# SIM2E: Benchmarking the Group Equivariant Capability of Correspondence Matching Algorithms

Shuai Su⊙, Zhongkai Zhao⊙, Yixin Fei⊙, Shuda Li⊙,
Qijun Chen⊙, and Rui Fan⊙

Tongji University, Shanghai 201804, China.
{sushuai, kanez, amyfei, qjchen}@tongji.edu.cn,
shuda.dexter.li@gmail.com, rui.fan@ieee.org

**Abstract.** Correspondence matching is a fundamental problem in computer vision and robotics applications. Solving correspondence matching problems using neural networks has been on the rise recently. Rotation-equivariance and scale-equivariance are both critical in correspondence matching applications. Classical correspondence matching approaches are designed to withstand scaling and rotation transformations. However, the features extracted using convolutional neural networks (CNNs) are only translation-equivariant to a certain extent. Recently, researchers have strived to improve the rotation-equivariance of CNNs based on group theories. Sim(2) is the group of similarity transformations in the 2D plane. This paper presents a specialized dataset dedicated to evaluating sim(2)-equivariant correspondence matching algorithms. We compare the performance of 16 state-of-the-art (SoTA) correspondence matching approaches. The experimental results demonstrate the importance of group equivariant algorithms for correspondence matching on various sim(2) transformation conditions. Since the subpixel accuracy achieved by CNN-based correspondence matching approaches is unsatisfactory, this specific area requires more attention in future works. Our dataset is publicly available at: mias.group/SIM2E.

**Keywords:** correspondence matching, computer vision, robotics, rotation-equivariance, scaling-equivariance, convolutional neural networks

## 1 Introduction

Correspondence matching is a key component in autonomous driving perception tasks, such as object tracking [1], simultaneous localization and mapping (SLAM) [2], multi-camera online calibration [3], 3D geometry reconstruction [4], panorama stitching [5], and camera pose estimation [6], as shown in Fig. 1. Sim(2) transformation consists of rotation, scaling, and translation. Sim(2)-equivariant correspondence matching is significantly important for autonomous driving, as vehicles often veer abruptly. Classical algorithms leverage a detector, a descriptor, and a matcher to determine correspondences. The detector and
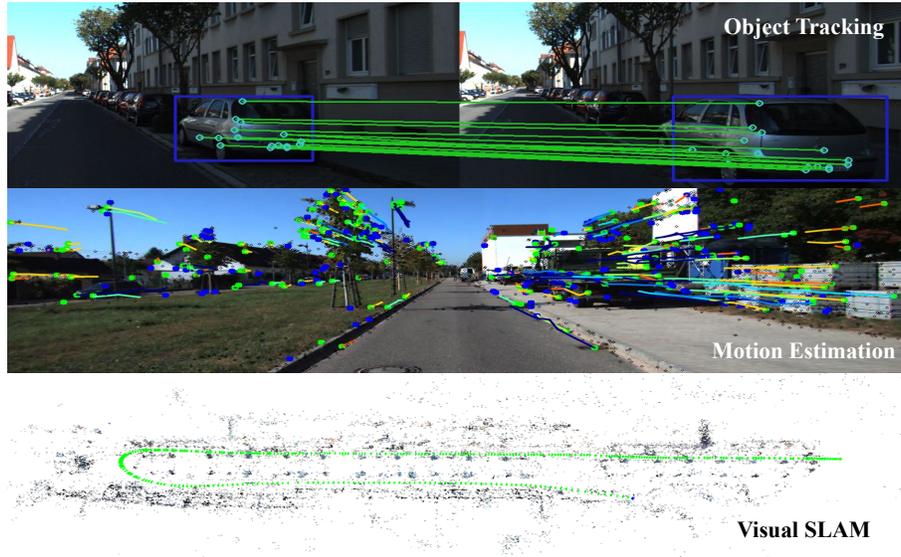
**Fig. 1.** Autonomous driving perception tasks involving correspondence matching.

descriptor provide the locations and the descriptions of interest points (point-like features in an image), and the matcher produces the final correspondences. Moravec *et al.* [7] presented the concept of interest points. Harris [8] judges whether the pixel is a corner based on the local image gradient changes. The scale-invariant feature transform (SIFT) [9] is a rotation-invariant and scale-invariant algorithm that consists of a detector and a descriptor. The distance among descriptors is computed using cosine distance. The correspondences of a given image pair are determined using the nearest neighbor matching algorithms. As a hand-crafted algorithm, SIFT [9] achieves rotation-invariance by computing the main directions of local features (in an image patch of 16x16 pixels.). ASIFT [10] aimed to improve the performance of SIFT [9] on affine transformation. Oriented FAST [11] and Rotated Binary Robust Independent Elementary Features (BRIEF) [12] (ORB) [13] greatly minimize the trade-off between accuracy and speed and have been widely used in visual SLAM [13–15].

Deep learning has been applied successfully in numerous computer vision tasks in recent years. LIFT [16] is an architecture of learning-based rotation-invariant feature detection and description approach. It consists of three modules: detector, orientation estimator, and descriptor. Similar to SIFT [9], a scale-space pyramid is used to obtain multi-scale correspondence detection results. SuperPoint [17] is a self-supervised framework for correspondence detection and description. It is a two-stage method: (1) in the first stage, a feature extractor is trained on a synthetic dataset generated by rendering patterns of corners; (2) in the second stage, the descriptor is trained on images from the COCO-dataset [18] that are augmented by random homography matrices including transformations

such as rotation, scaling, and translation. Unlike SuperPoint, D2Net [19] is a one-stage approach that jointly detects and describes correspondences. D2Net is trained using the correspondences obtained from large-scale structure from motion (SfM) reconstructions. R2D2 [20] proposes a framework to find more repeatable correspondences, and it can simultaneously estimate the reliability and repeatability of correspondences. DISK [21] uses reinforcement learning to realize end-to-end correspondence matching. It performs better at small angle changes but worse than SuperPoint [17] at large angle changes.

