

IMP: Iterative Matching and Pose Estimation with Adaptive Pooling

Fei Xue Ignas Budvytis Roberto Cipolla
 University of Cambridge
 {fx221, ib255, rc10001}@cam.ac.uk

Abstract

Previous methods solve feature matching and pose estimation using a two-stage process by first finding matches and then estimating the pose. As they ignore the geometric relationships between the two tasks, they focus on either improving the quality of matches or filtering potential outliers, leading to limited efficiency or accuracy. In contrast, we propose an *iterative matching and pose estimation framework (IMP)* leveraging the geometric connections between the two tasks: a few good matches are enough for a roughly accurate pose estimation; a roughly accurate pose can be used to guide the matching by providing geometric constraints. To this end, we implement a geometry-aware recurrent attention-based module which jointly outputs sparse matches and camera poses. Specifically, for each iteration, we first implicitly embed geometric information into the module via a pose-consistency loss, allowing it to predict geometry-aware matches progressively. Second, we introduce an efficient IMP, called *EIMP*, to dynamically discard keypoints without potential matches, avoiding redundant updating and significantly reducing the quadratic time complexity of attention computation in transformers. Experiments on YFCC100m, Scannet, and Aachen Day-Night datasets demonstrate that the proposed method outperforms previous approaches in terms of accuracy and efficiency. Code is available at <https://github.com/feixue94/imp-release>

1. Introduction

Feature matching and relative pose estimation are two fundamental tasks in computer vision and especially important to visual localization and 3D reconstruction. Traditionally, the two tasks are performed in two stages separately by first finding correspondences between keypoints extracted from two images with nearest neighbor (NN) matching and then estimating the relative pose from predicted matches with robust estimators, *e.g.* RANSAC [7, 8, 21, 34]. This pipeline has been the de-facto standard framework for decades [6]. However, due to repetitive textures/structures,

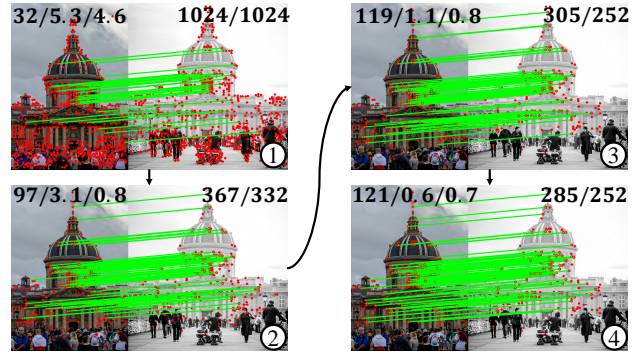


Figure 1. **Process of iterative matching and pose estimation.** For each image pair, we report the number inliers/rotation/translation errors (top-left) and retained keypoints in left/right images (top-right) at iterations from 1 to 4. In the iterative process, our method finds more inliers spanning almost the whole image, estimates increasingly precise pose and discards keypoints without true matches gradually.

changing appearances and viewpoint variations, matches given by NN often contain a large number of outliers, leading to poor pose accuracy [38, 39]. To mitigate this problem, some works [9, 11, 17, 29, 42, 53, 54, 57] filter potential outliers of predicted matches with neural networks to improve the pose accuracy. Although they report better results, their performance is limited by the quality of initial matches and require extra time for filtering at test time. Alternatively, advanced matchers such as SuperGlue [38] enhance the matching quality directly by using global information from all keypoints via transformers [47] with a fixed number (*e.g.* 9) of iterations. These methods have obtained remarkable performance. Yet, their quadratic time complexity for the attention computation degrades the efficiency in real applications. Some following works [13, 40, 43] explore more efficient variations, they run faster but are significantly less accurate (see Table 1 and 3).

In this paper, we aim to introduce an efficient and accurate framework for iterative matching and pose estimation. Our approach is built upon the following observations: (1) a few well distributed matches (*e.g.* 5) could give a roughly accurate pose (*e.g.* essential matrix); (2) in turn, a roughly

accurate pose could provide strong geometric constraints (*e.g.* epipolar line) to find more accurate matches at low cost; (3) the pose also reveals which keypoints have potential correspondences, preventing redundant operations. Based on the geometric connections of the two tasks, we propose an **iterative matching and pose estimation framework (IMP)**, to perform matching and pose estimation iteratively as opposed to in two separate stages. Specifically, we progressively augment descriptors with self and cross attention as [13,38,40,43], find matches and estimate the relative pose. As descriptors get gradually more discriminative, more correct matches can be found, leading to increasingly more precise pose, as shown in Fig. 1. However, due to the noise [28] and degeneration (*e.g.* co-planar keypoints) [14], not all inliers could give a good pose [5,19]. In addition to the classification loss mainly used by prior methods [13,38], we apply a pose-consistency loss [54] to the matching process, enabling the model to find matches which are not only accurate but also able to give a good pose.

Moreover, in order to avoid redundant operations on uninformative keypoints, we employ a sampling strategy by combining the matching and attention scores of keypoints and the uncertainty of predicted poses to adaptively remove useful keypoints, as shown in Fig. 1. Compared with prior sampling approaches [20,46] based mainly on attention scores, our adaptive strategy overcomes the over-sampling problem effectively. Our framework reduces the time cost from two aspects. First, in contrast to adopting a fixed number of iterations for all cases [13,38,40], it runs fewer iterations for easy cases with few viewpoint or appearance changes and more for challenging cases. Second, it reduces the cost of each iteration, significantly reducing the quadratic time complexity of attention computation. We also show that discarding potential outliers increases not only efficiency but also accuracy (see Sec. 5). The efficient version of IMP is called EIMP. Ours contributions are as follows:

- We propose to perform geometry-aware matching and pose estimation iteratively, allowing the two tasks to boost each other in an iterative manner.
- We adopt a robust sampling strategy to adaptively discard redundant keypoints in the iteration process, significantly decreasing the time complexity.
- We apply the pose uncertainty to the sampling strategy, which further improves the accuracy matching and pose estimation.

Our experiments on relative pose estimation and large-scale localization tasks demonstrate that our method outperforms previous competitors and is more efficient. We organize the rest of the paper as follows. In Sec. 2, we discuss related works. In Sec. 3, we give a detailed description

of our method. We test the performance of our model in Sec. 5 and conclude the paper in Sec. 6.

2. Related works

In this section, we discuss related work on local feature matching, efficient attention, and outlier filtering.

Local feature matching. Traditionally, NN matching and its variants, *e.g.*, mutual NN (MNN) and NN with ratio test (NN-RT) [30] are widely used for finding correspondences between two sets of keypoints [3, 18, 30, 36, 50]. Since they perform point-wise matching, they are fast but not robust to large viewpoint and appearance changes, *e.g.* illumination and season variations. Recently, SuperGlue [38] applies graph connections [51, 52] to keypoints and adopts with transformers [47] with self and cross attention to embed spatial information and achieves remarkable performance. Despite its excellent accuracy, its limitations are twofold. First, the complexity of attention is quadratic to the number of keypoints, degrading the efficiency in real applications. Second, it adopts a fixed number (9) of layers for message propagation for all input cases. This causes extra time cost especially for simple cases with few viewpoint and appearance changes, which require much fewer number of iterations to find potential inliers. Therefore, running a fixed number iterations for all cases leads to redundancy.

Some variants [13,40,43] try to solve the first problem by using a fixed number of seeded keypoints for message propagation [13], performing cluster-wise matching [40] or aggregating local and global information separately [43]. In spite of their high efficiency, they lose significant accuracy. Besides, like SuperGlue, they also update keypoints without correspondences redundantly in each iteration and adopt a fixed number of iterations for all cases constantly. Unlike these approaches, our model utilizes predicted poses to prevent extra iterations for easy cases. Additionally, as our method removes keypoints without true matches adaptively to reduce time complexity in each iteration, both efficiency and accuracy are guaranteed.

Efficient attention. Many works are proposed to mitigate the quadratic time complexity of the attention mechanism [47]. They reduce the complexity by learning a linear projection function [23, 48], a token selection classifier [20, 25, 35] or shared attention [12], to name a few. We refer the reader a comprehensive survey [44] for more details. Since most of these methods are designed to extract high-level features for downstream tasks *e.g.* image recognition [1, 10, 12, 20, 23, 46] usually with a fixed number of tokens, directly transferring these techniques to the matching task which has dynamic numbers of keypoints as input could cause performance loss. Instead, we utilize the geometric properties of the matching task to adaptively discard redundant keypoints. Our strategy depends mainly on the number of inliers in the input and is completely adaptive.

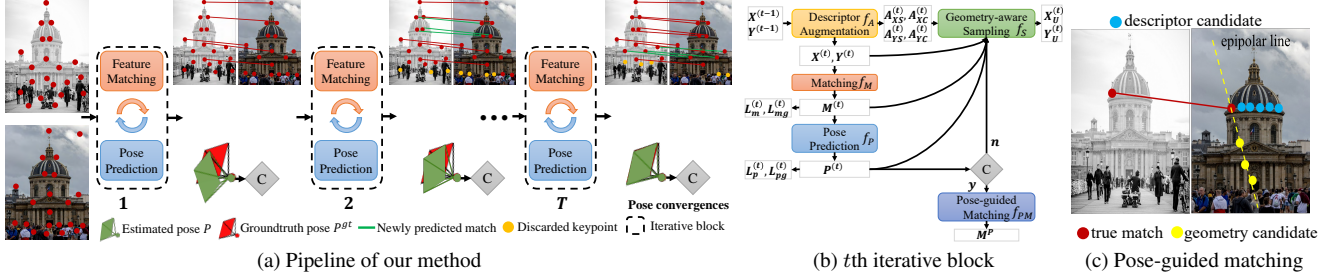


Figure 2. **Pipeline of our method.** In the iteration process, descriptors are gradually augmented by our recurrent attention-based module. Augmented descriptors are then used to compute matches $\mathcal{M}^{(t)}$, which are further used for estimating the pose $P^{(t)}$. As descriptors become more discriminative, more correct matches can be found, leading to more precise poses. The pose is utilized to provide geometric guidance to find more matches and discard redundant keypoints with geometry-aware sampling f_S , making next iterations faster. If the pose converges, iteration stops. Our recurrent attention-based module at t th iteration (b) and pose-guided matching (c) are also visualized.

