# DeepMatcher: A Deep Transformer-based Network for Robust and Accurate Local Feature Matching

Tao Xie[†], Kun Dai[†], Ke Wang, Ruifeng Li, Lijun Zhao

*Abstract*—Local feature matching between images remains a challenging task, especially in the presence of significant appearance variations, e.g., extreme viewpoint changes. In this work, we propose DeepMatcher, a deep Transformer-based network built upon our investigation of local feature matching in detector-free methods. The key insight is that local feature matcher with deep layers can capture more human-intuitive and simpler-to-match features. Based on this, we propose a Slimming Transformer (SlimFormer) dedicated for DeepMatcher, which leverages vector-based attention to model relevance among all keypoints and achieves long-range context aggregation in an efficient and effective manner. A relative position encoding is applied to each SlimFormer so as to explicitly disclose relative distance information, further improving the representation of keypoints. A layer-scale strategy is also employed in each Slim-Former to enable the network to assimilate message exchange from the residual block adaptively, thus allowing it to simulate the human behaviour that humans can acquire different matching cues each time they scan an image pair. To facilitate a better adaption of the SlimFormer, we introduce a Feature Transition Module (FTM) to ensure a smooth transition in feature scopes with different receptive fields. By interleaving the self- and cross-SlimFormer multiple times, DeepMatcher can easily establish pixel-wise dense matches at coarse level. Finally, we perceive the match refinement as a combination of classification and regression problems and design Fine Matches Module to predict confidence and offset concurrently, thereby generating robust and accurate matches. Experimentally, we show that DeepMatcher significantly outperforms the state-of-the-art methods on several benchmarks, demonstrating the superior matching capability of DeepMatcher. The code is available at https://github.com/XT-1997/DeepMatcher.

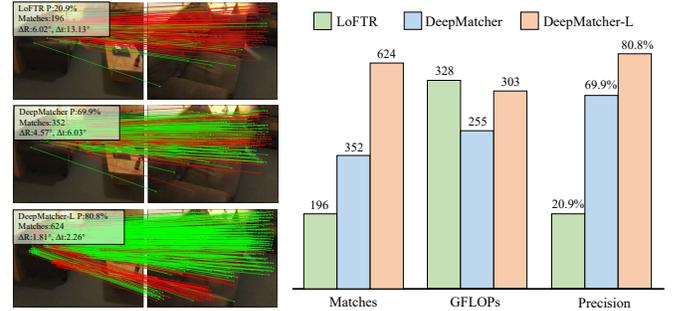*Index Terms*—Local feature matching, Pose Estimation, Trans-former.

Fig. 1. **The comparison between LoFTR and DeepMatcher under large viewpoint changes.** DeepMatcher families considerably outperform LoFTR with more dense and precise matches while using less GFLOPs.

## I. INTRODUCTION

LOCAL feature matching [1]–[4] is the prerequisite for a variety of geometric computer vision applications, including Simultaneous Localization and Mapping (SLAM) [5], [6] and Structure-from-Motion (SFM) [7], [8]. As a broadly acknowledged matching pipeline, detector-based matching [3], [9]–[18] is typically accomplished by (i) detecting and describing a set of sparse keypoints such as SIFT [9], ORB [10], and learning-based equivalents [13], [19], (2) instituting point-to-point correspondences via nearest neighbour search

or more advanced matching algorithms [20]. The use of a feature detector narrows the matching search space, revealing the general efficacy of such detector-based matching process. Nonetheless, when dealing with image pairs with severe viewpoint variations, such pipeline struggles to build reliable correspondences since the detectors are essentially incapable of extracting repeated keypoints in this case.

Parallel to the detector-based matching, another stream of research seeks to establish correspondences directly from original images by extracting visual descriptors on dense grids across an image, thus assuring that substantial repeating keypoints could well be captured [2], [21]–[28]. Earlier detector-free matching works [21]–[25] generally depend on iterative convolution based on correlation or cost volume to identify probable neighbourhood consensus. Transformer [29] has recently attracted considerable interest in computer vision due to its excellent model capability and superior potentials for capturing long-range relationships. On the basis of this insight, various studies base their modelling of long-range relationships on Transformer backbone [2], [26]–[28]. As a representative work, LoFTR [2] updates features by repeatedly interleaving the self- and cross-attention layers, and replace vanilla Transformer with linear Transformer [30] to achieve manageable computation cost. These detector-free methdos can generate repeatable keypoints in indistinct regions with poor textures or motion blur, thus yielding impressive results. After witnessing the success of detector-free methods, an intriguing issue arises: could we build a deeper yet compact local feature matcher to further improve performance while reducing computing costs? Intuitively, when individuals match images, they scan the images back and forth, and the more times they scan them, the easier it is for them to remember the easier-to-match features, which indicates that a deep local feature matcher could display superior matching ability.

However, as demonstrated in LoFTR, doubling the number of Transformer layers has little effect on the results, which is contrary to expectations. In this work, we argue that the following obstacles hinder us from developing a deep local feature matcher for detector-free methods:

(i) Typical detector-free methods begin with a convolution neural network (CNN) as the basic feature extractor, followed by Transformer layers to capture long-range relevance so that generating credible correspondences. In terms of context ranges, there is apparently a gap between the global receptive field of Transformer and the local neighborhood of CNN, which is detrimental to subsequent stages involving deep feature interaction.

(ii) The translation invariance of CNN causes ambiguity in scenes with recurring geometry patterns or symmetrical structures. Current detector-free methods utilize absolute position encodings before Transformer layers to tackle this issue, while the position information would disappear as the Transformer layers grow deeper. Moreover, humans naturally associate items across observations by referring to not only their absolute position but also their relative position.

(iii) Intuitively, the depth of the network is more prominent than width in the field of feature matching. However, as the linear Transformer layer in LoFTR goes deeper, the model fails to learn effective context aggregation from deeper layers since the linear Transformer uses a context-agnostic manner to approximate self-attention, which cannot efficaciously simulate relevance among all keypoints.

To this end, we propose DeepMatcher, a deep local feature matching network that can produce more human-intuitive and simpler-to-match features for accurate correspondence with less computational complexity, as shown in Fig. 1. Firstly, we utilize a CNN network to generate pixel tokens with enriched features. Secondly, an **Feature Transition Module (FTM)** is introduced to ensure a smooth transition from the locally aggregated features extracted by CNN to features with a global receptive field extracted by Transformer. Then, we propose a **Slimming Transformer (SlimFormer)** to build deep network that strengthens long-range global context modelling intra-/inter-images. Technically, SlimFormer leverages vector-based attention that efficiently handles pixel tokens with linear complexity for robust long-range global context aggregation. Besides, a relative position encoding is applied to each SlimFormer to clearly express relative distance information, boosting the network's capacity to convey information, particularly in deeper layers. Moreover, SlimFormer utilizes a layer-scale strategy that enables the network to assimilate message exchange from the residual block adaptively, thus allowing it to simulate the human behavior that human can receive different matching information each time they scan an image pair. By interleaving the self- and cross-SlimFormer multiple times, DeepMatcher learns the discriminative features to construct dense matches at the coarse level by **Coarse Matches Module (CMM)**. Ultimately, we view the match refinement as a combination of classification and regression problems and devise **Fine Matches Module (FMM)** to predict confidence and offset concurrently, obtaining robust and accurate matches.

To summarize, the main contributions of this work are as follows:

- We propose DeepMatcher, a deep Transformer-based network for local feature matching, achieving state-of-the-art results on various benchmarks.
- We propose a Feature Transition Module (FTM) to ensure a smooth transition from the locally aggregated features extracted by CNN to features with a global receptive field extracted by SlimFormer.
- We propose a Slimming Transformer (SlimFormer) that integrates long-range global context aggregation, relative position encoding, and layer-scale strategy to enable DeepMatcher to be extended into dozen layers.
- We propose Fine Matches Module (FMM) that views the match refinement as a combination of classification and regression problems to optimze coarse matches, deriving robust and accurate matches.

## II. RELATED WORK

### A. Detector-based Methods

The conventional pipeline of detector-based matching systems detects two sets of keypoints, describes them with high-dimensional vectors, and then implements a matching algorithm to generate matches between the two sets of keypoints [31], [32].

Regarding feature detection and description, there are numerous handcrafted methods that seek to strike a balance between accuracy and efficiency [9], [10], [33]. However, the handcrafted descriptors are fragile when coping with image pairs with extreme appearance variations. With the development of deep learning [34]–[37], numerous approaches leverage elaborate convolution neural network (CNN) to extract robust feature representations, hence achieving superior performance. SuperPoint [11] builds a large dataset of pseudo-ground truth interest point locations in real images, supervised by the interest point detector itself, as opposed to a large-scale human annotation. D2-Net [12] makes the collected keypoints more stable by delaying the detection to a later stage. Subsequently, the aforementioned methods utilize the nearest neighbor search, followed by a robust estimator, such as RANSAC or its variants [38]–[42], to find matches between the retrieved keypoints.

Recent researches [3], [16]–[18] has interpreted local feature matching as a graph matching problem involving two sets of features. These methods utilize keypoints as nodes to construct graph neural network (GNN), employ the self- and cross-attention layers in Transformer to exchange global visual and geometric messages across nodes, and then generate the matches in accordance with soft assignment matrixes. Typically, SuperGlue [3] utilizes self- and cross-attention in Transformer to integrate global context information, followed by the Sinkhorn algorithm to generate matches according to the soft assignment matrix. Nonetheless, the matrix multiplication in vanilla Transformer results in quadratic complexity with respect to the number of keypoints, making SuperGlue costly to deal with substantial keypoints. To tackle this problem, many approaches attempt to ameliorate the structure of

SuperGlue. SGMNet [16] exploits the sparsity of graph neural network to lower the computation complexity. ClusterGNN [18] employs a progressive clustering module adaptively to divide keypoints into different subgraphs to reduce computation. However, limited to inherent essence, detector-based approaches are incapable of extracting repeated keypoints when handling image pairs with large appearance variations.

