

VRS-NeRF: Visual Relocalization with Sparse Neural Radiance Field

Fei Xue¹, Ignas Budvytis¹, Daniel Olmeda Reino², Roberto Cipolla¹

¹University of Cambridge, ²Toyota Motor Europe
 {fx221,ib255,rc10001}@cam.ac.uk, daniel.olmeda.reino@toyota-europe.com

Abstract—Visual relocalization is a key technique to autonomous driving, robotics, and virtual/augmented reality. After decades of explorations, absolute pose regression (APR), scene coordinate regression (SCR), and hierarchical methods (HMs) have become the most popular frameworks. However, in spite of high efficiency, APRs and SCRs have limited accuracy especially in large-scale outdoor scenes; HMs are accurate but need to store a large number of 2D descriptors for matching, resulting in poor efficiency. In this paper, we propose an efficient and accurate framework, called VRS-NeRF, for visual relocalization with sparse neural radiance field. Precisely, we introduce an explicit geometric map (EGM) for 3D map representation and an implicit learning map (ILM) for sparse patches rendering. In this localization process, EGP provides priors of sparse 2D points and ILM utilizes these sparse points to render patches with sparse NeRFs for matching. This allows us to discard a large number of 2D descriptors so as to reduce the map size. Moreover, rendering patches only for useful points rather than all pixels in the whole image reduces the rendering time significantly. This framework inherits the accuracy of HMs and discards their low efficiency. Experiments on 7Scenes, CambridgeLandmarks, and Aachen datasets show that our method gives much better accuracy than APRs and SCRs, and close performance to HMs but is much more efficient. Source code is available at <https://github.com/feixue94/vrs-nerf>.

I. INTRODUCTION

Visual localization aims to estimate the rotation and position of a given image captured in a known environment. As a fundamental computer vision task, visual localization is the key technique to various applications such as virtual/augmented reality (VR/AR), robotics, and autonomous driving. After several decades of exploration, many excellent methods have been proposed [1], [2], [3], [4] and can be roughly categorized as absolute pose regression (APR) [4], [5], [6], [7], [8], [9], scene coordinate regression (SCR) [10], [11], [12], [13], and hierarchical methods (HM) [14], [3], [1]. APRs embed the map into high-level pose features and predict the 6-DoF pose with multi-layer perceptions (MLP); they are fast especially for large-scale scenes, but have limited accuracy due to implicit 3D information representation. Different with APRs, SCRs regress 3D coordinates for pixels to build 2D-3D matches directly and estimate the pose with PnP [15] and RANSAC [16]. Despite the high accuracy in indoor environments, SCRs can't scale up to outdoor large-scale scenes. Instead of using an end-to-end 2D-3D matches prediction, HMs adopt global features [17], [18], [19] to search reference images in the database and then build correspondences between keypoints extracted query and reference images; these 2D-2D matches are lifted to

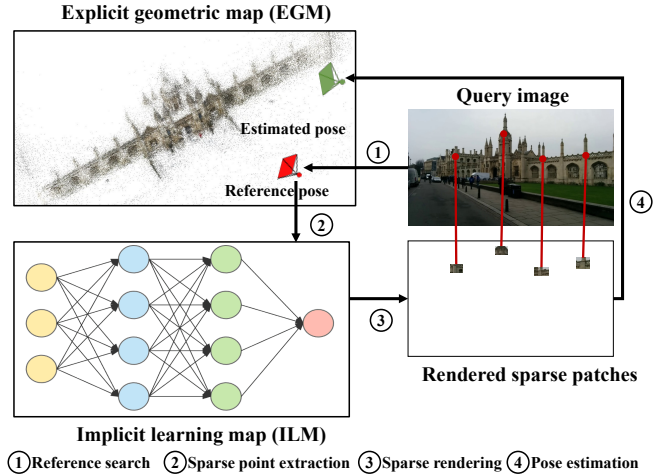


Fig. 1. **Overview of our framework.** We use an explicit geometric map (EGM) and an implicit learning map (ILM) to represent the environment. At test time, we first use the query image to find a reference image in EGM, then the pose and sparse points of this reference image are fed into ILM for sparse patch rendering. Finally the matches between the query and rendered patches are used for pose estimation with PnP [15] and RANSAC [16].

2D-3D matches and used for absolute pose estimation with PnP [15] and RANSAC [16] as SCRs. Because of high accuracy and flexibility, HMs are widely used recently. However, the huge memory cost of 2D keypoints storage impairs their efficiency in real applications.

In this paper, we aim to find an efficient and accurate solution to large-scale visual localization task. To achieve this, we seek help from neural radiance fields (NeRFs). NeRFs are first proposed for view synthesis [22]. Due to their powerful scene and object representation ability, NeRFs have been widely used for many other tasks including visual localization [23], [24]. Although LENS [24] and NeRF-loc [23] have applied NeRFs to APRs and SCRs respectively, their performance especially in outdoor scenes is still limited. Besides, direct usage of NeRFs for localization is inefficient as rendering all pixels of an image is slow.

Instead, we adopt a hybrid map of using NeRFs for efficient localization by rendering only useful sparse pixels. The hybrid map consists of two parts: explicit geometric map (EGM) and implicit learning map (ILM). EGM contains the sparse 3D points along with their 2D observations on reference images. ILM is the implicit map represented by NeRFs. At test time, 2D observations of reference images provide prior sparse pixel locations and camera poses as input to NeRFs. NeRFs return RGB values of each sparse

pixel. In order to improve the accuracy, we render a patch with constant size for each pixel. These rendered patches are further used to build 2D-3D matches for absolute pose estimation with PnP [15] and RANSAC [16].

An overview of our framework is shown in Fig. 1. With EGM and ILM, our method is able to render useful pixels online as opposed to relying on offline 2D descriptors for matching, making the localization system much more efficient. To allow current NeRFs work in large-scale scenes, we adopt a clustering-based strategy to divide the scene into smaller ones adaptively and automatically. The contributions of our method are summarized as follows:

- We propose a hybrid method combining explicit geometric map and implicit learning map for visual localization, making localization system efficient and accurate.
- Instead of rendering images, we render patches for useful sparse keypoints only, avoiding the time-consuming rendering process.
- We adopt a clustering-based strategy for scene division, enabling NeRFs to work in large-scale outdoor environments.