SuperGlue [22] achieves superior performance by modeling correspondence matching as a graph matching problem. The inputs of the graph neural network in SuperGlue [22] include descriptors, positions, and scores of keypoints. The Sinkhorn algorithm [23] is utilized to solve the optimal-transport problem. However, the graph edges in SuperGlue exponentially grow as the number of correspondences increases. SGMNet [24] uses a seeded graph to reduce computation and memory costs significantly. LoFTR [25] is a detector-free and end-to-end architecture. It uses convolutional neural networks (CNN) as feature extractors and a coarse-to-fine strategy to obtain more accurate pixel-level results. Similar to SuperGlue, it also fuses descriptors with the position information. Unlike LoFTR, another end-to-end correspondence matching network, referred to as MatchFormer [26], uses a hierarchical extract-and-match transformer. It is demonstrated that the correspondence matching operation can also be conducted in the encoder. RoRD [27] uses orthographic view generation to improve correspondence matching by increasing the visual overlap using orthographic projection. It also shows that rotation invariance can be improved by augmenting the training dataset with random rotation, scaling, and perspective transformations.

Group-equivariant convolutional neural networks (G-CNN) are equivariant under a specific transformation (*e.g.*, rotation, translation, *etc.*) which can also be represented by a special group. Researchers have designed G-CNNs using different mathematical approximations. Cohen *et al.* [28] proposed the first G-CNN. Li *et al.* [29] use the cyclic replacement to achieve P4-group equivariance. Cohen *et al.* [30] use the Fast Fourier Transform (FFT) to approximate the integral of a group. E2-CNN [31] is a general G-CNN framework that analyzes and models the orientation and symmetry of images. GIFT [32] is a rotation-equivariant and scaling-equivariant descriptor based on G-CNN. It uses E2-CNN [31] rather than conventional CNNs to describe local visual features. On the other hand, SEKD [33] is a group-equivariant correspondence detector based on G-CNN, which greatly improves the performance of rotation-equivariant correspondence matching. ReF [34] is a rotation-equivariant correspondence detection and description framework. It uses a G-CNN to extract group-equivariant feature maps and a group-pooling operation to get rotation-invariant descriptors. SE2-LoFTR [35] replaces the feature extractor of LoFTR with E2-CNN, achieving significantly better results on the rotated-HPatches dataset [27]. Furthermore, it also mentioned in [35] that the the position information is not rotation-equivariant while the descriptor is rotation-equivariant. The methods mentioned

above only consider the equivariance of local features. Unfortunately, the equivariance of position information is rarely discussed. Cieslewski *et al.* [36] presented an algorithm to match correspondences without descriptors, namely, only position information is used. This algorithm is evaluated on the KITTI [37] dataset (containing relatively ideal scenarios), demonstrating robust performance even without descriptors. Similar to [36], ZZ-Net [38] is an algorithm for matching two 2D point clouds. It demonstrates that correspondence matching without descriptors can work in rotation-only conditions. Therefore, the current research on the equivariance of position information needs to be further expanded.

## 2   SIM2E Dataset

### 2.1   Data Collection and Augmentation

To ensure the pixel-level accuracy of correspondence matching ground truth, we scrape frames from online time-lapse videos. The cameras used to capture such time-lapse videos are fixed. Our SIM2E dataset contains many challenging scenarios, such as moving clouds in the sky and changing illumination conditions. We choose the first frame of each video as the reference image and use the rest of the frames as target (query) images. We also publish our data augmentation code so that interested readers can conduct sim(2) transformations on our dataset according to their own needs.

### 2.2   SIM2E-SO2S, SIM2E-Sim2S, and SIM2E-PersS

The rotation and scaling operations produce many black backgrounds. To increase the difficulty of correspondence matching, we generate synthetic backgrounds to fill these black areas. Our dataset is split into three subsets: SIM2E-SO2S, SIM2E-Sim2S, and SIM2E-PersS.

- **SIM2E-SO2S (Rotation and Synthetic Background)**: The target images are rotated by random angles between $0°$ and $360°$. Scaling is not applied.
- **SIM2E-Sim2S (Rotation, Scaling, Translation and Synthetic Background)**: The target images are rotated by random angles between $0°$ and $360°$. Random scaling ranging between 0.4 and 1, and random translation ranging between 0 and 0.2 are also applied.
- **SIM2E-PersS: (Perspective Transformation and Synthetic Background)**: Random perspective transformations are applied to the target images, where the perspective parameters (the two elements on the 3rd row, the 1st and 2nd columns of the homography matrix, respectively) are random values between -0.0008 and 0.0008. The shear angle is randomly set to $[-10°, 10°]$. The target images are rotated by random angles between $0°$ and $360°$. Random scaling ranging between 0.4 and 1 is applied. Random translation ranging between 0 and 0.2 is applied.
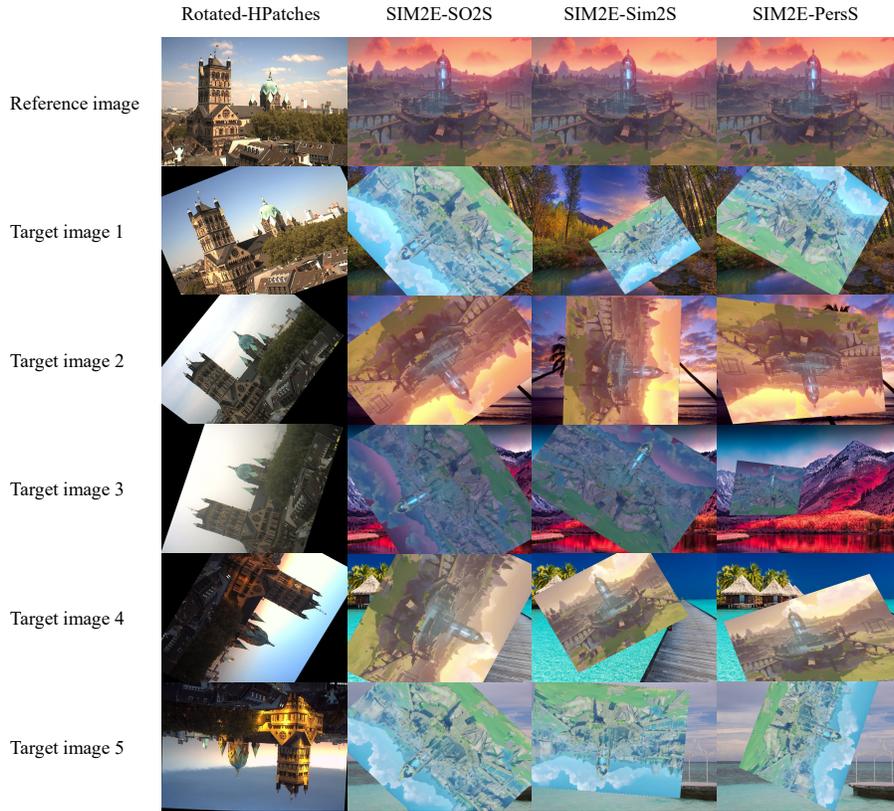
**Fig. 2.** Rotated-HPatches dataset and the three sub-sets of our created SIM2E dataset.

