This CVPR paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version; the final published version of the proceedings is available on IEEE Xplore.

EDM: Equirectangular Projection-Oriented Dense Kernelized Feature Matching

Dongki Jung ^{1,2} Jaehoon Choi ² Yonghan Lee ² Somi Jeong ¹ Taejae Lee ¹ Dinesh Manocha ² Suyong Yeon ¹ ¹NAVER LABS ²University of Maryland



Figure 1. (a) Previous state-of-the-art [15] struggles to achieve accurate dense matching in equirectangular projection (ERP) images due to inherent distortions. (b) The ERP image can be transformed into a cubemap image, which consists of six perspective images. However, this approach demands multiple independent iterations of inference for each pair of perspective images, increasing computational complexity and losing the global information in the ERP image. (c) Our proposed method, EDM, leverages the spherical camera model, rendering it robust against distortions. **Warp** refers to results obtained by multiplying the warped image with the predicted certainty map, demonstrating that our method yields more accurate dense matches.

Abstract

We introduce the first learning-based dense matching algorithm, termed Equirectangular Projection-Oriented Dense Kernelized Feature Matching (EDM), specifically designed for omnidirectional images. Equirectangular projection (ERP) images, with their large fields of view, are particularly suited for dense matching techniques that aim to establish comprehensive correspondences across images. However, ERP images are subject to significant distortions, which we address by leveraging the spherical camera model and geodesic flow refinement in the dense matching method. To further mitigate these distortions, we propose spherical positional embeddings based on 3D Cartesian coordinates of the feature grid. Additionally, our method incorporates bidirectional transformations between spherical and Cartesian coordinate systems during refinement, utilizing a unit sphere to improve matching performance. We demonstrate that our proposed method achieves notable performance enhancements, with improvements of +26.72 and +42.62 in AUC@5° on the Matterport3D and Stanford2D3D datasets. Project Page: https://jdk9405.github.io/EDM

1. Introduction

Omnidirectional images, also known as 360° images, provide significant advantages owing to their expansive fields of view, offering more contextual information and versatility [12, 21, 38, 63, 67]. These spherical images enable a comprehensive representation of environments, facilitating a deeper understanding of spatial information. Their utility extends to aiding robot navigation [40, 61] and autonomous vehicle driving [43] by minimizing blind spots. 360° images also can be utilized in a diverse range of applications, from creating immersive AR/VR experiences to practical uses in interior design [1], tourism [48], and real estate photography [5]. Integrating omnidirectional images into virtual house tours allows customers to experience an immersive view, enabling them to fully engage themselves in the service. Moreover, the adoption of omnidirectional images contributes to more efficient data collection. By replacing the need for multiple perspective images, omnidirectional images can reduce both the cost and time associated with data scanning. The large field of view provided by 360° images has also demonstrated superiority over narrower views in 3D motion estimation [18, 27, 42].

Feature matching plays a critical role in numerous 3D computer vision tasks, including mapping and localization. Traditionally, Structure from Motion (SfM) [49] leverages feature matching to estimate relative poses. Recent advancements have introduced semi-dense or dense approaches for feature matching such as LoFTR [55] and DKM [15], which demonstrate superior performance in repetitive or textureless environments compared to keypoint-based methods [13, 28, 36, 46, 47]. These methods have been mainly developed for perspective 2D images and videos, but encounter challenges when applied to omnidirectional images. For example, to adapt matching methods for spherical images, two prevalent approaches for sphere-to-plane projections are the equirectangular projection (ERP) and the cubemap projection [63]. ERP images exhibit significant distortions, particularly near the pole regions, which hinder the effective application of perspective methods. On the other hand, the cubemap format, consisting of six perspective images, can be processed independently without such distortions. However, this approach involves the costly computation of multiple inferences for each pair of perspective images, resulting in the loss of global information from a single spherical image and diminishing feature matching capabilities due to the reduced field of view in each perspective image. These challenges are shown in Fig. 1 (a) and (b).

Main Results In this paper, we propose EDM, a distortion-aware dense feature matching method for omnidirectional images, addressing challenges that existing detector-free approaches [15, 16, 55] struggle to overcome. To the best of our knowledge, EDM is the first learningbased method designed for dense matching and relative pose estimation between two omnidirectional images. As seen in Fig. 1, our method defines feature matching in 3D coordinates, specifically addressing the challenges posed by distortions of ERP images. We accomplish this based on the integration of two novel steps: a Spherical Spatial Alignment Module (SSAM) and specific enhancements in Geodesic Flow Refinement. The SSAM leverages spherical positional embeddings for ERP images and incorporates a decoder to generate the global matches. Furthermore, the Geodesic Flow Refinement step employs coordinate transformation to refine the residuals of correspondences. Compared to both recent sparse and dense feature matching methods [15, 16, 19, 69], our approach results in significant performance improvement of +26.72 and +42.62 AUC@5° in relative pose estimation for spherical images on the Matterport3D [5] and Stanford2D3D [2] datasets. Additionally, we evaluate our method qualitatively on the EgoNeRF [7] and OmniPhotos [4] datasets, demonstrating robust performance across diverse environments. The main contributions of this paper are summarized as follows:

• We introduce a novel approach for estimating dense matching across ERP images using geodesic flow on a unit sphere.