Our experiments on the popular indoor 7Scenes [25] and outdoor CambridgeLandmarks [4] and Aachen [26] datasets demonstrate that our method requires much less memory cost while preserving the accuracy. We hope this method could be a new baseline of applying NeRFs to visual localization task. We organize the rest of the paper as follows. In Sec. II, we discuss related works. In Sec. III, we introduce our method in detail. We test our approach in Sec. IV and conclude the paper in Sec. V.

II. RELATED WORK

In this section, we discuss related works about visual localization and NeRFs.

Visual localization. Visual localization methods can be roughly categorized as absolute pose regression (APR), scene coordinate regression (SCR), and hierarchical methods (HM). Posenet [4] is the first work introducing APR. Due to its simplicity, high memory and time efficiency, especially in large-scale scenes, a lot of variants have been proposed by introducing geometric loss [5], multi-view constraints [9], [6], [7], [27], [28], [29], view synthesis [30], feature selection [8], [28] and additional training data generation [24]. However, the accuracy is still limited because of the retrieval nature of APRs [31].

SCRs first regress 3D coordinates for each pixel in the query image and then estimate the pose with PnP [15] and RANSAC [16]. Initially, this is achieved via random forest with RGBD as input [10]. Later on, DSAC [11] and its variants [32], [20] extend it to RGB input and replace random forest with CNNs. More recently, hierarchical prediction [33], semantic-aware prediction [21], and the separation of backbone and prediction head [13], to name a few, are introduced for better accuracy and training efficiency. SCR compresses the map into a compact network and give very accurate poses in small scenes such as indoor

environments [25], [34], but have limited accuracy in large-scale scenes including CambridgeLandmarks [4] and Aachen dataset [26].

HMs [3], [1] estimate the pose of a query image by first finding reference images in the database, then building matches between keypoints in query and reference images, and finally compute the pose from associated 2D-3D matches lifted from 2D-2D ones with PnP [15] and RANSAC [16]. Traditionally, handcrafted SIFT [35] or ORB [36] features along with BoWs [37], [38] are widely used for the first two steps [14]. Due to the sensitivity of handcrafted features to illumination and seasonal changes, learned local features [39], [40], [41], [42], [43], [44], e.g., Lift [43], SuperPoint (SP) [39], SFD2 [42] and global features, e.g., NetVLAD (NV) [17], GeM [19] are used to replace classic ones. To further improve the accuracy, graph-based matchers, e.g., SuperGlue (SG) [45], SGMNet [46], IMP [47] are proposed to replace nearest matching for better 2D-2D matches. Nowadays, the combination of SPP+NV+SG [39], [17], [45] has set the new state-of-art on public datasets [25], [34], [4], [26]. Despite the outstanding accuracy, HMs need to store a large amount of local features to build 2D-3D correspondences, impairing the memory efficiency.

Our method absorbs advantages of HMs in terms of accuracy by introducing an explicit geometry map (EGM) to preserve geometric information. Additionally, we leverage the implicit learning map (ILM) for rendering useful pixels for high efficiency instead of storing redundant 2D descriptors for higher efficiency.

Neural radiance fields (NeRFs). NeRFs are first introduced for view synthesis [22]. After two-year exploration, many excellent works [48], [49], [50], [51] have been proposed to achieve higher quality and faster speed. Although some of them [51] report real-time rendering performance, they require to storage additional intermediate features, resulting in heavy memory cost. As NeRFs have impressive ability of scene representation, some works have applied them to visual localization task. LENS [24] uses NeRFs to render more images from different viewpoints for training absolute pose regression and obtains better accuracy. NeRF-loc [23] adopts a conditional NeRF to render 3D descriptors for matching. Despite its promising performance, it needs support images, features and depth maps as input, which leads to low memory efficiency. Moreover, the dense rendering of all pixels further impairs its time efficiency.

Different with LENS [24] and NeRF-loc [23], we use a hybrid map consisting of EGP and ILP for online sparse rendering to reduce time and memory cost and yield higher accuracy as well.

III. METHOD

In this section, we first introduce the explicit geometric map and implicit learning map in Sec. III-A and Sec. III-B, respectively. Then, we describe the process of localization with sparse NeRFs in Sec. III-C.

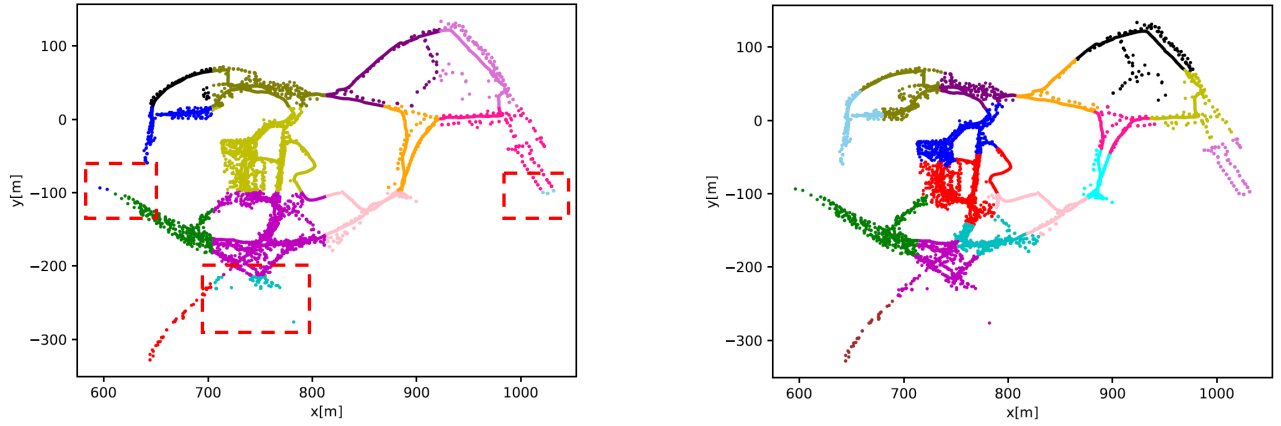


Fig. 2. **Visualization of scene division on Aachen dataset [26].** The uniform division of scene leads to imbalanced pieces (left) and our clustering on reference poses gives more balanced results (right).