### B. Detector-free Methods

Detector-free methods exclude the feature detector and generate dense matches directly from the original images. Earlier detector-free matching researches [21]–[25] generally utilize convolutional neural network (CNN) based on correlation or cost volume to identify probable neighbourhood consensus. DRC-Net [23] generates a 4D correlation tensor from the coarse-resolution features, which is refined by a learnable neighborhood consensus module to generate matches. Patch2Pix [24] proposes a weakly supervised approach to learn matches that are consistent with the epipolar geometry of image pairs. DFM [25] uses pre-trained VGG architecture as a feature extractor and captures matches without any additional training strategy. Although elevating the matching accuracy, these methods extract ambiguous feature representations and fail to discriminate incorrect matches owing to the limited receptive field of CNN.

To handle this issue, LoFTR [2], the pioneering detector-free GNN method, utlizes Transformer to realize global context information exchange and extracts matches in a coarse-to-fine manner. Matchformer [26] proposes a human-intuitive extract-and-match scheme that interleaves self- and cross-attention in each stage of the hierarchical encoder. Such a match-aware encoder releases the overloaded decoder and makes the model highly efficient. QuadTree [43] proposes a novel Transformer structure that builds token pyramids and computes attention in a coarse-to-fine manner. Then, the QuadTree Transformer is integrated into LoFTR and achieves superior matching performance. TopicFM [27] applies a topic-modeling strategy to encode high-level contexts in images, which improves the robustness of matching by focusing on the same semantic areas between the images. ASpanFormer [28] proposes a Transformer-based detector-free architecture, in which the flow maps are regressed in each cross-attention phase to perform local attention. These detector-free methods are capable of generating repeatable keypoints in indistinct regions with poor textures or motion blur, thus resulting in amazing results. However, the architecture of existing detector-free methods is designed as shallow-broad, and building a deeper and more compact local feature matcher to further improve performance while reducing computational costs has not been investigated.

### C. Efficient Transformer

In the vanilla Transformer, the memory cost is quadratic to the length of sequences due to the matrix multiplication, which has become a bottleneck for Transformer when dealing with long sequences. Recently, several approaches have been proposed to improve the efficiency of Transformer [30], [44], [45]. Linear Transformer [30] expresses self-attention as a
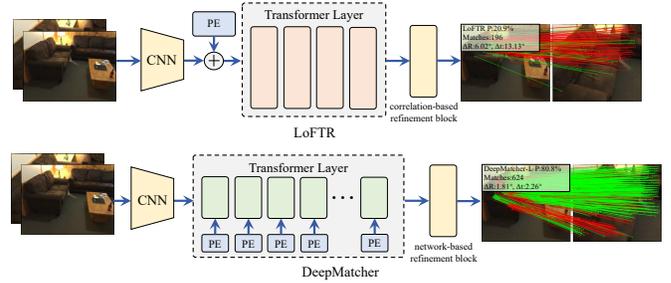


Fig. 2. **The comparison between LoFTR and DeepMatcher.** DeepMatcher designs a deep-narrow Transformer layers to capture more human-intuitive and simpler-to-match features. Besides, the position encoding (PE) is integrated to each Transformer layer to convey position information in deep layers. Moreover, a network-based refinement block is proposed to extract more precise matches.

linear dot product of kernel feature maps and makes use of the associativity property of matrix products to reduce the computational complexity. BigBird [44] combines local attention and global attention at certain positions and utilizes random attention on several randomly selected token pairs. FastFormer [45] uses additive attention mechanism to model global contexts, achieving effective context modeling with linear complexity.

## III. METHODOLOGY

### A. Overall

Intuitively, when humans match images, they scan the images back and forth, and the more times they scan them, the simpler it is for them to recall the easier-to-match features, which suggests that a deep local feature matcher can exhibit higher matching abilities. Thus, as shown in Fig. 2, we consider depth of the network is more prominent than width and present a deep Transformer-based network, namely DeepMatcher. As shown in Fig. 3, given the image pair $I_A$ and $I_B$, our network produces reliable and accurate matches across images in an end-to-end manner. The matching process starts with a CNN-based encoder to extract the fine-level features $\bar{F}_A, \bar{F}_B$ and coarse-level features $\hat{F}_A, \hat{F}_B$. Before feeding these features to Slimming Transformer (SlimFormer), we utilize the Feature Transition Module (FTM) to guarantee a smooth transition to SlimFormer. Then, we utilize SlimFormer to achieve long-range global context aggregation intra-/inter-images in an efficient and effective way. A relative position encoding is applied on each SlimFormer to explicitly model relative distance information, hence enhancing the DeepMatcher's ability to convey information, particularly in deeper layers. A layer-scale strategy is also leveraged in each SlimFormer to enables the network to assimilate message exchange from the feed-forward modules adaptively, thus allowing it to simulate the human behavior that human can receive different matching information each time they scan an image pair. After interleaving SlimFormer by $L$ times, the enhanced features $^L F_A^{seq}$ and $^L F_B^{seq}$ are utilized to establish coarse matches $H_c$, which are further optimized to fine matches $H_f$ using Correspondence Refine Module (CRM). In the following part, we introduce the details and underlying insights of each individual block.
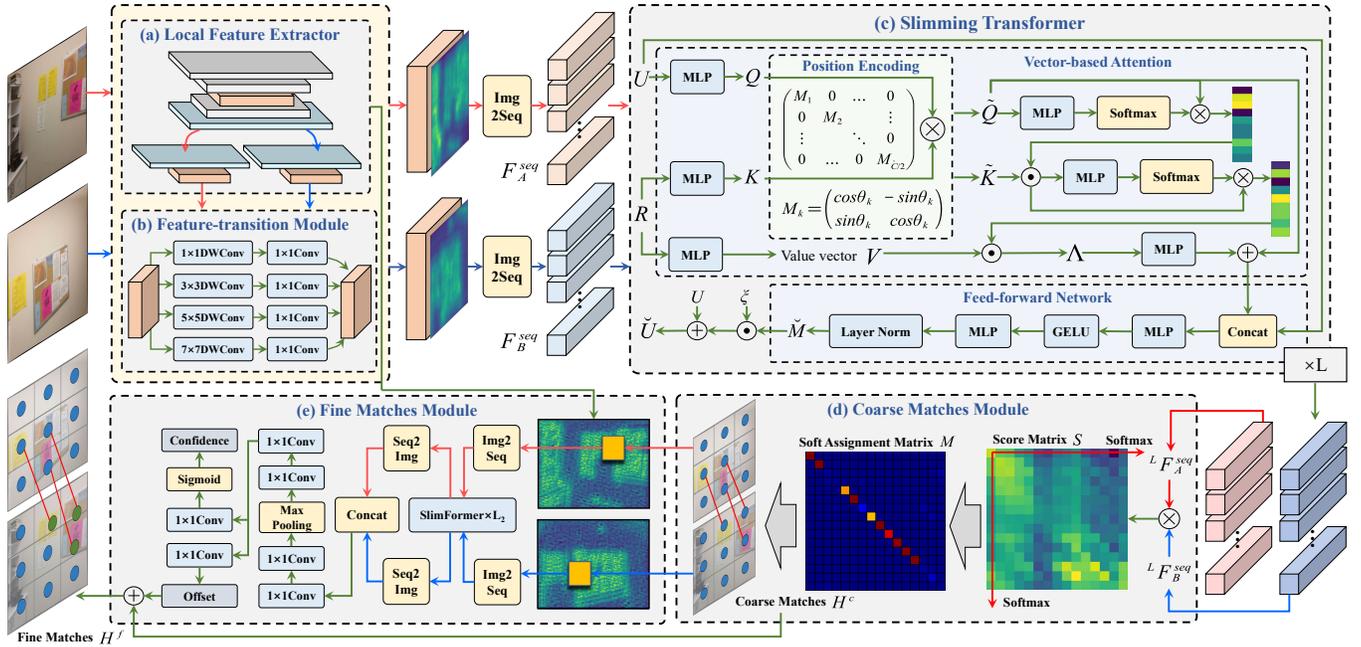
Fig. 3. **The network architecture of DeepMatcher.** DeepMatcher takes an image pair$(I_A, I_B)$ as input and generates transitional features from *Local Feature Extractor* and *Feature Transition Module*. Then, DeepMatcher interleaves *Slimming Transformer* by $L$ times to perform long-range context aggregation. *Coarse Matches Module* is utilized to establish coarse matches, which are optimize to fine matches by *Fine Matches Module*.

## B. Local Feature Extractor

As the first part of DeepMatcher, we use a standard convolutional neural network (CNN) with FPN [46] to extract the coarse-level features $\hat{F}_A$, $\hat{F}_B \in \mathbb{R}^{\hat{C} \times H/8 \times W/8}$, and the fine-level features $\bar{F}_A$, $\bar{F}_B \in \mathbb{R}^{\bar{C} \times H/2 \times W/2}$ for the image pair $I_A$ and $I_B$, where $H$ and $W$ are the height and width of the original images, $\hat{C}$, $\bar{C}$ denote feature dimension. For convenience, we denote $N = H/8 \times W/8$ as the number of pixel tokens. Since each pixel in $\hat{F}_A$, $\hat{F}_B$ represents an $8 \times 8$ grid in the original images $I_A, I_B$, we view the central position of all grids as the pixel coordinates $P_A, P_B \in \mathbb{R}^{N \times 2}$ of keypoints.

## C. Feature Transition Module (FTM)

In the subsequent steps, we construct graph neural network (GNN) and propose SlimFormer that leverages self-/cross-attention in Transformer to aggregate global context information intra-/inter-image. Nevertheless, there is apparently a gap between the feature extractor and SlimFormer in terms of context ranges, which is deleterious to subsequent steps involving deep feature interaction. Besides, representing features at multiple scales is so critical for discriminating objects or regions of varying sizes that it can ensure prominent features at various scales can be preserved for deep features aggregation. Thus, we propose a Feature Transition Module (FTM) inserted between the local feature extractor and SlimFormer to adjust the receptive fields of the extracted features, ensuring effective deep feature interaction in SlimFormer. Specifically, instead of directly using $(1 \times 1, 3 \times 3, 5 \times 5, 7 \times 7)$ convolution, FTM adopts $(1 \times 1, 3 \times 3, 5 \times 5, 7 \times 7)$ depth-wise convolution [47] followed by $1 \times 1$ point-wise convolution to reduce model parameters and computation, obtaining multi-scale feature representations. Then, we concatenate the features along channel

dimension, hence enlarging the receptive fields of $\hat{F}_A, \hat{F}_B$. Finally, we derive the updated $F_A^{ftm}, F_B^{ftm} \in \mathbb{R}^{\hat{C} \times H/8 \times W/8}$, which can be formulated as:

$$
\begin{aligned}
FTM(F) = [C_1^{1/4}(DW_1(F)) || C_1^{1/4}(DW_3(F)) || \\
C_1^{1/4}(DW_5(F)) || C_1^{1/4}(DW_7(F))], \\
F_A^{ftm} = FTM(\hat{F}_A), \quad F_B^{ftm} = FTM(\hat{F}_B),
\end{aligned}
\tag{1}
$$

where $C_1^{1/4}$ means using $1 \times 1$ convolution to squeeze the channel dimension to $\hat{C}/4$; $DW_1, DW_3, DW_5, DW_7$ mean depth-wise convolution with kernel size of 1, 3, 5, 7, respectively; $[\cdot||\cdot]$ means concatenation along the channel dimension.