## 2.3 Comparison with Other Public Correspondence Matching Datasets

As shown in Table. 1, the existing datasets for correspondence matching algorithm evaluation can be grouped into two types: 3D scenes [39–44] and planar scenes [27, 45].

Aachen Day-Night [39] is a public dataset designed to evaluate the performance of outdoor visual localization algorithms in changing illumination conditions (day-time and night-time). The dataset contains a scenario where images were taken with a hand-held camera at different times of the day. It is widely used to evaluate the performance of correspondence matching algorithms, especially when the illumination change is significant.

ScanNet [41] is a large-scale real-world dataset containing 2.5M RGB-D images (1513 scans acquired in 707 different places, such as offices, apartments, and bathrooms). All the scans are annotated with estimated calibration parameters, camera poses, reconstructed 3D surfaces, textured meshes, dense object-level semantic segmentations, and aligned computer-aided design (CAD) models.

**Table 1.** Comparison between our SIM2E dataset and other public datasets.

| Dataset | Type | Illumination Change | Rotation | Scaling | Dataset Size |
|---|---|---|---|---|---|
| AachenDayNight [39] | 3D | significant | small | medium | large |
| ScanNet [41] | 3D | slight | small | small | large |
| MegaDepth [42] | 3D | slight | small | large | large |
| Inloc [43] | 3D | medium | small | medium | large |
| TartanAir [44] | 3D | very significant | small | medium | large |
| Hpatches [45] | plane | significant | small | small | small |
| Rotated-HPatches [27] | plane | significant | large | medium | small |
| **SIM2E (ours)** | plane | very significant | large | large | small |

MegaDepth [42] is a large-scale dataset for the evaluation of depth estimation and/or correspondence matching algorithms. It uses SfM and multi-view stereo (MVS) techniques to acquire 3D point clouds, which can then be used to train and evaluate single-view depth estimation and/or correspondence matching networks. However, the 3D point clouds generated using SfM/MVS in the MegaDepth dataset are not sufficiently accurate and dense.

Compared to the Aachen Day-Night [39] and MegaDepth [42] datasets which were created in outdoor scenarios, the Inloc [43] dataset focuses on indoor localization problems. The Inloc dataset consists of a database of RGB-D images, geometrically registered to the floor maps and augmented with a separate set of RGB target images (annotated with manually verified ground-truth 6DoF camera poses in the global coordinate system of the 3D map).

Unlike the aforementioned datasets that are relatively ideal in terms of either motion or illumination conditions, TartanAir [44], a synthetic dataset used to evaluate visual SLAM algorithms, is collected using a photo-realistic simulator (with the presence of moving objects, changing illumination and weather conditions). Such a more challenging dataset fills the gap between synthetic and real-world datasets.

Hpatches [45] are created using other public datasets. It can be split into two subsets: illumination and viewpoint, which are two crucial aspects of correspondence matching. It can also be split into three subsets: EASY, HARD, and TOUGH, according to the sizes of the overlapping areas between reference and target images. Randomly rotating the target images in the Hpatches [45] dataset produces a new dataset, referred to as Rotated-HPatches [27].

Similar to the Hpatches and Rotated-HPatches datasets, our SIM2E dataset provides accurate subpixel correspondence matching ground truth. On the other hand, the illumination, rotation, and scaling changes are significant in our dataset. Therefore, compared to other existing public datasets, our SIM2E dataset can be used to evaluate the sim(2)-equivariant capability of correspondence matching algorithms more comprehensively. However, the size of the current version of our SIM2E dataset is small. We will therefore increase its size in our future work.

The rotation distributions of Rotated-HPatches and our SIM2E subsets are shown in Fig. 3. It can be observed that in the Rotated-HPatches dataset, slight
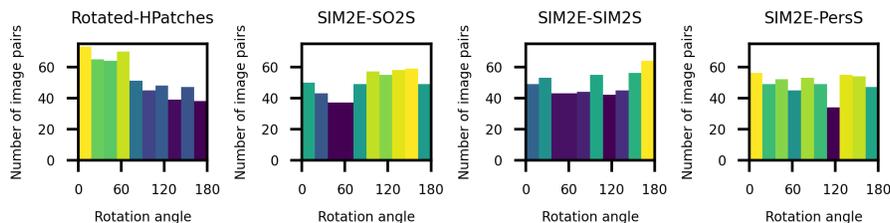
**Fig. 3.** Rotation distributions of the Rotated-HPatches and our SIM2E datasets.

rotations ($\leq 60°$) account for a large proportion. In contrast, the three subsets of our SIM2E dataset are uniformly distributed. Most existing learning-based correspondence matching approaches have poor rotation-equivariant capabilities, and their performances are satisfactory only when there are slight rotations. Therefore, our SIM2E dataset can provide more acceptable results when evaluating the rotation-equivariant capability of a given correspondence matching algorithm.

## 3 Experiments

### 3.1 Experimental Setup

The group-equivariant capabilities of six classical and ten learning-based correspondence matching approaches are evaluated on our SIM2E dataset.

For classical correspondence matching approaches, we use the OpenCV [46] implementations of AKAZE [47], BRISK [48], KAZE [49], ORB [13], FREAK [50], and SIFT [9] in our experiments. All these classical approaches use the nearest neighbor matching algorithm for correspondence matching. The ratio test technique (threshold is set to 0.7) is also used to improve the overall performance.