Outlier filtering. Ratio test [30] is often used to remove matches with high uncertainties. Recently, CNNs are used to leverage the spatial context of input keypoints to filter potential outliers [9, 11, 17, 29, 42, 53, 54, 57]. ACNe [42] introduces an attention mechanism to gather information from other keypoints for filtering. OANet [54] embeds geometric priors to further improve the performance. LMCNet [29] and CLNet [56] use local motion consistency to filter inconsistent matches. MS2DG-Net [17] infuses semantics into a graph network to enhance the accuracy. These methods have achieved better results than ratio test and could serve as an additional module to filter potential outliers. However, their performance relies heavily on the quality of initial matches and require extra cost for filtering at test time. In contrast, our method embed geometric information directly into the matching module, allowing the model to predict accurate and pose-ware correspondences end-to-end.

3. Method

We first give an introduction to performing iterative transformer-based matching in Sec. 3.1. Then, we describe our iterative pose estimation framework, the adaptive sampling strategy and application at test time in Sec. 3.2, Sec. 3.3 and Sec 3.4, respectively. We visualize the pipeline of our framework in Fig. 2a.

3.1. Iterative transformer-based matching

In this section, we introduce how to perform transformer-based matching iteratively.

Problem formulation. Given two sets of keypoints (e.g. SuperPoint [18] or SIFT [30]) $\mathcal{X} = \{x_1, x_2, \dots, x_m\}$, $\mathcal{Y} = \{y_1, y_2, \dots, y_n\}$ (m, n are the number of keypoints) extracted from two images, matchers predict matches between \mathcal{X} and \mathcal{Y} , denoted as $\mathcal{O} = \{o_1, o_2, \dots, o_k\}$, where $o_j = (x_k, y_l)$ is the matched pair. Each keypoint $x_i = (u_i, v_i, c_i, \mathbf{d}_i)$ comprises its 2D coordinates (u_i, v_i) , confidence c_i and descriptor $\mathbf{d}_i \in R^d$ (d is the descriptor size).

Descriptor augmentation. As [13, 38, 40], for each keypoint x_i , its coordinates (u_i, v_i) and confidence c_i are encoded into a high-dimension vector with a multi-layer perceptron (MLP) f_{enc} , which is then added to its descriptor \mathbf{d}_i , as: $\mathbf{d}'_i = \mathbf{d}_i + f_{enc}(u_i, v_i, c_i)$. \mathbf{d}'_i is used to replace \mathbf{d}_i as input for descriptor augmentation, as:

$$\mathcal{A}_{XS}^{(t)} = \text{softmax}\left(\frac{f_q^S(\mathcal{X}^{(t-1)})(f_k^S(\mathcal{X}^{(t-1)}))^T}{\sqrt{d}}\right), \quad (1)$$

$$\mathcal{A}_{XC}^{(t)} = \text{softmax}\left(\frac{f_q^C(\mathcal{X}^{(t-1)})(f_k^C(\mathcal{Y}^{(t-1)}))^T}{\sqrt{d}}\right), \quad (2)$$

$$\mathcal{X}^{(t)} = \mathcal{X}^{(t-1)} + f_{mlp}^S(f_p^S(\mathcal{A}_{XS}^{(t)})(f_v^S(\mathcal{X}^{(t-1)}))) \parallel \mathcal{X}^{(t-1)} + f_{mlp}^C(f_p^C(\mathcal{A}_{XC}^{(t)})(f_v^C(\mathcal{Y}^{(t-1)}))) \parallel \mathcal{X}^{(t-1)}. \quad (3)$$

$\mathcal{A}_{XS}^{(t)}$ and $\mathcal{A}_{XC}^{(t)}$ are self (S) and cross (C) attention matrices for $\mathcal{X}^{(t)}$. $f_{q/k/v/p}^{S/C}$ are FC layers. d is descriptor dimension. $f_{mlp}^{S/C}$ are 3-layer MLPs. \parallel is channel-wise concatenation. We use sharing attention mechanism [12] to further augment descriptors with pre-computed attention matrices:

$$\mathcal{X}^{(t)} = \mathcal{X}^{(t)} + f_{mlp}^S(f_p^S(\mathcal{A}_{XS}^{(t)})(\bar{f}_v^S(\mathcal{X}^{(t-1)}))) \parallel \mathcal{X}^{(t-1)} + f_{mlp}^C(f_p^C(\mathcal{A}_{XC}^{(t)})(\bar{f}_v^C(\mathcal{Y}^{(t-1)}))) \parallel \mathcal{X}^{(t-1)}. \quad (4)$$

$\bar{f}_v^{S/C}$ are FC layers as well. We conduct the same operations to $\mathcal{Y}^{(t-1)}$ to obtain augmented descriptors $\mathcal{Y}^{(t)}$, self $\mathcal{A}_{YS}^{(t)}$ and cross $\mathcal{A}_{YC}^{(t)}$ attention matrices. As shown in Fig. 2b, $\mathcal{X}^{(t)}$ and $\mathcal{Y}^{(t)}$ are augmented descriptors used for the next iteration. $\mathcal{A}_{XS}^{(t)}$, $\mathcal{A}_{XC}^{(t)}$, $\mathcal{A}_{YS}^{(t)}$ and $\mathcal{A}_{YC}^{(t)}$ are self and cross attention matrices used for message propagation in the augmentation function f_A . In prior works [13, 38], self and cross attention matrices are only used for augmenting descriptors. However, in our model, they are further utilized for adaptive sampling (Sec. 3.3).

Iterative matching prediction. Augmented descriptors $\mathcal{X}^{(t)}$, $\mathcal{Y}^{(t)}$ are used to compute the matching matrix $\mathcal{M}^{(t)} \in R^{m \times n}$ with function f_M . f_M first computes Euclidean distance $\mathcal{D}^{(t)} \in R^{m \times n}$ for all pairs in $\mathcal{X}_A^{(t)} \times \mathcal{Y}_A^{(t)}$ and then uses Sinkhorn algorithm [15, 41] to obtain $\mathcal{M}^{(t)}$ by optimizing $\mathcal{D}^{(t)}$ a transport problem [32]. Matches with

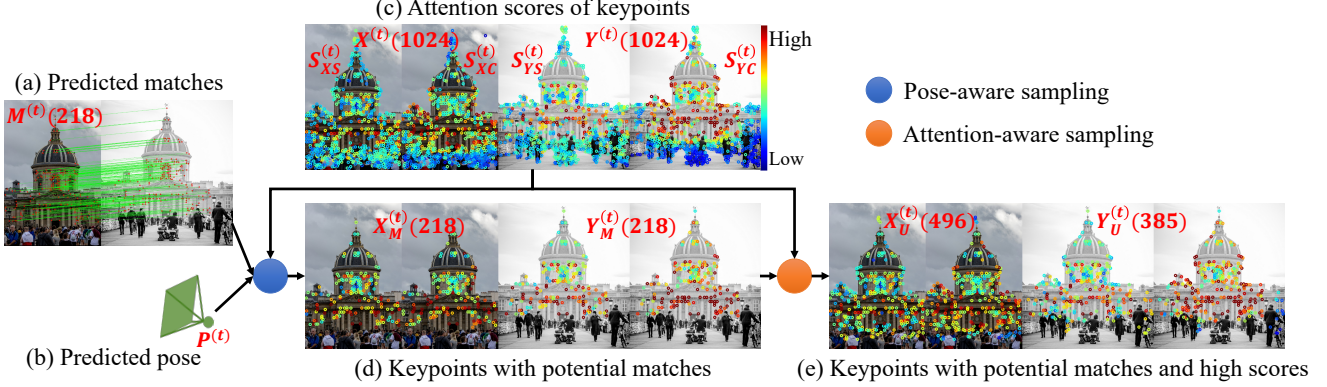


Figure 3. **Adaptive sampling.** At iteration t , predicted matches $\mathcal{M}^{(t)}$ and pose $P^{(t)}$ are used to select keypoints with potential correspondences. These selected keypoints are expanded by those with high self or cross attention scores. The finally preserved keypoints are those with potential matches and high contribution, guaranteeing both the efficiency and accuracy.

scores over than a certain threshold θ_m are deemed as predicted matches. In our framework, we predict matches at each iteration t for pose estimation and keypoint sampling as opposed to at only the last iteration [38, 40].

Iterative matching loss. Similar to SuperGlue [38], we adopt the classification loss of minimizing the negative log-likelihood of the matching matrix to enforce the network to predict correct matches for each iteration t , as:

$$L_m^{(t)} = - \sum_{(i,j) \in \mathcal{M}^{gt}} \log \bar{\mathcal{M}}_{i,j}^{(t)} - \sum_{i \in \mathcal{M}^{gt}} \log \bar{\mathcal{M}}_{i,n+1}^{(t)} - \sum_{j \in \mathcal{M}^{gt}} \log \bar{\mathcal{M}}_{m+1,j}^{(t)}. \quad (5)$$

$\bar{\mathcal{M}}^{(t)}$ is the expanded matrix of $\mathcal{M}^{(t)}$ with an additional row and column for unmatched keypoints [38] in $\mathcal{X}^{(t)}$ and $\mathcal{Y}^{(t)}$. \mathcal{M}^{gt} and $\bar{\mathcal{M}}^{gt}$ are the groundtruth matching matrix and its expansion. Applying the matching loss to each iteration t enables our model to predict correct matches at each iteration progressively.

3.2. Iterative pose estimation

In this section, we give a detailed description of how to predict poses and the usage of predicted poses for enhancing the matching process in each iteration, as shown in Fig. 2b.

Pose-consistency loss. Because of the noise [28] and degenerated keypoints [14], not all correct matches could give a good pose [5, 19]. Eq (5) only guarantees that keypoints with more discriminative descriptors have higher matching scores and are recognized first to attend pose estimation, but ignore the geometric requirements needed for robust pose estimation. Therefore, directly using all potential inliers with matching score over than θ_m for pose estimation could be inaccurate. Instead, we implicitly infuse geometric information into matching transformers, enforcing the matching module to focus first on matches which are not only correct but also have high probability to give a good pose.

