# Adaptive Spot-Guided Transformer for Consistent Local Feature Matching

Jiahuan Yu[*], Jiahao Chang[*], Jianfeng He, Tianzhu Zhang[†], Feng Wu

University of Science and Technology of China

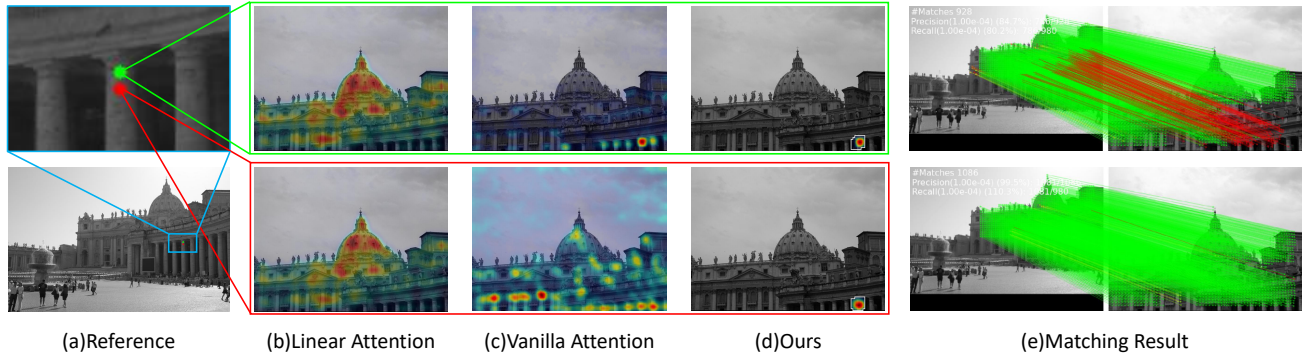(a)Reference      (b)Linear Attention      (c)Vanilla Attention      (d)Ours      (e)Matching Result

Figure 1. The visualization of the cross attention heatmaps and matching results. We sample two similar adjacent points in the reference image (a), marked with green and red. (b) are two heatmaps of the linear cross attention in LoFTR [50] when green and red pixels are queries. (c) are two heatmaps obtained from the vanilla cross attention. (d) are two heatmaps generated by our spot-guided attention. (e) are the comparison of the final matching results produced by LoFTR [50] (top) and our method (down).

## Abstract

*Local feature matching aims at finding correspondences between a pair of images. Although current detector-free methods leverage Transformer architecture to obtain an impressive performance, few works consider maintaining local consistency. Meanwhile, most methods struggle with large scale variations. To deal with the above issues, we propose Adaptive Spot-Guided Transformer (ASTR) for local feature matching, which jointly models the local consistency and scale variations in a unified coarse-to-fine architecture. The proposed ASTR enjoys several merits. First, we design a spot-guided aggregation module to avoid interfering with irrelevant areas during feature aggregation. Second, we design an adaptive scaling module to adjust the size of grids according to the calculated depth information at fine stage. Extensive experimental results on five standard benchmarks demonstrate that our ASTR performs favorably against state-of-the-art methods. Our code will be released on* https://astr2023.github.io.

## 1. Introduction

Local feature matching (LFM) is a fundamental task in computer vision, which aims to establish correspondence for local features across image pairs. As a basis for many 3D vision tasks, local feature matching can be applied in Structure-from-Motion (SfM) [49], 3D reconstruction [13], visual localization [48, 51], and pose estimation [18, 41]. Because of its broad applications, local feature matching has attracted substantial attention and facilitated the development of many researches [14, 27, 42, 44, 50]. However, finding consistent and accurate matches is still difficult due to various challenging factors such as illumination variations, scale changes, poor textures, and repetitive patterns.

To deal with the above challenges, numerous matching methods have been proposed, which can be generally categorized into two major groups, including detector-based matching methods [2, 14, 15, 39, 42, 47] and detector-free matching methods [9, 23, 27, 43, 44, 50]. Detector-based matching methods require to first design a keypoint detector to extract the keypoints between two images, and then establish matches between these extracted keypoints. The quality of detected keypoints will significantly affect the performance of detector-based matching methods. Therefore, many works aim to improve keypoint detection through multi-scale detection [36], repeatable and reliable

---
[*]Equal Contribution
[†]Corresponding Author

verification [42]. Thanks to the high-quality keypoints detected, these methods can achieve satisfactory performance while maintaining high computational and memory efficiency. However, these detector-based matching methods may have difficulty in finding reliable matches in textureless areas, where keypoints are challenging to detect. Differently, detector-free matching methods do not need to detect keypoints and try to establish pixel-level matches between local features. In this way, it is possible to establish matches in the texture-less areas. Due to the power of attention in capturing long-distance dependencies, many Transformer-based methods [9, 50, 52, 57] have emerged in recent years. As a representative work, considering the computation and memory costs, LoFTR [50] applies Linear Transformer [25] to aggregate global features at the coarse stage and then crops fixed-size grids for further refinement. To alleviate the problem caused by scale changes, COTR [24] calculate the co-visible area iteratively through attention mechanism. The promising performance of Transformer-based methods proves that attention mechanism is effective on local feature matching. Nevertheless, some recent works [28, 60] indicate Transformer lacks spatial inductive bias for continuous dense prediction tasks, which may cause inconsistent local matching results.

By studying the previous matching methods, we sum up two issues that are imperative for obtaining the dense correspondence between images. (1) **How to maintain local consistency.** The correct matching result usually satisfies the local matching consistency, i.e., for two similar adjacent pixels, their matching points are also extremely close to each other. Existing methods [24,50,57] utilize global attention in feature aggregation, introducing many irrelevant regions that affect feature updates. Some pixels are disturbed by noisy or similar areas and aggregate information from wrong regions, leading to false matching results. As shown in Figure 1 (b), for two adjacent similar pixels, highlighted regions of global linear attention are decentralized and inconsistent with each other. The inconsistency is also present in vanilla attention (see Figure 1 (c)). Therefore, it is necessary to utilize local consistency to focus the attention area on the correct place. (2) **How to handle scale variation.** In a coarse-to-fine architecture, since the attention mechanism at the coarse stage is not sensitive to scale variations, we should focus on the fine stage. Previous methods [9, 27, 50, 57] select fixed-size grids for matching at the fine stage. However, when the scale varies too much across images, the correct match point may be out of the range of the grid, resulting in matching failure. Hence, the scheme of cropping grids should be adaptively adjusted according to scale variation across views.

To deal with the above issues, we propose a novel Adaptive Spot-guided Transformer (ASTR) for consistent local feature matching, including a spot-guided aggregation mod-

ule and an adaptive scaling module. In the **spot-guided aggregation module**, towards the goal of maintaining local consistency, we design a novel attention mechanism called spot-guided attention: each point is guided by similar high-confidence points around it, focusing on a local candidate region at each layer. Here, we also adopt global features to enhance the matching ability of the network in the candidate regions. Specifically, for any point $p$, we pick the points with high feature similarity and matching confidence in the local area. Their corresponding matching regions are used for the next attention of point $p$. In addition, global features are applied to help the network to make judgments. The coarse feature maps are iteratively updated in the above way. With our spot-guided aggregation module, the red and green pixels are guided to the correct area, avoiding the interference of repetitive patterns (see Figure 1 (d)). In Figure 1 (e), our ASTR produces more accurate matching results, which maintains local matching consistency. In the **adaptive scaling module**, to fully account of possible scale variations, we attempt to adaptively crop different sizes of grids for alignment. In detail, we compute the corresponding depth map using the coarse matching result and leverage the depth information to crop adaptive size grids from the high-resolution feature maps for fine matching.

The contributions of our method could be summarized into three-fold: (1) We propose a novel Adaptive Spot-guided Transformer (ASTR) for local feature matching, including a spot-guided aggregation module and an adaptive scaling module. (2) We design a spot-guided aggregation module that can maintain local consistency and be unaffected by irrelevant regions while aggregating features. Our adaptive scaling module is able to leverage depth information to adaptively crop different size grids for refinement. (3) Extensive experimental results on five challenging benchmarks show that our proposed method performs favorably against state-of-the-art image matching methods.

## 2. Related Work

In this section, we briefly review several research lines that are related to sparse matching methods, dense matching methods, and vision Transformer.

