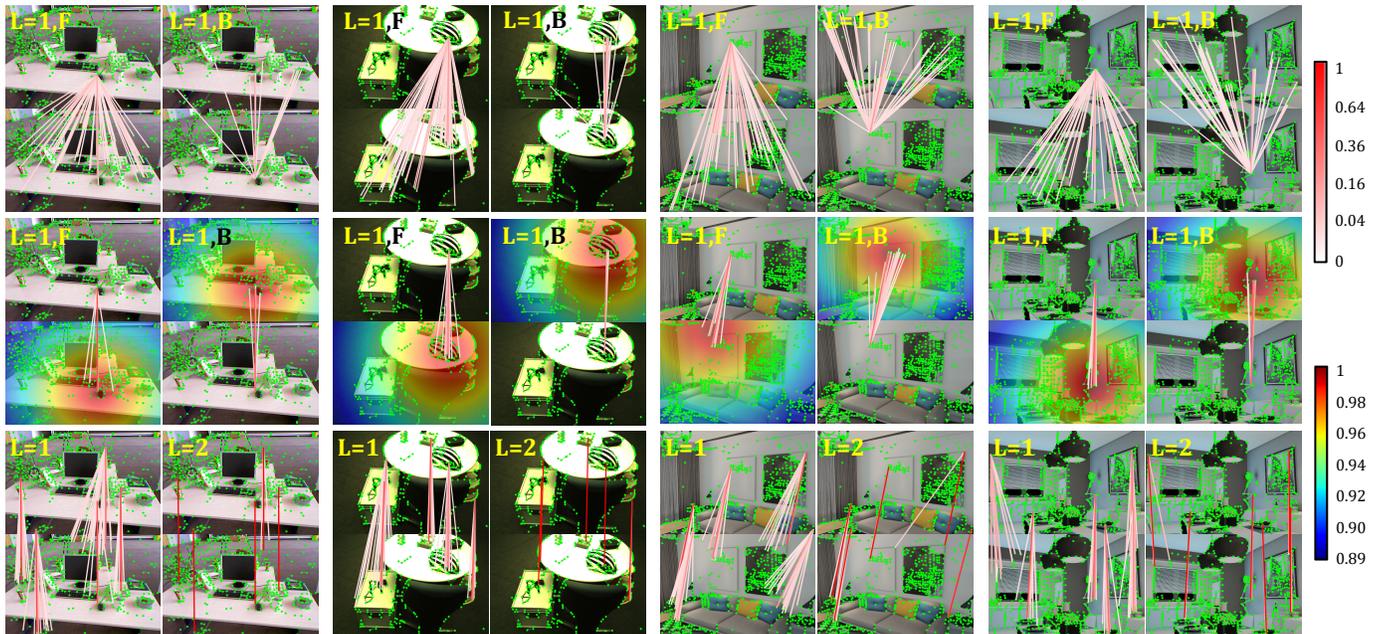
Fig. 8: Visualization of the spatial distribution prior and its effectiveness. The image pairs are sampled from the TUM-RGBD [34], ETH3D [35], and the InteriorNet [33] dataset. The cross-attention is drawn as lines across images, and the attentions lower than 0.01 are omitted for clarity. The "L=x" represents layer number "x", "F" and "B" denote the forward and backward cross-attention respectively. Top row: cross-attention of the first layer in SuperGlue [28]; Middle row: cross-attention of the first layer in the proposed model, and the spatial distribution priors are overlapped on images; Bottom row: cross-attention visualization of different keypoints in the proposed model. Notice that the prior assisted attentional GNN quickly can focus on the corresponding keypoints in the second attentional layer.
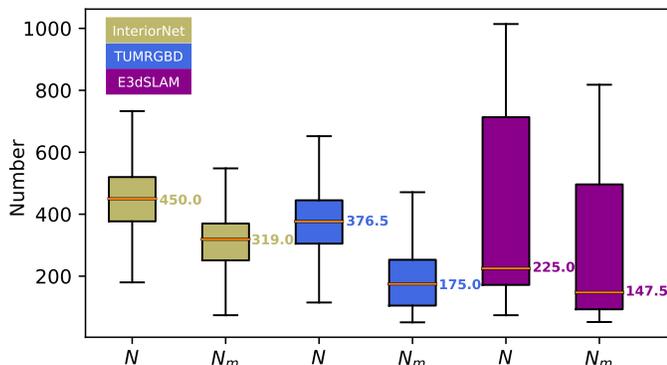


Fig. 9: The keypoints and matches distribution of the test split for InteriorNet [33], TUM-RGBD [34], and ETH3D [35] dataset. $N$ and $N_m$ denote the number of keypoints of images and the number of ground-truth matches of image pairs.

dynamic objects are excluded. This ends up with 37 training sequences and 12 testing sequences.

**HPatches dataset** [67] is a planar image pairs dataset with ground-truth homography matrix. It contains two subsets namely illumination and viewpoint, with 57 and 79 scenes respectively. This dataset is used to evaluate the generalization performance of the trained model.

**Oxford-Paris datatset** [68] has 6392 tourist pictures. We randomly selected 600 images and generated homography

image pairs with the following producers. All the sampled images are first cropped and resized to size $640 \times 480$. Then, a homography matrix is randomly sampled for each image in a way similar to [25]. At last, a corresponding image is generated with the homography transformation.

Since there is no motion prior for homography image pairs datasets, an initial homography matrix is obtained by adding noise to ground-truth homography transformation. As such, we can use the initial homography matrix to project the keypoints between images to attain the spatial distribution prior.

### B. Implementation details

To integrate the prior information to attentional GNN, we use single head attention rather than multi-head attention [63]. The position of keypoints is normalized by the image height and width. For direct prior integration, the $\sigma$ in Equation (2) is fixed to 0.1. And for probabilistic prior integration, the $\sigma$ is a trainable parameter, such that the network can determine the spatial region of each layer. In all the experiments, the ground-truth matches are obtained with Equation (16) with $th = 3$.