## D. Slimming Transformer (SlimFormer)

We flatten the updated enhanced features $F_A^{ftm}, F_B^{ftm}$ to be the input sequence for deep feature aggregation, obtaining $F_A^{seq}, F_B^{seq} \in \mathbb{R}^{N \times \hat{C}}$. Following [3], we view keypoints with features $F_A^{seq}, F_B^{seq}$ in image pairs as nodes to construct GNN, in which the global context aggregation intra-/inter-image is performed. Intuitively, more observations between images can result in more precise matches, indicating that deep feature interaction is essential for local features matching task. Nevertheless, the ablation study of LoFTR demonstrates that the matching performance has not been significantly improved with more Transformer layers. We attribute this phenomenon to the following reasons: (i) LoFTR only utilizes absolute position encoding before Transformer layers, where the position information would disappear when the Transformer layers grow deeper. Moreover, humans primarily associate objects by referring to their relative positions. (ii) The linear Transformer utilized in LoFTR uses a context-agnostic manner to approximate self-attention, which cannot

fully model relevance among all keypoints, especially in deep layers. To handle this dilemma, we propose SlimFormer that leverages relative position information and global context information to boost the capability of DeepMatcher to convey abundant information, hence extracting discriminative feature representations $^{L}F_{A}^{seq}, ^{L}F_{B}^{seq} \in \mathbb{R}^{\hat{C} \times H/8 \times W/8}$.

**Vector-based Attention (VAtt) Layer.** Instead of using a context-agnostic manner to approximate self-attention, we convert query vector to global query contexts and leverage element-wise product to model relevance among all keypoints. Technically, during each feature enhancement process, we utilize self-/cross-attention to aggregate long-range context information intra-/inter-images. For self-attention, the input features $U$ and $R$ are same (either $(F_{A}^{seq}, F_{A}^{seq})$ or $(F_{B}^{seq}, F_{B}^{seq})$). For cross-attention, the input features $U$ and $R$ are different (either $(F_{A}^{seq}, F_{B}^{seq})$ or $(F_{B}^{seq}, F_{A}^{seq})$). Firstly, SlimFormer transforms the input features $U$ and $R$ into the query, key, and value vectors $Q, K, V \in \mathbb{R}^{N \times \hat{C}}$.

$$Q = UW_Q, \quad K = RW_K, \quad V = RW_V, \quad (2)$$

where $W_Q, W_K, W_V \in \mathbb{R}^{\hat{C} \times \hat{C}}$ denote learnable weights for feature transformation. Then, we perform relative position encoding on query vector $Q$ and key $K$.

$$\tilde{Q} = DPE(Q), \quad \tilde{K} = DPE(K), \quad (3)$$

where $DPE(\cdot)$ means relative position encoding operation, described below.

Next, modeling the context information of the input features based on the interactions among $\tilde{Q}, \tilde{K}$, and $V$ is a critical problem for Transformer-like architectures. In the vanilla Transformer, dot-product attention mechanism leads to quadratic complexity, making it unrealistic to establish deep Transformer layers. A potential method to reduce the computational complexity is to summarize the attention matrices before modeling their interactions. Inspired by [45], we introduce vector-based attention that effectively models long-range interactions among pixel tokens to alleviate this bottleneck. Instead of computing a quadratic attention map $QK^T$ that encodes all possible interactions between candidate matches, we form a compact representation of query-key interactions via vector-based attention that computes the correlation between global query vector and each key vector. Specifically, we firstly leverage MLP to calculate the weight $\tilde{Q}_{imp} \in \mathbb{R}^{1 \times N}$ of each query vector:

$$\tilde{Q}_{imp} = Softmax(MLP(\tilde{Q})), \quad (4)$$

where $Softmax(\cdot)$ means softmax operation.

The global query vector $\breve{Q} \in \mathbb{R}^{1 \times \hat{C}}$ is set to be a linear combination of $\tilde{Q}$:

$$\breve{Q} = \tilde{Q}_{imp} \otimes \tilde{Q} \quad (5)$$

where $\otimes$ means matrix multiplication.

Then, we utilize the element-wise multiplication between the global query vector $\breve{Q}$ and each key vector to model their interaction, obtaining context-aware key vector $\tilde{K}_Q \in \mathbb{R}^{N \times \hat{C}}$:

$$\tilde{K}_Q = \breve{Q} \odot \tilde{K}, \quad (6)$$

where $\odot$ denotes element-wise multiplication.

We utilize a similar vector-based attention to extract global context-aware key vector $\breve{K}_Q$ and model the interaction between $\breve{K}_Q$ and $V$:

$$\tilde{K}_{Qimp} = Softmax(MLP(\tilde{K}_Q))$$
$$\breve{K}_Q = \tilde{K}_{Qimp} \otimes \tilde{K}_Q \quad (7)$$
$$\Lambda = \breve{K}_Q \odot V$$

Subsequently, we employ a MLP and short-cut structure to derive the global message $M \in \mathbb{R}^{N \times \hat{C}}$.

$$M = MLP(\Lambda) + \tilde{Q} \quad (8)$$

For convenience, we define the process of vector-based attention layer as:

$$M = VAtt(U, R) \quad (9)$$

**Feed-forward Network (FFN).** Inspired by conventional Transformers, we employ a feed-forward network applied to $M$ to extract discriminative features for effectively deep features aggregation. The feed-forward network consists of two fully-connected layers and a GELU activation function. The hidden dimension between the two fully-connected layers is extended by a scale rate $\gamma$ to learn abundant feature representation. This process can be formulated as:

$$FFN(U, M) = MLP_{1/\gamma}(GELU(MLP_{\gamma/2}([U||M]))), \quad (10)$$

where $MLP_{1/\gamma}, MLP_{\gamma/2}$ mean expand the channel dimension by $1/\gamma, \gamma/2$ times with a MLP, respectively; $[\cdot||\cdot]$ means concatenation along channel dimension; $GELU(\cdot)$ means GELU activation function. Ultimately, we obtain enhanced message $\breve{M} \in \mathbb{R}^{N \times \hat{C}}$.

**Layer Scale Strategy.** Intuitively, people obtain different message after observing images each time, which inspires us to propose a layer-scale strategy. Specifically, in accordance with ResNet [48], we utilize a shortcut structure to realize efficient training. Then, we design a learnable scaling factor $\xi$ to adaptively balance original features $U$ and enhanced message $\breve{M}$, which is formulated as.

$$\breve{U} = U + \xi \breve{M} \quad (11)$$

By incorporating $\xi$ into SlimFormer, SlimFormer can easily simulate the human behaviour that humans acquire different matching cues each time they scan an image pair.

**Relative Position Encoding (RPE).** The local feature extractor learns strict translation invariant features, which could cause ambiguity in scenes that have repetitive geometry texture or symmetric structures. Previous works [2], [3], [16], [26] attach a distinctive absolute positional embedding to each keypoint, thus alleviating such ambiguity. However, compared with absolute position, relative position is more conducive for humans to establish connections between objects. Therefore, we argue that incorporating the explicit relative position dependency during each deep feature aggregation is essential for distinguishing identical features. However, relative position is not applicable to transformers with linear complexity as they do not explicitly calculate the quadratic complexity attention

matrix. To this end, we employ rotary positional embedding (RoPE) [49] that leverages absolute position encoding to achieve relative position encoding without manipulating the attention matrix. Given a pixel token $T_i$ and its features $F_i \in \mathbb{R}^{\hat{C}}$, the rotary position encoding function is defined by:

$$
Pos(T_i, F_i) = \Theta(T_i)F_i = \begin{pmatrix} M_1 & 0 & \dots & 0 \\ 0 & M_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & M_{\hat{C}/2} \end{pmatrix} F_i,
$$
(12)

where $\Theta(T_i) \in \mathbb{R}^{\hat{C} \times \hat{C}}$ is a block diagonal matrix. Each block with size of $2 \times 2$ is defined by:

$$
M_k = \begin{pmatrix} cos\ i\theta_k & -sin\ i\theta_k \\ sin\ i\theta_k & cos\ i\theta_k \end{pmatrix}, \quad \theta_k = \frac{1}{10000^{2(k-1)/\hat{C}}} \quad (13)
$$

where $\theta_k$ encodes the index of the feature channel.

Compared to sinusoidal encoding [2], [3], [26], rotary positional embedding has two advantages: (i) $\Theta(\cdot)$ is an orthogonal function, the encoding only changes the feature's direction but not the feature's length, which could stabilize the learning process. (ii) The dot product of two encoded features $< Pos(T_i, F_i), Pos(T_j, F_j) >$ in self-attention of vanilla Transformer can be derived to:

$$
[\Theta(T_i)F_i]^T \Theta(T_j)F_j = (F_i)^T \Theta(T_j - T_i)F_j \quad (14)
$$

which means the relative 2D distance information can be explicitly revealed by the dot product.

Since RoPE injects position information by rotation, which maintains the norm of hidden representations unchanged, such positional encoding can be directly applied to linear complexity transformers as demonstrated in [49]. In SlimFormer, we implement this by employing rotary positional embedding into $Q, K$ to incorporate relative position information, as illustrated in Fig. 3 or Eq. (3). For more details about RoPE, we encourage readers to refer to original papers.