A. Explicit Geometric Map

We first build the explicit geometric map (EGM). Given reference images, we adopt current state-of-the-art SfM library colmap [52] to reconstruct the 3D environment. EGM consists of a set of 3D points $\mathcal{X} = \{X_1, \dots, X_m\}$ and reference images $\mathcal{I} = \{I_1, \dots, I_n\}$. Each reference image I_i has several 2D keypoints $\mathcal{P}_i = \{P_1^i, \dots, P_k^i\}$ corresponding to 3D points in \mathcal{X} . Each reference image I_i also has a global descriptor \mathbf{v}_i provided by NetVLAD [17] to search similar reference images for a given query image I_q . It is worthy to note that in EGM only global descriptors are retained and local descriptors of 2D keypoints are discarded due to their huge memory cost.

B. Implicit Learning Map

The implicit learning map (ILM) is a NeRF-based model compressing the scene in a single network implicitly. NeRF renders the RGB values of a single ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ where \mathbf{o} and \mathbf{d} are the origin and direction of the ray, and t is distance along the ray. Usually, the ray is split into a set of intervals $T_i = [t_i, t_{i+1})$. The features generated for T_i with function γ are fed into a MLP for the density τ and color \mathbf{c} prediction:

$$\tau_i, \mathbf{c}_i = MLP(\gamma(\mathbf{r}(T_i)), \mathbf{d}; \Theta), \mathbf{C}(\mathbf{r}, \mathbf{t}) = \sum_i \omega_i \mathbf{c}_i, \quad (1)$$

$$\omega_i = (1 - e^{-\tau_i(t_{i+1}-t_i)})e^{-\sum_{i' < i} \tau_{i'}(t_{i'+1}-t_{i'})}. \quad (2)$$

$\mathbf{C}(\mathbf{r}, \mathbf{t})$ is the predicted pixel value. For better performance, it is optimized from coarse to fine as:

$$\sum_{\mathbf{r} \in \mathcal{R}} 0.1 L_{rec}(\mathbf{C}(\mathbf{r}, \mathbf{t}^c), \mathbf{C}^*(\mathbf{r})) + L_{rec}(\mathbf{C}(\mathbf{r}, \mathbf{t}^f), \mathbf{C}^*(\mathbf{r})). \quad (3)$$

\mathcal{R} is the set of rays for training. \mathbf{t}^c and \mathbf{t}^f are predicted coarse and fine distances. $\mathbf{C}^*(\mathbf{r})$ is the ground-truth color. L_{rec} is the mean squared error. We adopt Mip-NeRF [49] as the base model in our ILM.

Scene division. A single NeRF may not be able to represent a scene of large scale [53], so we need to divide

the scene into several pieces. A straightforward method is to divide the scene uniformly. However, this usually makes sub-scenes imbalance because different areas with different density of objects. Instead, we propose to leverage to priors of reference images which reconstruct the scene. In detail, we perform clustering on the poses of reference images with K-means [54] with a given number of clusters N . This strategy allows us to segment the scene into more balanced pieces, as shown in Fig. 2.

In our experiment, we divide the large-scale Aachen dataset [26] into 16 sub-scenes and train 16 separate NeRFs $\Theta_1, \dots, \Theta_{16}$ independently to represent each of them, respectively.

C. Localization with Sparse NeRFs

Given a query image I_q , we first extract global descriptor \mathbf{v}_q and a set of 2D local keypoints $\mathcal{P}^q = \{P_1^q, \dots, P_k^q\}$.

Reference image search. As EGM contains global descriptors of all reference images, we perform nearest matching to find top k candidate images with the smallest distances between reference and query global descriptors. These candidate reference images are sorted as $\mathcal{I}^c = \{I_1^c, \dots, I_k^c\}$ in ascending order according to their distances to \mathbf{v}_q .

Matching with online renderings. As [3], we perform 2D-2D matching between the query and all candidate reference images independently. Since the EGM does not contain any 2D descriptors, we render the patches with ILM on the fly. For a keypoint P_{ij}^c from candidate reference image I_j^c , we generate $\lambda_r \times \lambda_r$ rays at center of P_{ij}^c denoted as \mathcal{R}_{ij}^c . \mathcal{R}_{ij}^c along with the pose of I_j^c is then fed into ILM for RGB values prediction. \mathcal{R}_{ij}^c is reshaped as a patch p_{ij}^c with size of $\lambda_r \times \lambda_r$. p_{ij}^c is finally used as the input to a local feature network [39], [42] as obtain the descriptor d_{ij}^c and score s_{ij}^c for 2D-2D matching. The whole process can be described as:

$$\mathcal{R}_{ij}^c = f_r(P_{ij}^c, \lambda_r), \mathbf{C}_{ij}^c = f_{nerf}(\mathcal{R}_{ij}^c, \Theta), \quad (4)$$

$$p_{ij}^c = f_p(\mathbf{C}_{ij}^c), (d_{ij}^c, s_{ij}^c) = f_l(p_{ij}^c), \quad (5)$$

where f_r, f_{nerf}, f_p, f_l are the ray generation, NeRF with weights of Θ , patch, and local feature extraction functions.

Group	Method	Chess	Fire	Heads	Office	Pumpkin	Kitchen	Stairs	Average (%)
APRs	Posenet [4]	13, 4.5	27, 11.3	17, 13.0	19, 5.6	26, 4.8	23, 5.4	35, 12.4	-
	MapNet [9]	8, 3.3	27, 11.7	18, 13.3	17, 5.2	22, 4.0	23, 4.9	30, 12.1	-
	LSG [6]	9, 3.3	26, 10.9	17, 12.7	18, 5.5	20, 3.7	23, 4.9	23, 11.3	-
	AtLoc [8]	10, 4.1	25, 11.4	16, 11.8	17, 5.3	21, 4.4	23, 5.4	26, 10.5	-
	GLNet [7]	8, 2.8	26, 8.9	17, 11.4	18, 5.1	15, 2.8	25, 4.5	23, 8.8	-
	SC-wLS [12]	3, 0.8	5, 1.1	3, 1.9	6, 0.9	8, 1.3	9, 1.4	12, 2.8	-
	PAEs [55]	12, 5.0	24, 9.3	14, 12.5	19, 5.8	18, 4.9	18, 6.2	25, 8.7	-
	LENS [24]	3, 1.3	10, 3.7	7, 5.8	7, 1.9	8, 2.2	9, 2.2	14, 3.6	-
SCRs	DSAC*	2, 1.1	2, 1.2	1, 1.8	3, 1.2	4, 1.3	4, 1.7	3, 1.2	96.0
	ACE [13]	2, 1.1	2, 1.8	2, 1.1	3, 1.4	3, 1.3	3, 1.3	3, 1.2	97.1
	NeRF-loc [23]	2, 1.1	2, 1.1	1, 1.9	2, 1.1	3, 1.3	3, 1.5	3, 1.3	89.5
HMs	SP+SG [39], [45]	0, 0.1	1, 0.2	0, 0.2	1, 0.2	1, 0.1	0, 0.1	2, 0.6	95.7
	SFD2+IMP [42], [47]	0, 0.1	1, 0.2	0, 0.2	1, 0.2	1, 0.2	0, 0.	2, 0.5	95.7
	Ours	0, 0.1	1, 0.2	0, 0.2	1, 0.2	1, 0.2	0, 0.1	3, 0.8	93.1