For learning-based correspondence matching approaches, we use the official weights of each model. These models were trained on different datasets, as detailed below:

- **SuperPoint** [17] is trained on the MS-COCO [18] dataset, a large-scale dataset for object detection and segmentation.
- **R2D2** [20] is trained on the Aachen Day-Night [39] dataset and a retrieval dataset [51].
- **ALIKE** [52] is trained on the MegaDepth [42] dataset.
- **GIFT** [32] is trained on the MS-COCO [18] dataset and finetuned on the GL3D [53] dataset (consisting of indoor and outdoor scenes).
- **RoRD** [27] is trained on the PhotoTourism [54] dataset, where the 3D structures of scenes are obtained using SfM.
- **SuperGlue** [22] is trained with the indoor models in the ScanNet [41] dataset and the outdoor models in the MegaDepth [42] dataset.

- **SGMNet** [24] is trained on the GL3D [53] dataset. Our experiments utilize the SIFT version of SGMNet, where the detector and descriptor are rotation-invariant.
- **LoFTR** [25] is trained with the same experimental setup as SuperGlue.
- **MatchFormer** [26] is trained with the same experimental setup as SuperGlue and LoFTR. Limited by our GPU memory, the lightweight version of MatchFormer is used in our experiments.
- **SE2-LoFTR** [35] is trained on the MegaDepth dataset.

Furthermore, the mean matching accuracy (MMA) is employed to quantify the performance of the aforementioned correspondence matching algorithms, which are run on a PC with an Intel Core i7-10870H CPU and an NVIDIA RTX3080-laptop GPU (having a 16GB DDR4 memory).

### 3.2   Comparison of the SoTA approaches on the Rotated-HPatches Dataset

The Rotated-HPatches [27] dataset is generated using the Hpatches [45] dataset to evaluate rotation-equivariant capability of correspondence matching methods. Each sub-folder of the Rotated-HPatches dataset contains one reference image and five target images. The target images are obtained by rotating the reference image at a random angle. The correspondence matching ground truth is acquired using the homography matrices between each pair of reference and target images. As illustrated in Fig. 4(a), the SoTA correspondence matching algorithms demonstrate significantly different performances on the Rotated-HPatches dataset.

The classical algorithms, such as AKAZE, BRISK, KAZE, and SIFT, achieve the best overall performances on the Rotated-HPatches dataset, as they consider both the scaling and rotation invariance of visual features. Benefiting from the higher dimensional feature descriptors, these four algorithms outperform ORB and FREAK.

On the other hand, SuperPoint, R2D2, and ALIKE are developed without considering rotation invariance. Therefore, their performances are relatively poor on the Rotated-HPatches dataset. GIFT [32] uses SuperPoint as the feature detector. Its feature descriptor is developed based on G-CNN to acquire the rotation-equivariant capability. As expected, GIFT significantly outperforms SuperPoint.

As can be seen from Table 2 and Fig. 4(a), when the tolerance $\delta$ exceeds 5, the learning-based methods demonstrate better performances than classical methods. For instance, SGMNet outperforms all classical methods when $\delta > 5$ and SE2-LoFTR shows similar performance to the classical methods when $\delta > 8$. Referring to [24], SGMNet is a lightweight version of SuperGlue and demonstrates slightly worse performance than SuperGlue. However, SuperGlue with SuperPoint performs much worse than SGMNet. This is probably because SGMNet uses SIFT as its detector and descriptor, which has the rotation-equivariant capability.

**Table 2.** The performance of SoTA correspondence matching approaches on the Rotated-HPatches dataset. $N$ denotes the average number of valid matches.

| Method | $\delta \leq 1$ | $\delta \leq 3$ | $\delta \leq 5$ | $\delta \leq 10$ | $N$ |
|---|---|---|---|---|---|
| AKAZE [47] | 0.426 | 0.744 | 0.812 | 0.841 | 203 |
| BRISK [48] | 0.419 | 0.745 | 0.816 | 0.841 | 282 |
| KAZE [49] | 0.423 | 0.753 | **0.825** | 0.858 | 505 |
| ORB [13] | 0.253 | 0.576 | 0.646 | 0.672 | 20 |
| SIFT [9] | **0.484** | **0.762** | 0.809 | 0.830 | 727 |
| FREAK [50] | 0.278 | 0.519 | 0.567 | 0.594 | 27 |
| SP [17] | 0.123 | 0.249 | 0.274 | 0.295 | 120 |
| R2D2 [20] | 0.057 | 0.129 | 0.144 | 0.153 | 158 |
| ALIKE [52] | 0.111 | 0.172 | 0.184 | 0.201 | 59 |
| GIFT [32] | 0.203 | 0.406 | 0.439 | 0.448 | 186 |
| RoRD [27] | 0.030 | 0.191 | 0.378 | 0.620 | 1077 |
| SPSG [22] | 0.185 | 0.425 | 0.491 | 0.552 | 479 |
| SGMNet [24] | 0.369 | 0.688 | 0.814 | **0.893** | 1278 |
| LoFTR [25] | 0.037 | 0.167 | 0.259 | 0.350 | 506 |
| MatchFormer [26] | 0.033 | 0.162 | 0.256 | 0.358 | 600 |
| SE2-LoFTR [35] | 0.197 | 0.556 | 0.720 | 0.842 | 1305 |

It can be observed that SoTA classical methods always perform better than learning-based methods when the tolerance is small. With the decrease in tolerance, the MMA scores achieved by learning-based methods drop considerably. This is probably because the learning-based methods are generally trained via self-supervised learning, where the tolerance to determine positive samples is typically set to 3. Furthermore, the model is difficult to converge in the training phase when the tolerance is too small, *e.g.*, less than 1, while the model's accuracy degrades dramatically when the tolerance is too large. Moreover, compared to classical methods, learning-based methods can always obtain more correspondences. Therefore, improving the subpixel accuracy of learning-based approaches without reducing the number of valid matches is a research area that requires more attention.

### 3.3   Comparison of the SoTA approaches on our SIM2E Dataset

In this paper, we quantify the group-equivalent capabilities of the aforementioned algorithms on the three subsets of our SIM2E Dataset.