**Self-/Cross-SlimFormer.** In summary, the SlimFormer is formatted as:

$$
Slim(U, R) = U + \xi FFN(U, VAtt(U, R)) \quad (15)
$$

We perform $L$ times of SlimFormer for feature enhancement. During the $l$-th feature enhancement, we use self-/cross-attention mechanism to integrate intra-/inter-image information, which can be formulated as:

$$
\begin{aligned}
{}^{l-1}F_A^{seq} &= Slim({}^{l-1}F_A^{seq}, {}^{l-1}F_A^{seq}), \\
{}^{l-1}F_B^{seq} &= Slim({}^{l-1}F_B^{seq}, {}^{l-1}F_B^{seq}), \\
{}^{l}F_A^{seq} &= Slim({}^{l-1}F_A^{seq}, {}^{l-1}F_B^{seq}), \\
{}^{l}F_B^{seq} &= Slim({}^{l-1}F_B^{seq}, {}^{l}F_A^{seq})
\end{aligned}
$$
(16)

Ultimately, we incorporate relative position information and global context message into enhanced features ${}^{L}F_A^{seq}, {}^{L}F_B^{seq}$.
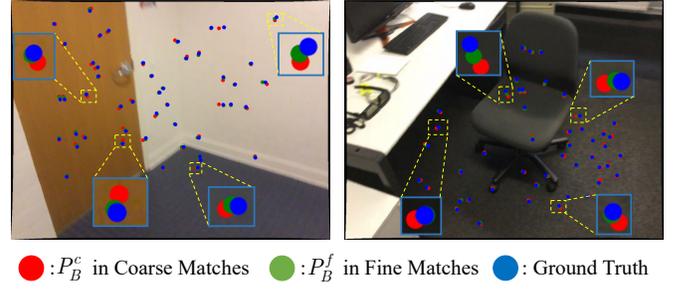
Fig. 4. **Visualization of refinement result.** The keypoints (green) in the fine matches approximate the ground-truth (blue).

### E. Coarse Matches Module (CMM)

Given ${}^{L}F_A^{seq}$ and ${}^{L}F_B^{seq}$, we utilize inner product of ${}^{L}F_A^{seq}, {}^{L}F_B^{seq}$ to calculate the score matrix $S \in \mathbb{R}^{N \times N}$.

$$
S(i, j) = \langle {}^{L}F_A^{seq}, {}^{L}F_B^{seq} \rangle, \quad (17)
$$

where $\langle \cdot, \cdot \rangle$ means the inner product. Subsequently, we apply softmax operator on both dimensions (denoted as dual-softmax operation) to convert the $S$ to soft assignment matrix $G \in \mathbb{R}^{N \times N}$:

$$
G = Softmax(S)_{col} \cdot Softmax(S)_{row}, \quad (18)
$$

where $Softmax(\cdot)_{col}$, $Softmax(\cdot)_{row}$ mean performing softmax on each column and row of $S$, respectively.

Then, for the $i$-th keypoint in $I_A$ and the $j$-th keypoint in $I_B$, we regard them as a pair of predicted coarse matches if they satisfy the following two conditions: (i) The soft assignment score is higher than a predefined threshold $\lambda$: $G(i, j) > \lambda$. (ii) They satisfy the mutual nearest neighbor (MNN) criteria, i.e., $G(i, j)$ is the maximum value in the corresponding row and column. Ultimately, we derive the index $D$ of anchor points in coarse matches:

$$
D = \{(i, j)|(i, j) \in \text{MNN}(M), M(i, j) > \lambda\} \quad (19)
$$

Given the index $D$ and the keypoints coordinates $P_A, P_B$, the coarse matches $H^c = \{(P_A^c, P_B^c)\}$ are formulated as:

$$
H^c = \{(P_A(i), P_B(j)) \mid \forall(i, j) \in D\} \quad (20)
$$

### F. Fine Matches Module (FMM)

After establishing coarse matches, a coarse-to-fine module is applied to refine these matches to the original picture resolution. However, the coarse-to-fine module in LoFTR only predicts the offset of the coarse matches without appraising whether the predicted matches are reliable. To tackle this issue, we view the match refinement as a combination of classification and regression problems and design Fine Matches Module to predict confidence and offset concurrently.

As shown in Fig. 3, for each coarse match, we locate its position at fine-level feature maps and crop two sets of local image patches with the size of $w \times w$, obtaining local features $\bar{F}_A^w, \bar{F}_B^w \in \mathbb{R}^{K \times \bar{C} \times w \times w}$, where $K$ is the number of coarse matches. Then, we flatten $\bar{F}_A^w, \bar{F}_B^w$ to be sequences, implement SlimFormer to perform $L_2$ times of global information passing, and rearrange the sequences into 2D feature maps,

obtaining $^{L_2}\bar{F}_A^w, {}^{L_2}\bar{F}_B^w$. The feature maps are concatenated along channel dimension and fed into a network, which is comprised of two convolution layers, a max pooling layer, and four convolution layers. The network predicts the offset $\Delta \in \mathbb{R}^{K \times 2}$ of the $P_B^c$ and the confidence $c \in \mathbb{R}^{K \times 1}$ of the predicted coarse matches:

$$
\begin{aligned}
\bar{F}_{mid} &= C_1(C_1(P_{max}(C_1(C_1([^{L_2}\bar{F}_A^w||^{L_2}\bar{F}_B^w]))))), \\
c &= Sig(C_1(\bar{F}_{mid})), \quad \Delta = C_1(\bar{F}_{mid}),
\end{aligned}
\tag{21}
$$

where $P_{max}$ means global max pooling operation; $C_1(\cdot)$ means $1 \times 1$ convolution; $[\cdot||\cdot]$ denotes concatenation along the channel dimension; $Sig(\cdot)$ means sigmoid function.

Ultimately, we obtain the fine matches $H^f = \{(P_A^f, P_B^f)\}$:

$$
H^f = \{(P_A^c(i), P_B^c(i) + \Delta(i)) \mid i \in \{1, 2, 3, ..., K\}\} \tag{22}
$$

### G. Loss

DeepMatcher generates final dense matches according to soft assignment matrix $G$ and offset $\Delta$. Therefore, the total loss $L^{all}$ of DeepMatcher comprises of matching loss $L^m$, regression loss $L^r$, and classification loss $L^c$.

$$
L^{all} = L^m + \beta L^r + \phi L^c, \tag{23}
$$

where $\beta$ and $\phi$ are weighting coefficient.

**Matching Loss.** Following [2], we calculate the index $E^{gt}$ of the ground truth matches, which are utilized in conjunction with soft assignment matrix $G$ to calculate matching loss $L^m$ defined as focal loss [50].

$$
\begin{aligned}
L^m = -[ & \frac{1}{|E^{gt}|} \sum_{(i,j)\in E^{gt}} \alpha(1 - G(i,j))^\eta log\ G(i,j) + \\
& \frac{1}{N - |E^{gt}|} \sum_{(i,j)\notin E^{gt}} (1 - \alpha)G(i,j)^\eta log\ (1 - G(i,j))],
\end{aligned}
\tag{24}
$$

where $\alpha$ is a weighting factor; $\eta$ is a focusing parameter; $|E^{gt}|$ means the number of ground truth matches.

**Regression Loss.** For predicted matches $\{(P_A^f, P_B^f)\}$, we project $P_A^f$ in the first image to second image, deriving $P_B^{gt}$. Then, the ground truth offset $\Delta^{gt}$ is formulated as:

$$
\Delta^{gt} = P_B^{gt} - P_B^f \tag{25}
$$

According to predicted offset $\Delta$ and ground truth offset $\Delta^{gt}$, we define the regression loss $L^r$ as:

$$
L^r = \frac{1}{K} \sum_{i=1}^{K} \|\Delta^{gt}(i) - \Delta(i)\|_2^2, \tag{26}
$$

where $K$ is the number of predicted matches. Notably, we ignore the predicted matches with $\Delta^{gt}$ larger than predefined threshold $\psi$.

**Classification Loss.** For the predicted matches with ground truth offset less than $\psi$, we regard them as positive and define the classification label as 1, while other matches are viewed as negative. Ultimately, we obtain the ground truth

confidence $c^{gt}$, while are utilized to calculate classification loss $L^c$ together with predicted confidence $c$.

$$
L^c = -\frac{1}{K} \sum_{i=1}^{K} \left[ c^{gt}(i)log\ c(i) + (1 - c^{gt}(i))log\ (1 - c(i)) \right]
\tag{27}
$$

## IV. EXPERIMENTS

### A. Implementation Details

**Architecture details.** We adopt a slightly modified ResNet-18 with FPN for local feature extraction. We use a width of 96 for the stem layer, followed by widths of [96, 128, 192] for the next three stages. We construct the FPN with levels $P_1$ through $P_3$ and take $P_3$ features as the coarse-level features, $P_1$ features as the fine-level features. Thus, the dimensions of fine-level and coarse-level feature maps are $\overline{C} = 96$, $\hat{C} = 192$, respectively. The scale rate $\gamma$ in feed-forward network is set to 4. Following SuperGlue, we set the confidence threshold $\lambda = 0.2$ to obtain coarse matches. Besides, we choose $w = 5$ to crop local windows in fine-level feature maps for matches refinement. To reconcile the coarse matching loss, regression loss, and classification loss, we set both weighting coefficients $\beta$ and $\phi$ to 0.2. For matching loss, we set the weighting factor $\alpha = 0.25$ and the focusing parameter $\eta = 2$. When making classification labels, we set $\psi$ to 8. In this work, we elaborately design two versions of DeepMatcher that interleave SlimFormer by $L = 6, 10$ times for feature enhancement, resulting in **DeepMatcher** and **DeepMatcher-L**.

**Training scheme for Scannet [51].** We train DeepMatcher on Scannet [51] dataset with 32 Tesla V100 GPUs for indoor local feature matching. In accordance with LoFTR, we sample 200 image pairs per scene at each epoch and balance scene variants over iterations. We employ the AdamW solver for optimization with a weight decay of 0.1. The initial learning rate is set to $6 \times 10^{-4}$ and will decrease by 0.5 every 3 epochs. We use gradient clipping that is set to 0.5 to avoid exploding gradients.