**Local Feature Matching.** Local feature matching can categorized into detector-based and detector-free methods. Detector-based methods can be divided into three stages: feature detection, feature description, and feature matching. SIFT [34] and ORB [46] are the most popular hand-crafted local features, while learning-based methods [2, 14, 15, 19, 42, 46, 63] also obtain good performance improvement compared to classical methods. There are also some works focusing on improving the feature matching stage. D2Net [15] fuses the detection and description stages. R2D2 [42] attempts to train a network to find reliable and repeatable local features. SuperGlue [47] pro-

poses an attention-based GNN network to update extracted local features in alternating self and cross attentions. However, detector-based methods rely on local feature extractors, which limits the performance in challenging scenarios such as repetitive textures, weak textures, and illumination changes. Unlike detector-based approaches, detector-free approaches do not require a local feature detector, but find dense feature matching between pixels directly. The classical methods [21, 35] exists, but few of them outperform detector-based methods. Learning-based methods change the game, which can be divided into cost-volume-based methods [27,44,53,54] and Transformer-based methods [9, 10, 22, 24, 50, 57]. Good performance have been achieved by cost-volume-based methods, but most of them are limited by the small receptive field of CNN, which is overcome by Transformer-based methods [50]. Detector-free methods attain better performance in local feature matching, so we adopt this paradigm as the baseline.

**Vision Transformer.** Transformer [56] has been proven to be better at capturing long-range correlations than CNN in vision tasks [7, 37, 38]. Despite the great success, the computational cost of vanilla attention at high resolution is unacceptable, so some approximations [25, 33, 52, 59] have been proposed, which inevitably leads to performance degradation. Linear Attention [25] approximates softmax with ELU [11] to reduce the computational complexity to linear but degrades the focusing ability of attention. Swin-Transformer [33] limits attention in local windows, which harms the ability to establish long-range associations. At the same time, QuadTree [52] calculates attention in a coarse-to-fine manner, and ASpanFormer [9] proposes an adaptive method for selecting attention spans, but few of them consider local consistency. Different from the existing attention mechanism, we explicitly model local consistency in our spot-guided attention without introducing excessive computation and memory costs.

**Local Feature Matching with Scale Invariance.** Scale variation is one of the main challenges faced by local feature matching. Many works have explored solutions. Handcrafted local features [5, 31, 45, 46] use Gaussian pyramid model to alleviate the problem. Following the hand-crafted methods, Some learning-based descriptors [2, 4, 32, 36, 42, 63] also use the multi-scale representation. ScaleNet [3] and Scale-Net [17], instead, try to directly estimate the scale ratio. Another popular paradigm is to perform a wrap or scaling operation to eliminate the distortion caused by the scale variance. GeoWrap [6] introduces a homography regression and warps images to increase overlap area. OETR [10] limits the keypoint detection in estimated overlap areas. COTR [24] estimates scale by finding co-visible regions, and then finds correspondence by recursively zooming. However, most of above methods require significant modifications to the network architecture, and introduce addi-

tional computation overhead. Therefore, we design a fully pluggable, lightweight and training-free module for coarse-to-fine architecture.

# 3. Our Approach

In this section, we present our proposed Adaptive Spot-guided Transformer (ASTR) for Consistent Local Feature Matching. The overall architecture is illustrated in Figure 2.

## 3.1. Overview

As shown in Figure 2, the proposed ASTR mainly consists of two modules, including a spot-guided aggregation module and an adaptive scaling module. Here we give a brief introduction to the entire process. Given an image pair $I_{Ref}$ and $I_{Src}$, to start with, we extract multi-scale feature maps of each image through a shared Feature Pyramid Network (FPN) [30]. We denote feature maps with the size of $1/i$ as $F^{1/i} = \{F_{ref}^{1/i}, F_{src}^{1/i}\}$. Then, $F^{1/32}$ and $F^{1/8}$ are fed into the spot-guided aggregation module for coarse-matching and depth maps. Here, the coarse matching result is acquired in three phases. First, we need to compute the similarity matrix, which can be given by $S(i, j) = \tau \langle F_{ref}^{1/8}(i), F_{src}^{1/8}(j) \rangle$ with flattened features, where $\tau$ is the temperature coefficient. Then we perform dural-softmax operator on $S$ to calculate matching matrix $P_c$:

$$P_c(i, j) = \text{softmax}(S(i, :))(i, j) \cdot \text{softmax}(S(:, j))(i, j). \quad (1)$$

Finally, we use the mutual nearest neighbor strategy and the threshold $\theta_c$ to filter out the coarse-matching result $M_c$. According to depth information and coarse-matching result, we can crop different size grids on the high-resolution feature map $F^{1/2}$. After linear self and cross attention layers, features of the cropped grids are used to produce the final fine-level matching result.

## 3.2. Spot-Guided Aggregation Module

Correct matching always satisfies the local matching consistency, i.e., the matching points of two similar adjacent pixels are also close to each other in the other image. When humans establish dense matches between two images, they will first scan through the two images quickly and keep in mind some landmarks that are easier to match correctly. For those trouble points similar to surrounding landmarks, it is not easy to obtain correct matches in the beginning. But now, they can focus attention around the matching points of landmarks to revisit trouble points' matches. In this way, more correctly matched landmarks are obtained. After several iterations of the above process, eventually, they will get the matching result for the whole image. Inspired by this idea, we design a spot-guided aggregation module. Section 3.2.1 introduces the preliminaries of vanilla attention and linear attention. Section 3.2.2 describes our spot-guided attention mechanism. Section 3.2.3 demonstrates the design of the entire spot-guided aggregation module.
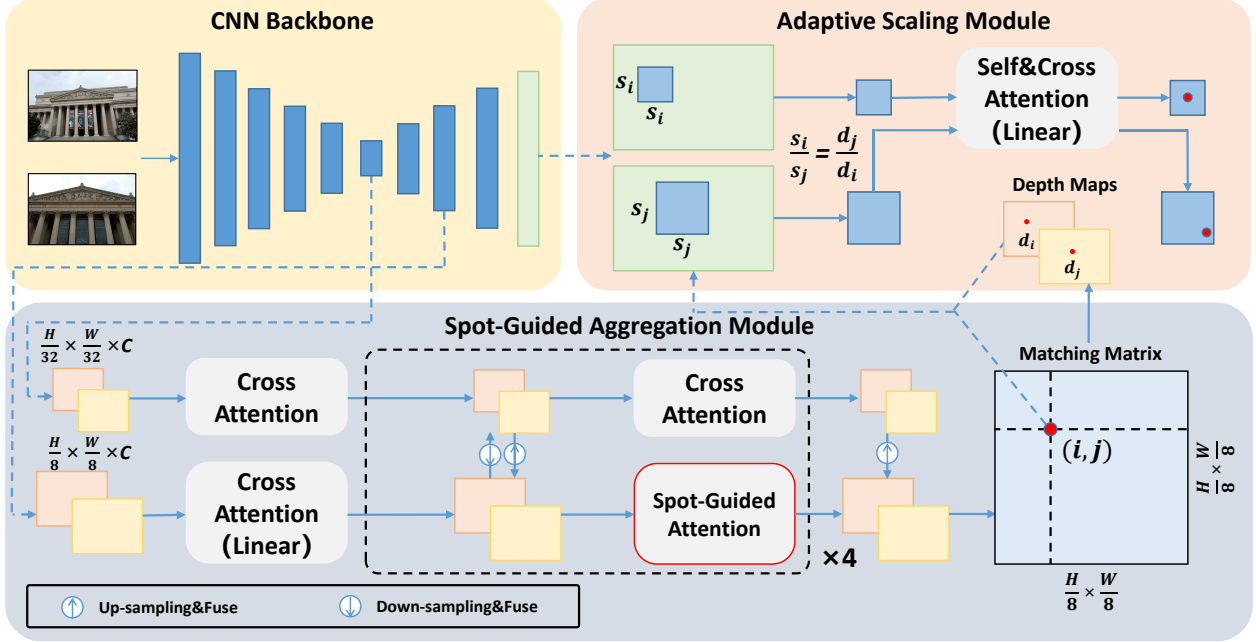
Figure 2. The architecture of ASTR. Our ASTR consists of two major components: spot-guided aggregation module and the adaptive scaling module. "Cross Attention" means vanilla cross attention, unless otherwise stated. Please refer to the text for detailed architecture.