To train the proposed models, we sample image pairs from InteriorNet [33], TUM-RGBD [34], and ETH3D [35] dataset with a maximum of 300/300/200 image pairs, a maximum frame interval of 10/10/8, and a minimum overlap score of 0.3. The image pairs with less than 50 ground-truth matches are removed. An initial pose is integrated with the IMU measurements or accelerator readings between two sampled
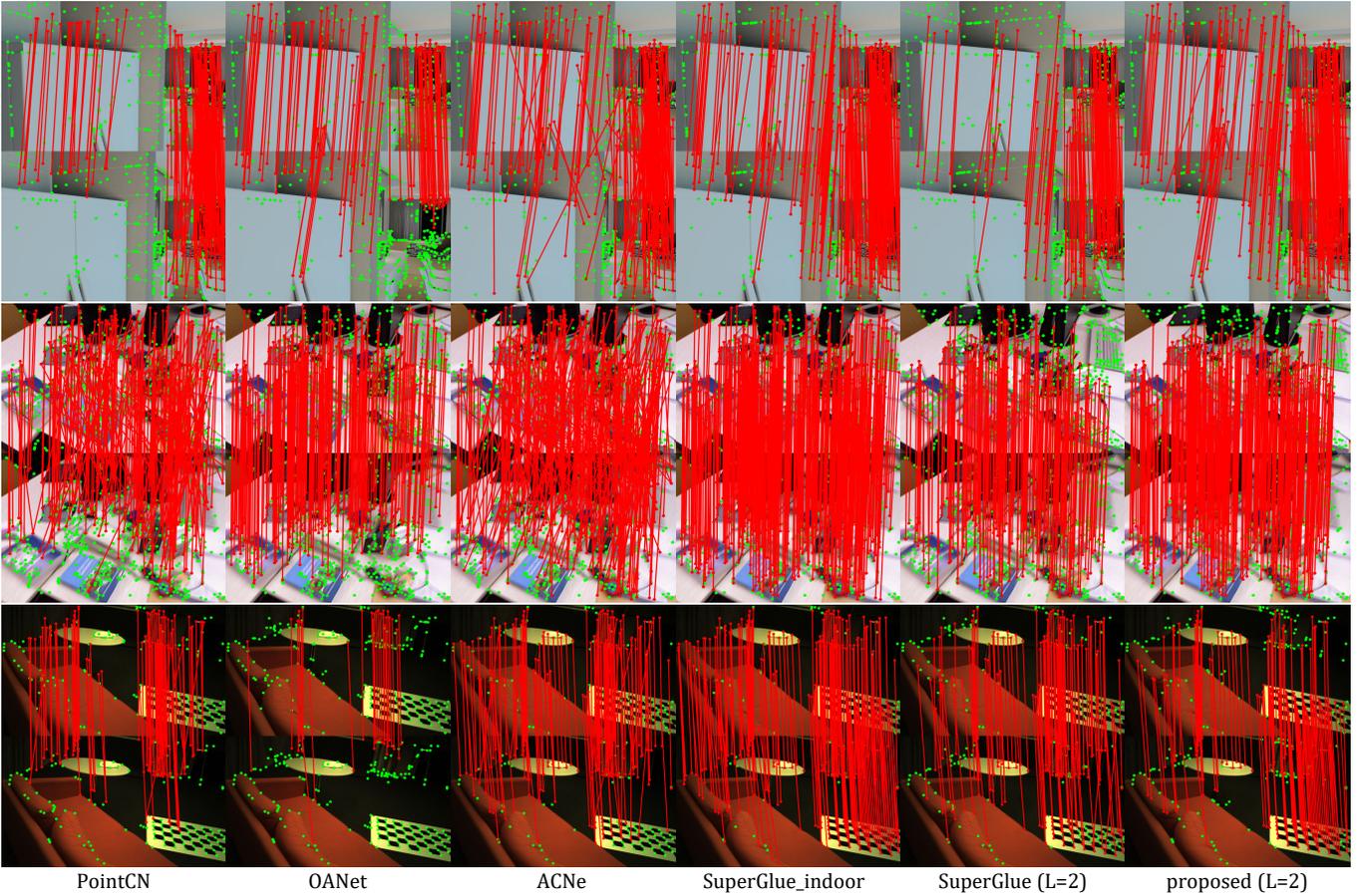
Fig. 10: Matching examples of different learned methods on InteriorNet [33] (top), TUM-RGBD [34] (middle), and ETH3D [35] (bottom) test set.
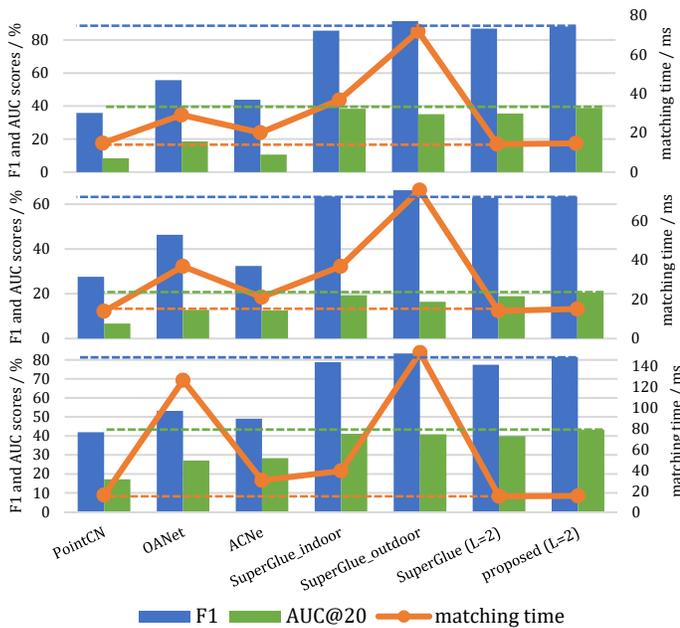


Fig. 11: The F1 score, AUC@20° and the matching time of different learned methods on InteriorNet [33] (top), TUM-RGBD [34] (middle), and ETH3D [35] (bottom) test set.

images [31], [32]. If the translation error between the initial pose and the ground-truth pose is greater than 0.1 meters or the rotation error is greater than 8 degrees, then this image pair is discarded. For those image pairs in TUM-RGBD dataset [34] without valid IMU measurements, we generated a synthetic pose by sampling translation and rotation from a zero-mean Gaussian distribution with $\sigma = 0.01\Delta t$.