**Training scheme for MegaDepth [52].** We train Deep-Matcher on MegaDepth [52] datasets with 32 Tesla V100 GPUs for outdoor local feature matching. Following LoFTR, we randomly sample 100 pairs from each sub-scene during each epoch of training. We train DeepMatcher for 30 epochs in total. We also employ the AdamW solver for optimization with a weight decay of 0.1. The initial learning rate is set to $8 \times 10^{-4}$, with a linear learning rate warm-up in 3 epochs from 0.1 to the initial learning rate. We decay the learning rate by 0.5 every 4 epochs starting from the 4-th epoch.

### B. Indoor Pose Estimation

Typically, indoor pose estimation task is hampered by motion blur and significant viewpoint shifts. There are commonly extensive regions of low textures in indoor scenes. To evaluate the performance of DeepMatcher in such situations, we conducted indoor pose estimation experiments.

**Dataset.** We use ScanNet [51] dataset to validate the effectiveness of DeepMatcher on indoor pose estimation task. ScanNet consists of 1513 RGB-D sequences with RGB images

TABLE I
**INDOOR POSE ESTIMATION EVALUATION** ON SCANNET DATASET. THE AUC@$(5°, 10°, 20°)$ IS REPORTED.

| Local features | Matcher | Pose estimation AUC | | |
|---|---|---|---|---|
| | | @5° | @10° | @20° |
| Detector-based Methods | | | | |
| D2-Net [12] | NN | 5.25 | 14.53 | 27.96 |
| ContextDesc [53] | ratio test [9] | 6.64 | 15.01 | 25.75 |
| | NN | 9.43 | 21.53 | 36.40 |
| | NN + OANet [54] | 11.76 | 26.90 | 43.85 |
| SuperPoint [11] | SuperGlue [3] | 16.16 | 33.81 | 51.84 |
| | SGMNet [16] | 15.40 | 32.06 | 48.32 |
| | DenseGAP [17] | 17.01 | 36.07 | 55.66 |
| | HTMatch [55] | 15.11 | 31.42 | 48.23 |
| Detector-free Methods | | | | |
| | LoFTR [2] | 22.06 | 40.80 | 57.62 |
| | QuadTree [43] | 24.90 | 44.70 | 61.80 |
| | MatchFormer [26] | 24.31 | 43.90 | 61.41 |
| ___ | ASpanFormer [28] | 25.60 | 46.00 | **63.30** |
| | DeepMatcher | 25.38 | 44.38 | 60.35 |
| | DeepMatcher-L | **27.32** | **46.25** | 62.49 |

TABLE II
**OUTDOOR POSE ESTIMATION EVALUATION** ON MEGADEPTH DATASET. THE AUC@$(5°, 10°, 20°)$ IS REPORTED.

| Local features | Matcher | Pose estimation AUC | | |
|---|---|---|---|---|
| | | @5° | @10° | @20° |
| Detector-based Methods | | | | |
| | SuperGlue [3] | 42.18 | 61.16 | 75.96 |
| SuperPoint [11] | DenseGAP [17] | 41.17 | 56.87 | 70.22 |
| | ClusterGNN [18] | 44.19 | 58.54 | 70.33 |
| Detector-free Methods | | | | |
| | DRC-Net [23] | 27.01 | 42.96 | 58.31 |
| | LoFTR [2] | 52.80 | 69.19 | 81.18 |
| | QuadTree [43] | 54.60 | 70.50 | 82.20 |
| ___ | TopicFM [27] | 54.10 | 70.10 | 81.60 |
| | MatchFormer [26] | 52.91 | 69.74 | 82.00 |
| | ASpanFormer [28] | 55.30 | 71.50 | 83.10 |
| | DeepMatcher | 55.71 | 72.25 | 83.49 |
| | DeepMatcher-L | **56.98** | **73.11** | **84.15** |

and corresponding ground-truth poses in indoor environments. Following [3], we select 230M image pairs with overlap values ranging from 0.4 to 0.8 as the training set and 1500 image pairs as the test set. All images are resized to $640 \times 480$.

**Evaluation Protocol.** In accordance with [2], [3], we report the area under the cumulative curve (AUC) of pose errors at the thresholds $(5°, 10°, 20°)$, where pose errors are defined as the maximum of translational and rotational errors between ground-truth poses and predicted poses by DeepMatcher. Specifically, given the predicted dense matches, we utilize OPENCV to calculate the essential matrix $E$ and relative pose $\tilde{T}$ of image pairs. Then, the pose errors $\Delta T$ are defined as the maximum of translational and rotational errors between ground-truth relative pose $T = [R|t]$ and estimated relative pose $\tilde{T} = [\tilde{R}|\tilde{t}]$:

$$\Delta T = max(\Delta t, \Delta R),$$
$$\Delta t = arccos(\frac{\tilde{t} \cdot t}{||\tilde{t}||_2 \cdot ||t||_2}), \quad \Delta R = arccos(\frac{tr(\tilde{R}^T R) - 1}{2}),$$

(28)

where $\Delta t$ and $\Delta R$ denote the translational error and rotational error, respectively; $R, t$ is the ground-truth rotation matrix and translation vector; $\tilde{R}, \tilde{t}$ mean the predicted rotation matrix and translation vector; $tr(\cdot)$ means the trace of a matrix.

Given the pose errors of all image pairs, we plot the cumulative error distribution curve, whose area at three thresholds $(5°, 10°, 20°)$ are computed as AUC@$(5°, 10°, 20°)$.

**Results.** As illustrated in Table I, we observe that the detector-free methods achieve superior performance than detector-based methods since the detector struggles to extract repeatable keypoints when handling image pairs with significant viewpoint change. Wherein, Deep-Matcher and DeepMatcher-L outperform all cutting-edge detector-based and detector-free methods by a great margin. More specifically, DeepMatcher-L surpasses detector-based method DenseGAP by $(10.31\%, 10.18\%, 6.83\%)$ in

terms of AUC@$(5°, 10°, 20°)$, demonstrating the superiority of detector-free structure. Compared with the pioneering method LoFTR, DeepMatcher-L realizes superior performance with the improvement of $(5.26\%, 5.45\%, 4.87\%)$. Furthermore, DeepMatcher-L outperforms the state-of-the-art detector-free method ASpanFormer by $(1.72\%, 0.25\%)$ in terms of AUC@$(5°, 10°)$, proving the deep Transformer architecture is essential to extract more human-intuitive and easier-to-match features. Additionally, DeepMatcher-L only consumes $77.65\%$ GFLOPs and achieves $26.95\%$ inference speed boost compared with ASpanFormer, as demonstrated in Table IV.

### C. Outdoor Pose Estimation

Outdoor pose estimation remains a challenging task owing to the intricate 3D geometry, extreme illumination and viewpoint changes. To demonstrate the efficacy of DeepMatcher in overcoming these obstacles, an outdoor pose estimation experiment is conducted.

**Dataset.** We utilize MegaDepth [52] to conduct the outdoor pose estimation experiment. MegaDepth contains 1M internet images from 196 different scenes. These images come from photo-tourism and contain challenging conditions, including large viewpoint and illumination variations. Following [2], [14], we select 100 image pairs each scene for training and 1500 image pairs for testing. Images are resized such that their longer dimensions are equal to 840.

**Evaluation Protocol.** We use the same evaluation metrics AUC@$(5°, 10°, 20°)$ as the indoor pose estimation task.

**Results.** As shown in Table II, we can observe that DeepMatcher families surpass other methods in all evaluation metrics. Specifically, DeepMatcher-L noticeably outperforms the cutting-edge detector-based method ClusterGNN by $(12.79\%, 14.57\%, 13.82\%)$ in AUC@$(5°, 10°, 20°)$ since the detector struggles to extract repeatable keypoints in image pairs with extreme viewpoint change. Besides, compared with the baseline approach LoFTR, DeepMatcher-L achieves superior performance with the improvement of
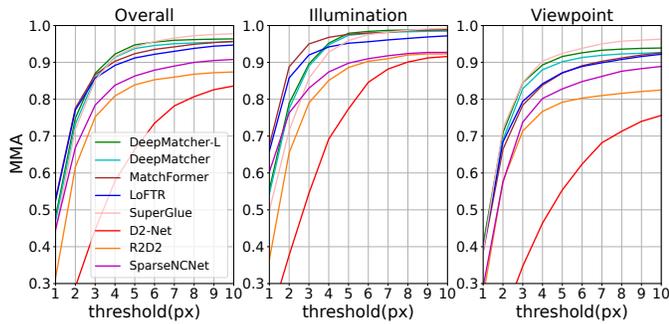
Fig. 5. **Image matching evaluation** on HPatches dataset.



Fig. 6. **Visualization of the predicted matches**. The mismatches, whose reprojection errors are larger than 5px, are colored red.

$(4.18\%, 3.92\%, 2.97\%)$. Moreover, DeepMatcher-L also surpasses the state-of-the-art detector-free method ASpanFormer by $(1.68\%, 1.61\%, 1.05\%)$, further validating the rationality of the deep Transformer structures.

### D. Image Matching

As a fundamental visual task, image matching plays an important role in several applications. Therefore, we conduct an image matching experiment to validate the performance of DeepMatcher.

**Dataset.** We conduct homography estimation experiments on the HPatches dataset [56]. Following [12], we select 108 sequences from HPatches. Each sequence consists of a ground-truth homography matrix and 6 images of progressively larger illumination (52 sequences with illumination changes) or viewpoint changes (56 sequences with viewpoint changes).

**Evaluation Protocol.** We adopt the generally employed mean matching accuracy (MMA) as metric, i.e., the average proportion of correct correspondences per image pair [12]. Specifically, the keypoints from the $i$-th query image are projected into the reference image by using the provided homography matrix $H_i$. Then, the matches with reprojection errors that are lower than a predefined threshold $t$ are deemed correct. Finally, we compute the average percentage of correct matches across all image pairs and define MMA as:

$$MMA(t) = \frac{1}{HP} \sum_{i=1}^{HP} \left( \frac{\sum_{j=1}^{N^f} \mathbb{1}(t - ||H_i(P_{A,i,j}^f) - P_{B,i,j}^f||_2)}{N^f} \right),$$
(29)

where $HP$ means the number of image pairs in HPatches; $N^f$ means the number of predicted matches; $\mathbb{1}(\cdot)$ is a binary indicator function whose output is 1 for non-negative value and 0 otherwise; $t$ is the threshold of reprojection error, varying from 1 to 10 pixels; $H_i(\cdot)$ means warping the keypoints in the $i$-th query image to reference image by ground-truth homography matrix; $(P_{A,i,j}^f, P_{B,i,j}^f)$ means the pixel coordinates of the $j$-th match in the $i$-th image pair.