TABLE I

LOCALIZATION ACCURACY ON 7SCENE DATASET [25]. WE REPORT THE MEDIAN TRANSLATION (CM) AND ROTATION ($^{\circ}$) ERRORS AND THE AVERAGE SUCCESS RATIO OF POSES WITHIN ERROR OF $5cm, 5^{\circ}$.

Compared with image-wise rendering, our sparse rendering focuses only on the useful keypoints with corresponding 3D points in the map, reducing the time significantly. For example, an image with size of 480×640 from 7Scenes dataset [25] has 307,200 rays in total. However, for 500 keypoints with patch size of 15×15 , the number of rays is 112,500 which is $2.7 \times$ fewer. As patches have overlap, for each unique ray, we only render it once, which further makes the rendering efficient. The online rendering of patches introduce additional time cost of feature extraction, because they are only patches with limited size, the additional time cost is limited.

Pose estimation. The 2D-2D matches between the query I^q and candidate reference images \mathcal{I}^c are lifted to 2D-3D matches between I^q and EGM. Then, we use RANSAC [16] and EPnP [15] to recover the absolute pose.

IV. EXPERIMENT

In this section, we compare the performance of our method with previous approaches for localization in terms both accuracy and efficiency.

A. Implementation

We implement our model on PyTorch [56]. We use Mip-NeRF [48] as our basic NeRF model for scene representation. Each model is trained with Adam optimizer [57], batch size of 1024 on a single RTX 3090 GPU for 50,000 iterations in total. The initial learning rate is set to 0.001 and is decayed at the ratio of $1e-8$ after 5,000 iterations. The patch size λ_r is set to 15 in our experiments.

B. Datasets, Metrics and Baselines

Dataset. We evaluate our method on three public datasets including the indoor 7Scenes [25] and outdoor Cambridge-Landmarks [4] and Aachen [26] datasets. For 7Scenes and CambridgeLandmarks, we assign each sub-scene a NeRF model. As Aachen dataset is a large-scale outdoor scene, we divide it into 16 sub-scenes with the clustering strategy as

introduced in Sec. III-B. Following previous methods [3], [1], [42], NetVLAD (NV) [17] is adopted to provide 20, 20 and 50 candidate reference images for query images in 7Scenes, CambridgeLandmarks, and Aachen datasets, respectively. We use SFD2 [42] and IMP [47] as the local feature and matcher because SFD2+IMP gives better performance than SP+SG on Aachen dataset and has smaller map size.

Metrics. As [32], [3], [14], we report the median rotation and translation errors on 7Scenes [25] and Cambridge-Landmarks [4]. Besides, the success ratio of pose error within $5cm, 5^{\circ}$ and $25cm, 5^{\circ}$ are also provided. For Aachen dataset [26], we use the official metric of success ratio at error thresholds of $0.25m, 2^{\circ}$, $0.5m, 5^{\circ}$ and $5m, 10^{\circ}$.

By introducing sparse NeRFs, we aim to reduce the high memory cost of HMs and the time cost of image-wise rendering. Therefore, we additionally analyze the memory and time efficiency.

Baselines. We compare our system with previous APRs [4], [7], [6], [12], [8], [55], [9], SCRs [11], [32], [13] and hierarchical methods [14], [3], [1], [42]. Besides, we also compare our approach with previous NeRF-based methods including LENS [24] and NeRF-loc [23].

C. Pose Accuracy Analysis

7Scenes. As shown in Table I, we compare our approach with previous APRs [4], [9], [6], [7], [8], [12], [55], [24], SCRs [32], [13], [23] and HMs [39], [45], [42], [47]. APRs give the largest errors because their similar behaviors to image retrieval in the localization process [31], resulting in limited pose accuracy. Since most APRs only report median errors, their success ratios are not available. SCRs obtain much higher accuracy than APRs due to their explicit 3D coordinates regression. HMs achieve the best accuracy in terms of median errors. However, they are less robust to textureless areas due to the reliance on sparse keypoints, so they report slightly worse accuracy ratio than some SCRs such as DSAC* [32] and ACE [13]. Although our approach renders sparse patches for localization, it yields close performance to

Group	Method	Kings College	Great Court	Old Hospital $t(cm), R(^{\circ}), Percent(25cm, 2^{\circ})$	Shop Facade	St Mary Church	Average (%)
APRs	Posenet [4]	88, 1.0	683, 3.5	88, 3.8	157, 3.3	320, 3.3	-
	MapNet [9]	107, 1.9	785, 3.8	149, 4.2	200, 4.5	194, 3.9	-
	GLNet [7]	59, 0.7	667, 3.0	50, 2.9	190, 3.3	188, 2.8	-
	PAEs [55]	90, 1.5	-	207, 2.6	99, 3.9	164, 4.2	-
	LENS [24]	33, 0.5	-	44, 0.9	27, 1.6	53, 1.6	-
SCRs	DSAC* [32]	13, 0.4	40, 0.2	20, 0.3	6, 0.3	13, 0.4	64.88
	ACE [13]	28, 0.4	42, 0.2	31, 0.6	5, 0.3	19, 0.6	54.68
	NeRF-loc [23]	7, 0.2	25, 0.1	18, 0.4	11, 0.2	4, 0.2	-
HMs	SP+SG [39], [45]	7, 0.1	12, 0.1	9, 0.2	2, 0.1	4, 0.1	89.4
	SFD2+IMP [42], [47]	7, 0.1	11, 0.1	10, 0.2	2, 0.1	4, 0.1	89.1
	Ours	9, 0.1	-	11, 0.2	2, 0.1	5, 0.2	89.3

TABLE II

LOCALIZATION ACCURACY ON CAMBRIDGELANDMARKS DATASET [4]. WE REPORT THE MEDIAN TRANSLATION (CM), ROTATION ($^{\circ}$) ERRORS AND THE LOCALIZATION PRECISION WITH TRANSLATION AND ROTATION ERRORS WITHIN $25cm, 2^{\circ}$. - INDICATES NO VALUES AVAILABLE. DUE TO POOR IMAGE QUALITY, WE FAILED TO TRAIN A NERF MODEL FOR GREAT COURT.