To this end, at each iteration t , matches with score over

θ_m in $\mathcal{M}^{(t)}$ along with their scores are used for fundamental matrix estimation with weighted 8 points function f_P , as:

$$P^{(t)} = f_P(\{(x_i^{(t)}, y_i^{(t)}, s_i^{(t)})\}). \quad (6)$$

$P^{(t)} \in R^{3 \times 3}$ is the predicted fundamental matrix and $\{(x_i^{(t)}, y_i^{(t)}, s_i^{(t)})\}$ are predicted matches with matching score $s_i^{(t)} = M^{(t)}[idx(x_i^{(t)}), idx(y_i^{(t)})]$ ($idx(\cdot)$ is the index function). We enforce the geometric consistency between $P^{(t)}$ and the groundtruth P^{gt} by minimizing the pose and epipolar errors jointly, as:

$$L_p^{(t)} = \|P^{(t)} - P^{gt}\|_2, \quad (7)$$

$$L_{pg}^{(t)} = \frac{1}{N_{gt}} \sum_k f_{epipolar}(P^{(t)}, x_k^{gt}, y_k^{gt}), \quad (8)$$

$$L_{mg}^{(t)} = \frac{1}{N_p^{(t)}} \sum_k f_{epipolar}(P^{gt}, x_k^{(t)}, y_k^{(t)}). \quad (9)$$

$f_{epipolar}$ is Sampson distance [22]. (x_k^{gt}, y_k^{gt}) and $(x_k^{(t)}, y_k^{(t)})$ are the groundtruth and predicted matches. $L_p^{(t)}$ and $L_{pg}^{(t)}$ enforce the predicted pose $P^{(t)}$ to be correct with constraints provided by the groundtruth pose P^{gt} and matches (x_k^{gt}, y_k^{gt}) , respectively. $L_{mg}^{(t)}$ additionally ensures the correctness of predicted matches $(x_k^{(t)}, y_k^{(t)})$ with groundtruth pose P^{gt} . N_{gt} and $N_p^{(t)}$ are the number of groundtruth and predicted matches.

For each iteration, our final loss combines $L_m^{(t)}$, $L_{pg}^{(t)}$ and $L_{mg}^{(t)}$ with weights of α_m , α_p and α_g , as:

$$L^{(t)} = \alpha_m L_m^{(t)} + \alpha_p L_p^{(t)} + \alpha_g (L_{pg}^{(t)} + L_{mg}^{(t)}). \quad (10)$$

We apply $L^{(t)}$ to each iteration and compute the total loss over T iterations, as $L_{total} = \frac{1}{T} \sum_t L^{(t)}$.

3.3. Adaptive geometry-aware sampling

In fact, lots of keypoints are uninformative and a large number of keypoints don't have correspondences in the other image (about 70% of 2k keypoints in YFCC100m dataset [45]). Updating these keypoints leads to extra time

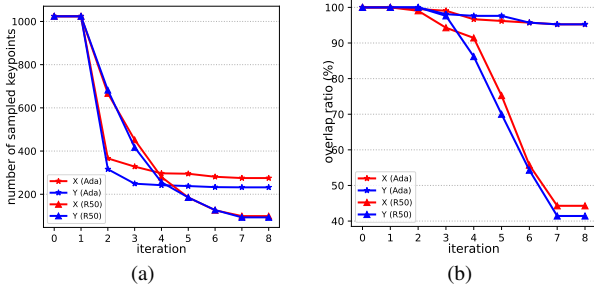


Figure 4. **Number of Retained keypoints and overlap ratio.** We show (a) the number of retained keypoints and (b) ratio of inliers in two sets at each iteration. Compared to R50, the adaptive sampling reduces redundant keypoints more effectively at early stages (a) and preserves much more inliers at latter iterations (b).

cost, so we propose an effective strategy to remove these keypoints, as shown in Fig. 3.

Informative keypoints. The information contained by each keypoint is defined by its contribution to others in the attention matrix $\mathcal{A}_{XS}^{(t)} \in R^{m \times n \times h}$ (m, n are the number of keypoints in the query and key and h is the number of heads). As [44, 46], we compute the score of each keypoint by averaging values along the head and key dimension, as $S_{XS}^{(t)} = \frac{1}{n \cdot h} \sum_{(i,j)} \mathcal{A}_{XS,i,j}$. $S_{XS} \in R^m$ is normalized with sum of 1. We also compute scores from $\mathcal{A}_{XS}^{(t)}$, $\mathcal{A}_{YS}^{(t)}$, $\mathcal{A}_{YC}^{(t)}$ as $S_{XC}^{(t)}$, $S_{YS}^{(t)}$, $S_{YC}^{(t)}$. Self and cross attention scores of 1024 keypoints in $\mathcal{X}^{(t)}$ and $\mathcal{Y}^{(t)}$ are visualized in Fig. 3c. Although attention scores can be directly used for sampling keypoints by keeping a certain ratio of keypoints with the highest scores [25, 46], this strategy is prone to over-pruning. Instead, we use keypoints with potential matches as guidance to mitigate this problem.

Adaptive sampling. At iteration t , the matching matrix $\mathcal{M}^{(t)}$ containing matching confidence of all pairs, reveals which keypoints potentially have true correspondences. We visualize the potential matches and scores of self and cross attention in Fig. 3a and Fig. 3d.

Based on matching matrix $\mathcal{M}^{(t)}$, we generate two subsets $\mathcal{X}_{\mathcal{M}}^{(t)} \subseteq \mathcal{X}^{(t)}$ and $\mathcal{Y}_{\mathcal{M}}^{(t)} \subseteq \mathcal{Y}^{(t)}$ which contain matched keypoints (matching score over θ_m). Since keypoints in $\mathcal{X}_{\mathcal{M}}^{(t)}$ and $\mathcal{Y}_{\mathcal{M}}^{(t)}$ are potential inliers, they can provide guidance to find more informative ones. In detail, let $S_{MSX}^{(t)} = \{S_{XS}^{(t)}[idx(x)], s.t. x \in \mathcal{X}_{\mathcal{M}}^{(t)}\}$ be self-attention scores of keypoints in $\mathcal{X}_{\mathcal{M}}^{(t)}$. We generate another set of keypoints with high self-attention scores as $\mathcal{X}_S^{(t)} = \{x, s.t. x \in \mathcal{X}^{(t)}, S_{XS}^{(t)}[idx(x)] \geq f_{md}(S_{MSX}^{(t)})\}$ (f_{md} returns the median value). By repeating this process, we obtain another subset from $\mathcal{X}^{(t)}$ with high cross-attention scores as $\mathcal{X}_C^{(t)}$ and two sets $\mathcal{Y}_S^{(t)}$, $\mathcal{Y}_C^{(t)}$ from $\mathcal{Y}^{(t)}$ with high self and cross attention scores as well. The final sets are the union of informative keypoints and those with matches, as $\mathcal{X}_U^{(t)} =$

$\mathcal{X}_{\mathcal{M}}^{(t)} \cup \mathcal{X}_S^{(t)} \cup \mathcal{X}_C^{(t)}$, $\mathcal{Y}_U^{(t)} = \mathcal{Y}_{\mathcal{M}}^{(t)} \cup \mathcal{Y}_S^{(t)} \cup \mathcal{Y}_C^{(t)}$. As shown in Fig. 3e, $\mathcal{X}_U^{(t)}$, $\mathcal{Y}_U^{(t)}$ will replace $\mathcal{X}^{(t)}$, $\mathcal{Y}^{(t)}$ to join the next iteration. The sampling reduces the number of keypoints from 1024 to 496 and 358 in two sets, reducing the time cost significantly for next iterations.

Fig. 4a visualizes the number of keypoints in $\mathcal{X}_U^{(t)}$ and $\mathcal{Y}_U^{(t)}$. Fig. 4b shows the value of $\frac{|\mathcal{X}_U^{(t)} \cap \mathcal{X}_F|}{|\mathcal{X}_F|}$ and $\frac{|\mathcal{Y}_U^{(t)} \cap \mathcal{Y}_F|}{|\mathcal{Y}_F|}$ (\mathcal{X}_F and \mathcal{Y}_F are two sets of matched keypoints without any sampling; $|\cdot|$ indicates the number of elements), which measure the ability of retaining inliers. For comparison, we also generate two sets of keypoints based on sampling ratio by choosing the top 50% keypoints with highest self and cross attention scores as $\mathcal{X}_{R50S}^{(t)}$, $\mathcal{X}_{R50C}^{(t)}$, $\mathcal{Y}_{R50S}^{(t)}$, and $\mathcal{Y}_{R50C}^{(t)}$. The final sampled sets are $\mathcal{X}_{R50}^{(t)} = \mathcal{X}_{R50S}^{(t)} \cup \mathcal{X}_{R50C}^{(t)}$ and $\mathcal{Y}_{R50}^{(t)} = \mathcal{Y}_{R50S}^{(t)} \cup \mathcal{Y}_{R50C}^{(t)}$. The number of retained keypoints and overlap ratios are also visualized.

Fig. 4a shows the adaptive strategy performs more effectively at discarding keypoints at early stages and preserving keypoints at later iterations. Fig. 4b illustrates that although R50 has the close ratio to the adaptive strategy at the 2nd and 3rd iterations by keeping more redundant keypoints, the former loses overlap ratio dramatically after 4 iterations while the latter still keeps the ratio over 90%.

Adaptive sampling with pose uncertainty. The matching matrix $\mathcal{M}^{(t)}$ might not be very accurate at the first several iterations when descriptors are not discriminative, impairing the accuracy. To mitigate the problem, we make use of the predicted pose. Note, if most of the predicted matches are correct, the pose is more precise with higher inlier ratio. However, if predicted matches contain many outliers, the pose is probably not accurate and has low inlier ratio. Consequently, we define the uncertainty of the pose on its consistency with matches, as $r^{(t)} = \frac{|\{(x_i, y_i), s.t. f_{epipolar}(P^{(t)}, x_i, y_i) \leq \theta_e\}|}{|\{(x_i, y_i, s_i), s.t. s_i \geq \theta_m\}|}$ (θ_e is threshold determining inliers). We use $r^{(t)}$ to adjust the sampling threshold θ_m , as $\theta_m^{(t)} = \theta_m * r^{(t)}$, allowing the model to dynamically sample fewer keypoints when both matches and pose are accurate and more when they are not.