**Experimental results on the SIM2E-SO2S subset** Our SIM2E-SO2S subset is created in a similar fashion to the Rotated-HPatches dataset. Since we select time-lapse videos with more challenging illumination conditions, apply more uniformly distributed random rotations, and add synthetic backgrounds to the target image, the SIM2E-SO2S subset is expected to reflect the correspondence matching algorithms' group-equivariant capabilities more comprehensively. As can be observed from Fig. 4(b), all the SoTA methods perform
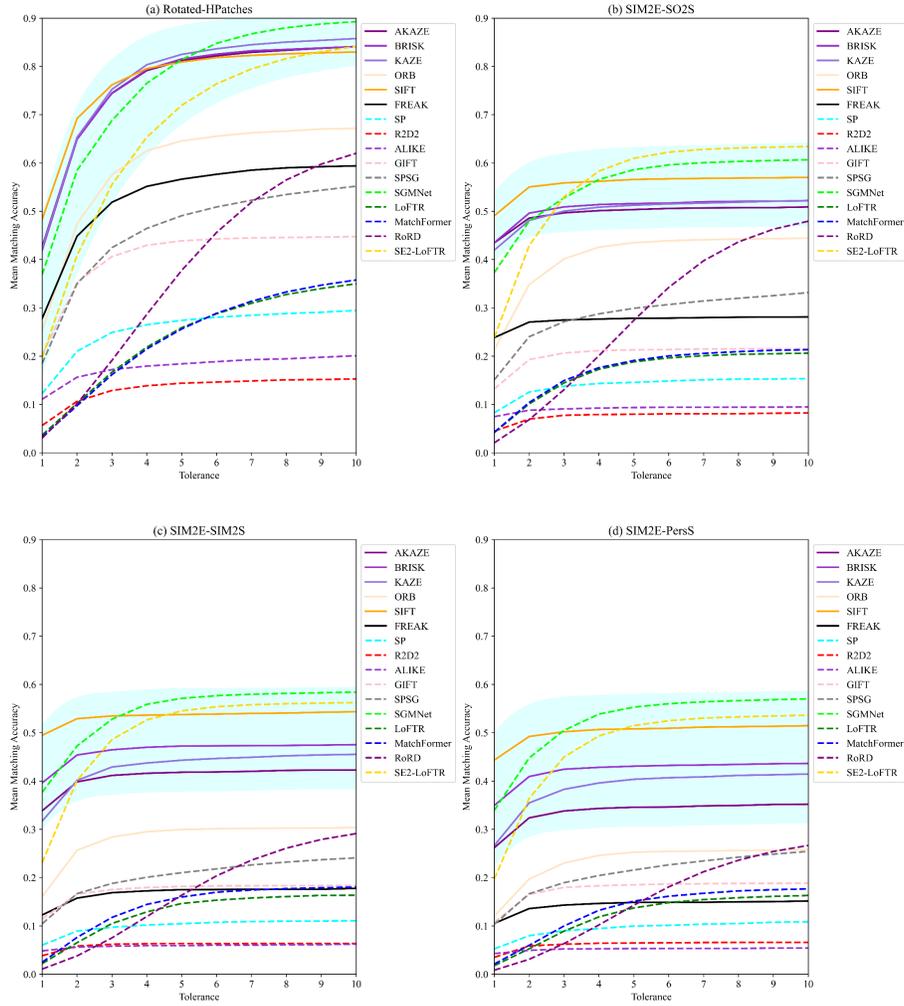
**Fig. 4.** MMA results of six classical and ten learning-based correspondence matching approaches on the Rotated-HPatches dataset and our SIM2E-SO2S, SIM2E-SIM2S, and SIM2E-PersS subsets. MMA results for the classical methods are shown as solid lines, while MMA results for the learning-based methods are shown as dashed lines.

much worse on our SIM2E-SO2S subset because it is more challenging than the Rotated-HPatches dataset. Furthermore, similar to the experimental results in the Rotated-HPatches experiments (see Fig. 4(a)), SE2-LoFTR, SGMNet, AKAZE, BRISK, KAZE, and SIFT also achieve the best group-equivariant capabilities on our SIM2E-SO2S subset (see Fig. 4(b)). This validates the effectiveness of our SIM2E-SO2S subset in terms of evaluating a correspondence matching algorithm's group-equivariant capability. Moreover, the MMA scores

**Table 3.** The performance of SoTA correspondence matching approaches on the SIM2E-SO2S subset. $N$ denotes the average number of valid matches.

| Method | $\delta \leq 1$ | $\delta \leq 3$ | $\delta \leq 5$ | $\delta \leq 10$ | $N$ |
|---|---|---|---|---|---|
| AKAZE [47] | 0.435 | 0.497 | 0.504 | 0.509 | 72 |
| BRISK [48] | 0.435 | 0.509 | 0.516 | 0.522 | 113 |
| KAZE [49] | 0.420 | 0.500 | 0.513 | 0.522 | 129 |
| ORB [13] | 0.213 | 0.402 | 0.435 | 0.444 | 24 |
| SIFT [9] | **0.491** | **0.559** | 0.566 | 0.570 | 264 |
| FREAK [50] | 0.238 | 0.275 | 0.278 | 0.281 | 9 |
| SP [17] | 0.083 | 0.138 | 0.146 | 0.154 | 40 |
| R2D2 [20] | 0.045 | 0.078 | 0.080 | 0.082 | 53 |
| ALIKE [52] | 0.075 | 0.091 | 0.094 | 0.095 | 22 |
| GIFT [32] | 0.132 | 0.207 | 0.213 | 0.215 | 51 |
| RoRD [27] | 0.021 | 0.131 | 0.274 | 0.480 | 630 |
| SPSG [22] | 0.151 | 0.271 | 0.299 | 0.332 | 192 |
| SGMNet [24] | 0.373 | 0.528 | 0.586 | 0.607 | 1226 |
| LoFTR [25] | 0.042 | 0.145 | 0.188 | 0.206 | 329 |
| MatchFormer [26] | 0.042 | 0.150 | 0.191 | 0.213 | 451 |
| SE2-LoFTR [35] | 0.239 | 0.530 | **0.610** | **0.634** | 1640 |