HMs and outperforms APRs and SCRs significantly in terms median errors. Similar to HMs, our model is also sensitive to textureless regions. As our EGM inherits advantages of HMs, it outperforms previous approaches LENS [24] and NeRF-loc [23] which introduce NeRFs to APRs and SCRs, respectively.

CambridgeLandmarks. Table II demonstrates the results of previous and our methods on the CambridgeLandmarks dataset [4]. We report the median translation (cm) and rotation ($^{\circ}$) errors and the success ratio of poses within error threshold of $25cm, 2^{\circ}$. Similar to results on 7Scenes dataset [25], APRs give over $2\times$ larger errors than SCRs due to the missing 3D information embedding.

SCRs report promising accuracy in terms of median translation and rotation errors. However, their success ratios within the error threshold of $25cm, 2^{\circ}$ are much worse than HMs. Even the state-of-the-art DSAC* [32] and ACE [13] fail to achieve comparable accuracy to HMs. These comparisons reveal that SCRs are not that accurate as they are expected to be in outdoor scenes. Note that each sub-scene in CambridgeLandmarks dataset such as Kings College is of size about $500m \times 50m$, far from large-scale environments. We hope the future methods especially SCRs could also include the success ratio metric so as to provide a deeper understanding of their performance.

HMs are still the most accurate methods in terms of both median errors and success ratios. As our approach also preserves the explicit geometric information as explicit geometric map, its results are as accurate as HMs' and are much more accurate than APRs and SCRs.

Compared with prior NeRF-based LENS [24] and NeRF-loc [23], our system also achieves the significantly better accuracy.

Aachen. We finally show the results on Aachen dataset [26]. As Aachen dataset is a large-scale dataset consisting of images with illumination, season and large-viewpoint changes, few APRs and SCRs report numbers on this dataset. We follow the official metric by showing

Method	Day $(2^{\circ}, 0.25m)/(5^{\circ}, 0.5m)/(10^{\circ}, 5m)$	Night $(2^{\circ}, 0.25m)/(5^{\circ}, 0.5m)/(10^{\circ}, 5m)$
ESAC [58]	42.6 / 59.6 / 75.5	3.1 / 9.2 / 11.2
HSCNet [33]	71.1 / 81.9 / 91.7	32.7 / 43.9 / 65.3
AS [14]	85.3 / 92.2 / 97.9	39.8 / 49.0 / 64.3
LBR [1]	88.3 / 95.6 / 98.8	84.7 / 93.9 / 100.0
SIFT [35]	82.8 / 88.1 / 93.1	30.6 / 43.9 / 58.2
SP+SPG [39], [45]	89.6 / 95.4 / 98.8	86.7 / 93.9 / 100.0
SFD2+IMP [42], [47]	89.7 / 96.5 / 98.9	84.7 / 94.9 / 100.0
Ours (15)	60.8 / 67.8 / 73.1	19.4 / 22.4 / 25.5
Ours (31)	70.1 / 76.9 / 80.9	44.9 / 51.0 / 62.2

TABLE III

RESULTS ON AACHEN DATASET [26], [59]. WE REPORT THE POSE RATIOS WITHIN ERROR THRESHOLDS OF $2^{\circ}, 0.25m, 5^{\circ}, 0.5m$ AND $10^{\circ}, 5m$.

the success ratio of poses at error thresholds of $0.25m, 2^{\circ}, 0.5m, 5^{\circ}$ and $5m, 10^{\circ}$. SCRs including ESAC [58] and HSCNet [33] give relatively worse accuracy especially on night images due to severe illumination changes. Among HMs, the combination of SFD2 and IMP [42], [47] achieves the state-of-the-art performance.

Due to poor rendering quality of insufficient observations and large illumination changes, our method gives poor performance when the patch size is 15. When using larger patch size (31), our model gives promising accuracy. However, a large gap between the performance of our method and state-of-the-art SFD2+IMP. As our approach is the first the apply NeRFs on large-scale Aachen dataset, we hope more works can make it better in the future.

D. Map size and Time Analysis

In this section, we analyze the map size and running time.

Map size. In table IV, we show the map size of APRs, SCRs, HMs and our method. For APRs and SCRs, the map size is the model size. For HMs, the map size is the sum of local descriptors, global descriptors, and 3D points. As our method discards local descriptors and introduces NeRFs, the