3.4. Iterative process at test time

At test time, after each iteration, we compute matches $\mathcal{M}^{(t)}$ and estimate the pose from matches with RANSAC [21]. We adopt the relative error between consecutively predicted pose $P^{(t)}$ and $P^{(t-1)}$ as stop criteria to determine if continuing iteration. Specifically, if the maximum angular errors of rotations and translations is less than θ_P , the iteration stops. Benefiting from the embedded geometry information, for about 55% of the cases in YFCC100m dataset [45], 6 iterations rather than 9 are enough to find a good pose. The predicted pose is further used for adaptive sampling at test time (Sec. 3.3).

Once a good pose is predicted, more matches can be found from the guidance of the pose with function f_{PM} . f_{MP} first computes the epipolar distance for all pairs in $\mathcal{X}^{(t)}, \mathcal{Y}^{(t)}$ as $\mathcal{M}_e^{(t)}$ where $\mathcal{M}_{eij}^{(t)} = \frac{f_{epipolar}(\mathcal{X}_i^{(t)}, \mathcal{Y}_j^{(t)})}{(y^T P^{(t)} x)^2} = \frac{f_{epipolar}(x, y)}{(P^{(t)}x)_1^2 + (P^{(t)}x)_2^2 + (P^{(t)T}y)_1^2 + (P^{(t)T}y)_2^2}$. Then $\mathcal{M}_e^{(t)}$ is binarized with errors smaller than $12px$ as 1 otherwise 0 as $\bar{\mathcal{M}}_e^{(t)}$. Finally, the pose-guided matching matrix is obtained as $\mathcal{M}^P = \mathcal{M}^{(t)} \odot \bar{\mathcal{M}}_e^{(t)}$ (\odot is element-wise multiplication). As shown in Fig. 2c, matches with high uncertainties caused by similar descriptors at regions with repetitive textures can be effectively mitigated by pose constraints. These matches usually require more number iterations to be found.

4. Experiment Setup

Training. We implement the model in PyTorch [31] and train it on MegaDepth dataset [27] for 900k iterations with Adam optimizer [24], batch size of 16, initial learning rate of 0.0001. The learning rate decays with ratio of 0.99996 after 200k iterations and is fixed at 0.00001. We use 1024 SuperPoint [18] and RootSIFT [3, 30] keypoints for training. As with SuperGlue [38], our network has T (9) iterative blocks. $\alpha_m, \alpha_p, \alpha_g, \theta_m, \theta_P, \theta_e$ are set to 0.6, 0.2, 0.2, 0.2, 1.5 and 0.005, respectively. Note that unlike SuperGlue which is first trained on Oxford-Paris dataset [33] and then fine-tuned on Scannet [16] and MegaDepth [27] datasets for indoor and outdoor models respectively, our models are trained only on the MegaDepth dataset from scratch.

Datasets and metrics. We first test our method on YFCC100m [45] and Scannet [16] datasets to evaluate the performance on relative pose estimation. YFCC100m is a large-scale outdoor dataset consisting of images with large illumination and season changes and viewpoint variations. Scannet is an indoor dataset widely used for depth prediction [4] and pose estimation [13, 38]. For relative pose estimation, we report the accurate cumulative error curve (AUC) [38] at thresholds of $5^\circ, 10^\circ,$ and 20° . The error is the maximum angular errors of rotations and normalized translations. We also report the mean matching score (M.S.) and mean precision (Prec.) of matches.

We additionally test our method on large-scale localization task at Aachen Day-Night (v1.0 and v1.1) datasets [39, 55]. Aachen v1.0 dataset comprises of 4,328 reference and 922 query (824 day, 98 night) images. Aachen v1.1 extends v1.0 by adding additional 2,369 reference and 93 night query images. We use HLoc [37] pipeline for mapping and localization and report the accuracy at error thresholds of $0.25m/2^\circ, 0.5m/5^\circ,$ and $5m/10^\circ$.

Baselines. The baselines comprise of simple matchers such as MNN and NN-RT [30]. Also, filtering-based methods including OANet [54], AdaLAM [11], CLNet [56] and LMCNet [29] are also included. The final part is

Feature	Matcher	@5°	@10°	@20°	M.S.(%)	Prec.(%)
RootSIFT [30]	NN-RT [30]	26.7	43.2	59.4	4.4	56.4
	AdaLAM (4k) [11]	27.5	44.5	60.5	6.3	84.3
	OANet [54] [54]	22.4	36.3	50.3	5.6	53.7
	CLNet [56]	33.0	52.1	68.5	7.8	75.2
	LMCNet [29]	35.8	55.6	72.2	-	86.9
	SuperGlue* [38]	35.3	56.1	73.6	-	-
	SuperGlue [13, 38]	35.1	54.2	70.9	16.6	81.7
	SGMNet [13]	34.8	54.1	70.9	17.1	86.1
	ClusterGNN [40]	32.8	50.3	65.9	-	-
	IMP	36.7	56.6	72.9	18.0	87.6
EIMP	36.8	56.3	72.8	13.7	89.8	
SuperPoint [18]	MNN	6.5	15.4	28.5	16.2	16.2
	AdaLAM (2k) [11]	20.8	36.5	51.9	10.9	72.0
	OANet [54]	19.2	34.5	50.3	9.4	62.1
	CLNet [56]	27.8	46.4	63.8	11.9	75.1
	LMCNet [29]	17.4	33.2	51.1	-	88.9
	SuperGlue* [38]	37.1	57.2	73.6	21.7	88.5
	SuperGlue [13, 38]	33.2	53.5	70.8	19.7	78.7
	SGMNet [13]	33.0	53.0	70.0	<u>22.3</u>	85.1
	ClusterGNN [40]	35.3	56.1	73.6	-	-
	IMP	39.4	59.4	75.2	23.0	84.9
EIMP	37.9	57.9	74.0	19.9	88.4	

Table 1. **Results on YFCC100m dataset [45].** We report the AUC of relative poses at error thresholds of $5^\circ, 10^\circ,$ and 20° . The mean matching score (M.S.) and mean matching accuracy (Prec.) are also reported. The **best** and second best results are highlighted.

transformer-based matchers such as SuperGlue [38], SGMNet [13], and ClusterGNN [40]. As it is difficult to reproduce the results of SuperGlue from custom training (official training code is not released), we follow SGMNet and show results of original SuperGlue (SuperGlue*) and the model trained by SGMNet (SuperGlue).

5. Experiment Results

In this section, we first show results on relative pose estimation and localization tasks in Sec. 5.1 and Sec. 5.2. Next, we analyze the computational cost in Sec. 5.3. Finally, we conduct a full ablation study to test each component in our framework in Sec. 5.4.

5.1. Relative pose estimation

Comparison with filtering-based methods. Table 1 shows our IMP and efficient IMP (EIMP) give higher accuracy than previous filtering-based approaches such as OANet [54] and LMCNet [29] for both RootSIFT [3, 30] and SuperPoint [18] features. Because these filtering-based methods rely purely on the geometric information to filter potentially wrong correspondences, their performance is influenced heavily by the quality of initial matches. In contrast, we use both the geometric information and descriptors for matching, which are able to give more accurate matches, resulting in more precise poses.

Comparison with matchers. Augmented with spatial information, SuperGlue [38] and its variants [13, 40] outperform MNN and NN-RT obviously. Although SGMNet [13] and ClusterGNN [40] are faster, they give worse

Feature	Matcher	@5°	@10°	@20°	M.S.(%)	Prec.(%)
RootSIFT [30]	NN+RT [30]	9.1	19.8	32.7	2.3	28.8
	AdaLAM (4k) [11]	8.2	18.6	31.0	3.1	47.6
	OANet [54]	10.7	23.1	37.4	3.2	36.9
	CLNet [56]	5.8	15.1	26.8	2.0	43.9
	LMCNet [29]	8.5	19.3	32.4	-	47.0
	SuperGlue* [38]	-	-	-	-	-
	SuperGlue [13, 38]	14.7	29.4	45.6	8.4	42.2
	SGMNet [13]	14.4	29.9	46.0	8.8	45.6
	IMP	15.6	30.9	47.4	5.8	42.0
	EIMP	15.3	30.8	46.6	7.7	45.7
SuperPoint [18]	MNN	9.4	21.6	36.4	13.3	30.2
	AdaLAM (2k) [11]	6.7	15.8	27.4	13.2	44.2
	OANet [54]	10.0	25.1	38.0	10.6	44.6
	CLNet [56]	4.1	11.0	21.6	8.6	44.2
	LMCNet [29]	14.6	33.6	53.6	-	36.8
	SuperGlue* [38]	16.2	32.6	49.3	16.8	47.9
	SuperGlue [13, 38]	12.0	26.3	42.4	15.0	45.9
	SGMNet [13]	16.4	32.1	48.7	17.0	48.1
	IMP	16.6	33.1	49.4	15.8	42.0
	EIMP	15.9	32.4	48.9	15.9	46.2

Table 2. **Results on Scannet dataset [16].** We show the AUC of relative poses at error thresholds of 5°, 10°, and 20°. The mean matching score (M.S.) and mean matching accuracy (Prec.) are also reported. The **best** and second best results are highlighted.

accuracy than SuperGlue* due to the loss of information in the message propagation process. With implicitly embedded geometric information, our IMP gives more pose-aware matches and hence outperforms SuperGlue* especially for very precise pose estimation at error threshold of 5°. Enhanced by geometric information, our EIMP even gives slightly better results than SuperGlue* at error thresholds of 5° and 10°. Our EIMP outperforms previous efficient methods such as SGMNet and ClusterGNN at all error thresholds but has close even higher efficiency.

Table 2 shows results on the Scannet dataset [16]. Our IMP and EIMP also give higher accuracy than filtering-based approaches and other matchers on RootSIFT keypoints. When using SuperPoint, our IMP gives slightly better performance than SuperGlue* and SGMNet, which yield close performance to our EIMP. LMCNet [29] reports the best accuracy at error thresholds of 10° and 20° with SuperPoint because LMCNet uses SuperGlue* to provide initial matches and is further trained on the indoor dataset [49].

Qualitative comparison. As shown in Fig. 5, in the iteration process, IMP and EIMP produce more inliers and more accurate poses progressively. Besides, EIMP dynamically discards useless keypoints after each iteration. We also show the results of SuperGlue* [38] and SGMNet [13]. Due to large viewpoint changes, both SuperGlue and SGMNet give much fewer inliers from only a small area and have larger errors compared to our models. In the iterative process, inliers span almost the whole meaningful regions, resulting in more robust pose estimation.