**Results.** As illustrated in Fig. 5, we can observe that Deep-Matcher families achieve superior performance than detector-free methods (i.e. MatchFormer, LoFTR, SparseNCNet). Under varying illumination conditions, DeepMatcher yields inferior performance at low thresholds, while achieving outstanding performance when the threshold is larger than 5. Moreover, when handling image pairs with viewpoint changes, DeepMatcher exhibits extremely superior robotness compared
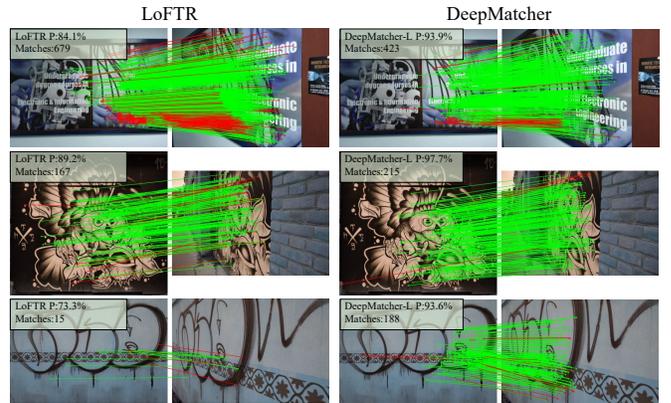
TABLE III
**HOMOGRAPHY ESTIMATION EVALUATION** ON HPATCHES DATASET.

| Local features | Matcher | Overall | Illumination | Viewpoint |
|---|---|---|---|---|
| | | CCM ($\varepsilon < 1/3/5$ pixels) | | |
| Detector-based Methods | | | | |
| D2-Net [12] | NN | 0.38/0.71/0.82 | 0.66/**0.95**/0.98 | 0.12/0.49/0.67 |
| R2D2 [13] | NN | 0.47/0.77/0.82 | 0.63/0.93/**0.98** | 0.32/0.64/0.70 |
| ASLFeat [57] | NN | 0.48/0.81/0.88 | 0.62/0.94/**0.98** | 0.34/0.69/0.78 |
| | NN | 0.46/0.78/0.85 | 0.57/0.92/0.97 | 0.35/0.65/0.74 |
| SuperPoint [11] | SuperGlue [3] | 0.51/0.82/0.89 | 0.60/0.92/**0.98** | 0.42/0.71/0.81 |
| | SGMNet [16] | 0.52/**0.85/0.91** | 0.59/0.94/**0.98** | **0.46/0.74/0.84** |
| | ClusterGNN [18] | 0.52/0.84/0.90 | 0.61/0.93/**0.98** | 0.44/**0.74**/0.81 |
| Detector-free Methods | | | | |
| | SparseNCNet [22] | 0.36/0.65/0.76 | 0.62/0.92/0.97 | 0.13/0.40/0.58 |
| | Patch2Pix [24] | 0.50/0.79/0.87 | 0.71/**0.95**/0.98 | 0.30/0.64/0.76 |
| —— | LoFTR [2] | **0.55**/0.81/0.86 | 0.74/**0.95**/0.98 | 0.38/0.69/0.76 |
| | MatchFormer [26] | **0.55**/0.81/0.87 | **0.75/0.95**/0.98 | 0.37/0.68/0.78 |
| | DeepMatcher | 0.50/0.81/0.90 | 0.62/0.93/**0.98** | 0.38/0.70/0.81 |
| | DeepMatcher-L | 0.51/0.83/**0.91** | 0.64/0.94/**0.98** | 0.39/0.72/**0.84** |

with other detector-free methods. As shown in Fig. 6, we select three pairs of image pairs from HPatches dataset and visualize the matches predicted by LoFTR and DeepMatcher-L to further validate the robustness of DeepMatcher-L to viewpoint variations.

### E. Homography Estimation

Since the distribution and number of matches are essential to estimate reliable geometry relationship between image pairs, we conduct a homography estimation experiment to comprehensively evaluate the performance of DeepMatcher.

**Dataset.** We assess DeepMatcher on HPatches dataset, which is widely used for homography estimation task.

**Evaluation Protocol.** Following the corner correctness metric (CCM) utilized in [24], we report the percentage of image pairs with average corner errors $\varepsilon$ smaller than 1/3/5 pixels. Specifically, based on the predicted dense matches, we use OPENCV to calculate the homography matrix $\tilde{H}_i$ for the $i$-th image pair. Subsequently, four corners in the query image are projected into the reference image by using the ground-truth homography matrix $H_i$ and the predicted homography matrix
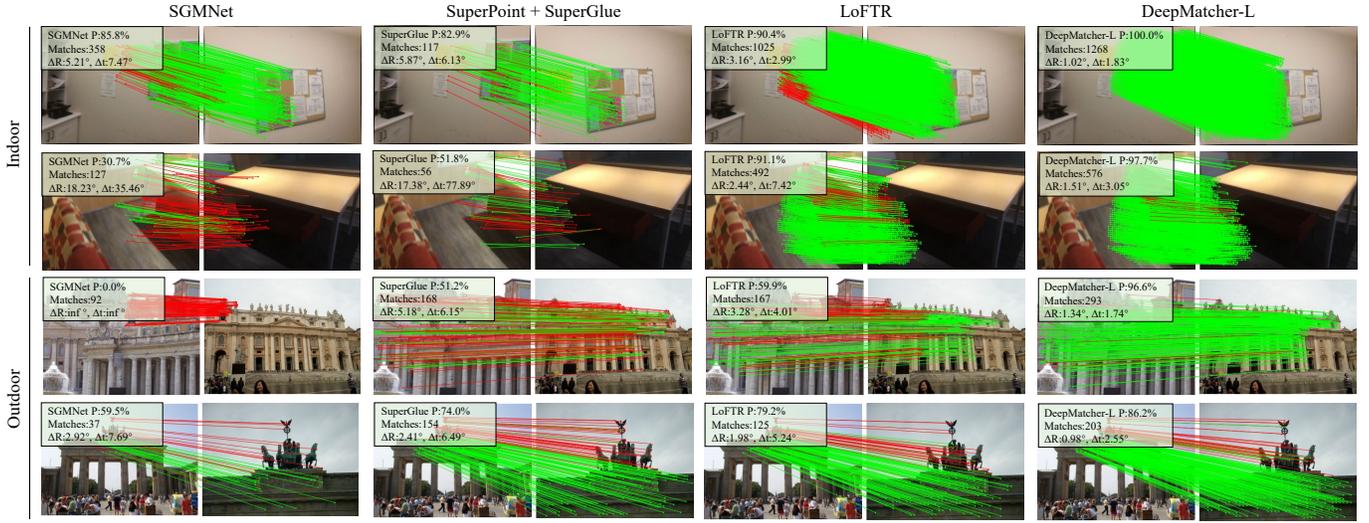
Fig. 7. **Visualization of the predicted matches.** The matches are colored by their reprojection errors (green indicates correct matches, and red indicates mismatches). We set the error threshold to 10 and 15 pixels for indoor and outdoor scenes.
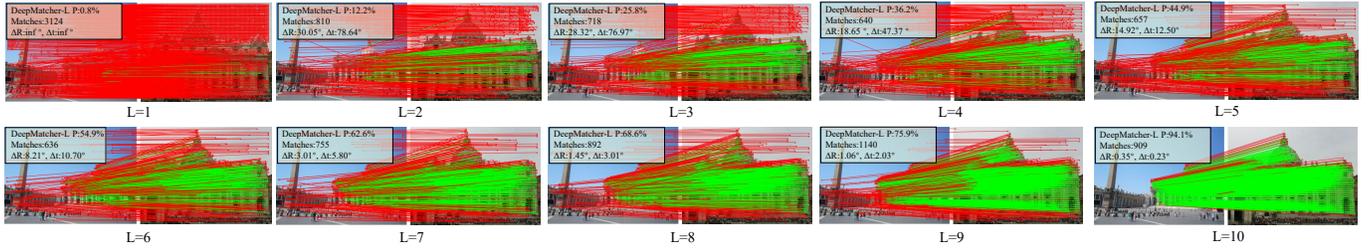


Fig. 8. The predicted matches of DeepMatcher-L after each SlimFormer. The threshold $\lambda$ of the soft assignment matrix is set to 0.

$\tilde{H}_i$, respectively. Ultimately, we calculate average reprojection error as corner error $\varepsilon_i$, thereby obtaining the CCM:

$$
\begin{aligned}
\varepsilon_i &= \frac{\sum_{p \in P_{co}} ||H_i(p) - \tilde{H}_i(p)||_2}{4}, \\
CCM(t) &= \frac{\sum_{i=1}^{HP} \mathbb{1}(t - \varepsilon_i)}{HP},
\end{aligned}
\tag{30}
$$

where $P_{co} = \{(0,0), (W_o - 1, 0), (0, H_o - 1), (W_o - 1, H_o - 1)\}$ means the four corners coordinates of the query image; $H_i(\cdot)$ and $\tilde{H}_i(\cdot)$ mean warping the corners in the $i$-th query image to reference image by ground-truth homography matrix and predicted homography matrix, respectively; $t \in \{1, 3, 5\}$ means the predefined threshold.

**Results.** As shown in Table III, DeepMatcher achieves the best performance among the detector-free methods under extreme viewpoint changes. Specifically, DeepMatcher outperforms LoFTR and MatchFormer with the improvement of $(1\%, 5\%)$ and $(2\%, 3\%)$ when thresholds are set to $3, 5$ pixels, repsectively. Furthermore, DeepMatcher-L surpasses LoFTR and MatchFormer with the improvement of $(1\%, 3\%, 8\%)$ and $(2\%, 4\%, 6\%)$. Besides, the detector-based methods are more robust to viewpoint variations, while the detector-free methods realize superior performance when handling image pairs with extreme illumination changes. In comparison, DeepMatcher strikes a decent balance when handling image pairs with various viewpoint and illumination changes.