**Table 4.** The performance of SoTA correspondence matching approaches on the SIM2E-SIM2S subset. $N$ denotes the average number of valid matches.

| Method | $\delta \leq 1$ | $\delta \leq 3$ | $\delta \leq 5$ | $\delta \leq 10$ | $N$ |
|---|---|---|---|---|---|
| AKAZE [47] | 0.338 | 0.412 | 0.418 | 0.423 | 22 |
| BRISK [48] | 0.396 | 0.465 | 0.472 | 0.475 | 51 |
| KAZE [49] | 0.317 | 0.429 | 0.443 | 0.455 | 54 |
| ORB [13] | 0.161 | 0.284 | 0.300 | 0.303 | 10 |
| SIFT [9] | **0.495** | **0.535** | 0.538 | 0.544 | 165 |
| FREAK [50] | 0.123 | 0.169 | 0.175 | 0.178 | 3 |
| SP [17] | 0.060 | 0.097 | 0.104 | 0.110 | 20 |
| R2D2 [20] | 0.038 | 0.062 | 0.063 | 0.064 | 23 |
| ALIKE [52] | 0.048 | 0.058 | 0.059 | 0.062 | 6 |
| GIFT [32] | 0.114 | 0.175 | 0.182 | 0.184 | 32 |
| RoRD [27] | 0.011 | 0.075 | 0.164 | 0.291 | 287 |
| SPSG [22] | 0.105 | 0.188 | 0.210 | 0.241 | 118 |
| SGMNet [24] | 0.377 | 0.528 | **0.571** | **0.584** | 802 |
| LoFTR [25] | 0.022 | 0.105 | 0.146 | 0.164 | 155 |
| MatchFormer [26] | 0.025 | 0.118 | 0.160 | 0.180 | 232 |
| SE2-LoFTR [35] | 0.231 | 0.486 | 0.545 | 0.563 | 847 |

achieved by these six methods differ more significantly in the SIM2E-SO2S experiments. Therefore, our SIM2E-SO2S subset more comprehensively quantifies the group-equivariant capability of a given correspondence matching algorithm.

**Table 5.** The performance of SoTA correspondence matching approaches on the SIM2E-PersS subset. $N$ denotes the average number of valid matches.

| Method | $\delta \leq 1$ | $\delta \leq 3$ | $\delta \leq 5$ | $\delta \leq 10$ | $N$ |
|---|---|---|---|---|---|
| AKAZE [47] | 0.262 | 0.338 | 0.346 | 0.352 | 14 |
| BRISK [48] | 0.350 | 0.424 | 0.431 | 0.436 | 34 |
| KAZE [49] | 0.267 | 0.383 | 0.404 | 0.414 | 42 |
| ORB [13] | 0.120 | 0.230 | 0.253 | 0.258 | 7 |
| SIFT [9] | **0.443** | 0.502 | 0.508 | 0.515 | 130 |
| FREAK [50] | 0.105 | 0.143 | 0.149 | 0.152 | 3 |
| SP [17] | 0.052 | 0.090 | 0.100 | 0.108 | 13 |
| R2D2 [20] | 0.035 | 0.062 | 0.065 | 0.066 | 16 |
| ALIKE [52] | 0.043 | 0.052 | 0.053 | 0.054 | 5 |
| GIFT [32] | 0.107 | 0.180 | 0.185 | 0.188 | 24 |
| RoRD [27] | 0.008 | 0.063 | 0.143 | 0.267 | 257 |
| SPSG [22] | 0.103 | 0.189 | 0.216 | 0.254 | 120 |
| SGMNet [24] | 0.339 | **0.505** | **0.553** | **0.570** | 758 |
| LoFTR [25] | 0.018 | 0.089 | 0.137 | 0.163 | 120 |
| MatchFormer [26] | 0.021 | 0.100 | 0.151 | 0.177 | 161 |
| SE2-LoFTR [35] | 0.197 | 0.450 | 0.515 | 0.536 | 772 |

**Experimental results on the SIM2E-SIM2S subset**  Compared to the SIM2E-SO2S subset, the SIM2E-SIM2S subset contains scaling and translation transformations that shrink the size of the overlapping area between image pairs. As illustrated in Fig. 4(c) and Table 4, the performances of SE2-LoFTR, SGM-Net, SIFT, and BRISK remain stable in the SIM2E-SIM2S experiments, while other models' performances degrade dramatically. Therefore, our SIM2E-SIM2S subset can be used to quantify not only the rotation-equivariant capability but also the scaling-equivariant capability of correspondence matching algorithms.

**Experimental results on the SIM2E-PersS subset**  Compared to the SIM2E-SO2S and SIM2E-SIM2S subsets, the SIM2E-PersS subset contains random perspective transformations. Therefore, correspondence matching on the SIM2E-PersS subset is more challenging. Similarly, in the SIM2E-PersS experiments (see Fig. 4(d) and Table 5), the performances of SE2-LoFTR, SGMNet, and SIFT remain stable, while the performances of BRISK, AKAZE, and KAZE degrade more significantly. Therefore, our SIM2E-PersS subset can be utilized to quantify sim(2)-equivariant capability of correspondence matching algorithms in a more in-depth manner.