During the training phase, we extracted 512 keypoints with SuperPoint [25] for each image. The non-maximum suppression radius is four pixels and keypoint threshold is 0.005. If there are no enough keypoints, additional keypoints were sampled from a uniform distribution to ensure that each image has 512 keypoints. To train the proposed models, we initialized them with the pre-trained position encoder and attentional GNN in SuperGlue [28]. And the proposed models were trained on InteriorNet [33], TUM-RGBD [34], and ETH3D [35] datasets together for a maximum of 300 epochs with adam optimizer [69] whose learning rate is $10^{-4}$ and batch size is 64. The validation loss was monitored, and if the validation loss was not reduced within 20 epochs, the training steps were early stopped. At last, the best models were selected and saved based on the validation loss.

While in the test and evaluation stages, the keypoints extraction configuration is the same as the training phase, except
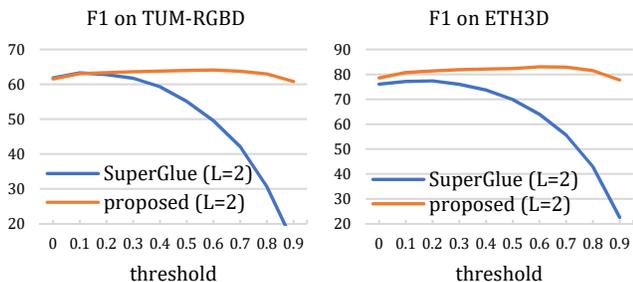
Fig. 12: The F1 scores with respective to the different matching confidence thresholds of SueprGlue with two attentional GNN layers and the proposed method on TUM-RGBD [34] and ETH3D [35] test sets.



Fig. 13: The running time distribution of the learned methods on ETH3D [35] test set. The outliers are not shown in this figure for clarity.

there is no limit on the number of keypoints. The distribution of keypoints and matches on the test set is shown in Fig. 9. The median numbers of keypoints are around 450, 377, and 225 for InteriorNet [33], TUM-RGBD [34], and ETH3D [35], respectively. The overall number of feature points ranges from 200 to 800, 100 to 700, and 80 to 1000 for the three datasets, respectively. Moreover, the number of ground truth matches of these three test sets is generally more than 100. The above conditions ensure that the image pairs of the test set have sufficient keypoints and matches, and can obtain accurate and comparable test results.

### C. Visualization of the spatial distribution prior

The visualization of the cross-attention and spatial distribution prior is shown in Fig. 6 and Fig. 8. Without the assistance of any prior information, the SuperGlue [28] has to model the attention from scratch. The model has to focus on almost all the keypoints in another image in the first few layers to find a potential match for one keypoint (Fig. 6 and the top row of Fig. 8). And the proposed spatial distribution prior of keypoints across images is shown in the middle row of Fig. 8. We can see that the spatial distribution prior gives a region of potential matches in another image. With the assistance of spatial distribution prior, the proposed model can focus on the potential keypoints in the first cross-attention layer (middle row of Fig. 8), and quickly focus on the correct matches in the second cross-attention layer (bottom row of Fig. 8).

### D. Comparison to related works

*1) Setups:* We compared the proposed method with the handcrafted and learned methods. The proposed method of two attentional GNN layers was trained with projection loss (will be discussed in Section IV-E) for comparisons. For the handcrafted methods, two variations of nearest neighbor (NN) matching, the mutual NN with PyTorch implementation, and the FLANN with OpenCV implementation when the test ratio is 0.7, were evaluated. For the learned methods, the PointCN [16], OANet [18], ACNe [19], and SuperGlue [28] were assessed. During testing, the official released code and model weights were used. Specifically, we used the model trained on GL3D-v2 for OANet [18]. For ACNe [19], the indoor model weights were used, and the keypoint matches with a combined
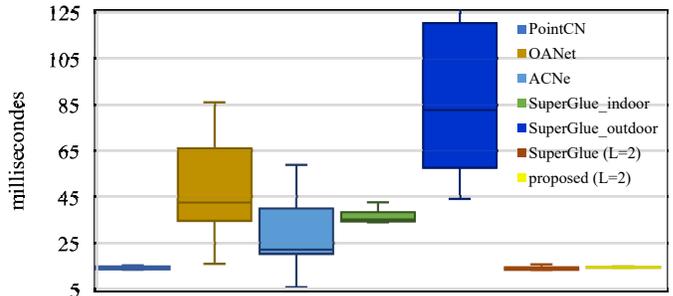
weight is greater than $10^{-7}$ were taken as correct matches. For SuperGlue [28], we assessed the official "indoor" and "outdoor" model weights, as well as the SueprGlue of two attentional GNN layers trained on our training set with the same training configures as for the proposed methods. The matching confidence thresholds of the assignment matrix of SuperGlue and the proposed models were both set to 0.2 as in [28], and we found this configuration gives the best trade-off between precision and recall. All the models were evaluated on a desktop platform with Intel Core i5-4590 and GeForce GTX TITAN X (PASCAL).

*2) Qualitative comparisons:* The matching results of different learned methods are shown in Fig. 10 and Fig. 11. As can be seen from these figures, the PointNet-like methods generate some false matches (results of PointCN [16] and ACNe [19] in Fig. 10) and obtain fewer matches (results of PointCN [16] and OANet [18] in Fig. 10) than the feature-based models. The possible reason is that the PointNet-like methods only take advantage of the position of keypoints and cannot model the appearance or other hidden features. The full SuperGlue [28] model, on the other hand, achieves the best matching results on the test images. Nevertheless, the full SuperGlue model consumes a lot of computation time (Fig. 11), and is not suitable for real-time SLAM systems. Simplifying SuperGlue by reducing the number of GNN layers can decrease its computational efforts, but it also degrades the matching performance (as shown in Fig. 11). With the assistance of prior information, the proposed method fills this performance decline without increasing the computation time (as shown in Fig. 11).

*3) Matching results:* The ground-truth matches are first obtained according to Equation (16) with $th = 3$. Then the matching precision $P_m$, recall $R_m$, and F1-score is given as

$$
\begin{aligned}
P_m &= TP/(TP + FP) \\
R_m &= TP/(TP + FN) \\
F1 &= 2P_m R_m/(P_m + R_m)
\end{aligned}
\tag{22}
$$

where $TP$, $FP$, and $FN$ denote the true positives, false positives, and false negatives respectively. The $TP$, $TP + FP$, and $TP + FN$ represent the number of correct matches, estimated matches, and ground-truth matches in the keypoints matching case respectively.