- We propose a Spherical Spatial Alignment Module that utilizes Gaussian Process regression and spherical positional embeddings to establish 3D correspondences between omnidirectional images. In addition, we use Geodesic Flow Refinement by enabling conversions between coordinates to refine the displacement on the surface of the sphere.
- With azimuth rotation for data augmentation, we achieve state-of-the-art performance in dense matching and relative pose estimation between two omnidirectional images.

2. Related Work

Omnidirectional Images The popularity of consumerlevel 360° cameras has led to increased interest in spherical images, which offer comprehensive coverage of the field of view from a single vantage point. These images are often represented using equirectangular projection (ERP) [63], facilitating their utilization in various computer vision tasks. Recent advancements in computer vision have leveraged ERP images for diverse tasks such as object detection [11, 53], semantic segmentation [24, 66], depth estimation [25, 32, 33, 45, 50, 60, 65], omnidirectional Simultaneous Localization and Mapping [62], scene understanding [54], and neural rendering [8, 26, 29, 37].

Despite the utility of ERP images, their unique geometry presents several challenges in visual representation. As ERP images are obtained through projecting a sphere onto a plane, a single spherical image can be expressed by multiple distinct ERP images. Additionally, ensuring perfect alignment of their left and right extremities is essential. While some research methods have introduced rotation-equivariant convolutions [9, 17] to address these issues, their implementation often demands increased computational resources. To mitigate this constraint, we propose an azimuth rotation approach for data augmentation, under the assumption that maintaining the downward orientation of scanned omnidirectional images parallel to gravity offers benefits [3].

Feature Matching Local feature matching has relied on detector-based methods, encompassing both traditional hand-crafted techniques [36, 46] and learning-based approaches [13, 28, 34, 44, 59]. These methods typically involve detecting keypoints, computing descriptor distances between paired keypoints, and performing matching via mutual nearest neighbor search. SuperGlue [47] introduces a learning-based paradigm, optimizing visual descriptors using an attentional graph neural network and an optimal matching layer. However, detector-based methods face limitations in terms of accurately detecting keypoints, particularly in repetitive or indiscriminative regions. In contrast, detector-free or dense methods [15, 16, 39, 55, 57, 58] offer a solution to the keypoint detection issue, providing dense feature matches at the pixel level.

While the aforementioned methods are tailored for perspective images, they often fail to address the unique challenges of spherical cameras. SPHORB [69], an extension of ORB [46], mitigates distortion in ERP images using a geodesic grid and local planar approximation [14]. Similarly, learning-based matching methods such as SphereGlue [19, 20] and PanoPoint [68] adapt keypoint matching techniques for spherical imagery. CoVisPose [23, 41] explores layout features for estimating camera poses over large baselines yet remains constrained by detected feature information. Therefore, we propose a novel dense matching method that extracts all matches without keypoint detection in spherical images.

3. Preliminaries

3.1. Spherical and Cartesian Coordinate



Figure 2. Coordinate system.

$$\begin{cases} S^{x} = \sin(\theta)\cos(\phi) \\ S^{y} = \sin(\phi) \\ S^{z} = \cos(\theta)\cos(\phi) \end{cases} \begin{cases} \theta = \arctan(\frac{S^{x}}{S^{z}}) \\ \phi = \arctan(\frac{S^{y}}{|\mathbf{S}|}) \end{cases}$$
(1)

Although ERP images are displayed in 2D space, they actually represent a collection of flattened rays normalized to a unit scale within a spherical camera model. Thus, we can express the coordinate conversion equation $\mathbf{u} = \boldsymbol{\pi}(\mathbf{S})$ between the spherical coordinates $\mathbf{u} = (\theta, \phi)$ and the 3D Cartesian coordinates $\mathbf{S} = (S^x, S^y, S^z)$ as shown in Fig. 2. Each value of $\theta \in [-\pi, \pi]$ and $\phi \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ indicates the longitude and latitude. We utilize this coordinate transformation $\boldsymbol{\pi}(\cdot)$ in Section 4.1 and Section 4.2 to handle the spherical camera model effectively.