Feature	Matcher	Day (0.25m, 2°)/(0.5m, 5°)		Night (5m, 10°)	
		Day	Night	Day	Night
SuperPoint [18]	MNN	85.4 / 93.3 / 97.2	75.5 / 86.7 / 92.9	-	-
	OANet [54]	-	77.6 / 86.7 / 98.0	-	-
	AdaLAM [11]	-	78.6 / 86.7 / 98.0	-	-
	ELAM [43]	-	78.6 / 87.8 / 96.9	-	-
	SuperGlue [13, 38]	-	76.5 / 88.8 / 99.0	-	-
	SuperGlue* [38]	89.6 / 95.4 / 98.8	86.7 / 93.9 / 100.0	-	-
	SGMNet [13]	86.8 / 94.2 / 97.7	83.7 / 91.8 / 99.0	-	-
	ClusterGNN [40]	89.4 / 95.5 / 98.5	81.6 / 93.9 / 100.0	-	-
	IMP	89.1 / 95.4 / 99.0	86.7 / 94.9 / 100.0	-	-
	EIMP	90.0 / 96.5 / 99.2	84.7 / 94.9 / 100.0	-	-
SuperPoint [18]	MNN	87.9 / 93.6 / 96.8	70.2 / 84.8 / 93.7	-	-
	AdaLAM [11]	-	73.3 / 86.9 / 97.9	-	-
	SuperGlue* [38]	89.8 / 96.6 / 99.4	75.9 / 90.1 / 100.0	-	-
	SGMNet [13]	88.7 / 96.2 / 98.9	75.9 / 89.0 / 99.0	-	-
	IMP	89.1 / 95.4 / 99.0	75.9 / 92.7 / 99.5	-	-
	EIMP	90.0 / 96.5 / 99.0	77.0 / 91.6 / 99.5	-	-

Table 3. **Results on Aachen v1.0 (top) and v1.1 (bottom) dataset [39, 55].** The **best** and second best results are highlighted.

5.2. Visual localization

As most filtering methods don't provide results on visual localization task, we only show results of OANet [54] and AdaLAM [11]. OANet and AdaLAM give close accuracy to MNN at error thresholds of 0.25m, 2° and 0.5m, 5°, but much better performance at threshold of 5m, 10°, because test images with larger viewpoint changes are more sensitive to outliers which can be partially filtered by OANet and AdaLAM. Our IMP obtains similar performance to SuperGlue* [38] and outperforms SGMNet and ClusterGNN especially for night images. Note that our models are only trained on the MegaDepth dataset [27] while SuperGlue* is additionally pretrained on Oxford and Pairs dataset [33].

EIMP slightly outperforms IMP and SuperGlue*. That is because in the long-term large-scale localization task, query and reference images usually have larger viewpoint and illumination changes with more keypoints without true matches. EIMP effectively discards these keypoints, guaranteeing the quality of matches and thus improving the localization accuracy.

5.3. Running time

As the source code of ClusterGNN [40] is not released, we mainly compare our method with SuperGlue* [38] and SGMNet [13]. The total time for each method is the summary of the time used at each iteration. IMP reduces the time by reducing the number of iterations and EIMP further decreases the time of each iteration.

Fig. 6a shows the number of iterations for relative pose estimation on YFCC100m dataset [45]. SuperGlue* and SGMNet adopt a fixed number of layers and only give matches at the last layer, so they don't report any results before the last iteration. In contrast, IMP needs only 4 iterations for 30% cases and 5 iterations for 55% cases to find a good pose (relative pose error less than 1.5), avoiding ex-



Figure 5. **Qualitative comparison.** We visualize the iterative process (left→right) of IMP (top) and EIMP (bottom). For each pair, we report the number inliers/matches and rotation/translation errors (bottom) and the number of preserved keypoints in two sets (top) at each iteration. **Inliers** between two images are visualized. In the iterative process, our model not only finds more inliers but enforces inliers to span wider regions. However, both SuperGlue* and SGMNet predict fewer matches from a small area, resulting higher pose errors.

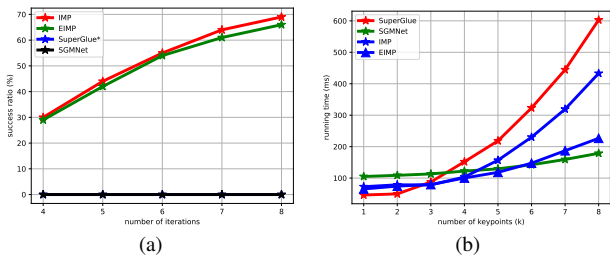


Figure 6. **Running time analysis.** We show (a) the number of iterations required for each pair YFCC100m dataset [45] and (b) the mean running time of 1k pairs on RTX 3090.

tra computation. Although EIMP loses some matches by discarding useless keypoints, it obtains close results to IMP. Fig. 6b shows the running time of SuperGlue*, SGMNet, IMP and EIMP. When using 1k and 2k keypoints, IMP and EIMP run slower than SuperGlue* because of additional in shared attention. However, when using more than 3k keypoints, IMP runs faster because it uses fewer number of iterations (10) for sinkhorn optimization [15, 41]. Benefiting from adaptive pooling strategy, EIMP works even faster than SGMNet when using less than 6k keypoints but gives better accuracy, as demonstrated in Table 1.

5.4. Ablation study

We conduct a full ablation study for relative pose estimation on YFCC100m dataset [45] with SuperPoint [18] keypoints to verify all components in our model including the shared attention (S), pose-consistency loss (C), pose-aware iteration (P), adaptive pooling (A), and adaptive pooling with pose uncertainty (U).

Table 4 shows that shared attention (S) improves both the pose and matching accuracy by propagating more message to each keypoint. By infusing the geometric information in the iterative process, pose-consistency loss (C) and pose-aware iteration (P) also enhance the performance of IMP and EIMP. Comparisons between IMP-SC and EIMP-SCA indicate that our adaptive pooling loses 5% matches

Model	S	C	P	A	U	@5°	@10°	20°	M.S.(%)	Prec.(%)
IMP	X	X	X	X	X	36.9	57.1	73.7	23.7	85.2
IMP-S (base)	✓	X	X	X	X	38.6	58.3	74.6	23.6	87.2
IMP-SC	✓	✓	X	X	X	<u>39.1</u>	<u>59.0</u>	<u>75.0</u>	23.7	87.7
IMP-SCP (full)	✓	✓	✓	X	X	39.4	59.4	75.2	23.0	84.9
EIMP-S	✓	X	X	X	X	35.4	55.5	72.1	7.8	91.2
EIMP-SA (base)	✓	X	X	✓	X	36.4	56.5	73.0	<u>18.4</u>	88.8
EIMP-SCA	✓	✓	X	✓	X	36.9	56.6	<u>73.3</u>	<u>18.4</u>	88.7
EIMP-SCPA	✓	✓	✓	X	X	<u>37.2</u>	<u>57.2</u>	<u>73.3</u>	18.3	86.5
EIMP-SCPAU (full)	✓	✓	✓	✓	✓	37.9	57.9	74.0	19.9	88.4

Table 4. **Ablation study.** We test the efficacy of shared attention (S), pose-consistency loss (C), adaptive pooling (A), pose-aware iteration (P), and adaptive sampling with pose uncertainty (U) on YFCC100m dataset [45] with SuperPoint [18] features. The **best** and **second best** results are highlighted.

(M.S. 23.7 vs. 18.4) and gains 1% inliers (Prec. 87.7 vs. 88.7). The pose uncertainty (U) further increases the number of matches (M.S. 18.3 vs. 19.9) and inliers (Prec. 86.5 vs. 88.4) and enhances the pose accuracy about 0.7% by dynamically adjusting the sampling ratio to avoid over-pruning at each iteration. EIMP-S, which is a ratio-based sampling, gives promising pose accuracy but loses over 10% matches (M.S. 7.8 vs. 18.4) than adaptive pooling, EIMP-SA, which effectively mitigates this problem.

6. Conclusions

In this paper, we propose the iterative matching and pose estimation framework, allowing the two tasks to boost each other and thus improving the accuracy and efficiency. Particularly, we embed the geometric information into the matching module, enforcing the model to predict matches which are not only accurate but also able to give a good pose. Moreover, in each iteration, we utilize the predicted matches, relative pose, and attention scores to remove keypoints without potential true matches adaptively at each iteration, improving the efficiency and preserving the accuracy. Experiments demonstrate that our method achieves better

performance than previous approaches on relative pose estimation and large-scale localization tasks and has high efficiency as well.

Acknowledgment

This project is supported by Toyota Motor Europe.

Supplementary Material

In the supplementary material, we provide additional implementation details, qualitative results and analysis in Sec. A and Sec. B, respectively.

A. Implementation

A.1. Training

We train our model on the MegaDepth dataset [27] from scratch. Following SuperGlue [38], we use 153 scenes with 130k images in total for training and 36 for validation. For each category, we first detect 4096 keypoints for all images and then build correspondences for image pairs with overlap ratio from 0.3-1.0. Matches with re-projection errors less than 5px are deemed as inliers, resulting in different number of inliers for different pairs. In the training process, for each epoch, we randomly choose 80 pairs for each scene. For each pair, we randomly choose 1024 keypoints with inliers ranging from 32 to 512 between two images, respectively. We observe that samples with high inlier ratios boost the convergence and those with low inlier ratios enhance the ability of models for finding matches for tough cases at test time.

Our matching loss is identical to the assignment loss of SGMNet [13], which is more stable than the original version in SuperGlue [38], as analyzed in the SGMNet paper. These modifications enable us to train the model on the MegaDepth dataset [27] from scratch without requiring any pretraining.

A.2. Architecture

As [13, 38, 40], we use self and cross attention to gather global information for each keypoint in two sets. Considering the similarities of attention matrices in two consecutive iterations, we adopt the shared attention mechanism [12] to speed up the message propagation process at low cost. A detailed architecture of self and cross attention is shown in Fig. 7. We adopt the identical position encoder as SuperGlue [38].