### F. Understanding DeepMatcher

**Qualitative Results Visualization.** To further exhibit the capability of DeepMatcher to handle image pairs with extreme appearance settings, e.g., sparse texture, motion blur, large viewpoint and illumination changes, we visualize the matches predicted by SGMNet, SuperGlue, LoFTR, and DeepMatcher-L. As shown in Fig. 7, we can observe that DeepMatcher-L achieves dense and accurate matching performance.

**Visual Descriptors Enhancement Efficacy Analysis.** To validate the effectiveness of performing $L$ times of Slim-Former for feature enhancment, we visualize the matching results of DeepMatcher after each SlimFormer. As illustrated in Fig. 8, we can observe that the matching precision is promoted consistently, demonstrating interleaving SlimFormer can effectively integrate intra-/inter-image information, hence extracting easier-to-match features.

**Efficiency Analysis.** To validate the efficiency of Deep-Matcher, we compare several cutting-edge detector-free methods in terms of parameters, flops, and inference speed to determine their computational cost and storage consumption. We resize the input images to $640 \times 480$ and conduct all experiments on a single NVIDIA TITAN RTX GPU. When counting runtime, we run the test code 500 times and report the average time to eliminate occasionality. Notably, we only compare DeepMatcher families with other detector-free methods. As shown in Table IV, we can observe that DeepMatcher families realize competitive inference speed with

TABLE IV
**EFFICIENCY ANALYSIS.** SEVERAL APPROACHES ARE COMPARED IN TERMS OF PARAMETERS (MB), GFLOPS, AND RUNTIME (S). WE ALSO RECORD THE COMPUTATIONAL COMPLEXITY (TC) OF THE ATTENTION LAYER. $N$ DENOTES THE PIXLE TOKEN NUMBER, $C$ DENOTES THE FEATURE DIMENSION, $k$ DENOTES SELECTED TOKEN NUMBER, $r$ DENOTES DOWN-SCALE RATIO, $w$ DENOTES THE SAMPLE NUMBER.

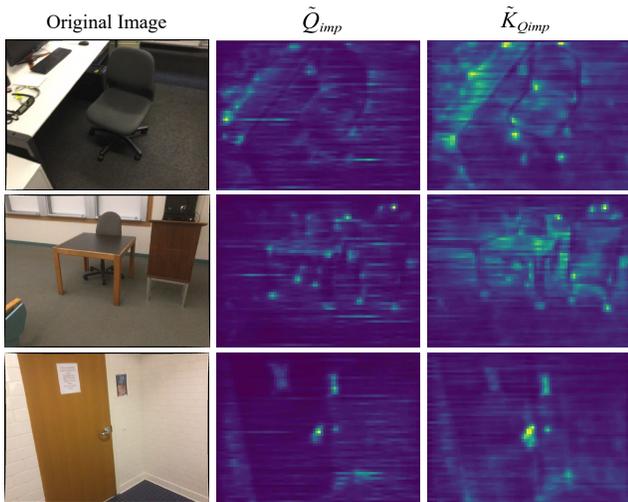| Methods | Params | GFLOPs | Runtime | TC |
|---|---|---|---|---|
| LoFTR [2] | 11.06 | 328.67 | **0.079** | $\mathcal{O}(NC^2)$ |
| QuadTree [43] | 13.21 | 382.01 | 0.152 | $\mathcal{O}(kNC)$ |
| MatchFormer [26] | 22.37 | 396.95 | 0.357 | $\mathcal{O}(N^2C/r)$ |
| AspanFormer [28] | 15.05 | 391.38 | 0.141 | $\mathcal{O}((N/r)^2C + NwC)$ |
| DeepMatcher | **10.96** | **255.99** | 0.089 | $\mathcal{O}(NC)$ |
| DeepMatcher-L | 15.51 | 303.90 | 0.103 | $\mathcal{O}(NC)$ |



Fig. 9. **Visualization of weights in SlimFormer structure**. SlimFormer emphasises keypoints at object boundaries to incorporate global context.

less GFLOPs since SlimFormer leverages element-wise product to model relevance among all keypoints. Compared with the baseline LoFTR, DeepMatcher and DeepMatcher-L only consume $(77.89\%, 92.46\%)$ GFLOPs. Moreover, compared with cutting-edge detector-free methods QuadTree, MatchFormer and AspanFormer, DeepMatcher-L exhibits more efficient matching performance with $(20.45\%, 23.44\%, 22.35\%)$ less GFLOPs and $(32.24\%, 71.15\%, 26.95\%)$ inference speed boost.

Furthermore, we also record the dominant computational complexity of the attention layers in various methods. As shown in Table IV, DeepMatcher achieves the minimum computational complexity of the attention layer. Specifically, compared with the baseline LoFTR, DeepMatcher reduces the dominant computational complexity from $\mathcal{O}(NC^2)$ to $\mathcal{O}(NC)$.

**Weight Analysis.** To explore which keypoints SlimFormer pays attention to when extracting global vectors, we visualize the weight $\tilde{Q}_{imp}$ and $\tilde{K}_{Qimp}$ in Eq. (4) and Eq. (7), respectively. As shown in Fig. 9, we can observe that SlimFormer primarily pays attention to the prominent keypoints at object boundaries that involve tremendous visual and geometry information. Consequently, SlimFormer exhibits puissant capability to aggregate global context information effectively.
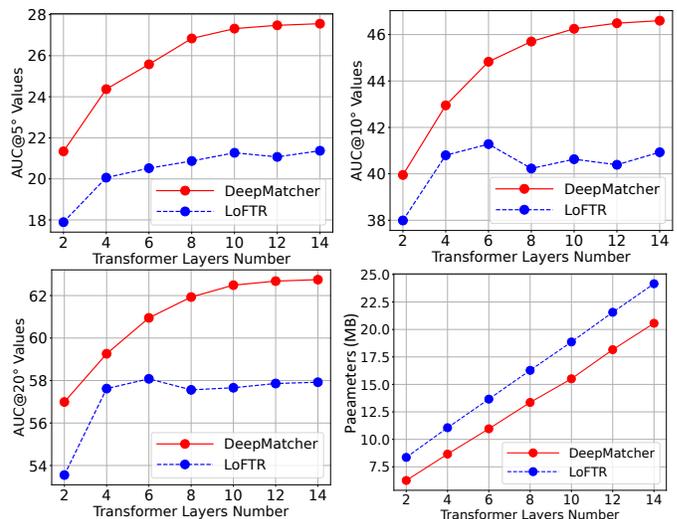


Fig. 10. **The AUC values and parameters** of DeepMatcher and LoFTR with the Transformer layers increasing.

**Deep Transformer Architecture Analysis.** To verify the opinion that deep Transformer architecture is essential to extract more human-intuitive and easier-to-match features, we record the indoor pose estimation precision with the number of Transformer layers increasing. As shown in Fig. 10, since SlimFormer leverages vector-based attention for robust long-range global context aggregation and utilizes layer-scale strategy and relative position encoding to enhance the representation of keypoints, the AUC values of DeepMatcher are promoted consistently with the Transformer layers going deep, while the accuracy of LoFTR is constantly fluctuating. Besides, the matching performance of DeepMatcher is significantly superior to LoFTR. Moreover, the parameters of DeepMatcher and LoFTR increase linearly with the number of Transformer layers, while DeepMatcher occupies fewer parameters.

*G. Ablation Study*

**Effect of the Proposed Modules.** To thoroughly validate the rationality of each module, we conduct indoor pose estimation experiments using different variants of DeepMatcher. As illustrated in Table V, we can observe that all components contribute to the outstanding performance of DeepMatcher.

(i), (ii) Using only self- and cross-SlimFormer layers leads to a severe decrease in matching performance, demonstrating interleaving the self- and cross-SlimFormer layers can effectively integrate intra-/inter-image message. (iii) Removing the Feature Transition Module results in a much lower accuracy $(-0.66\%, -0.47\%, -0.33\%)$, proving the effectiveness of ensuring smooth transition between feature extractor and SlimFormer in terms of context ranges. (iv) Removing Relative Position Encoding spawns a large drop in pose estimation accuracy $(-1.36\%, -1.67\%, -1.38\%)$, proving the relative position information is crucial to distinguish similar features. (v) Removing the Fine Matches Module results in lower AUC values $(-8.34\%, -12.56\%, -14.86\%)$, indicating the effectiveness of refining coarse matches.

**Effect of the Learnable Scale Factor** $\xi$**.** To validate that layer-scale strategy can simulate the human behaviour that

TABLE V
**ABLATION STUDY WITH DIFFERENT VARIANTS OF DEEPMATCHER** ON SCANNET DATASET.

| Methods | Pose estimation AUC | | |
|---|---|---|---|
| | @5° | @10° | @20° |
| (i) w only self-SlimFormer layers | 20.38 | 37.71 | 53.89 |
| (ii) w only cross-SlimFormer layers | 22.33 | 40.27 | 56.64 |
| (iii) w/o Feature Transition Module | 24.72 | 43.91 | 60.02 |
| (iv) w/o Relative Potition Encoding | 24.02 | 42.71 | 58.97 |
| (v) w/o Fine Matches Module | 17.04 | 31.82 | 45.49 |
| DeepMatcher full | **25.38** | **44.38** | **60.35** |

TABLE VI
**ABLATION STUDY WITH DIFFERENT LEARNABLE SCALE FACTORS** $\xi$ ON SCANNET DATASET.

| Methods | Pose estimation AUC | | |
|---|---|---|---|
| | @5° | @10° | @20° |
| w residual scale factor $\xi$ | **25.38** | **44.38** | **60.35** |
| w/o residual scale factor $\xi$ | 24.50 | 42.91 | 59.32 |

TABLE VII
**ABLATION STUDY WITH DIFFERENT POSITION ENCODING** ON SCANNET DATASET.

| Methods | Pose estimation AUC | | |
|---|---|---|---|
| | @5° | @10° | @20° |
| w/o position encoding | 24.02 | 42.71 | 58.97 |
| w absolute position encoding | 24.52 | 43.14 | 59.51 |
| w relative position encoding | **25.38** | **44.38** | **60.35** |

TABLE VIII
**ABLATION STUDY WITH DIFFERENT COARSE-TO-FINE MODULES** ON SCANNET DATASET.

| Methods | Pose estimation AUC | | |
|---|---|---|---|
| | @5° | @10° | @20° |
| Coarse-to-fine Module in LoFTR | 23.74 | 42.26 | 58.86 |
| Fine Matches Module | **25.38** | **44.38** | **60.35** |

humans can acquire different matching cues each time they scan an image pair to further improve matching performance, we remove the learnable scale factor $\xi$ and conduct an ablation experiment. As shown in Table VI, we can observe that introducing the residual scaling factor leads to superior performance.