3.2. Dense Kernelized Feature Matching

Dense matching is the task of finding dense correspondence and estimating 3D geometry from two images (I_A, I_B) . Recently, DKM [15] introduced a kernelized global matcher and warp refinement, formulating this problem as finding a mapping $f \rightarrow \mathbf{u}$ where \mathbf{u} are 2D spatial coordinates. First, DKM extracts multi-scale features using a ResNet50 encoder [22],

$$\{f_{\mathcal{A}}^{l}\}_{l=1}^{L} = \text{Encoder}(I_{\mathcal{A}}), \quad \{f_{\mathcal{B}}^{l}\}_{l=1}^{L} = \text{Encoder}(I_{\mathcal{B}}),$$
(2)

where the strides are defined as elements of the set $l \in \{2^0, ..., 2^{L-1}\}$. Coarse features are associated with stride $\{32, 16\}$, and fine features correspond to $\{8, 4, 2, 1\}$.

At the coarse level, it consists of a kernelized regression to estimate the posterior mean $\mu_{\mathcal{A}|\mathcal{B}}$ using a Gaussian Process (GP) formulation. GP regression generates a probabilistic distribution using the feature information conditioned on frame \mathcal{B} to estimate coarse global matches. The normalized 2D feature grid $f_B^{\text{grid}} \in \mathbb{R}^{h \times w \times 2}$, where h and w denote the resolution of the feature grid, is embedded into $\chi_{\mathcal{B}}$ with an additional cosine embedding [51] to induce multimodality in GP. The embedded coordinates are processed by an exponential cosine similarity kernel K to calculate $\mu_{\mathcal{A}|\mathcal{B}}$,

$$\mu_{\mathcal{A}|\mathcal{B}} = K_{\mathcal{A}\mathcal{B}} (K_{\mathcal{B}\mathcal{B}} + \sigma_n^2 I)^{-1} \chi_{\mathcal{B}}^{\text{coarse}}, \qquad (3)$$

$$\begin{cases} K_{mn} = \exp\left(\tau\left(\frac{f_m \cdot f_n}{\sqrt{(f_m \cdot f_m)(f_n \cdot f_n) + \varepsilon}} - 1\right)\right),\\ \chi_{\mathcal{B}}^{\text{coarse}} = \cos(W f_{\mathcal{B}}^{\text{grid}} + b), \end{cases}$$
(4)

where $\tau = 5$, $\epsilon = 10^{-6}$, and the standard deviation of the measurement noise $\sigma_n = 0.1$ in the experiments. Wand b are the weights and biases of a 1×1 convolution layer. Then, CNN embedding decoder [64] yields the initial global matches $\hat{u}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}$ and confidence of matches $\hat{c}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}$ from the concatenation of the reshaped estimated posterior mean $\mu_{\mathcal{A}|\mathcal{B}}^{\text{grid}}$ and the coarse features,

$$(\hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}, \hat{c}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}) = \text{Decoder}(\mu_{\mathcal{A}|\mathcal{B}}^{\text{grid}} \oplus f_{\mathcal{A}}^{\text{coarse}}).$$
 (5)

At the fine level, the warp refiners estimate the residual displacement using the previous matches and feature information. The process is described as follows,

$$(\triangle \hat{\mathbf{u}}_{\mathcal{A} \to \mathcal{B}}^{l+1}, \ \triangle \hat{c}_{\mathcal{A} \to \mathcal{B}}^{l+1})$$

$$= \operatorname{Refiner}^{l+1} \left(f_{\mathcal{A}}^{l+1} \oplus f_{\mathcal{B} \to \mathcal{A}}^{l+1} \oplus \operatorname{Corr}_{\Omega_{k}}^{l+1} \oplus \hat{\mathbf{u}}_{\mathcal{A} \to \mathcal{B}}^{l+1} - \mathbf{u}_{\mathcal{A}}^{l+1} \right),$$

$$(6)$$

$$\begin{cases} f_{\mathcal{B}\to\mathcal{A}}^{l+1} = f_{\mathcal{B}} \langle \hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l+1} \rangle, \\ f_{\mathcal{B}\to\mathcal{A}, \ \Omega_{k}}^{l+1} = f_{\mathcal{B}} \langle \Omega_{k}, (\hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l+1}) \rangle, \\ Corr_{\Omega_{k}}^{l+1} = \sum_{\text{channel}} f_{\mathcal{A}}^{l+1} f_{\mathcal{B}\to\mathcal{A}, \ \Omega_{k}}^{l+1}, \end{cases}$$
(7)

where $\Omega_k(\mathbf{u}) = \mathbf{u} + \mathbf{p} (\|\mathbf{p}\|_{\infty} \leq k)$ is the patch sized $k, \langle \cdot \rangle$ means the bilinear interpolation function, $Corr_{\Omega_k}^{l+1}$ represents local correlation between the features, and $\mathbf{u}_{\mathcal{A}}^{l+1}$ indicates the grid in $f_{\mathcal{A}}^{l+1}$. Finally, it recursively updates the matching points and confidence by adding the residuals to