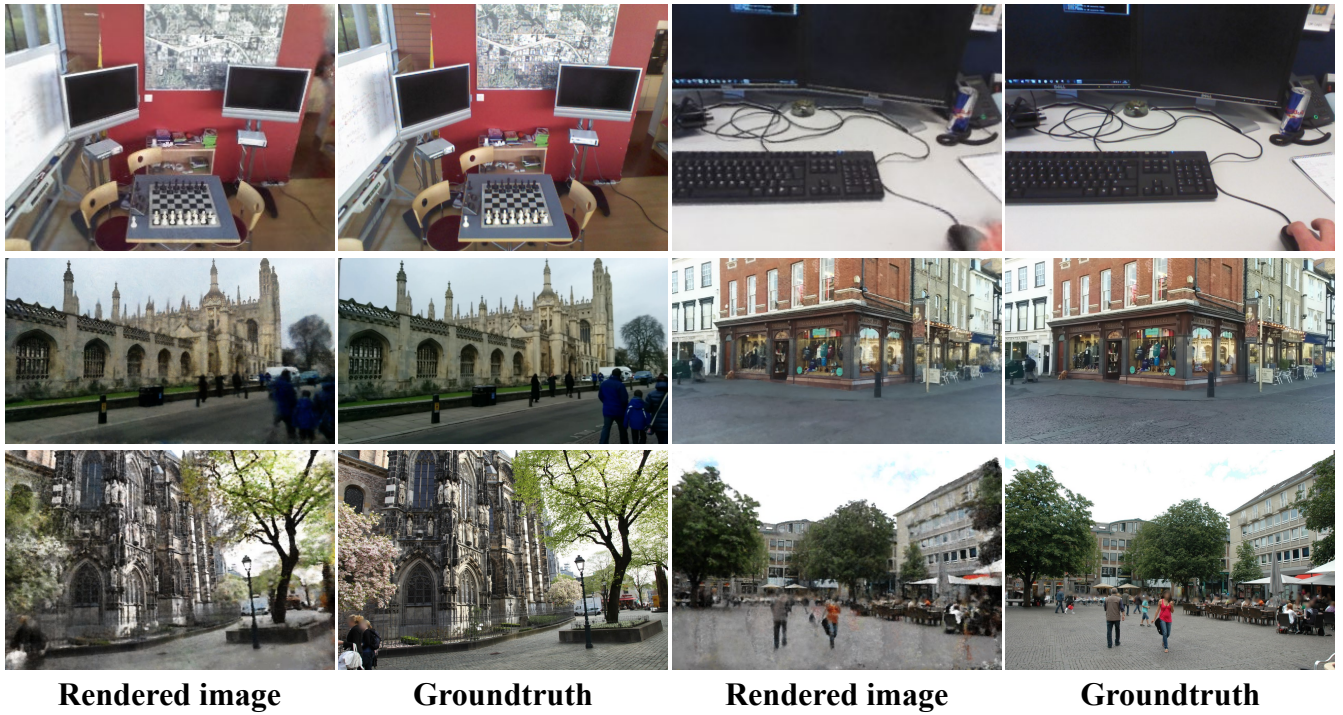


Fig. 3. **Visualization of rendered image.** We visualize the rendered and groundtruth images from 7Scenes [25] (top), CambridgeLandmarks [4] (middle) and Aachen [26] (bottom) datasets.

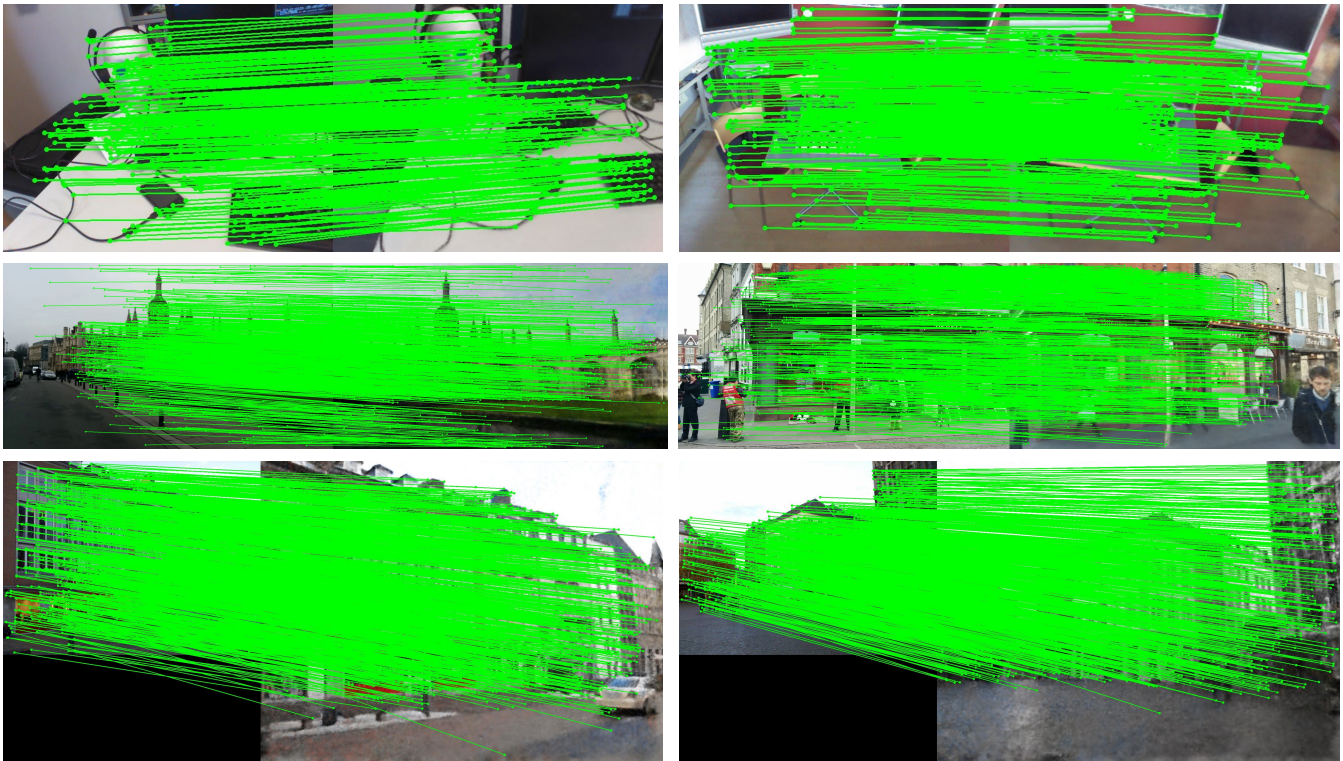


Fig. 4. **Visualization of matches.** We visualize the matches between query (left) and reference (right) images from 7Scenes [25] (top), CambridgeLandmarks [4] (middle) and Aachen [26] (bottom) datasets. Note that we perform matching the patches with size of 15×15 and show the whole rendered reference images for visualization only.

map size of our model is the sum of global descriptors, 3D points and NeRFs. For simplicity, we report the numbers on a

sub-scene of 7Scenes and CambridgeLandmarks and Aachen dataset.

Method	Chess [25]	Kings College [4]	Aachen [26]
Posenet [4]	50MB	50MB	-
DSAC* [32]	28MB	28MB	-
ACE [13]	4MB	4MB	-
SP+SG [39], [45]	8.9GB	13.3GB	53.5GB
SFD2+IMP [42], [47]	4.0GB	6.6GB	29.5GB
Ours	1.0GB	0.76GB	4.5GB

TABLE IV

MAP SIZE. WE REPORT THE MAP SIZE OF POSENET [4], DSAC* [32], ACE [13], SP+SG [39], [45], SFD2+IMP [42], [47] AND OUR METHOD.