TABLE I: The matching performance of different methods on InteriorNet [33], TUM-RGBD [34], and ETH3D [35] test sets. $P_m$, $R_m$, and $F1$ denote the matching precision, recall and F1 score respectively. L=2 and L=3 denote there are 2 and 3 self- and cross- attentional GNN layers in the model respectively. The best and the second best are marked as **bold** and **blue**.

| | InteriorNet [33] | | | TUM-RGBD [34] | | | ETH3D [35] | | |
|---|---|---|---|---|---|---|---|---|---|
| | $P_m$ | $R_m$ | $F1$ | $P_m$ | $R_m$ | $F1$ | $P_m$ | $R_m$ | $F1$ |
| mNN | 81.49 | 93.43 | 86.93 | 52.05 | 84.13 | 62.58 | 68.69 | 87.08 | 76.05 |
| FLANN | 82.26 | 87.96 | 84.90 | 52.52 | 76.71 | 60.61 | 69.02 | 76.36 | 71.51 |
| PointCN [16] | 52.36 | 27.70 | 35.77 | 33.14 | 25.68 | 27.55 | 53.91 | 35.87 | 41.90 |
| OANet [18] | 78.74 | 44.62 | 55.63 | 53.86 | 46.71 | 46.33 | 67.91 | 48.16 | 53.13 |
| ACNe [19] | 48.16 | 40.41 | 43.83 | 30.86 | 36.78 | 32.38 | 50.16 | 49.26 | 49.01 |
| SuperGlue_indoor [28] | 78.89 | 94.06 | 85.65 | 52.08 | **87.03** | **63.37** | 68.93 | **94.00** | 78.76 |
| SuperGlue_outdoor [28] | **88.00** | 95.33 | 91.40 | **56.53** | 84.71 | 66.15 | **75.27** | 95.06 | 83.37 |
| SuperGlue (L=2) [28] | **84.70** | 89.43 | 86.83 | 55.10 | 77.01 | 62.79 | 71.21 | 85.90 | 77.39 |
| proposed (L=2) | 82.84 | **94.37** | **88.10** | 60.98 | **84.47** | 63.32 | **72.94** | 93.95 | **81.37** |

TABLE II: The pose estimation AUC under 5°, 10° and 20° of the learned methods on InteriorNet [33], TUM-RGBD [34] and ETH3D [35] test sets. The best and the second best are marked as **bold** and **blue**.

| | InteriorNet [33] | | | TUM-RGBD [34] | | | ETH3D [35] | | |
|---|---|---|---|---|---|---|---|---|---|
| | AUC@5 | AUC@10 | AUC@20 | AUC@5 | AUC@10 | AUC@20 | AUC@5 | AUC@10 | AUC@20 |
| PointCN [16] | 2.18 | 3.89 | 8.28 | 1.62 | 3.43 | 6.67 | 2.37 | 7.25 | 17.15 |
| OANet [18] | 3.48 | 8.39 | 18.52 | 1.43 | 4.26 | 12.75 | 7.14 | 14.93 | 27.04 |
| ACNe [19] | 4.69 | 6.46 | 10.52 | **4.49** | 7.29 | 12.46 | 8.24 | 15.26 | 28.29 |
| SuperGlue_indoor [28] | **13.72** | **24.01** | **38.26** | 3.19 | 9.12 | **19.29** | **11.03** | **23.40** | **41.04** |
| SuperGlue_outdoor [28] | **10.27** | 20.43 | 34.99 | 3.74 | 6.86 | 16.37 | 9.33 | 22.41 | 40.85 |
| SuperGlue (L=2) [28] | 9.86 | 19.36 | 35.37 | 3.45 | **9.42** | 18.90 | 6.97 | 22.98 | 39.94 |
| proposed (L=2) | 9.91 | **21.76** | **38.80** | **6.37** | **11.35** | **20.49** | **11.47** | **25.53** | **43.23** |

The detailed matching results are shown in Table I. It is noticeable that the mNN and FLANN both have good matching precision ($P_m$) and recall ($R_m$) on our test sets, which is only slightly lower than the SuperGlue [28] and the proposed method. While the PointNet-like models PointCN [16], OANet [18], and ACNe [19] do not perform well on our test sets. This could due to the fact that they follow the putative filtering strategy, so they require massive putative keypoint pairs (usually > 2000) for each image pair. However, we only extracted the most distinctive keypoints as in most SLAM systems, resulting in fewer keypoints in the image and fewer putative matches in the image pairs. The full SuperGlue [28] model with "outdoor" weights (nine attentional GNN layers) achieves the best matching F1 scores, with an F1 score of 91.40%, 66.15%, and 83.37% on InteriorNet [33], TUM-RGBD [34], and ETH3D [35] respectively. Nonetheless, the SuperGlue of two attentional GNN layers drop the matching performance by 4.57%, 3.36%, and 5.98% for the F1 score on InteriorNet, TUM-RGBD, and ETH3D test sets respectively. The proposed methods, which are assisted by the spatial distribution prior of keypoints, can regain the performance drop of reduction of GNN layers as in Table I. Specifically, the proposed method promotes the recall $R_m$ by 4.94%, 7.46%, and 8.05%, F1 scores by 1.27%, 0.53%, and 3.98% on InteriorNet, TUM-RGBD, and ETH3D test sets receptively. Furthermore, as shown in Fig. 12, without the assistance of the spatial distribution prior, SuperGlue of two attentional GNN layers fails to recover correct matches with higher matching confidence (lower F1 score than the proposed method). Some matching examples are shown in Fig. 10.

*4) Pose estimation accuracy:* The purpose of keypoints matching is to estimate the relative pose of two images. So the pose accuracy is also evaluated. Given the predicted keypoint matches of an image pair, the essential matrix is

TABLE III: The matching precision $P_m$, recall $R_m$, $F1$ scores, and homography accuracy $Acc_H$ on homographic datasets. HPatches_i and HPatches_v denote the illumination and viewpoint subset of HPatches [67]. The OxfordParis is the manually generated homography image pairs from OxfordParis dataset [68] as described in Section IV-B. The matching confidence of the model is 0.2.

| Dataset | model | $P_m$ | $R_m$ | $F1$ | $Acc_H$ |
|---|---|---|---|---|---|
| HPatches_i | SuperGlue_outdoor | 80.75 | 84.92 | 82.52 | 94.39 |
| | SuperGlue (L=2) | 61.82 | 60.92 | 60.77 | 94.74 |
| | proposed (L=2) | 62.42 | 67.89 | 64.73 | 92.63 |
| HPatches_v | SuperGlue_outdoor | 82.08 | 83.67 | 82.73 | 56.95 |
| | SuperGlue (L=2) | 55.64 | 51.23 | 52.13 | 45.76 |
| | proposed (L=2) | 61.01 | 63.73 | 62.18 | 50.17 |
| OxfordParis | SuperGlue_outdoor | 76.12 | 76.54 | 75.51 | 76.50 |
| | SuperGlue (L=2) | 52.49 | 42.77 | 45.58 | 55.33 |
| | proposed (L=2) | 52.26 | 52.73 | 51.42 | 61.83 |

obtained with OpenCV function `findEssentialMat`, and then the relative pose is recovered with `recoverPose`. As in previous works [16], [18], [28], we calculate the pose angular differences between ground truth and estimated pose, and report the area under curve (AUC) with a maximum error of threshold 5°, 10°, and 20°.