A.3. Inference

At test time, for relative pose estimation on YFCC100m [45] and Scannet [16] datasets, we use the identical testing pairs and number of keypoints for evaluation as SuperGlue [38] and SGMNet [13]. As for the

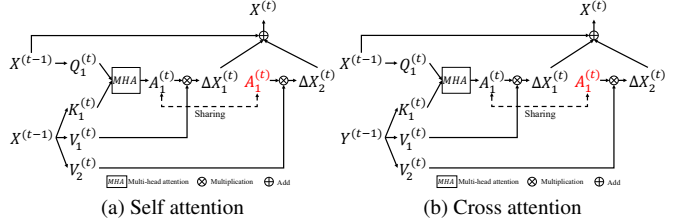


Figure 7. **Architecture.** Detailed architecture of self (a) and cross (b) attention with sharing attention matrix.

metrics, we report the **exact** cumulative error curve (AUC) rather than the **approximate** one, because the former measures the error based on groundtruth poses while the latter does it based on predicted poses.

When evaluating our model on Aachen Day-Night v1.0 and v1.1 datasets [39, 55], we adopt HLoc [37] pipeline for mapping and localization, as previous methods [13, 38, 40]. We leverage NetVLAD [2] to provide 50 reference images for each query image in the localization process.

In our adaptive sampling process, in order to avoid losing too many potential inliers or informative keypoints, we set the minimum number of preserved keypoints for each image to 256. Once the number of keypoints is smaller or equal to 256, we don't perform any sampling. This strategy is also applied to the ratio-based sampling (R50).

B. Results

In this section, we first provide more ablation studies in Sec. B.1. Then, we show and discuss more qualitative results of SuperGlue* (official SuperGlue) [38], SGMNet [13] and our IMP and EIMP on Scannet [16], YFCC100m [45], and Aachen Day-Night datasets [39, 55] in Sec. B.2, Sec. B.3, and Sec. B.4, respectively.

B.1. Ablation study

Pose-consistency loss (C). Pose-consistency loss forces our model to predict matches which are not only correct but also able to give a good pose by implicitly embedding geometric information into the matching module. Fig. 8a shows the influence of pose-consistency loss to the number of iterations required to find a good pose. With pose-consistency loss, as expected, when using the same number of iterations, IMP gives higher success ratios than IMP without this loss (IMP w/o C) because pose-consistency loss allows pose-aware matches pump out first to join the pose estimation, reducing the number of iterations. EIMP adopts the sampling strategy which filters some unreliable matches, therefore the pose-consistency has little influence on the success ratios of EIMP.

Number of iterations. A our model is able to predict matches at each iteration. We test its performance on rela-

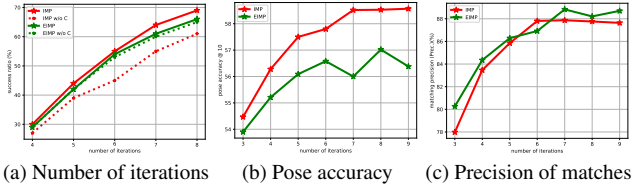


Figure 8. We report the number of iterations for IMP and EIMP with and without pose-consistency (w/o C) loss (a), relative pose accuracy @ 10° (b), and precision of matches (c) of IMP and EIMP using different number of iterations on YFCC100m dataset [45].

tive pose estimation by progressively increasing the number of iterations from 3 to 9. As the number of iterations increases, keypoints become more discriminative with more geometric information embedded, so both IMP and EIMP achieve more precise poses, as shown in Fig. 8b. Additionally, their matching precision (Prec.) also gets higher gradually, indicating that they find more inliers, as shown in Fig. 8c. Fig. 8b and Fig. 8c show that both IMP and EIMP report almost the best performance when using 7 or 8 iterations, which means IMP and EIMP could even run faster by reducing the maximum number iterations from 9 to 7 or with marginal performance loss. Note that both IMP and EIMP are trained to predict matches at each iteration, so we don’t need to retrain or fine-tune IMP and EIMP to achieve this.

B.2. Qualitative results on Scannet

In Fig. 9, we visualize predicted matches and relative poses of SuperGlue* [38], SGMNet [13], our IMP and EIMP. We observe that for simple cases (Fig. 9 (1)), all matchers give similar numbers of inliers. However, both our IMP and EIMP obtain smaller rotation and translation errors than SuperGlue* and SGMNet because the embedded geometric information in our matching module. In the iteration process, rather than finding inliers from a cluster, both IMP and EIMP expand the areas with inliers, which allows our models to find more potential inliers and make the pose estimation more stable. When testing images become more difficult (Fig. 9 (2) and (3)), the behavior of our models in expanding inliers over the whole meaningful regions of images can be observed more clearly. Especially for Fig. 9 (3) where all keypoints are extracted from regions with repetitive textures, matching methods based on pure descriptors can hardly discriminate these keypoints, so SuperGlue* and SGMNet fail to give enough inliers. At this time, geometric constraints play an indispensable role at finding correct matches, which explains the success of IMP and EIMP.

B.3. Qualitative results on YFCC100m

Fig. 10 shows the predicted matches and relative poses reported by SuperGlue* [38], SGMNet [13], our IMP and EIMP. For simple cases (Fig. 10 (1)), all methods give a large number of inliers while IMP and EIMP report smaller pose errors in the iteration process even using fewer inliers. Note that instead of keeping all keypoints, our EIMP with adaptive sampling effectively reduces the number of keypoints from 2,000 to 865 and 255 in the iteration process, which significantly decreases the time complexity for self and cross attention computation 7. We also observe that for tough cases (Fig. 10 (2) and (3)), due to large viewpoint and illumination changes, SuperGlue* and SGMNet fail to report comparable number of inliers to our models, resulting in higher rotation and translation errors. In contrast, our models still progressively increase the number of inliers from different regions in the iteration process. These well distributed inliers lead to smaller pose errors. By comparing the results of EIMP in Fig. 10 (2) and Fig. 10 (3), we see that the number of preserved keypoints are based on the number of potential inliers in the image pair: more potential inliers result in more retrained keypoints. That is because our sampling strategy is fully adaptive.

B.4. Qualitative results on Aachen

In Fig. 11, we show the inliers between query and reference images in the large-scale localization task on Aachen v1.1 dataset [39, 55]. Different with image pairs in Scannet [16] and YFCC100m [45] datasets, query and reference images in Aachen dataset are captured under totally different conditions usually with extremely large viewpoint (Fig. 11 (1)-(4)) and illumination (Fig. 11 (5)-(8)) changes, making finding matches difficult. As the groundtruth poses of query images are not available, we use HLoc [37] framework to visualize the inliers between query and reference images after the PnP [26] + RANSAC [21] for all methods.

Fig. 11 ((1)-(4)) shows that when image pairs have large viewpoint changes, both SuperGlue* and SGMNet fail to find enough correct matches. That is because geometric constraints are more useful for finding matches in two images with large viewpoint changes and both SuperGlue* and SGMNet ignore this information. However, we embed the geometric information into the matching module, so our IMP and EIMP work much better, guaranteeing the localization success.

When query images have large illumination changes with reference images, corresponding keypoints from two images are less discriminative, so SuperGlue* and SGMNet only give slightly more inliers than for images with large viewpoint, as shown in Fig. 11((5)-(8)). As our model additionally leverages geometric constraints to find matches, both IMP and EIMP successfully obtain a large number of inliers. Note that compared to SuperGlue* and SGMNet,

both IMP and EIMP find inliers from the almost the whole overlap regions of the two images as opposed to some clusters (Fig. 11 (2)-(5) and (8)).

References

- [1] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. In *ECCV*, 2022. [2](#)
- [2] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. In *CVPR*, 2016. [9](#)
- [3] Relja Arandjelović and Andrew Zisserman. Three things everyone should know to improve object retrieval. In *CVPR*, 2012. [2, 6](#)
- [4] Gwangbin Bae, Ignas Budvytis, and Roberto Cipolla. Multi-View Depth Estimation by Fusing Single-View Depth Probability with Multi-View Geometry. In *CVPR*, 2022. [6](#)
- [5] Daniel Barath, Luca Cavalli, and Marc Pollefeys. Learning To Find Good Models in RANSAC. In *CVPR*, 2022. [2, 4](#)
- [6] Daniel Barath and Zuzana Kukelova. Relative Pose from SIFT Features. In *ECCV*, 2022. [1](#)
- [7] Daniel Barath, Jiri Matas, and Jana Noskova. MAGSAC: marginalizing sample consensus. In *CVPR*, 2019. [1](#)
- [8] Daniel Barath, Jana Noskova, Maksym Ivashechkin, and Jiri Matas. MAGSAC++, a fast, reliable and accurate robust estimator. In *CVPR*, 2020. [1](#)
- [9] Eric Brachmann and Carsten Rother. Neural-Guided RANSAC: Learning where to sample model hypotheses. In *ICCV*, 2019. [1, 3](#)
- [10] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. [2](#)
- [11] Luca Cavalli, Viktor Larsson, Martin Ralf Oswald, Torsten Sattler, and Marc Pollefeys. Handcrafted outlier detection revisited. In *ECCV*, 2020. [1, 3, 6, 7](#)
- [12] Boyu Chen, Peixia Li, Baopu Li, Chuming Li, Lei Bai, Chen Lin, Ming Sun, Junjie Yan, and Wanli Ouyang. PSViT: Better vision transformer via token pooling and attention sharing. *arXiv preprint arXiv:2108.03428*, 2021. [2, 3, 9](#)
- [13] Hongkai Chen, Zixin Luo, Jiahui Zhang, Lei Zhou, Xuyang Bai, Zeyu Hu, Chiew-Lan Tai, and Long Quan. Learning to match features with seeded graph matching network. In *CVPR*, 2021. [1, 2, 3, 6, 7, 9, 10, 12, 14, 15](#)
- [14] Ondrej Chum, Tomas Werner, and Jiri Matas. Two-view geometry estimation unaffected by a dominant plane. In *CVPR*, 2005. [2, 4](#)
- [15] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *NIPS*, 2013. [3, 8](#)
- [16] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM ToG*, 2017. [6, 7, 9, 10, 12](#)
- [17] Luanyuan Dai, Yizhang Liu, Jiayi Ma, Lifang Wei, Taotao Lai, Changcai Yang, and Riqing Chen. MS2DG-Net: Progressive Correspondence Learning via Multiple Sparse Semantics Dynamic Graph. In *CVPR*, 2022. [1, 3](#)
- [18] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *CVPRW*, 2018. [2, 3, 6, 7, 8](#)
- [19] Hongyi Fan, Joe Kileel, and Benjamin Kimia. On the Instability of Relative Pose Estimation and RANSAC’s Role. In *CVPR*, 2022. [2, 4](#)
- [20] Mohsen Fayyaz, Soroush Abbasi Koohpayegani, Farnoush Rezaei, and Sommerlade1 Hamed Pirsiavash2 Juergen Gall. Adaptive Token Sampling For Efficient Vision Transformers. In *ECCV*, 2022. [2](#)
- [21] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 1981. [1, 5, 10, 15](#)
- [22] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. [4](#)
- [23] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *ICML*, 2020. [2](#)
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. [6](#)
- [25] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosior, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *ICML*, 2019. [2, 5](#)
- [26] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. EPnP: An Accurate O(n) Solution to the PnP Problem. *IJCV*, 81:155–166, 2009. [10, 15](#)
- [27] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *CVPR*, 2018. [6, 7, 9](#)
- [28] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-Perfect Structure-from-Motion with Featuremetric Refinement. In *ICCV*, 2021. [2, 4](#)
- [29] Yuan Liu, Lingjie Liu, Cheng Lin, Zhen Dong, and Wenping Wang. Learnable motion coherence for correspondence pruning. In *CVPR*, 2021. [1, 3, 6, 7](#)
- [30] David G Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 2004. [2, 3, 6, 7](#)
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. [6](#)
- [32] Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019. [3](#)
- [33] Filip Radenović, Ahmet Iscen, Giorgos Toliás, Yannis Avrithis, and Ondřej Chum. Revisiting oxford and paris: Large-scale image retrieval benchmarking. In *CVPR*, 2018. [6, 7](#)
- [34] Rahul Raguram, Ondrej Chum, Marc Pollefeys, Jiri Matas, and Jan-Michael Frahm. USAC: A universal framework for random sample consensus. *TPAMI*, 2012. [1](#)
- [35] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision

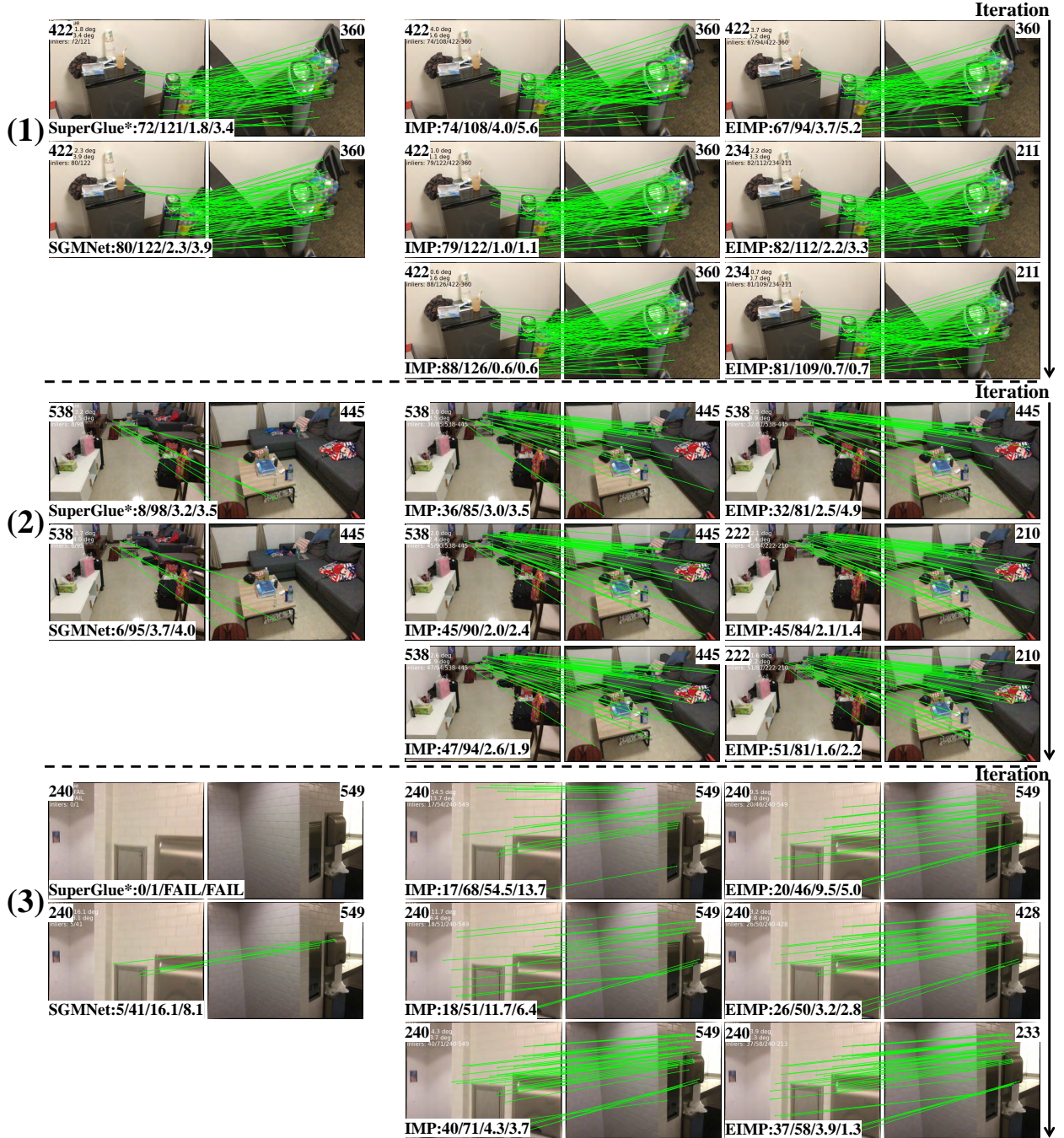


Figure 9. **Qualitative results on Scannet dataset [16].** At the left-bottom of each image pair, we report the number of inliers/matches and rotation/translation errors of SuperGlue* [38] (official SuperGlue), SGMNet [13], our IMP and EIMP. Besides, the number of keypoints in each image are shown at the top of each pair. For simple case (1), all methods give similar numbers of inliers, but IMP and EIMP obtain smaller rotation and translation errors than SuperGlue* and SGMNet because of the embedded geometric information. In the iteration process, rather than finding inliers from a cluster, both IMP and EIMP expand the areas with inliers, which allows our models to find more potential inliers and make the pose estimation more stable. When testing images become more difficult ((2), (3)), the behavior of our models in expanding inliers over the whole meaningful regions of images can be observed more clearly. Especially for (3) where all keypoints are extracted from regions with repetitive textures, matching methods based on descriptors can hardly discriminate these keypoints, so SuperGlue* and SGMNet fail to give enough inliers. At this time, geometric constraints play an indispensable role at finding correct matches, which explains the success of IMP and EIMP.

- transformers with dynamic token sparsification. In *NeurIPS*, 2021. 2
- [36] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. ORB: An efficient alternative to SIFT or SURF. In *ICCV*, 2011. 2
- [37] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From Coarse to Fine: Robust Hierarchical Localization at Large Scale. In *CVPR*, 2019. 6, 9, 10, 15
- [38] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *CVPR*, 2020. 1, 2, 3, 4, 6, 7, 9, 10, 12, 14, 15
- [39] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, et al. Benchmarking 6dof outdoor visual localization in changing conditions. In *CVPR*, 2018. 1, 6, 7, 9, 10, 15
- [40] Yan Shi, Jun-Xiong Cai, Yoli Shavit, Tai-Jiang Mu, Wensen Feng, and Kai Zhang. ClusterGNN: Cluster-based Coarse-to-Fine Graph Neural Network for Efficient Feature Matching. In *CVPR*, 2022. 1, 2, 3, 4, 6, 7, 9
- [41] Richard Sinkhorn and Paul Knopp. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific Journal of Mathematics*, 21(2):343–348, 1967. 3, 8
- [42] Weiwei Sun, Wei Jiang, Eduard Trulls, Andrea Tagliasacchi, and Kwang Moo Yi. Acne: Attentive context normalization for robust permutation-equivariant learning. In *CVPR*, 2020. 1, 3
- [43] Suwichaya Suwanwimolkul and Satoshi Komorita. Efficient linear attention for fast and accurate keypoint matching. In *ICMR*, 2022. 1, 2, 7
- [44] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *ACM Computing Surveys*, 2020. 2, 5
- [45] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. *Communications of the ACM*, 2016. 4, 5, 6, 7, 8, 9, 10, 14
- [46] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 2, 5
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 2
- [48] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. 2
- [49] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *ICCV*, 2013. 7
- [50] Fei Xue, Ignas Budvytis, and Roberto Cipolla. SFD2: Semantic-guided Feature Detection and Description. In *CVPR*, 2023. 2
- [51] Fei Xue, Xin Wang, Zike Yan, Qiuyuan Wang, Junqiu Wang, and Hongbin Zha. Local supports global: Deep camera relocalization with sequence enhancement. In *ICCV*, 2019. 2
- [52] Fei Xue, Xin Wu, Shaojun Cai, and Junqiu Wang. Learning multi-view camera relocalization with graph neural networks. In *CVPR*, 2020. 2
- [53] Kwang Moo Yi, Eduard Trulls, Yuki Ono, Vincent Lepetit, Mathieu Salzmann, and Pascal Fua. Learning to find good correspondences. In *CVPR*, 2018. 1, 3
- [54] Jiahui Zhang, Dawei Sun, Zixin Luo, Anbang Yao, Lei Zhou, Tianwei Shen, Yurong Chen, Long Quan, and Honggen Liao. Learning two-view correspondences and geometry using order-aware network. In *ICCV*, 2019. 1, 2, 3, 6, 7
- [55] Zichao Zhang, Torsten Sattler, and Davide Scaramuzza. Reference pose generation for long-term visual localization via learned features and view synthesis. *IJCV*, 2021. 6, 7, 9, 10, 15
- [56] Chen Zhao, Yixiao Ge, Feng Zhu, Rui Zhao, Hongsheng Li, and Mathieu Salzmann. Progressive correspondence pruning by consensus learning. In *ICCV*, 2021. 3, 6, 7
- [57] Zhen Zhong, Guobao Xiao, Linxin Zheng, Yan Lu, and Jiayi Ma. T-Net: Effective Permutation-Equivariant Network for Two-View Correspondence Learning. In *ICCV*, 2021. 1, 3