### 3.2.1 Preliminaries

The calculation of vanilla attention requires three inputs: query $Q$, key $K$, and value $V$. The output of vanilla attention is a weighted sum of the value, where the weight matrix is determined by the query and its corresponding key. The process can be described as

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T)V. \quad (2)$$

However, in vision tasks, the size of the weight matrix $\text{softmax}(QK^T)$ increases quadratically as the image resolution grows. When the image resolution is large, the memory and computational cost of vanilla attention is unacceptable. To solve this problem, Linear attention [25] is proposed to replace the softmax operator with the product of two kernel functions:

$$\text{Linear\_attention}(Q, K, V) = \phi(Q)(\phi(K^T)V), \quad (3)$$

where $\phi(\cdot) = \text{elu}(\cdot) + 1$. Since the number of feature channels is much smaller than the number of pixels, the computational complexity is reduced from quadratic to linear.

### 3.2.2 Spot-Guided Attention

It is known from the local matching consistency that the matching points of similar adjacent pixels are also close to each other. In Figure 10, we illustrate the case that the reference image as query aggregates features from the source image. Given reference and source feature maps $F^{1/8} = \{F_{ref}^{1/8}, F_{src}^{1/8}\}$, we compute a matching matrix $P_s$
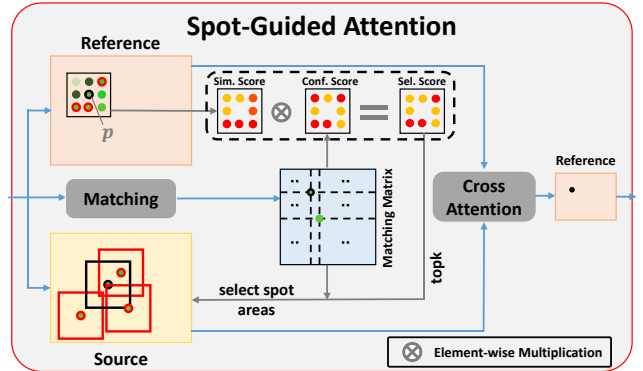


Figure 3. The illustration of our spot-guided attention.

across images. For any pixel $p$ in Figure 10, we first compute the similarity score $\text{S}_{\text{sim}}(p) \in R^{l^2-1}$ between $p$ and other pixels in the $l \times l$ area around $p$. Specifically, the similarity score can be obtained as

$$\text{S}_{\text{sim}}(p) = \text{softmax}_i \left( \left\{ \langle F_{ref}^{1/8}(p), F_{src}^{1/8}(p_i) \rangle \right\}_{p_i \in N(p)} \right), \quad (4)$$

where $\langle \cdot, \cdot \rangle$ is the inner product, and $N(p)$ is the set of pixels in the $l \times l$ field around pixel $p$. In addition, we should also consider the reliability of points in $N(p)$. For each $p_i \in N(p)$, confidence can be viewed as the highest similarity to all pixels on the source images. Meanwhile, we can also get the matching point position of $p_i$, denoted as $\text{Loc}(p_i)$. Hence, $\text{Loc}(p_i)$ and confidence score $\text{S}_{\text{conf}}(p) \in R^{l^2-1}$ can
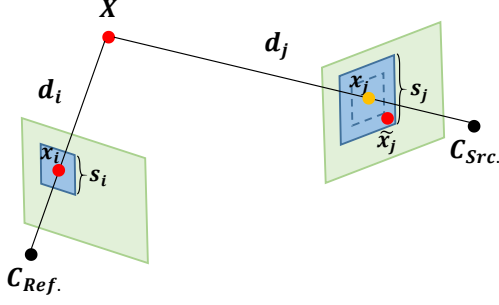
4

Figure 4. The illustration of our adaptive scaling module. On the left is the reference image, whose optical center is $C_{Ref}$. On the right is the source image, whose optical center is $C_{Src}$. $x_i$ and $x_j$ are the projections of the real-world point $X$.

be computed in the following way:

$$\begin{aligned} \text{S}_{\text{conf}}(p) &= \{\max(\text{P}_s(p_i,:))\}_{p_i \in N(p)} \cdot \\ \text{Loc}(p_i) &= \text{argmax}(\text{P}_s(p_i,:)), \end{aligned} \quad (5)$$

Combining two scores, we select $p$ and top-k points Topk$(p)$ whose matching points are used as seed points Seed$(p)$:

$$\begin{aligned} \text{Topk}(p) &= \{p\} \cup \text{topk}\{\text{S}_{\text{sim}}(p) \cdot \text{S}_{\text{conf}}(p)\}, \\ \text{Seed}(p) &= \{\text{Loc}(q)\}_{q \in \text{Topk}(p)}, \end{aligned} \quad (6)$$

Following that, we extend $l \times l$ regions centered on these seed points Seed$(p)$ on $I_{src}$, which are the spot areas of $p$. Finally, cross attention is performed between $p$ and corresponding spot areas. After exchanging the source image and the reference image, the source feature map is updated in the same way.

### 3.2.3 Spot-Guided Feature Aggregation

For the input features $F^{1/32}$ and $F^{1/8}$, $F^{1/32}$ is updated by vanilla cross attention, and $F^{1/8}$ is updated by linear cross attention for initialization. Then, two features of different resolutions are fed into the spot-guided aggregation blocks. In each block, $F^{1/32}$ and $F^{1/8}$ are first fused into each other in the following way:

$$\begin{aligned} \hat{F}^{1/32} &= F^{1/32} + \text{Conv}_{1\times1}(\text{Down}(F^{1/8})), \\ \hat{F}^{1/8} &= F^{1/8} + \text{Conv}_{1\times1}(\text{Up}(F^{1/32})), \end{aligned} \quad (7)$$

where $\hat{F}^{1/32}$ and $\hat{F}^{1/8}$ are features after fusion. $\text{Down}(\cdot)$ and $\text{Up}(\cdot)$ are downsampling and upsampling. And then, $\hat{F}^{1/32}$ aggregate features across images by vanilla attention. In the meantime, $\hat{F}^{1/8}$ aggregate features across images by spot-guided attention. After four spot-guided aggregation blocks, 1/32-resolution features are fused into 1/8-resolution features, which are used to obtain the coarse-matching result $M_c$.

## 3.3. Adaptive Scaling Module

At the fine stage, previous methods usually crop fixed-size grids based on the coarse matching result. When there is a large scale variation, fine matching may fail since the ground-truth matching points are out of grids. Thus, we refer to depth information to adaptively crop grids of different sizes between images. Section 3.3.1 describes the way to obtain depth information from the coarse-matching result. Section 3.3.2 demonstrates the process of adaptively cropping grids.

### 3.3.1 Depth Information

With the coarse-level matching result, we can obtain the relative pose of two images $\{R, T\}$ through RANSAC [16]. It should be noted that the $T$ calculated here has a scale uncertainty, i.e., $T_{real} = \alpha T$, where $\alpha$ is the scale factor. Given the image coordinates of any pair of matching points $\{x_i, x_j\}$ from coarse-level matching result, they satisfy the following equation:

$$d_j K_j^{-1}(x_j, 1)^T = d_i R K_i^{-1}(x_i, 1)^T + \alpha T, \quad (8)$$

where $d_i$ and $d_j$ are the depth values of $x_i$ and $x_j$. $K_i$ and $K_j$ are corresponding camera intrinsics. We let $p_i = R K_i^{-1}(x_i, 1)^T$ and $p_j = K_j^{-1}(x_j, 1)^T$. From Equation (8) it can be deduced that:

$$\begin{aligned} & d_j p_j = d_i p_i + \alpha T, \\ \Rightarrow & \begin{cases} (d_j/\alpha)p_j \wedge p_i = 0 + T \wedge p_i, \\ 0 = (d_i/\alpha)p_i \wedge p_j + T \wedge p_j, \end{cases} \\ \Rightarrow & \begin{cases} d_j/\alpha = \text{mean}(\text{div}(T \wedge p_i, p_j \wedge p_i)), \\ d_i/\alpha = \text{mean}(\text{div}(-T \wedge p_j, p_i \wedge p_j)), \end{cases} \end{aligned} \quad (9)$$

where $\wedge$ indicates outer product. $\text{div}(\cdot, \cdot)$ denotes element-wise division between two vectors. $\text{mean}(\cdot)$ is the scalar mean of each component of a vector. In this way, we have obtained depth information of $x_i$ and $x_j$ with scale uncertainty.