**Effect of the Relative Position Encoding.** Humans leverage the relative position information to establish the connection between objects. Therefore, the relative position encoding is more conducive to realize elaborate scene parsing. To prove this opinion, we implement an ablation experiment using three structures: (i) Removing all position encoding in all SlimFormer. (ii) Using absolute position encoding proposed in LoFTR. (iii) Using relative position encoding. As shown in Table VII, we can observe that both absolute and relative position encoding boost AUC values, in which the relative position encoding exhibits more superior performance than absolute position encoding with the improvement of $(0.86\%, 1.24\%, 0.84\%)$.

**Effect of the Fine Matches Module.** Compared with the coarse-to-fine module used in LoFTR, FMM views the match refinement as a combination of classification and regression tasks. To validate the availability of FMM, we conduct an ablation experiment. As shown in Table VIII, we can observe that using FMM significantly achieves superior performance with the improvement of $(1.64\%, 2.12\%, 1.49\%)$, proving the rationality of predicting offset and confidence concurrently using a network.

## V. CONCLUSION

In this work, we propose DeepMatcher, a deep Transformer-based network for local feature matching. DeepMatcher simulates human behaviors when humans match image pairs, including: (1) Deep SlimFormer layers of the network to aggregate information intra-/inter-images; (2) Layer-scale strategy to assimilate message exchange from each layer adaptively. Besides, relative position encoding is applied to each layer so as to explicitly disclose relative distance information, hence improving the representation of DeepMatcher. We also propose Fine Matches Module to refine the coarse matches, thus generating robust and accurate matches. Extensive experiments demonstrate that DeepMatcher surpasses state-of-the-art approaches on several benchmarks, confirming the superior matching capability of DeepMatcher.

## REFERENCES

[1] J. Chen, S. Chen, X. Chen, Y. Dai, and Y. Yang, "Csr-net: Learning adaptive context structure representation for robust feature correspondence," *IEEE Transactions on Image Processing*, vol. 31, pp. 3197–3210, 2022.

[2] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, "Loftr: Detector-free local feature matching with transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8922–8931.

[3] 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*, 2020, pp. 4938–4947.

[4] L. Zheng, G. Xiao, Z. Shi, S. Wang, and J. Ma, "Msa-net: Establishing reliable correspondences by multiscale attention network," *IEEE Transactions on Image Processing*, vol. 31, pp. 4598–4608, 2022.

[5] 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.

[6] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós, "Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam," *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021.

[7] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.

[8] H. Cui, D. Tu, F. Tang, P. Xu, H. Liu, and S. Shen, "Vidsfm: Robust and accurate structure-from-motion for monocular videos," *IEEE Transactions on Image Processing*, vol. 31, pp. 2449–2462, 2022.

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

[10] 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.

[11] 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 (CVPR) Workshops*, June 2018.

[12] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[13] J. Revaud, P. Weinzaepfel, C. De Souza, N. Pion, G. Csurka, Y. Cabon, and M. Humenberger, "R2d2: repeatable and reliable detector and descriptor," in *NeurIPS*, 2019.

[14] M. J. Tyszkiewicz, P. Fua, and E. Trulls, "Disk: Learning local features with policy gradient," in *NeurIPS*, 2020.

[15] Y. Xia and J. Ma, "Locality-guided global-preserving optimization for robust feature matching," *IEEE Transactions on Image Processing*, vol. 31, pp. 5093–5108, 2022.

[16] H. Chen, Z. Luo, J. Zhang, L. Zhou, X. Bai, Z. Hu, C.-L. Tai, and L. Quan, "Learning to match features with seeded graph matching network," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6301–6310.

[17] Z. Kuang, J. Li, M. He, T. Wang, and Y. Zhao, "Densegap: Graph-structured dense correspondence learning with anchor points," *arXiv preprint arXiv:2112.06910*, 2021.

[18] Y. Shi, J.-X. Cai, Y. Shavit, T.-J. Mu, W. Feng, and K. Zhang, "Clustergnn: Cluster-based coarse-to-fine graph neural network for efficient feature matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 517–12 526.

[19] 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.

[20] 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.

[21] I. Rocco, M. Cimpoi, R. Arandjelović, A. Torii, T. Pajdla, and J. Sivic, "Neighbourhood consensus networks," *Advances in neural information processing systems*, vol. 31, 2018.

[22] I. Rocco, R. Arandjelović, and J. Sivic, "Efficient neighbourhood consensus networks via submanifold sparse convolutions," in *European Conference on Computer Vision*. Springer, 2020, pp. 605–621.

[23] X. Li, K. Han, S. Li, and V. Prisacariu, "Dual-resolution correspondence networks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 346–17 357, 2020.

[24] Q. Zhou, T. Sattler, and L. Leal-Taixe, "Patch2pix: Epipolar-guided pixel-level correspondences," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4669–4678.

[25] U. Efe, K. G. Ince, and A. Alatan, "Dfm: A performance baseline for deep feature matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4284–4293.

[26] Q. Wang, J. Zhang, K. Yang, K. Peng, and R. Stiefelhagen, "Matchformer: Interleaving attention in transformers for feature matching," *arXiv preprint arXiv:2203.09645*, 2022.

[27] K. Truong Giang, S. Song, and S. Jo, "Topicfm: Robust and interpretable feature matching with topic-assisted," *arXiv e-prints*, pp. arXiv–2207, 2022.

[28] H. Chen, Z. Luo, L. Zhou, Y. Tian, M. Zhen, T. Fang, D. McKinnon, Y. Tsin, and L. Quan, "Aspanformer: Detector-free image matching with adaptive span transformer," in *European Conference on Computer Vision*. Springer, 2022, pp. 20–36.

[29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[30] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are rnns: Fast autoregressive transformers with linear attention," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5156–5165.

[31] M. Karpushin, G. Valenzise, and F. Dufaux, "Keypoint detection in rgbd images based on an anisotropic scale space," *IEEE Transactions on Multimedia*, vol. 18, no. 9, pp. 1762–1771, 2016.

[32] V. Balntas, E. Riba, D. Ponsa, and K. Mikolajczyk, "Learning local feature descriptors with triplets and shallow convolutional neural networks." in *Bmvc*, vol. 1, no. 2, 2016, p. 3.

[33] X. Jiang, J. Ma, J. Jiang, and X. Guo, "Robust feature matching using spatial clustering with heavy outliers," *IEEE Transactions on Image Processing*, vol. 29, pp. 736–746, 2019.

[34] M. Sun, W. Suo, P. Wang, Y. Zhang, and Q. Wu, "A proposal-free one-stage framework for referring expression comprehension and generation via dense cross-attention," *IEEE Transactions on Multimedia*, 2022.

[35] K. Fu, Q. Zhao, and I. Y.-H. Gu, "Refinet: A deep segmentation assisted refinement network for salient object detection," *IEEE Transactions on Multimedia*, vol. 21, no. 2, pp. 457–469, 2018.

[36] Y. Song, X. Chen, X. Wang, Y. Zhang, and J. Li, "6-dof image localization from massive geo-tagged reference images," *IEEE Transactions on Multimedia*, vol. 18, no. 8, pp. 1542–1554, 2016.

[37] B. Fan, Y. Yang, W. Feng, F. Wu, J. Lu, and H. Liu, "Seeing through darkness: Visual localization at night via weakly supervised learning of domain invariant features," *IEEE Transactions on Multimedia*, 2022.

[38] 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.

[39] V. Fragoso, P. Sen, S. Rodriguez, and M. Turk, "Evsac: accelerating hypotheses generation by modeling matching scores with extreme value theory," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 2472–2479.

[40] D. Barath, J. Matas, and J. Noskova, "Magsac: marginalizing sample consensus," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 197–10 205.

[41] V. Mousavi, M. Varshosaz, F. Remondino, S. Pirasteh, and J. Li, "A two-step descriptor-based keypoint filtering algorithm for robust image matching," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–21, 2022.

[42] V. Fragoso, C. Sweeney, P. Sen, and M. Turk, "Ansac: Adaptive non-minimal sample and consensus," *arXiv preprint arXiv:1709.09559*, 2017.

[43] S. Tang, J. Zhang, S. Zhu, and P. Tan, "Quadtree attention for vision transformers," *arXiv preprint arXiv:2201.02767*, 2022.

[44] M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontanon, P. Pham, A. Ravula, Q. Wang, L. Yang *et al.*, "Big bird: Transformers for longer sequences," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 283–17 297, 2020.

[45] C. Wu, F. Wu, T. Qi, Y. Huang, and X. Xie, "Fastformer: Additive attention can be all you need," *arXiv preprint arXiv:2108.09084*, 2021.

[46] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.

[47] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[48] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[49] J. Su, Y. Lu, S. Pan, B. Wen, and Y. Liu, "Roformer: Enhanced transformer with rotary position embedding," *arXiv preprint arXiv:2104.09864*, 2021.

[50] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.

[51] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "Scannet: Richly-annotated 3d reconstructions of indoor scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5828–5839.

[52] Z. Li and N. Snavely, "Megadepth: Learning single-view depth prediction from internet photos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2041–2050.

[53] Z. Luo, T. Shen, L. Zhou, J. Zhang, Y. Yao, S. Li, T. Fang, and L. Quan, "Contextdesc: Local descriptor augmentation with cross-modality context," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2527–2536.

[54] J. Zhang, D. Sun, Z. Luo, A. Yao, L. Zhou, T. Shen, Y. Chen, L. Quan, and H. Liao, "Learning two-view correspondences and geometry using order-aware network," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5845–5854.

[55] Y. Cai, L. Li, D. Wang, X. Li, and X. Liu, "Htmatch: An efficient hybrid transformer based graph neural network for local feature matching," *Signal Processing*, p. 108859, 2022.

[56] 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.

[57] Z. Luo, L. Zhou, X. Bai, H. Chen, J. Zhang, Y. Yao, S. Li, T. Fang, and L. Quan, "Aslfeat: Learning local features of accurate shape and localization," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 6589–6598.