Five learned models are assessed in Table II. The Pointnet-like models PointCN [16], OANet [18], and ACNe [19] yield lower AUCs. The possible reason is that the number of keypoints is not sufficient for them to recover good matches. The visualization in Fig. 10 suggests they recall fewer matches and produces many false matches. The AUCs of the full SuperGlue [28] model, are generally higher than the PointNet-like models. And the SuperGlue with "indoor" weights obtains the best AUC performance among previously learned methods, yields 38.26%, 19.29%, and 41.04% of AUC@20° on InteriorNet [33], TUM-RGBD [34], and ETH3D [35] test sets respectively.
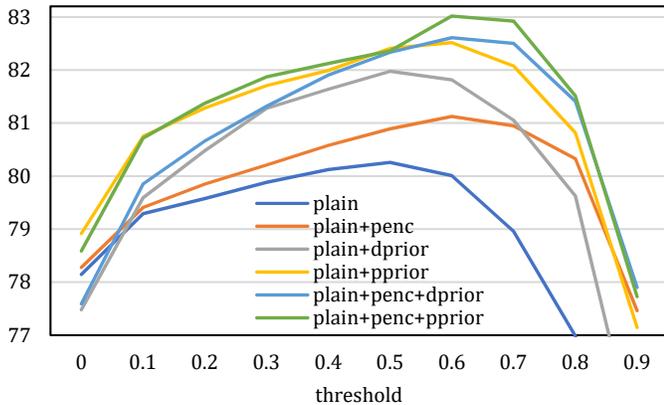
Fig. 14: The matching F1 scores with respect to matching confidence of different network configurations on ETH3D dataset [35].
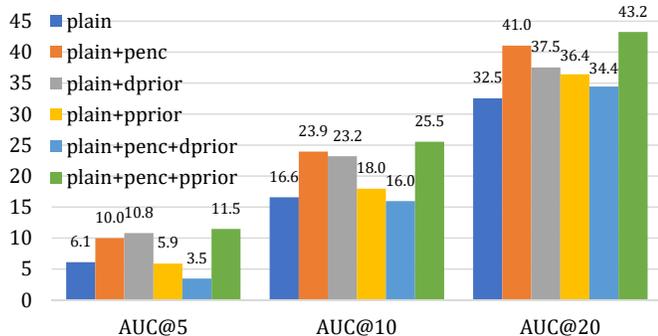


Fig. 15: The pose estimation AUC of different network configurations on ETH3D dataset [35].

As the SuperGlue of two attentional GNN layers was trained on our training set, it also gives AUCs on par with the full SuperGlue model, with 35.37%, 18.09%, and 39.94% of AUC@20° on InteriorNet, TUM-RGBD, and ETH3D test sets respectively. Due to the assistance of spatial distribution prior, the proposed method produces higher AUCs than SuperGlue of two attentional GNN layers by 3.43%, 1.59%, and 3.29% of AUC@20° on InteriorNet, TUM-RGBD, and ETH3D test sets respectively, and outperforms the other learned methods.

*5) Running time:* We tested all the learned models with the same agenda to measure the matching time of an image pair. The results are shown in Fig. 13. Among the learned matching networks, the model sizes of ACNe [19], OANet [18], and SuperGlue [28] are larger, which results in longer running times. Since the running time is influenced by the number of points, the larger model in turn amplifies the dispersion of the running time. While the PointCN [16], SuperGlue [28] of two GNN layers and the proposed method have the shortest matching time of about 15 milliseconds. Considering the matching metrics and pose estimation accuracy of the proposed method are much better than those of PointCN [16] and SuperGlue [28] of two GNN layers, the effectiveness of the proposed prior integration and projection loss is verified.

*6) Generalization to homographic image pairs:* To demonstrate the generalization of the proposed method, experiments

TABLE IV: Ablation studies on ETH3D dataset [35]. The matching precision $P_m$, recall $R_m$ and $F1$ score are computed when the matching confidence is 0.2. The plain network contains two attentional GNN layers and an optimization layer, the "penc" is the MLP position encoder in Equation (3), and the "dprior" and "pprior" denote the direct and probabilistic prior integration method respectively.

| plain | penc | dprior | pprior | $P_m$ | $R_m$ | $F1$ |
|---|---|---|---|---|---|---|
| ✓ | | | | 70.71 | 92.71 | 79.58 |
| ✓ | ✓ | | | 70.73 | 93.35 | 79.85 |
| ✓ | | ✓ | | 72.48 | 92.04 | 80.48 |
| ✓ | | | ✓ | **73.29** | 93.09 | 81.28 |
| ✓ | ✓ | ✓ | | 72.13 | 93.25 | 80.66 |
| ✓ | ✓ | | ✓ | 72.94 | **93.95** | **81.37** |