### 3.3.2 Adaptive Scaling Strategy

As shown in Figure 4, $x_i$ and $x_j$ are a pair of matching points at the coarse stage. $d_i$ and $d_j$ are depth values of $x_i$ and $x_j$. To begin with, we crop a $s_i \times s_i$ region centered on $x_i$. When the scale changes too much, the correct matching point $\widetilde{x_j}$ may be beyond the $s_i \times s_i$ region around $x_j$. Because everything looks small in the distance and big on the contrary, the size of cropped grid $s_j$ should satisfy:

$$\frac{s_j}{s_i} = \frac{d_i}{d_j} = (\frac{d_i}{\alpha})(\frac{d_j}{\alpha})^{-1}, \quad (10)$$

Following the above approach, we can crop different sizes of grids adaptively according to the scale variation. After the same refinement as LoFTR [50], we get the final matching position $\widetilde{x_j}$ of $x_i$.

### 3.4. Loss Function

Our loss function mainly consists of three parts, spot matching loss, coarse matching loss, and fine matching loss. Spot matching loss is the cross entropy loss to supervise the matching matrix during spot-guided attention:

$$L_s = -\frac{1}{|M_c^{gt}|} \sum_{(i,j) \in M_c^{gt}} \log P_s(i,j), \qquad (11)$$

where $M_c^{gt}$ is the ground truth matches at coarse resolution. Coarse matching loss is also the cross entropy loss to supervise the coarse matching matrix:

$$L_c = -\frac{1}{|M_c^{gt}|} \sum_{(i,j) \in M_c^{gt}} \log P_c(i,j). \qquad (12)$$

Fine matching loss $L_f$ is a weighted $L_2$ loss same as LoFTR [50]. Therefore, our total loss is:

$$L_{total} = L_s + L_c + L_f. \qquad (13)$$

## 4. Experiments

In this section, we evaluate our ASTR with extensive experiments. First of all, we introduce implementation details, followed by experiments on five benchmarks and some visualizations. Finally, we conduct a series of ablation studies to verify the effectiveness of each component.

### 4.1. Implementation Details

We implement the proposed model in Pytorch [40]. Our ASTR is trained on the MegaDepth dataset [29]. In the training phase, we input images with the size of $832 \times 832$ for training. The CNN extractor is a deepened ResNet-18 [20] with features at $1/32$ resolution. In spot-guided attention, we set the kernel size of local region $l$ to 5 and $k$ to 4 in topk. Threshold $\theta_c$ in coarse matching is chosen to 0.2. At the fine stage, window size $s_i$ in the reference image is fixed to 5, and window size $s_j$ in the source image will be adaptively calculated according to the depth information. In particular, $s_j/s_i$ is clamped into $[1,3]$. Our network is trained for 15 epochs with a batch size of 8 by Adam [26] optimizer. The initial learning rate is $1 \times 10^{-3}$. In order to establish spot-guided attention efficiently, we implement a highly optimized general sparse attention operator based on CUDA. Please refer to the Supplementary Material for more details about the operator.

### 4.2. Homography Estimation

**Dataset and Metric.** HPatches [1] is a popular benchmark for image matching. Following [15] , we choose 56 sequences under significant viewpoint changes and 52 sequences with large illumination variation to evaluate the

Table 1. Evaluation on HPatches [1] for homography estimation.

| Category | Method | Homography est. AUC | | | matches |
|---|---|---|---|---|---|
| | | @3px | @5px | @10px | |
| Detector-based | D2Net [15]+NN | 23.2 | 35.9 | 53.6 | 0.2K |
| | R2D2 [42]+NN | 50.6 | 63.9 | 76.8 | 0.5K |
| | DISK [55]+NN | 52.3 | 64.9 | 78.9 | 1.1K |
| | SP [14]+SuperGlue [47] | 53.9 | 68.3 | 81.7 | 0.6K |
| | Patch2Pix [64] | 46.4 | 59.2 | 73.1 | 1.0K |
| Detector-free | Sparse-NCNet [43] | 48.9 | 54.2 | 67.1 | 1.0K |
| | COTR [24] | 41.9 | 57.7 | 74.0 | 1.0K |
| | DRC-Net [27] | 50.6 | 56.2 | 68.3 | 1.0K |
| | LoFTR [50] | 65.9 | 75.6 | 84.6 | 1.0K |
| | PDC-Net+ [54] | 66.7 | 76.8 | 85.8 | 1.0k |
| | **ASTR(ours)** | **71.7** | **80.3** | **88.0** | 1.0K |

performance of our ASTR trained on MegaDepth [29]. We use the same evaluation protocol as LoFTR [50]. We report the area under the cumulative curve (AUC) of the corner error distance up to 3, 5, and 10 pixels, respectively. We limit the maximum number of output matches to 1k.

**Results.** In Table 1, we can see that our ASTR achieves new state-of-the-art performance on HPatches [1] under all error thresholds, which strongly proves the effectiveness of our method. ASTR outperforms the best method before (PDC-net+ [54]), achieving a large margin of **4.4%** under 3 pixels, **3.5%** under 5 pixels, and **2.5%** under 10 pixels. Thanks to the proposed spot-guided aggregation module and adaptive scaling module, our method can yield more accurate matches under extreme viewpoint and illumination variations.

### 4.3. Relative Pose Estimation

**Dataset and Metric.** We use MegaDepth [29] and Scan-Net [12] to demonstrate the performance of our ASTR in relative pose estimation. MegaDepth [29] is a large-scale outdoor dataset that contains 1 million internet images of 196 different outdoor scenes. Each scene is reconstructed by COLMAP [49]. Depth maps as intermediate results can be converted to ground truth matches. We sample the same 1500 pairs as [50] for testing. All test images are resized such that their longer dimensions are 1216. ScanNet [12] is usually used to validate the performance of indoor pose estimation. It is composed of monocular sequences with ground truth poses and depth maps. Wide baselines and extensive textureless regions in image pairs make Scan-Net [12] challenging. For a fair comparison, we follow the same testing pairs and evaluation protocol as [50]. And all test images are resized to $640 \times 480$. Note that we use our ASTR trained on MegaDepth [29] to evaluate its performance on ScanNet [12]. We report the AUC of the pose error at thresholds $(5°, 10°, 20°)$, where pose error is the maximum angular error in rotation and translation. The angular error is computed between the ground truth pose and the predicted pose.

**Results.** As shown in Table 2, our ASTR outperforms other state-of-the-art methods on MegaDepth [29]. In par-

Table 2. Evaluation on MegaDepth [29] for outdoor relative position estimation.

| Category | Method | Pose estimation AUC | | |
|---|---|---|---|---|
| | | @5° | @10° | @20° |
| Detector-based | SP [14]+SuperGlue [47] | 42.2 | 59.0 | 73.6 |
| | SP [14]+SGMNet [8] | 40.5 | 59.0 | 73.6 |
| Detector-free | DRC-Net [27] | 27.0 | 42.9 | 58.3 |
| | PDC-Net+(H) [54] | 43.1 | 61.9 | 76.1 |
| | LoFTR [50] | 52.8 | 69.2 | 81.2 |
| | MatchFormer [57] | 53.3 | 69.7 | 81.8 |
| | QuadTree [52] | 54.6 | 70.5 | 82.2 |
| | ASpanFormer [9] | 55.3 | 71.5 | 83.1 |
| | **ASTR(ours)** | **58.4** | **73.1** | **83.8** |

Table 3. Evaluation on ScanNet [12] for indoor relative position estimation. * indicates models trained on MegaDepth [29].