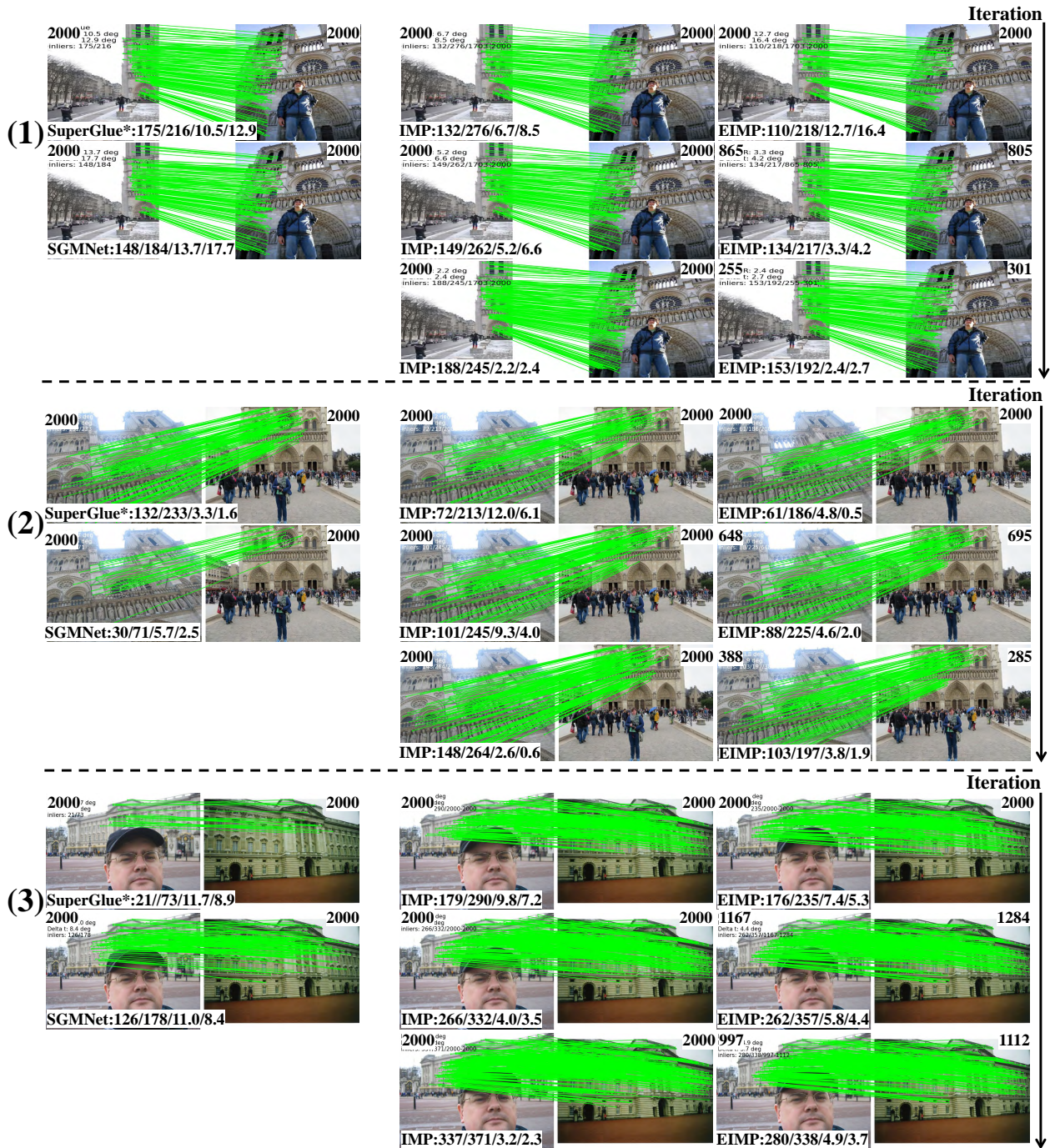


Figure 10. **Qualitative results on YFCC100m dataset** [45] At the left-bottom of each image pair, we report the number of inliers/matches and rotation/translation errors of SuperGlue* [38] (official SuperGlue), SGMNet [13], our IMP and EIMP. Besides, the number of keypoints in each image are shown at the top of each pair. For simple case (1), all methods give a large number of inliers, but IMP and EIMP report smaller pose errors even when using fewer inliers. Note that instead of keeping all keypoints, EIMP effectively reduces the number of keypoints from 2,000 to 865 and 255 in the iteration process, significantly decreasing the time complexity for self and cross attention computation 7. For tough cases ((2), (3)), due to large viewpoint and illumination changes, SuperGlue* and SGMNet fail to report comparable number of inliers to our models, resulting in higher rotation and translation errors. In contrast, our models still progressively increase the number of inliers from different regions in the iteration process. These well distributed inliers lead to smaller pose errors. By comparing the results of EIMP in (2) and (3), we see that the number of preserved keypoints are based on the number of potential inliers in the image pair: more potential inliers result in more retained keypoints. That is because our sampling strategy is fully adaptive.

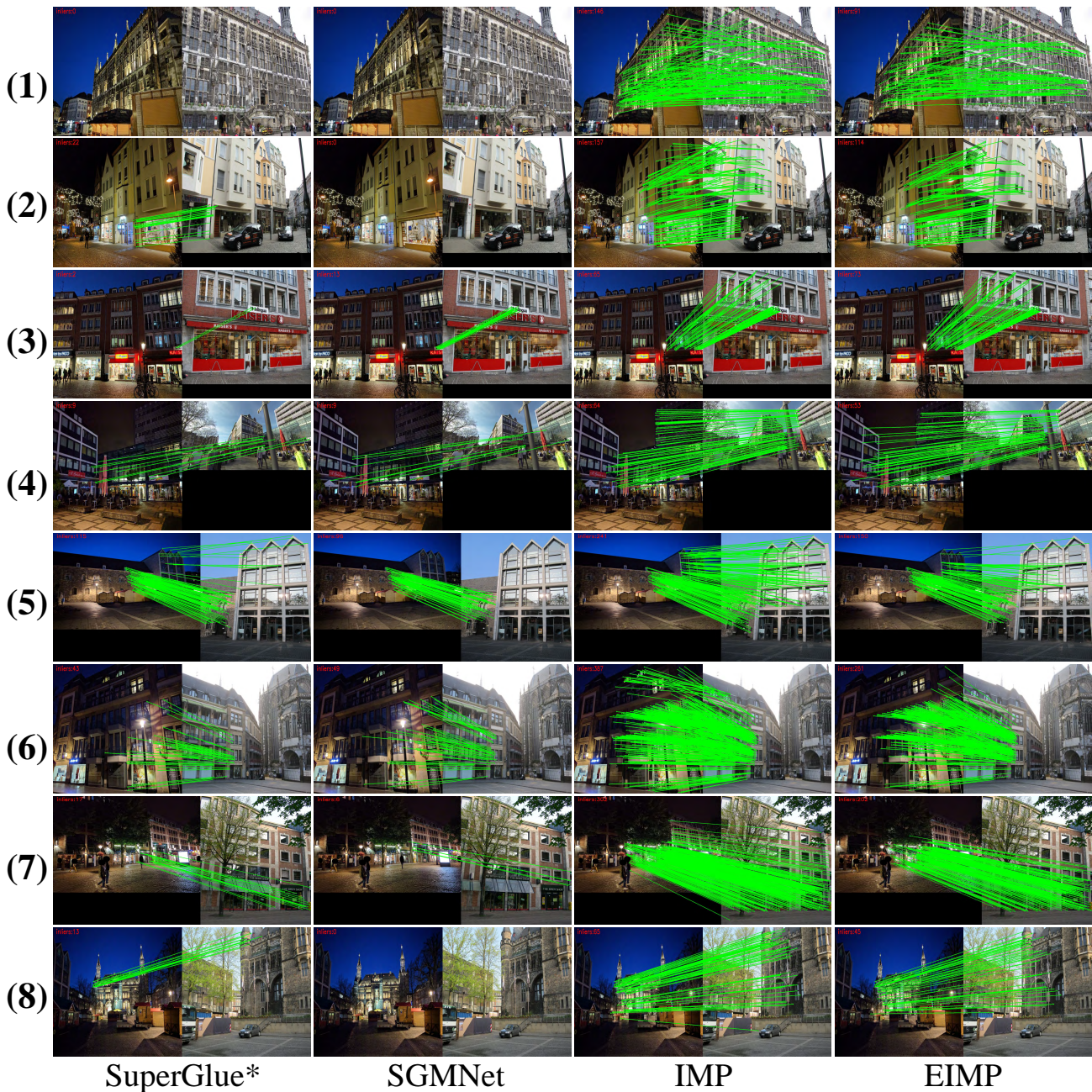


Figure 11. **Qualitative results on Aachen v1.1 dataset [39, 55].** We visualize inliers of query (left) and reference (right) images under larger viewpoint ((1)-(4)) and illumination ((5)-(8)) changes of SuperGlue* [38] (official SuperGlue), SGMNet [13], our IMP and EIMP. As the groundtruth poses of each query images are not available, we utilize HLoc [37] framework to visualize inliers given by PnP [26] + RANSAC [21]. Testing pairs (1)-(4) show that when image pairs have large viewpoint changes, both SuperGlue* and SGMNet fail to find enough correct matches. That is because geometric constraints are more useful for finding matches in two images with large viewpoint changes and both SuperGlue* and SGMNet ignore it. However, we embed the geometric information into the matching module, so our IMP and EIMP work much better, guaranteeing the localization success. For cases (5)-(8), when query images have large illumination changes with reference images, corresponding keypoints from two images are less discriminative, so SuperGlue* and SGMNet only give slightly more inliers than for images with large viewpoint. As our model additionally leverages geometric constraints to find matches, both IMP and EIMP successfully obtain a large number of inliers. Note that compared to SuperGlue* and SGMNet, both IMP and EIMP find inliers from the almost the whole overlap regions of the two images as opposed to some clusters ((2)-(5), (8)).