Table IV shows that both APRs [4] and SCRs [32], [13] are memory-efficient as they compress the map into neural networks at the cost of accuracy loss. HMs have a larger size of map due to the storage of 2D descriptors. SFD2+IMP [42], [47] has smaller size than SP+SG [39], [45] because SFD2 has smaller dimension of 2D descriptors. By discarding 2D descriptors, our method reduces the map size significantly. Our method has larger size of map on Aachen dataset as we use 16 NeRFs to represent the whole scene.

Although our method has much smaller size than HMs, its size is still much larger than that of APRs and SCRs. Therefore, how to design an accurate and efficient localization system is still a challenging task and is worth of further exploration in the future.

E. Ablation study

We conduct an ablation study to explore the influence of different patch sizes to pose accuracy. Table V shows as the patch size increases from 8×9 to 15×15 , the pose accuracy also increases. It is more obvious on Kings College as it is an outdoor scene and the query and reference images have larger viewpoint and illumination changes. However, for indoor heads, due to little changes between the query and reference images, the improvement of increasing patch size is not obvious.

Moreover, as the patch size increases, it takes longer to render a patch. Therefore, the final solution is the balance between accuracy and efficiency. For indoor scenes without large changes between query and reference images, we suggest using smaller patch sizes to high efficiency. For outdoor scenes with large viewpoint, illumination changes between the query and reference images, a larger patch size can bring better accuracy.

F. Qualitative results

Rendered image. Fig. 3 shows the rendered images in 7Scenes [25], CambridgeLandmarks [4] and Aachen [26] datasets. As 7Sevens dataset consists of only indoor scenes, the rendered images have higher quality. CambridgeLandmarks and Aachen datasets are both outdoor datasets, so there are still some artifacts on rendered images due to dynamic objects (e.g. pedestrian, trees), large depths and insufficient observations. Since our method relies mainly on

Patch size	heads	Kings College	Rendering time (ms)
9×9	0.3, 0.4, 97.2%	10.6, 0.2, 82.8%	0.65
11×1	0.2, 0.3, 97.8	9.5, 0.1, 88.9%	0.96
13×13	0.2, 0.3, 98.6%	9.0, 0.1, 90.1%	1.35
15×15	0.2, 0.3, 99.6%	9.1, 0.1, 90.7%	1.79
15×15 (GT)	0.2, 0.3, 99.5%	8.2, 0.1, 94.2%	1.79

TABLE V

ABLATION STUDY OF PATCH SIZE. WE VERIFY THE EFFICACY OF DIFFERENT SIZES OF PATCHES TO THE POSE ACCURACY. FOR CHESS, WE REPORT THE MEDIAN TRANSLATION (CM), ROTATION ($^{\circ}$) AND POSE ACCURACY AT THE ERROR THRESHOLD OF 5cm, 5° . FOR KINGS COLLEGE, WE REPORT THE MEDIAN TRANSLATION (CM), ROTATION ($^{\circ}$) AND POSE ACCURACY AT THE ERROR THRESHOLD OF 25cm, 2° . WE ALSO PROVIDE THE TIME OF RENDERING A PATCH.



Fig. 5. **Failed cases.** We visualize the failed cases of rendered images (left) and their corresponding groundtruth images (right).

sparse patches with corresponding 3D points the map, the influence of artifacts to localization is partially mitigated.

Matching. We also visualize the matches between query and rendered reference images in Fig. 4. Although we show the whole image, in matching process, only patches with size of 15×15 at the center of sparse keypoints are used for matching. Due to the high quality of rendered images, we are able to obtain good matches between the query and rendered images, guaranteeing the success of localization even in outdoor scenes.

Failed cases. Fig. 5 shows failed cases of rendered images. We failed to render high quality images mainly because of insufficient 2D observations of physical 3D points and dynamic objects. The inaccurate camera poses also have negative influence to the rendering quality. This happens mainly on Aachen dataset probably because of insufficient images for some areas.

V. CONCLUSION

In this paper, we propose a new method of applying NeRFs to visual localization task. In detail, we introduce the explicit geometric map (EGM) and implicit learned map (ILM) to provide sparse keypoints and rendered patches to build sparse matches between query and rendered images. By adopting sparse rendering from sparse points provided by EGM, our approach avoids time-consuming full image rendering. With ILM represented by NeRFs, our method discards the memory-consuming 2D descriptors. Therefore, our system is more efficient. However, the accuracy on large-scale Aachen dataset is still limited compared to state-of-the-art methods. We hope this work could be a baseline and the more researchers can make it better in the future.