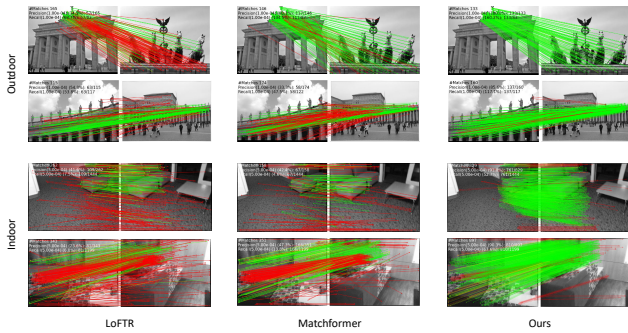| Category | Method | Pose estimation AUC | | |
|---|---|---|---|---|
| | | @5° | @10° | @20° |
| Detector-based | D2-Net [15]+NN | 5.3 | 14.5 | 28.0 |
| | SP [14]+OANet [61] | 11.8 | 26.9 | 43.9 |
| | SP [14]+SuperGlue [47] | 16.2 | 33.8 | 51.8 |
| Detector-free | DRC-Net [27]* | 7.7 | 17.9 | 30.5 |
| | MatchFormer [57]* | 15.8 | 32.0 | 48.0 |
| | LoFTR-OT [50]* | 16.9 | 33.6 | 50.6 |
| | Quadtree [52]* | 19.0 | 37.3 | 53.5 |
| | **ASTR(ours)*** | **19.4** | **37.6** | **54.4** |



Figure 5. Qualitative results of dense matching on MegaDepth [29] and ScanNet [12].

ticular, our ASTR improves by **3.1%** in AUC@5° and **1.6%** in AUC@10°. Table 3 summarizes the performance comparison between the proposed ASTR and state-of-the-art methods on ScanNet [12]. Our ASTR ranks first when only considering models not trained on ScanNet [12], indicating the impressive generalization of our method. Thanks to the proposed spot-guided aggregation module and adaptive scaling module, our method can yield more correct matches, resulting in more accurate pose estimation. In order to further demonstrate the effectiveness of our ASTR, in Figure 5, we visually demonstrate the comparison with other methods on the matching result. Notably, our methods can better handle the challenges such as textureless areas, repetitive patterns, and scale variations.

## 4.4. Visual Localization

**Dataset and Metric.** In this experiment, InLoc [51] and Aachen Day-Night v1.1 [62] are used to verify the ability of our ASTR in visual localization. InLoc [51] is an in-

Table 4. Visual localization evaluation on the InLoc [51] benchmark.

| Method | DUC1 | DUC2 |
|---|---|---|
| | (0.25m, 10°) / (0.5m, 10°) / (1m, 10°) | |
| Patch2Pix [64](w.SP [47]+CAPS [58]) | 42.4 / 62.6 / 76.3 | 43.5 / 61.1 / 71.0 |
| LoFTR [50] | 47.5 / 72.2 / 84.8 | 54.2 / 74.8 / **85.5** |
| MatchFormer [57] | 46.5 / 73.2 / 85.9 | **55.7** / 71.8 / 81.7 |
| ASpanFormer [9] | 51.5 / **73.7** / 86.4 | 55.0 / 74.0 / 81.7 |
| **ASTR(ours)** | **53.0 / 73.7 / 87.4** | 52.7 / **76.3** / 84.0 |

Table 5. Visual localization evaluation on the Aachen Day-Night benchmark v1.1 [62].

| Method | Day | Night |
|---|---|---|
| | (0.25m, 2°) / (0.5m, 5°) / (1m, 10°) | |
| **Localization with matching pairs provided in dataset** | | |
| R2D2 [42]+NN | - | 71.2 / 86.9 / 98.9 |
| ASLFeat [36]+NN | - | 72.3 / 86.4 / 97.9 |
| SP [14]+SuperGlue [47] | - | 73.3 / 88.0 / 98.4 |
| SP [14]+SGMNet [8] | - | 72.3 / 85.3 / 97.9 |
| **Localization with matching pairs generated by HLoc** | | |
| LoFTR [50] | 88.7 / 95.6 / 99.0 | 78.5 / 90.6 / 99.0 |
| ASpanFormer [9] | 89.4 / 95.6 / 99.0 | 77.5 / 91.6 / 99.0 |
| AdaMatcher [22] | 89.2 / **95.9 / 99.2** | **79.1** / 92.1 / 99.5 |
| **ASTR(ours)** | **89.9** / 95.6 / 99.2 | 76.4 / **92.1 / 99.5** |

Table 6. Ablation Study of each component on MegaDepth [29].

| Index | Multi-Level | Spot-Guided (l = 5, k = 4) | Scaling | Pose estimation AUC | | |
|---|---|---|---|---|---|---|
| | | | | @5° | @10° | @20° |
| 1 | | | | 45.6 | 62.2 | 75.3 |
| 2 | ✓ | | | 46.7 | 63.1 | 76.3 |
| 3 | ✓ | ✓ | | 47.7 | 64.5 | 77.4 |
| 4 | ✓ | ✓ | ✓ | **48.3** | **65.0** | **77.7** |

door dataset with 9972 RGBD images, of which 329 RGB images are employed as queries for visual localization. The challenge of InLoc [51] mainly comes from texture-less regions and repetitive patterns under large viewpoint changes. In Aachen Day-Night v1.1 [62], 824 day-time images and 191 night-time images are chosen as queries for outdoor visual localization. Large illumination and view-point changes pose challenges for Aachen [62]. For both benchmarks, we evaluate the performance of our ASTR trained on MegaDepth [29] in the same way as [50]. The metrics of Inloc [51] and Aachen [62] are the same, which measure the percentage of images registered within given error thresholds.

**Results.** For InLoc [51] benchmark, our method achieves the best performance on DUC1 and is on par with state-of-the-art methods on DUC2 (in Tabel 4). For Aachen [62] benchmark, our ASTR performs comparative with others on Day and Night scenes (in Tabel 5). Overall, our method exhibits strong generalization ability in visual localization.

## 4.5. Ablation Study

To deeply analyze the proposed method, we perform detailed ablation studies on MegaDepth [29] to evaluate the effectiveness of each component in ASTR. Here, we use images with a size of 544 for training and evaluation. As shown in Table 6, we intend to gradually add these components to the baseline. The baseline (Index-1) we used is
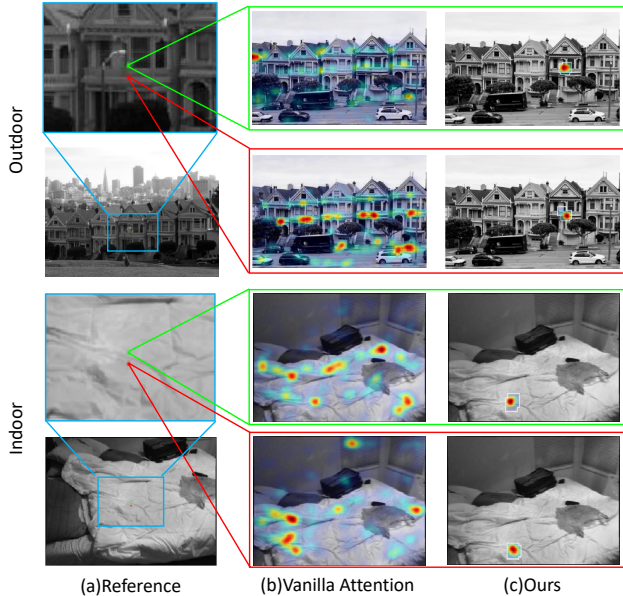
Figure 6. Visualization of vanilla and spot-guided cross attention maps on MegaDepth [29] (outdoor) and ScanNet [12] (indoor).

Table 7. Ablation Study with different $k$ and $l$ in spot-guided attention on MegaDepth [29].

| $k(l=5)$ | Pose estimation AUC | | | $l(k=4)$ | Pose estimation AUC | | |
|---|---|---|---|---|---|---|---|
| | @5° | @10° | @20° | | @5° | @10° | @20° |
| 1 | 46.0 | 62.7 | 76.2 | | | | |
| 2 | 47.5 | 64.0 | 77.1 | 3 | 46.7 | 63.2 | 76.1 |
| 3 | 47.3 | 63.8 | 76.7 | 5 | **47.7** | **64.5** | **77.4** |
| 4 | **47.7** | **64.5** | **77.4** | 7 | 47.2 | 63.4 | 76.8 |
| 5 | 47.1 | 63.7 | 77.0 | 9 | 43.0 | 60.5 | 74.8 |
| 6 | 46.9 | 63.6 | 76.6 | | | | |

slightly different from LoFTR [50]. More details can be found in Supplementary Material.