are also carried out on homography image pairs. Table III illustrates the matching results of full SuperGlue [28] (with "outdoor" weights), SuperGlue of two attentional GNN layers, and the proposed method. Note that the latter two models were trained only on the RGBD datasets, no homography image pairs are included in the training process. As in Table III, the full SuperGlue [28] with "outdoor" weights achieves the best matching performance, giving F1 scores of 82.52%, 82.73%, and 75.51% on HPatches illumination, viewpoint, and OxfordParis, respectively. However, F1 scores of the two GNN layers version of SuperGlue drops to 60.77%, 51.12%, and 45.58% on these datasets. While the F1 scores of the proposed method significantly outperform the two GNN layers version of SuperGlue, which are 64.73%, 62.18%, and 51.42% respectively. Moreover, as the viewpoint differences on HPatches viewpoint and OxfordParis are larger, the assistance of the spatial distribution prior is more significant. Therefore, the improvement of the F1 score is more obvious, which increases by 10.05% and 5.84% respectively. In addition to the matching performance, the homography estimation accuracy $Acc_H$ is also evaluated. The mean reprojection error of the four corners of the image is first computed, and then the homography estimation accuracy is defined by the accuracy of the corner error under a threshold of three pixels. As shown in Table III, the homography accuracy $Acc_H$ of the proposed method is slightly lower than SuperGlue [28] on the Hpatches illumination dataset. The proposed method, on the other hand, significantly outperforms SueprGlue of two attentional GNN layers on HPatches viewpoint and OxfordParis datasets. The reason for this phenomenon is that for the Hpatches illumination dataset, the homography matrix is identity, and the noise added during generating the priors interferes with the matching process. For HPatches viewpoint and OxfordParis datasets, there are significant viewpoint differences in a pair of images. Thus the noise is trivial and the prior plays a vital role in improving the matching performance.

### E. Ablation studies

To investigate the effectiveness of each model part, this subsection studies the position encoder, the direct and probabilistic prior integration, the matching loss and projection loss, as well as the number of attentional GNN layers to the matching performance.

*1) Position encoder:* The MLP in Equation (3) makes use of position by embedding it to feature space, while the proposed prior integration utilizes the position by propagating contextual keypoint features through the attentional GNN. The two methods of position utilization are compared in this section. As shown in Table IV, the position encoder (penc) promotes the F1 score on the plain network by 0.27%. The proposed direct prior (dprior) and probabilistic prior integration boost the F1 score by 0.9% and 2.7% respectively. The same results can also be found in Fig. 14, where the F1 score is plotted with respect to the matching confidence. It could be because the embedding of positional encoding changes the feature distribution of descriptors, and it is difficult for the network to learn such embedding. While the proposed prior integration does not change the distribution of any features, but only aggregates the nearby features in the original feature space, so the proposed prior integration methods significantly improve the matching performance.

*2) Direct and probabilistic prior integration:* As shown in Table IV, compared to the plain network, after integrating the direct and probabilistic prior into attentional GNN, the precision is improved by 1.77% and 2.58% respectively, and the F1 score is improved 0.9% and 1.7% respectively. With the keypoint position encoder module, the total improvement of the F1 score is 1.08% and 1.79% for the direct and probabilistic prior integration respectively. The same promotion can also be found in Fig. 14. These primary matching results indicate the probabilistic prior integration is more efficient than direct prior integration.

In terms of the accuracy of the estimated pose, inconsistent results are observed for different approaches with the matching results as in Fig. 15. The improvements of AUCs by the direct prior integration are larger than those of probabilistic prior integration, and the improvements of both prior integrations are lower than the positional encoder except that of AUC under 5°. Another interesting result is that the combination of positional encoding and direct prior integration actually deteriorates the pose estimation accuracy. The possible reason could be that the positional embedding feature could not be processed correctly by the direct prior integration strategy. While the combination of positional encoding and probabilistic prior integration can promote each other, and improve the pose accuracy significantly.

*3) Matching and projection loss:* We evaluated the matching loss of Equation (17) and the projection loss of Equation (21) by training the plain network (with two attentional GNN layers and an optimization layer) on the same train dataset with a fixed random seed. As expected, the projection loss recalls more matches than matching loss due to its relaxation of the matching threshold, thus yields higher $F1$ scores as shown in Fig. 16(a). Thus, we use the projection loss for all experiments.

*4) Number of GNN layers:* Since the introduction of the spatial distribution prior can assist the model to utilize the contextual features, so that the number of attentional GNN layers is reduced for efficiency. To verify the impact of the number of GNN layers on the performance, we trained and assessed the models of one to five attention GNN layers. The
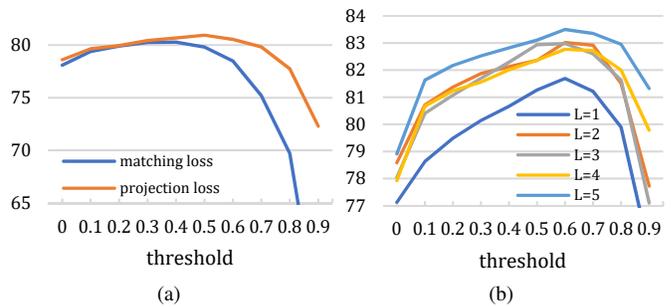


Fig. 16: Ablation studies of loss functions and the number of attentional GNN layers. (a): comparison of the matching loss and projection loss of F1 scores with respect to the matching confidence. (b): the matching F1 scores of the proposed model with different numbers of attentional GNN layers.

TABLE V: Matching performance for the different number of GNN layers on ETH3D [35] test set with matching confidence of 0.2.

|  | L=1 | L=2 | L=3 | L=4 | L=5 |
|---|---|---|---|---|---|
| $Pm$ | 70.45 | 72.94 | 72.63 | 73.08 | 74.17 |
| $Rm$ | 92.96 | 93.95 | 93.70 | 93.34 | 94.03 |
| $F1$ | 79.49 | 81.37 | 81.07 | 81.24 | 82.17 |
| AUC@10 | 24.03 | 25.53 | 18.92 | 24.03 | 21.85 |
| time (ms) | 12.09 | 15.36 | 18.57 | 22.36 | 26.02 |

matching results are shown in Fig. 16(b). Clearly, the model with more GNN layers has a better matching performance. However, when the number of GNN layers is larger than one, there is no significant performance improvement, but the running time is greatly increased. As shown in Table V, the pose estimation accuracy even decreases with the increase of the number of GNN layers, as a deeper model is harder to train. To balance the performance and efficiency, we use two attentional GNN layers.

## V. CONCLUSIONS

In this paper, motion prior from other sensors such as the IMU is adopted to obtain the spatial distribution prior of keypoints, which is exploited to streamline the attentional keypoint matching network. Specifically, the spatial distribution prior is naturally integrated into the attentional GNN network with the probabilistic perspective of attention. Thus, the number of GNN layers can be reduced, and the running time is decreased to about 15ms while keeping the matching performance. Besides, a loss using the pixel projection errors is proposed to train the network to achieve better matching performance. The experiments on SLAM datasets InteriorNet, TUM-RGBD, and ETH3D validate the effectiveness and efficiency of the proposed method. A similar idea can be adopted to study other image processing problems such as motion de-blurring, visual tracking for autonomous systems. All these problems will be studied in our future research.