REFERENCES

- [1] F. Xue, I. Budvytis, D. O. Reino, and R. Cipolla, "Efficient Large-scale Localization by Global Instance Recognition," in *CVPR*, 2022.
- [2] C. Toft, E. Stenborg, L. Hammarstrand, L. Brynte, M. Pollefeys, T. Sattler, and F. Kahl, "Semantic match consistency for long-term visual localization," in *ECCV*, 2018.
- [3] P.-E. Sarlin, C. Cadena, R. Siegwart, and M. Dymczyk, "From Coarse to Fine: Robust Hierarchical Localization at Large Scale," in *CVPR*, 2019.
- [4] A. Kendall, M. Grimes, and R. Cipolla, "Posenet: A convolutional network for real-time 6-dof camera relocalization," in *IJCV*, 2015.
- [5] A. Kendall and R. Cipolla, "Geometric loss functions for camera pose regression with deep learning," in *CVPR*, 2017.
- [6] F. Xue, X. Wang, Z. Yan, Q. Wang, J. Wang, and H. Zha, "Local supports global: Deep camera relocalization with sequence enhancement," in *ICCV*, 2019.
- [7] F. Xue, X. Wu, S. Cai, and J. Wang, "Learning multi-view camera relocalization with graph neural networks," in *CVPR*, 2020.
- [8] B. Wang, C. Chen, C. X. Lu, P. Zhao, N. Trigoni, and A. Markham, "Atloc: Attention guided camera localization," in *AAAI*, 2020.
- [9] S. Brahmabhatt, J. Gu, K. Kim, J. Hays, and J. Kautz, "MapNet: Geometry-aware learning of maps for camera localization," in *CVPR*, 2018.
- [10] J. Shotton, B. Glocker, C. Zach, S. Izadi, A. Criminisi, and A. Fitzgibbon, "Scene coordinate regression forests for camera relocalization in rgb-d images," in *CVPR*, 2013.
- [11] E. Brachmann, A. Krull, S. Nowozin, J. Shotton, F. Michel, S. Gumhold, and C. Rother, "DSAC-differentiable RANSAC for camera localization," in *CVPR*, 2017.
- [12] X. Wu, H. Zhao, S. Li, Y. Cao, and H. Zha, "Sc-wls: Towards interpretable feed-forward camera re-localization," in *ECCV*, 2022.
- [13] E. Brachmann, T. Cavallari, and V. A. Prisacariu, "Accelerated coordinate encoding: Learning to relocalize in minutes using rgb and poses," in *CVPR*, 2023.
- [14] T. Sattler, B. Leibe, and L. Kobbelt, "Efficient & effective prioritized for large-scale image-based localization," *TPAMI*, 2016.
- [15] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Epnnp: An accurate o (n) solution to the pnp problem," *IJCV*, 2009.
- [16] 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.
- [17] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," in *CVPR*, 2016.
- [18] S. Hausler, S. Garg, M. Xu, M. Milford, and T. Fischer, "Patch-netvlad: Multi-scale fusion of locally-global descriptors for place recognition," in *CVPR*, 2021.
- [19] F. Radenović, G. Toliás, and O. Chum, "Fine-tuning cnn image retrieval with no human annotation," *TPAMI*, vol. 41, no. 7, pp. 1655–1668, 2018.
- [20] E. Brachmann and C. Rother, "Learning less is more-6d camera localization via 3d surface regression," in *CVPR*, 2018.
- [21] I. Budvytis, M. Teichmann, T. Vojir, and R. Cipolla, "Large scale joint semantic re-localisation and scene understanding via globally unique instance coordinate regression," in *BMVC*, 2019.
- [22] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, 2021.
- [23] J. Liu, Q. Nie, Y. Liu, and C. Wang, "Nerf-loc: Visual localization with conditional neural radiance field," in *ICRA*, 2023.
- [24] A. Moreau, N. Piasco, D. Tsishkou, B. Stanculescu, and A. de La Fortelle, "Lens: Localization enhanced by nerf synthesis," in *CoRL*, 2022.
- [25] B. Glocker, S. Izadi, J. Shotton, and A. Criminisi, "Real-time rgb-d camera relocalization," in *ISMAR*, 2013, pp. 173–179.
- [26] T. Sattler, T. Weyand, B. Leibe, and L. Kobbelt, "Image retrieval for image-based localization revisited," in *BMVC*, 2012.
- [27] X. Li and H. Ling, "Pogo-net: pose graph optimization with graph neural networks," in *ICCV*, 2021.
- [28] X. Li and H. Ling, "Gtcar: Graph transformer for camera relocalization," in *ECCV*, 2022.
- [29] M. O. Turkoglu, E. Brachmann, K. Schindler, G. J. Brostow, and A. Monszpart, "Visual camera re-localization using graph neural networks and relative pose supervision," in *3DV*, 2021.
- [30] H. Li, P. Xiong, H. Fan, and J. Sun, "Dfanet: Deep feature aggregation for real-time semantic segmentation," in *CVPR*, 2019.
- [31] T. Sattler, Q. Zhou, M. Pollefeys, and L. Leal-Taixe, "Understanding the limitations of cnn-based absolute camera pose regression," in *CVPR*, 2019.
- [32] E. Brachmann and C. Rother, "Visual camera re-localization from rgb and rgb-d images using dsac," *TPAMI*, vol. 44, no. 9, pp. 5847–5865, 2022.
- [33] X. Li, S. Wang, Y. Zhao, J. Verbeek, and J. Kannala, "Hierarchical scene coordinate classification and regression for visual localization," in *CVPR*, 2020.
- [34] J. Valentin, A. Dai, M. Nießner, P. Kohli, P. Torr, S. Izadi, and C. Keskin, "Learning to navigate the energy landscape," in *3DV*, 2016.
- [35] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *IJCV*, 2004.
- [36] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," in *ICCV*, 2011.
- [37] J. Sivic and A. Zisserman, "Efficient visual search of videos cast as text retrieval," *TPAMI*, vol. 31, no. 4, pp. 591–606, 2008.
- [38] D. Gálvez-López and J. D. Tardos, "Bags of binary words for fast place recognition in image sequences," *T-RO*, vol. 28, no. 5, pp. 1188–1197, 2012.
- [39] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *CVPRW*, 2018.
- [40] J. Revaud, P. Weinzaepfel, C. R. de Souza, and M. Humenberger, "R2D2: Repeatable and reliable detector and descriptor," in *NeurIPS*, 2019.
- [41] M. J. Tyszkiewicz, P. Fua, and E. Trulls, "DISK: Learning local features with policy gradient," in *NeurIPS*, 2020.
- [42] F. Xue, I. Budvytis, and R. Cipolla, "SFD2: Semantic-guided Feature Detection and Description," in *CVPR*, 2023.
- [43] K. M. Yi, E. Trulls, V. Lepetit, and P. Fua, "Lift: Learned invariant feature transform," in *ECCV*, 2016.
- [44] 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 *CVPR*, 2019.
- [45] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superglue: Learning feature matching with graph neural networks," in *CVPR*, 2020.
- [46] 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 *ICCV*, 2021.
- [47] F. Xue, I. Budvytis, and R. Cipolla, "Imp: Iterative matching and pose estimation with adaptive pooling," in *CVPR*, 2023.
- [48] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-nerf 360: Unbounded anti-aliased neural radiance fields," in *CVPR*, 2022.
- [49] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Zip-nerf: Anti-aliased grid-based neural radiance fields," in *ICCV*, 2023.
- [50] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.

- [51] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics (ToG)*, vol. 42, no. 4, pp. 1–14, 2023.
- [52] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *CVPR*, 2016.
- [53] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar, "Block-nerf: Scalable large scene neural view synthesis," in *CVPR*, 2022.
- [54] D. Arthur and S. Vassilvitskii, "K-means++ the advantages of careful seeding," in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 2007.
- [55] Y. Shavit and Y. Keller, "Camera pose auto-encoders for improving pose regression," in *ECCV*, 2022.
- [56] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *NeurIPS*, 2019.
- [57] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.
- [58] E. Brachmann and C. Rother, "Expert sample consensus applied to camera re-localization," in *CVPR*, 2019.
- [59] T. Sattler, W. Maddern, C. Toft, A. Torii, L. Hammarstrand, E. Stenborg, D. Safari, M. Okutomi, M. Pollefeys, J. Sivic, *et al.*, "Benchmarking 6dof outdoor visual localization in changing conditions," in *CVPR*, 2018.