**Effectiveness of Spot-Guided Aggregation Module.** We divide the spot-guided aggregation module into multi-level cross attention and spot-guided attention for ablation studies. We first add vanilla cross attention layers at 1/32 resolution to the baseline (Index-2 in Table 6). Comparing the results of Index-2 and Index-1, we conclude that 1/32 resolution global interaction across images is beneficial for image matching. Then, in Index-3, linear attention layers at 1/8 resolution are substituted for the spot-guided attention layers. The performance of Index-3 is improved compared with Index-2, which verifies the effectiveness of our spot-guided attention. In Figure 6, we visualize vanilla and our spot-guided cross attention maps for contrast, showing that spot-guided attention can indeed avoid interference from unrelated areas.

To maximize the effectiveness of our spot-guided attention, we explore how to set suitable parameters $l$ and $k$. First, in the setting of Index-3, we fix $l = 5$ and vary $k$ from 1 to 6. After observing the results in Table 7, the performance drops when $k$ is smaller than 4 or larger than 4.
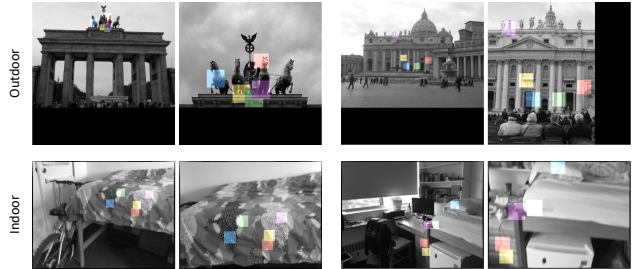


Figure 7. Visualization of grids from adaptive scaling module on MegaDepth [29] (outdoor) and ScanNet [12] (indoor).

Then, we fix $k = 4$ and vary $l$ from 3 to 9. As shown in Table 7, we find that the model achieves the best performance at $l = 5$. The reason may be that the spot area is too small to provide sufficient information from another image when using small $k$ or $l$. With large $k$ or $l$, for a certain pixel, some matching areas of low confidence or dissimilar points will damage its feature aggregation.

**Effectiveness of Adaptive Scaling Module.** As shown in Table 6, comparing the results of Index-4 and Index-3, we can see that the performance is improved, which indicates that coarse-level matching results are better refined with adaptive scaling module. In Figure 7, we visualize the cropped grids from adaptive scaling module, indicating that our adaptive scaling module can adaptively crop grids of different sizes according to scale variations.

## 5. Conclusion

In this paper, we propose a novel Adaptive Spot-guided Transformer (ASTR) for consistent local feature matching. To model local matching consistency, we design a spot-guided aggregation module to make most pixels avoid the impact of irrelevant areas, such as noisy and repetitive regions. To better handle large scale variation, we use the calculated depth information to adaptively adjust the size of grids at the fine stage. Extensive experimental results on five benchmarks demonstrate the effectiveness of the proposed method.

**Limitation.** Although our adaptive scaling module is lightweight and pluggable, it demands camera pose estimation in the coarse stage, which requires the camera intrinsic parameters. While camera intrinsic parameters are obtainable in standard datasets and most real-world scenarios, there are still some images from wild that lack them, rendering the adaptive scaling module disabled in such cases.

# Adaptive Spot-Guided Transformer for Consistent Local Feature Matching
## ——Supplementary Material——

In this supplementary material, we first introduce the general sparse attention operator in Section 6. In Section 7, we provide some details about our experiment. In Section 8, we show additional visualizations about the spot-guided attention and adaptive scaling modules.
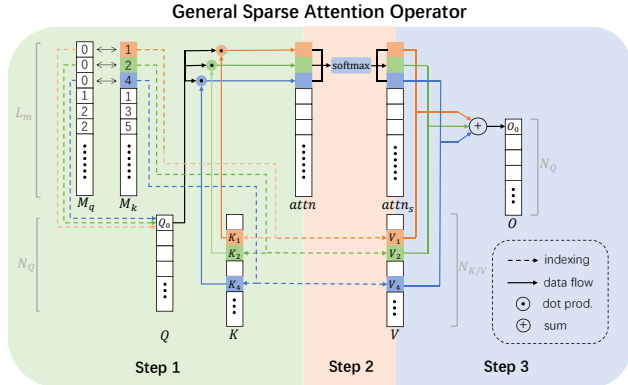
**General Sparse Attention Operator**



Figure 8. The illustration of our general sparse attention operator.

## 6. General Sparse Attention Operator

Due to irregular key/value token number for each query in Spot Attention, the naive implementation by PyTorch [40] is not efficient for memory and computation, which uses a mask to set unwanted values in the attention map to 0. More generally, the same problem also exists when the numbers of key corresponding to queries are not the same. Inspired by PointNet [?] and Stratified Transformer [?], we implement a general sparse attention operator using CUDA that is efficient in terms of memory and computation. We attempt to only compute the necessary attention between much less query/key tokens.

We can divide a vanilla attention operator into 3 steps. Inputs are grouped as query $Q$, key $K$ and value $V$. First, the attention map $A$ is computed by dot production as $A = QK^T$. Then, a softmax operator is performed on the attention map: $A_s = \mathrm{softmax}(A/\sqrt{d_k})$. Finally, the updated query $O$ can be obtained by $O = A_s V$. We optimize these three steps separately.

In the step 1, because only a few results in $A$ are useful for sparse attention, we do not need to compute the full $A$. Instead, we compute the dot productions between $L_m$ pairs of query and key. $M_q$ and $M_k$ record the indexes of query and key tokens whose dot productions are needed. The length of $M_q$ and $M_k$ are both $L_m$. Here, we denote
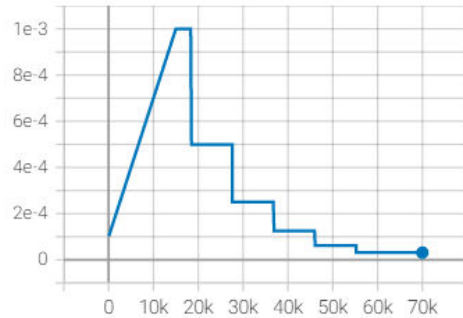


Figure 9. Learning rate curve while training on MegaDepth [29].

the sparse attention map as $attn$, which is calculated by

$$attn[i] = Q[M_q[i]]K[M_k[i]]^T, \ i = 0, 1, \cdots, L_m - 1. \tag{14}$$

In the step 2, we group the elements in $attn$ with the same query index and apply $\mathrm{softmax}$ on each group. The result is denoted as $attn_s$.

In the step 3, we compute the updated query

$$O[q] = \sum_{M_q[i]=q} attn_s[i] \cdot V[M_k[i]]. \tag{15}$$

All of three steps are implemented in CUDA.

Compared with the naive implementation using PyTorch [40], our highly optimized implementation reduces the memory and time complexity from $\mathcal{O}(N_q \cdot N_k \cdot N_h \cdot N_d^2)$ to $\mathcal{O}(L_m \cdot N_h \cdot N_d^2)$, where $N_q$, $N_k$ and $N_h$ are separately the numbers of query tokens, key tokens and attention heads, and $N_d$ is the dimension of each head. Considering $L_m \ll N_q \cdot N_k$, our implementation is much more efficient than the naive implementation.

In particular, we also calculate the matching matrix in spot-guided attention in this way and set the probability of unrelated pixels to 0, which can greatly reduce the memory and computation cost.