## ACKNOWLEDGMENT

REFERENCES

[1] C. G. Harris and M. Stephens, "A combined corner and edge detector." in *Alvey vision conference*, vol. 15. Citeseer, 1988, pp. 10–5244.

[2] J. Shi, "Good features to track," in *1994 Proceedings of IEEE conference on computer vision and pattern recognition*. IEEE, 1994, pp. 593–600.

[3] S. Alkaabi and F. Deravi, "Candidate pruning for fast corner detection," *Electronics Letters*, vol. 40, no. 1, pp. 18–19, 2004.

[4] A. Barroso-Laguna, E. Riba, D. Ponsa, and K. Mikolajczyk, "Key. net: Keypoint detection by handcrafted and learned cnn filters," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 5836–5844.

[5] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.

[6] C. Ding and Z. Ma, "Multi-Camera Color Correction via Hybrid Histogram Matching," *IEEE Transactions on Circuits and Systems for Video Technology*, Dec. 2020, doi: 10.1109/TCSVT.2020.3038484.

[7] Q. Wang, W. Chen, X. Wu, and Z. Li, "Detail-enhanced multi-scale exposure fusion in YUV color space," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 8, pp. 2418–2429, Aug. 2020.

[8] C. Zheng, Z. Li, Y. Yang, and S. Wu, "Single Image Brightening via Multi-Scale Exposure Fusion with Hybrid Learning," *IEEE Transactions on Circuits and Systems for Video Technology*, Jun. 2020, doi: 10.1109/TCSVT.2020.3009235.

[9] J. Ma, J. Zhao, J. Jiang, H. Zhou, and X. Guo, "Locality Preserving Matching," *International Journal of Computer Vision*, vol. 127, no. 5, pp. 512–531, May 2019.

[10] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.

[11] O. Pele and M. Werman, "A linear time histogram metric for improved sift matching," in *European conference on computer vision*. Springer, 2008, pp. 495–508.

[12] X. Li and Z. Hu, "Rejecting mismatches by correspondence function," *International Journal of Computer Vision*, vol. 89, no. 1, pp. 1–17, 2010.

[13] Y. Lipman, S. Yagev, R. Poranne, D. W. Jacobs, and R. Basri, "Feature matching with bounded distortion," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 3, pp. 1–14, 2014.

[14] J. Ma, J. Zhao, J. Tian, A. L. Yuille, and Z. Tu, "Robust Point Matching via Vector Field Consensus," *IEEE Transactions on Image Processing*, vol. 23, no. 4, pp. 1706–1721, Apr. 2014.

[15] W.-Y. Lin, F. Wang, M.-M. Cheng, S.-K. Yeung, P. H. Torr, M. N. Do, and J. Lu, "CODE: Coherence based decision boundaries for feature correspondence," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 1, pp. 34–47, 2017.

[16] K. M. Yi, E. Trulls, Y. Ono, V. Lepetit, M. Salzmann, and P. Fua, "Learning to Find Good Correspondences," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT, USA: IEEE, Jun. 2018, pp. 2666–2674.

[17] C. Zhao, Z. Cao, C. Li, X. Li, and J. Yang, "NM-Net: Mining Reliable Neighbors for Robust Feature Correspondences," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[18] J. Zhang, D. Sun, Z. Luo, A. Yao, L. Zhou, T. Shen, Y. Chen, H. Liao, and L. Quan, "Learning Two-View Correspondences and Geometry Using Order-Aware Network," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. Seoul, Korea (South): IEEE, Oct. 2019, pp. 5844–5853.

[19] W. Sun, W. Jiang, E. Trulls, A. Tagliasacchi, and K. M. Yi, "ACNe: Attentive Context Normalization for Robust Permutation-Equivariant Learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 286–11 295.

[20] H. Bay, T. Tuytelaars, and L. Van Gool, "Surf: Speeded up robust features," in *European conference on computer vision*. Springer, 2006, pp. 404–417.

[21] K. M. Yi, E. Trulls, V. Lepetit, and P. Fua, "LIFT: Learned Invariant Feature Transform," in *European Conference on Computer Vision*, vol. 9910. Cham: Springer, 2016, pp. 467–483.

[22] X. Han, T. Leung, Y. Jia, R. Sukthankar, and A. C. Berg, "Matchnet: Unifying feature and metric learning for patch-based matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3279–3286.

[23] A. Mishchuk, D. Mishkin, F. Radenovic, and J. Matas, "Working hard to know your neighbor's margins: Local descriptor learning loss," in *Advances in Neural Information Processing Systems*, Jan. 2018.

[24] Y. Tian, X. Yu, B. Fan, F. Wu, H. Heijnen, and V. Balntas, "SOSNet: Second Order Similarity Regularization for Local Descriptor Learning," in *Conference on Computer Vision and Pattern Recognition*, Dec. 2019.

[25] 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.

[26] M. Dusmanu, I. Rocco, T. Pajdla, M. Pollefeys, J. Sivic, A. Torii, and T. Sattler, "D2-Net: A Trainable CNN for Joint Description and Detection of Local Features," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, Jun. 2019, pp. 8084–8093.

[27] J. Revaud, P. Weinzaepfel, C. D. Souza, N. Pion, G. Csurka, Y. Cabon, and M. Humenberger, "R2D2: Repeatable and Reliable Detector and Descriptor," in *NeurIPS*, 2019, p. 12.

[28] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "SuperGlue: Learning Feature Matching with Graph Neural Networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Mar. 2020, pp. 4938–4947.

[29] 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.

[30] J. Zhang and S. Singh, "LOAM: Lidar Odometry and Mapping in Real-time." Robotics: Science and Systems Foundation, Jul. 2014.

[31] J. Henawy, Z. Li, W. Y. Yau, G. Seet, and K. W. Wan, "Accurate IMU Preintegration Using Switched Linear Systems For Autonomous Systems," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. Auckland, New Zealand: IEEE, Oct. 2019, pp. 3839–3844.

[32] J. Henawy, Z. Li, W. Y. Yau, and G. Seet, "Accurate IMU Factor Using Switched Linear Systems For VIO," *IEEE Transactions on Industrial Electronics*, vol. 68, no. 9, pp. 1–10, Sep. 2021.