### 3.4   Discussion

The experimental results presented above show that SGMNet, SE2-LoFTR, SIFT, BRISK, KAZE, and AKAZE demonstrate similar group-equivariant capabilities. As can be seen in Table 2, the performances of these algorithms are
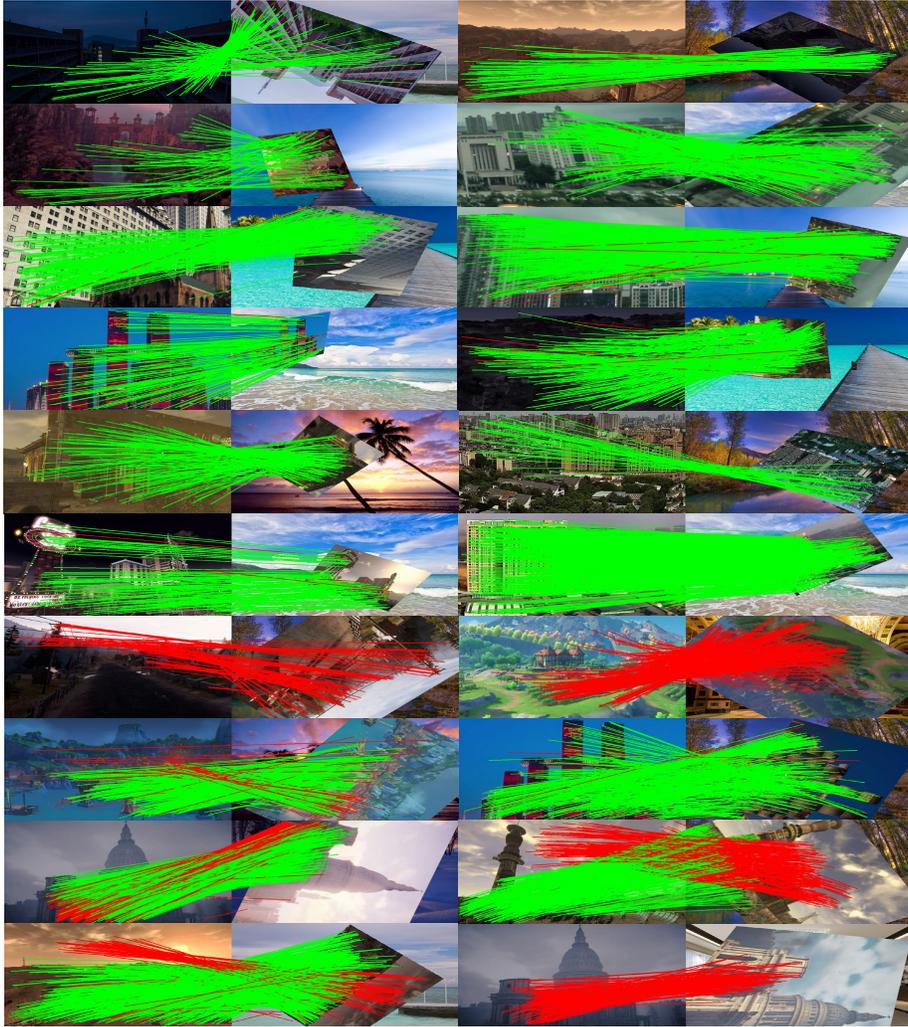
**Fig. 5.** Correspondence matching results on our SIM2E dataset.

very similar, while SE2-LoFTR achieves the worst performance on the Rotated-HPatches dataset. As can be observed in Fig. 4(b)-(d), our created SIM2E dataset can reflect the group-equivariant capabilities of the existing correspondence matching algorithms more comprehensively. The compared SoTA methods achieve the best performances on the SIM2E-SO2S subset and the worst overall performances on the SIM2E-PersS subset. Therefore, we believe that our created SIM2E dataset can help users obtain more objective and in-depth evaluation results of their developed correspondence matching algorithms' group-equivariant capabilities.

## 4   Conclusion

This paper presented a benchmark dataset for the evaluation of sim(2)-equivariant capability of correspondence matching approaches. We first discussed the classical and learning-based methods and the mainstream of developing group-equivariant network architectures. We qualitatively and quantitatively evaluated sixteen SoTA correspondence matching algorithms on the Rotated-HPatches dataset and three subsets of our created SIM2E dataset. These results suggest that our SIM2E dataset is much more challenging than public correspondence matching datasets, and it can comprehensively reflect the group-equivariant capability of SoTA correspondence matching approaches. In summary, group-equivariant detection, group-equivariant description, and group-equivariant position information are vital for group-equivariant correspondence matching. SuperGlue, LoFTR, and SGMNet use neural networks to fuse global position information and local feature descriptors, and achieve superior performances over others. However, obtaining group equivariance of position information is still challenging, as discussed in [35]. The scaling-equivariant and rotation-equivariant capabilities of learning-based approaches are close to classical approaches. However, the sub-pixel accuracy achieved by the former is still unsatisfactory.

## 5   Acknowledgements

## References

1. Huiyu Zhou et al. Object tracking using SIFT features and mean shift. *Computer vision and image understanding*, 113(3):345–352, 2009. 1
2. Yang Yu et al. Accurate and robust visual localization system in large-scale appearance-changing environments. *IEEE/ASME Transactions on Mechatronics*, 2022. DOI: 10.1109/TMECH.2022.3177237. 1
3. Yonggen Ling and Shaojie Shen. High-precision online markerless stereo extrinsic calibration. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1771–1778. IEEE, 2016. 1
4. Rui Fan et al. Road surface 3D reconstruction based on dense subpixel disparity map estimation. *IEEE Transactions on Image Processing*, 27(6):3025–3035, 2018. 1
5. Matthew Brown and David G Lowe. Automatic panoramic image stitching using invariant features. *International journal of computer vision*, 74(1):59–73, 2007. 1
6. Rui Fan and Ming Liu. Road damage detection based on unsupervised disparity map segmentation. *IEEE Transactions on Intelligent Transportation Systems*, 21(11):4906–4911, 2019. 1

7. Hans P Moravec. Techniques towards automatic visual obstacle avoidance. 1977. 2

8. Chris Harris, Mike Stephens, et al. A combined corner and edge detector. In *Alvey vision conference*, volume 15, pages 10–5244. Citeseer, 1988. 2

9. David G Lowe. Object recognition from local scale-invariant features. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1150–1157. Ieee, 1999. 2, 7, 9, 11, 12

10. Jean-Michel Morel and Guoshen Yu. ASIFT: A new framework for fully affine invariant image comparison. *SIAM journal on imaging sciences*, 2(2):438–469, 2009. 2

11. Edward Rosten and Tom Drummond. Machine learning for high-speed corner detection. In *European conference on computer vision*, pages 430–443. Springer, 2006. 2

12. Michael Calonder, Vincent Lepetit, Christoph Strecha, and Pascal Fua. BRIEF: Binary robust independent elementary features. In *European conference on computer vision*, pages 778–792. Springer, 2010. 2

13. 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, 7, 9, 11, 12

14. Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. ORB-SLAM: A versatile and accurate monocular SLAM system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015. 2