## 7. Experimental Details

### 7.1. Training Details

To reduce the GPU memory, we randomly sample $50\%$ of ground truth matches to supervise the matching matrix at the coarse stage. And we sample $20\%$ of the maximum number of coarse-level possible matches at the fine stage. We train ASTR on MegaDepth [29] for 15 epochs. The initial learning rate is $1 \times 10^{-3}$, with a linear learning rate warm-up for 15000 iterations. The learning rate curve is shown in Figure 9.
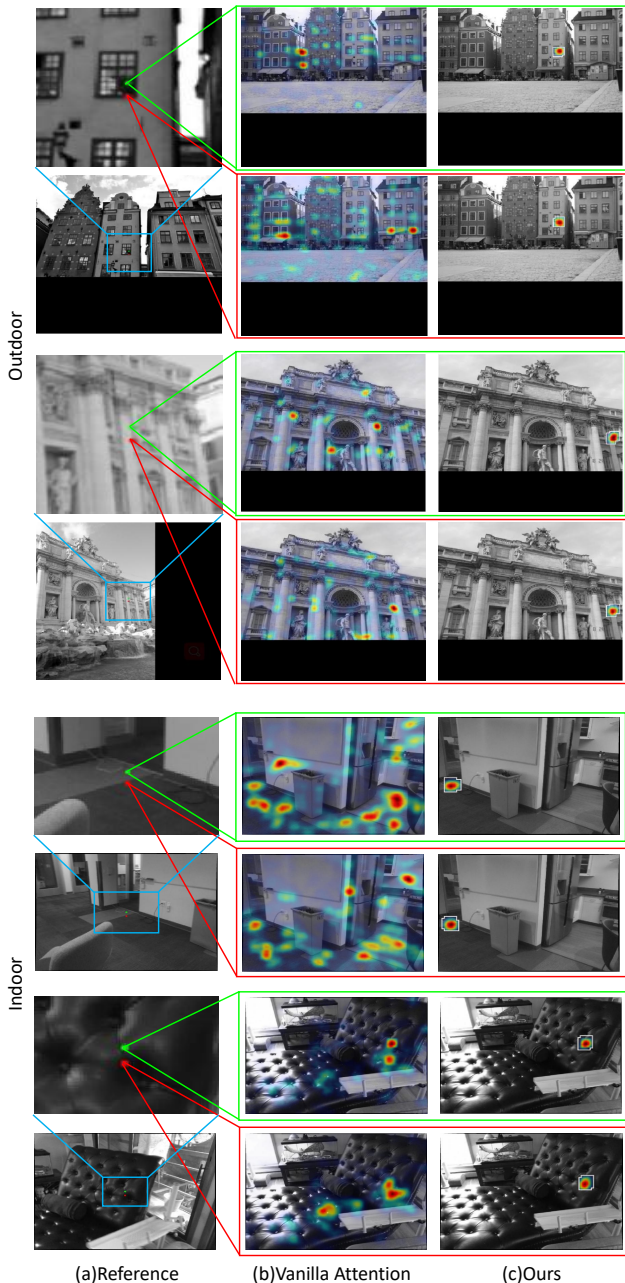
Figure 10. Visualizations of vanilla and spot-guided attention maps on MegaDepth [29] (outdoor) and ScanNet [12] (indoor).
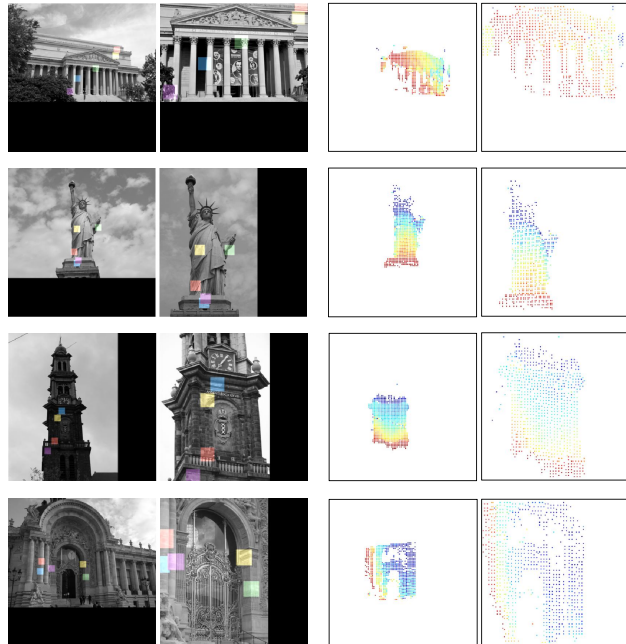


Figure 11. Visualizations of grids from adaptive scaling module and corresponding depth maps on MegaDepth [29]. Note that we use depth values with scale uncertainty to compose the depth maps.

## 7.2. Differences between Baseline and LoFTR

There are two main differences between our baseline and LoFTR [50].

**(1) Normalized Positional Encoding.** LoFTR [50] adopts the absolute sinusoidal positional encoding by fol-

lowing [7]:

$$
\text{PE}_i(x,y) = \begin{cases} \sin(w_k \cdot x), & i = 4k \\ \cos(w_k \cdot x), & i = 4k+1 \\ \sin(w_k \cdot y), & i = 4k+2 \\ \cos(w_k \cdot y), & i = 4k+3 \end{cases}, \quad (16)
$$

where $w_k = \frac{1}{10000^{2k/d}}$, $d$ denotes the number of feature channels and $i$ is the index for feature channels. Considering the gap in image resolution between training and testing, we utilize the normalized positional encoding as [9], which is proven to mitigate the impact of image resolution changes in [9]. The normalized positional encoding $\text{NPE}_i(\cdot, \cdot)$ can be expressed as

$$
\text{NPE}_i(x,y) = \text{PE}_i(x * \frac{W_{train}}{W_{test}}, y * \frac{H_{train}}{H_{test}}), \quad (17)
$$

where $W_{train/test}$ and $H_{train/test}$ are width and height of training/testing images.

**(2) Convolution in Attention.** Chen et al. [9] find that replacing the self attention with convolution can improve the performance. Hence, we deprecate self attention and MLP, and utilize a $3 \times 3$ convolution in our ASTR.

## 7.3. CNN Backbone

Here we leverage a deepened version of Feature Pyramid Network (FPN) [30], which achieves a minimum resolution

of 1/32. The initial dimension for the stem is still 128 as LoFTR [50], and the number of feature channels for subsequent stages is [128, 196, 256, 256, 256].

## 8. Visualization Results

In Figure 10, we pick up two similar adjacent pixels as queries and visualize the corresponding attention maps of vanilla and our spot-guided attention for comparison. The vanilla attention mechanism is vulnerable to repetitive textures, while our spot-guided attention can focus on the correct areas in these repeated texture regions. Because large scale variation occurs frequently on outdoor datasets, we mainly visualize the grids from the adaptive scaling module and corresponding depth maps on MegaDepth [29]. As shown in Figure 11, our adaptive scaling module can adjust the size of grids according to depth information.

## References

[1] Vassileios Balntas, Karel Lenc, Andrea Vedaldi, and Krystian Mikolajczyk. Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5173–5182, 2017. 6

[2] Axel Barroso-Laguna, Edgar Riba, Daniel Ponsa, and Krystian Mikolajczyk. Key. net: Keypoint detection by handcrafted and learned cnn filters. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5836–5844, 2019. 1, 2, 3

[3] Axel Barroso-Laguna, Yurun Tian, and Krystian Mikolajczyk. Scalenet: A shallow architecture for scale estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12808–12818, 2022. 3

[4] Axel Barroso-Laguna, Yannick Verdie, Benjamin Busam, and Krystian Mikolajczyk. Hdd-net: Hybrid detector descriptor with mutual interactive learning. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 3

[5] Herbert Bay, Andreas Ess, Tinne Tuytelaars, and Luc Van Gool. Speeded-up robust features (surf). *Computer vision and image understanding*, 110(3):346–359, 2008. 3

[6] Gabriele Berton, Carlo Masone, Valerio Paolicelli, and Barbara Caputo. Viewpoint invariant dense matching for visual geolocalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12169–12178, 2021. 3

[7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 3, 10

[8] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6301–6310, 2021. 7

[9] Hongkai Chen, Zixin Luo, Lei Zhou, Yurun Tian, Mingmin Zhen, Tian Fang, David Mckinnon, Yanghai Tsin, and Long Quan. Aspanformer: Detector-free image matching with adaptive span transformer. *arXiv preprint arXiv:2208.14201*, 2022. 1, 2, 3, 7, 10

[10] Ying Chen, Dihe Huang, Shang Xu, Jianlin Liu, and Yong Liu. Guide local feature matching by overlap estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 365–373, 2022. 3