[33] W. Li, S. Saeedi, J. McCormac, R. Clark, D. Tzoumanikas, Q. Ye, Y. Huang, R. Tang, and S. Leutenegger, "InteriorNet: Mega-scale Multi-sensor Photo-realistic Indoor Scenes Dataset," in *British Machine Vision Conference (BMVC)*, 2018.

[34] 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*, Oct. 2012, pp. 573–580.

[35] T. Schops, T. Sattler, and M. Pollefeys, "BAD SLAM: Bundle Adjusted Direct RGB-D SLAM," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, Jun. 2019, pp. 134–144.

[36] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.

[37] V. Balntas, E. Riba, D. Ponsa, and K. Mikolajczyk, "Learning local feature descriptors with triplets and shallow convolutional neural networks." in *BMVC*, vol. 1, 2016, p. 3.

[38] Y. Tian, B. Fan, and F. Wu, "L2-Net: Deep Learning of Discriminative Patch Descriptor in Euclidean Space," in *2017 IEEE Conference on Computer Vision and Pattern Recognition*. Honolulu, HI: IEEE, Jul. 2017, pp. 6128–6136.

[39] Y. Ono, E. Trulls, P. Fua, and K. M. Yi, "LF-Net: Learning Local Features from Images," in *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., 2018, pp. 6234–6244.

[40] Y. Verdie, K. Yi, P. Fua, and V. Lepetit, "Tilde: A temporally invariant learned detector," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5279–5288.

[41] N. Savinov, A. Seki, L. Ladicky, T. Sattler, and M. Pollefeys, "Quad-networks: unsupervised learning to rank for interest point detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1822–1830.

[42] X. Guo and X. Cao, "Good match exploration using triangle constraint," *Pattern Recognition Letters*, vol. 33, no. 7, pp. 872–881, 2012.

[43] J. Ma, J. Jiang, H. Zhou, J. Zhao, and X. Guo, "Guided locality preserving feature matching for remote sensing image registration," *IEEE transactions on geoscience and remote sensing*, vol. 56, no. 8, pp. 4435–4447, 2018.

[44] J. Bian, W.-Y. Lin, Y. Matsushita, S.-K. Yeung, T.-D. Nguyen, and M.-M. Cheng, "Gms: Grid-based motion statistics for fast, ultra-robust feature correspondence," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4181–4190.

[45] Y.-T. Hu, Y.-Y. Lin, H.-Y. Chen, K.-J. Hsu, and B.-Y. Chen, "Matching images with multiple descriptors: An unsupervised approach for locally adaptive descriptor selection," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 5995–6010, 2015.

[46] J. Maier, M. Humenberger, M. Murschitz, O. Zendel, and M. Vincze, "Guided matching based on statistical optical flow for fast and robust correspondence analysis," in *European Conference on Computer Vision*. Springer, 2016, pp. 101–117.

[47] Y. Liu, L. De Dominicis, B. Wei, L. Chen, and R. R. Martin, "Regularization based iterative point match weighting for accurate rigid transformation estimation," *IEEE transactions on visualization and computer graphics*, vol. 21, no. 9, pp. 1058–1071, 2015.

[48] L. Torresani, V. Kolmogorov, and C. Rother, "Feature Correspondence via Graph Matching: Models and Global Optimization," in *European Conference on Computer Vision*, 2008.

[49] S. Liu, H. Wang, Y. Wei, and C. Pan, "Bb-homography: Joint binary features and bipartite graph matching for homography estimation," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 25, no. 2, pp. 239–250, 2015.

[50] R. Zhang and W. Wang, "Second- and high-order graph matching for correspondence problems," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 10, pp. 2978–2992, 2018.

[51] Y. F. Yu, G. Xu, K. K. Huang, H. Zhu, L. Chen, and H. Wang, "Dual calibration mechanism based l2,p-norm for graph matching," *IEEE Transactions on Circuits and Systems for Video Technology*, 2020, doi:10.1109/TCSVT.2020.3023781.

[52] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[53] N. Ufer and B. Ommer, "Deep semantic feature matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6914–6923.

[54] W. Yu, X. Sun, K. Yang, Y. Rui, and H. Yao, "Hierarchical semantic image matching using CNN feature pyramid," *Computer Vision and Image Understanding*, vol. 169, pp. 40–51, 2018.

[55] S. Khan, M. Nawaz, X. Guoxia, and H. Yan, "Image correspondence with cur decomposition-based graph completion and matching," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 9, pp. 3054–3067, 2020.

[56] F. Kou, Z. Li, C. Wen, and W. Chen, "Multi-scale exposure fusion via gradient domain guided image filtering," in *2017 IEEE International Conference on Multimedia and Expo (ICME)*, 2017, pp. 1105–1110.

[57] J. Zheng, Z. Li, Z. Zhu, S. Wu, and S. Rahardja, "Hybrid patching for a sequence of differently exposed images with moving objects," *IEEE Transactions on Image Processing*, vol. 22, no. 12, pp. 5190–5201, 2013.

[58] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.

[59] M. Cuturi, "Sinkhorn Distances: Lightspeed Computation of Optimal Transport," in *Advances in Neural Information Processing Systems 26*. Curran Associates, Inc., 2013, pp. 2292–2300.

[60] R. Mur-Artal and J. D. Tardós, "Visual-inertial monocular SLAM with map reuse," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 796–803, 2017.

[61] T. Qin, P. Li, and S. Shen, "VINS-Mono: A Robust and Versatile Monocular Visual-Inertial State Estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, Aug. 2018.

[62] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-Manifold Preintegration for Real-Time Visual–Inertial Odometry," *IEEE Transactions on Robotics*, vol. 33, no. 1, pp. 1–21, 2017.

[63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is All you Need," in *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017, pp. 5998–6008.

[64] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7794–7803.

[65] G. Peyré and M. Cuturi, "Computational Optimal Transport," *arXiv:1803.00567 [stat]*, Mar. 2020.

[66] J. Munkres, "Algorithms for the assignment and transportation problems," *Journal of the society for industrial and applied mathematics*, vol. 5, no. 1, pp. 32–38, 1957.

[67] 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.

[68] F. Radenović, A. Iscen, G. Tolias, Y. Avrithis, and O. Chum, "Revisiting oxford and paris: Large-scale image retrieval benchmarking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5706–5715.

[69] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," in *ICLR*, 2015.