15. Raul Mur-Artal and Juan D Tardós. ORB-SLAM2: An open-source SLAM system for monocular, stereo, and RGB-D cameras. *IEEE transactions on robotics*, 33(5):1255–1262, 2017. 2

16. Kwang Moo Yi, Eduard Trulls, Vincent Lepetit, and Pascal Fua. LIFT: Learned invariant feature transform. In *European conference on computer vision*, pages 467–483. Springer, 2016. 2

17. 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. 2, 3, 7, 9, 11, 12

18. Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2, 7

19. Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-Net: A trainable cnn for joint description and detection of local features. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 8092–8101, 2019. 3

20. Jerome Revaud, Philippe Weinzaepfel, César De Souza, Noe Pion, Gabriela Csurka, Yohann Cabon, and Martin Humenberger. R2D2: Repeatable and reliable detector and descriptor. *arXiv preprint arXiv:1906.06195*, 2019. 3, 7, 9, 11, 12

21. Michał Tyszkiewicz, Pascal Fua, and Eduard Trulls. DISK: Learning local features with policy gradient. *Advances in Neural Information Processing Systems*, 33:14254–14265, 2020. 3

22. Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020. 3, 7, 9, 11, 12

23. Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013. 3

24. 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. 3, 8, 9, 11, 12

25. 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. 3, 8, 9, 11, 12

26. 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. 3, 8, 9, 11, 12

27. Udit Singh Parihar, Aniket Gujarathi, Kinal Mehta, Satyajit Tourani, Sourav Garg, Michael Milford, and K Madhava Krishna. RoRD: Rotation-robust descriptors and orthographic views for local feature matching. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1593–1600. IEEE, 2021. 3, 5, 6, 7, 8, 9, 11, 12

28. Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999. PMLR, 2016. 3

29. Junying Li, Zichen Yang, Haifeng Liu, and Deng Cai. Deep rotation equivariant network. *Neurocomputing*, 290:26–33, 2018. 3

30. Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical CNNs. *arXiv preprint arXiv:1801.10130*, 2018. 3

31. Maurice Weiler and Gabriele Cesa. General E(2)-equivariant steerable CNNs. *Advances in Neural Information Processing Systems*, 32, 2019. 3

32. Yuan Liu, Zehong Shen, Zhixuan Lin, Sida Peng, Hujun Bao, and Xiaowei Zhou. GIFT: Learning transformation-invariant dense visual descriptors via group CNNs. *Advances in Neural Information Processing Systems*, 32, 2019. 3, 7, 8, 9, 11, 12

33. Jongmin Lee, Byungjin Kim, and Minsu Cho. Self-supervised equivariant learning for oriented keypoint detection. *arXiv preprint arXiv:2204.08613*, 2022. 3

34. Abhishek Peri, Kinal Mehta, Avneesh Mishra, Michael Milford, Sourav Garg, and K Madhava Krishna. ReF-rotation equivariant features for local feature matching. *arXiv preprint arXiv:2203.05206*, 2022. 3

35. Georg Bökman and Fredrik Kahl. A case for using rotation invariant features in state of the art feature matchers. *arXiv preprint arXiv:2204.10144*, 2022. 3, 8, 9, 11, 12, 14

36. Titus Cieslewski, Michael Bloesch, and Davide Scaramuzza. Matching features without descriptors: implicitly matched interest points. *arXiv preprint arXiv:1811.10681*, 2018. 4

37. Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 4

38. Georg Bökman, Fredrik Kahl, and Axel Flinth. ZZ-Net: A universal rotation equivariant architecture for 2D point clouds. *arXiv preprint arXiv:2111.15341*, 2021. 4

39. Torsten Sattler, Tobias Weyand, Bastian Leibe, and Leif Kobbelt. Image retrieval for image-based localization revisited. In *BMVC*, volume 1, page 4, 2012. 5, 6, 7

40. 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. 5

41. 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. 5, 6, 7

42. 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. 5, 6, 7

43. 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. 5, 6

44. Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. TartanAir: A dataset to push the limits of visual SLAM. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4909–4916. IEEE, 2020. 5, 6

45. 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. 5, 6, 8

46. Gary Bradski. The OpenCV library. *Dr. Dobb's Journal: Software Tools for the Professional Programmer*, 25(11):120–123, 2000. 7

47. Pablo F Alcantarilla and T Solutions. Fast explicit diffusion for accelerated features in nonlinear scale spaces. *IEEE Trans. Patt. Anal. Mach. Intell*, 34(7):1281–1298, 2011. 7, 9, 11, 12

48. Stefan Leutenegger, Margarita Chli, and Roland Y Siegwart. BRISK: Binary robust invariant scalable keypoints. In *2011 International conference on computer vision*, pages 2548–2555. Ieee, 2011. 7, 9, 11, 12

49. Pablo Fernández Alcantarilla, Adrien Bartoli, and Andrew J Davison. KAZE features. In *European conference on computer vision*, pages 214–227. Springer, 2012. 7, 9, 11, 12

50. Alexandre Alahi, Raphael Ortiz, and Pierre Vandergheynst. FREAK: Fast retina keypoint. In *2012 IEEE conference on computer vision and pattern recognition*, pages 510–517. Ieee, 2012. 7, 9, 11, 12

51. Filip Radenović, Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondřej Chum. Revisiting oxford and paris: Large-scale image retrieval benchmarking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5706–5715, 2018. 7

52. Xiaoming Zhao, Xingming Wu, Jinyu Miao, Weihai Chen, Peter CY Chen, and Zhengguo Li. ALIKE: Accurate and lightweight keypoint detection and descriptor extraction. *IEEE Transactions on Multimedia*, 2022. 7, 9, 11, 12

53. Tianwei Shen, Zixin Luo, Lei Zhou, Runze Zhang, Siyu Zhu, Tian Fang, and Long Quan. Matchable image retrieval by learning from surface reconstruction. In *Asian conference on computer vision*, pages 415–431. Springer, 2018. 7, 8

54. Noah Snavely, Steven M Seitz, and Richard Szeliski. Photo tourism: Exploring photo collections in 3D. In *ACM siggraph 2006 papers*, pages 835–846. 2006. 7