[11] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015. 3

[12] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 6, 7, 8, 10

[13] 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 Transactions on Graphics (ToG)*, 36(4):1, 2017. 1

[14] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 224–236, 2018. 1, 2, 6, 7

[15] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-net: A trainable cnn for joint detection and description of local features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 6, 7

[16] 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*, 24(6):381–395, 1981. 5

[17] Yujie Fu and Yihong Wu. Scale-net: Learning to reduce scale differences for large-scale invariant image matching. *arXiv preprint arXiv:2112.10485*, 2021. 3

[18] Alexander Grabner, Peter M Roth, and Vincent Lepetit. 3d pose estimation and 3d model retrieval for objects in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3022–3031, 2018. 1

[19] Jianfeng He, Tianzhu Zhang, Yuhui Zheng, Mingliang Xu, Yongdong Zhang, and Feng Wu. Consistency graph modeling for semantic correspondence. *IEEE Transactions on Image Processing*, 30:4932–4946, 2021. 2

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 6

[21] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 3

[22] Dihe Huang, Ying Chen, Shang Xu, Yong Liu, Wenlong Wu, Yikang Ding, Chengjie Wang, and Fan Tang. Adaptive assignment for geometry aware local feature matching. *arXiv preprint arXiv:2207.08427*, 2022. 3, 7

[23] Shuaiyi Huang, Qiuyue Wang, Songyang Zhang, Shipeng Yan, and Xuming He. Dynamic context correspondence network for semantic alignment. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2010–2019, 2019. 1

[24] Wei Jiang, Eduard Trulls, Jan Hosang, Andrea Tagliasacchi, and Kwang Moo Yi. Cotr: Correspondence transformer for matching across images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6207–6217, 2021. 2, 3, 6

[25] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165. PMLR, 2020. 2, 3, 4

[26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[27] Xinghui Li, Kai Han, Shuda Li, and Victor Prisacariu. Dual-resolution correspondence networks. *Advances in Neural Information Processing Systems*, 33, 2020. 1, 2, 3, 6, 7

[28] Zhenyu Li, Zehui Chen, Xianming Liu, and Junjun Jiang. Depthformer: Exploiting long-range correlation and local information for accurate monocular depth estimation. *arXiv preprint arXiv:2203.14211*, 2022. 2

[29] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2041–2050, 2018. 6, 7, 8, 9, 10, 11

[30] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017. 3, 10

[31] Ce Liu, Jenny Yuen, and Antonio Torralba. Sift flow: Dense correspondence across scenes and its applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):978–994, 2010. 3

[32] Dongfang Liu, Yiming Cui, Liqi Yan, Christos Mousas, Baijian Yang, and Yingjie Chen. Densernet: Weakly supervised visual localization using multi-scale feature aggregation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6101–6109, 2021. 3

[33] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 3

[34] David G Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91–110, 2004. 2

[35] Bruce D Lucas, Takeo Kanade, et al. *An iterative image registration technique with an application to stereo vision*, volume 81. Vancouver, 1981. 3

[36] Zixin Luo, Lei Zhou, Xuyang Bai, Hongkai Chen, Jiahui Zhang, Yao Yao, Shiwei Li, Tian Fang, and Long Quan. Aslfeat: Learning local features of accurate shape and localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6589–6598, 2020. 1, 3, 7

[37] Meng Meng, Tianzhu Zhang, Zhe Zhang, Yongdong Zhang, and Feng Wu. Adversarial transformers for weakly supervised object localization. *IEEE Transactions on Image Processing*, 31:7130–7143, 2022. 3

[38] Meng Meng, Tianzhu Zhang, Zhe Zhang, Yongdong Zhang, and Feng Wu. Task-aware weakly supervised object localization with transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 3

[39] Yuki Ono, Eduard Trulls, Pascal Fua, and Kwang Moo Yi. Lf-net: learning local features from images. In *Advances in Neural Information Processing Systems*, pages 6237–6247, 2018. 1

[40] 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. *arXiv preprint arXiv:1912.01703*, 2019. 6, 9

[41] Mikael Persson and Klas Nordberg. Lambda twist: An accurate fast robust perspective three point (p3p) solver. In *Proceedings of the European Conference on Computer Vision*, pages 318–332, 2018. 1

[42] Jerome Revaud, Philippe Weinzaepfel, César Roberto de Souza, and Martin Humenberger. R2d2: repeatable and reliable detector and descriptor. In *Advances in Neural Information Processing Systems*, 2019. 1, 2, 3, 6, 7

[43] Ignacio Rocco, Relja Arandjelović, and Josef Sivic. Efficient neighbourhood consensus networks via submanifold sparse convolutions. In *Proceedings of the European Conference on Computer Vision*, pages 605–621, 2020. 1, 6

[44] Ignacio Rocco, Mircea Cimpoi, Relja Arandjelović, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Neighbourhood consensus networks. In *Advances in Neural Information Processing Systems*, pages 1658–1669, 2018. 1, 3

[45] Edward Rosten and Tom Drummond. Machine learning for high-speed corner detection. In *Computer Vision–ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7-13, 2006. Proceedings, Part I 9*, pages 430–443. Springer, 2006. 3

[46] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *2011 International conference on computer vision*, pages 2564–2571. Ieee, 2011. 2, 3

[47] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4938–4947, 2020. 1, 2, 6, 7

[48] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8601–8610, 2018. 1

[49] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE Con-*

*ference on Computer Vision and Pattern Recognition*, pages 4104–4113, 2016. 1, 6

[50] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8922–8931, 2021. 1, 2, 3, 5, 6, 7, 8, 10, 11

[51] Hajime Taira, Masatoshi Okutomi, Torsten Sattler, Mircea Cimpoi, Marc Pollefeys, Josef Sivic, Tomas Pajdla, and Akihiko Torii. Inloc: Indoor visual localization with dense matching and view synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7199–7209, 2018. 1, 7

[52] Shitao Tang, Jiahui Zhang, Siyu Zhu, and Ping Tan. Quadtree attention for vision transformers. *arXiv preprint arXiv:2201.02767*, 2022. 2, 3, 7

[53] Prune Truong, Martin Danelljan, and Radu Timofte. Glunet: Global-local universal network for dense flow and correspondences. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6258–6268, 2020. 3

[54] Prune Truong, Martin Danelljan, Radu Timofte, and Luc Van Gool. Pdc-net+: Enhanced probabilistic dense correspondence network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 3, 6, 7

[55] Michał Tyszkiewicz, Pascal Fua, and Eduard Trulls. Disk: Learning local features with policy gradient. *Advances in Neural Information Processing Systems*, 33:14254–14265, 2020. 6

[56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3

[57] Qing Wang, Jiaming Zhang, Kailun Yang, Kunyu Peng, and Rainer Stiefelhagen. Matchformer: Interleaving attention in transformers for feature matching. *arXiv preprint arXiv:2203.09645*, 2022. 2, 3, 7

[58] Qianqian Wang, Xiaowei Zhou, Bharath Hariharan, and Noah Snavely. Learning feature descriptors using camera pose supervision. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 757–774. Springer, 2020. 7

[59] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. 3

[60] Guanglei Yang, Hao Tang, Mingli Ding, Nicu Sebe, and Elisa Ricci. Transformer-based attention networks for continuous pixel-wise prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16269–16279, 2021. 2

[61] Jiahui Zhang, Dawei Sun, Zixin Luo, Anbang Yao, Lei Zhou, Tianwei Shen, Yurong Chen, Long Quan, and Hongen Liao. Learning two-view correspondences and geometry using order-aware network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5845–5854, 2019. 7

[62] Zichao Zhang, Torsten Sattler, and Davide Scaramuzza. Reference pose generation for long-term visual localization via learned features and view synthesis. *International Journal of Computer Vision*, 129(4):821–844, 2021. 7

[63] Lei Zhou, Siyu Zhu, Tianwei Shen, Jinglu Wang, Tian Fang, and Long Quan. Progressive large scale-invariant image matching in scale space. In *Proceedings of the IEEE international conference on computer vision*, pages 2362–2371, 2017. 2, 3

[64] Qunjie Zhou, Torsten Sattler, and Laura Leal-Taixe. Patch2pix: Epipolar-guided pixel-level correspondences. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4669–4678, 2021. 6, 7