Figure 3. Overview of our approach. It consists of three steps: Multi-scale Feature Extraction, Spherical Spatial Alignment Module (Sec. 4.1), and Geodesic Flow Refinement (Sec. 4.2).

the previous information and upsampling until reaching the same resolution as the input images,

$$\hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l} = \hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l+1} + \triangle \hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l+1}, \\
\hat{c}_{\mathcal{A}\to\mathcal{B}}^{l} = \hat{c}_{\mathcal{A}\to\mathcal{B}}^{l+1} + \triangle \hat{c}_{\mathcal{A}\to\mathcal{B}}^{l+1}.$$
(8)

4. Our Proposed Method

The overall process is illustrated in Fig. 3. Following the approach outlined in Section 3.2, we first utilize ERP images I_A and I_B as input and extract multi-scale features f_A and f_B . Different from [15], we reformulate the problem as finding a mapping $f \rightarrow \mathbf{S}$ using 3D Cartesian coordinates. We introduce the Spherical Spatial Alignment Module, a global matcher utilizing a spherical camera system to compensate for distortions caused by sphere-to-plane projection in ERP images. We then formalize the geodesic flow on a unit sphere and establish projections between equirectangular and spherical spaces to refine matches. In addition, to enhance the robust accuracy of our method, we leverage randomized azimuth rotation during the training process.

4.1. Spherical Spatial Alignment Module

Our Spherical Spatial Alignment Module (SSAM) conducts global matching at a coarse level through Gaussian Process (GP) regression, depicted in Fig. 4. GP predicts the posterior mean $\mu_{\mathcal{A}|\mathcal{B}}$ from the embeddings as in Eq. 3. Due to the pronounced distortions in the polar regions of ERP images, spherical positional embedding/encoding is frequently employed to mitigate this challenge [6, 30, 31]. Here, we explicitly apply positional embeddings with 3D Cartesian coordinates, derived from the 2D spherical feature grid and the inverse transformation function $\pi^{-1}(\cdot)$,

$$\chi_{\mathcal{B}}^{\text{coarse}} = \cos(W\pi^{-1}(f_{\mathcal{B}}^{\text{grid}}) + b).$$
(9)

Our proposed positional embedding facilitates the utilization of embedded coordinates $\chi_{\mathcal{B}}^{\text{coarse}}$ to promote distortion awareness within the ERP images. Additionally, this embedding ensures structural consistency along the boundaries



Figure 4. Our Spherical Spatial Alignment Module. We present Spherical Positional Embedding (red dotted box). The embedding decoder generates the global matches $\hat{\mathbf{S}}_{A\to B}^{\text{coarse}}$. Here, the gray curved lines represent the geodesic flow between \mathbf{S}_A and \mathbf{S}_B . \oplus denotes concatenation, \otimes means reshape and matrix multiplication. We provide the matrix dimensions of intermediate features for reference.



Figure 5. Our proposed Geodesic Flow Refinement. Refining the displacement along curved lines on the spherical surface presents significant challenges. To address this, we project the displacement into the ERP space for refinement (Cartesian to spherical) and subsequently unproject it back onto the spherical surface for further refinement (spherical to Cartesian).

of ERP images by leveraging relative spatial information within the 3D Cartesian grid. The outputs of the subsequent embedding decoder provide the initial global matches $\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}$ on the unit sphere and the ERP certainty map $\hat{c}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}$.

$$\left(\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}, \hat{c}_{\mathcal{A}\to\mathcal{B}}^{\text{coarse}}\right) = \text{Decoder}(\mu_{\mathcal{A}|\mathcal{B}} \oplus f_{\mathcal{A}}^{\text{coarse}}).$$
(10)

4.2. Geodesic Flow Refinement

In our SSAM approach, as the geodesic flow must reside on the unit sphere, directly defining warp refinement on the surface of the sphere makes it impossible to update the residuals linearly. Thus, we circumvent this problem by enabling a conversion between the 3D Cartesian coordinates and the 2D equirectangular space, as illustrated in Fig. 5,

$$\hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l+1} = \boldsymbol{\pi}(\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{l+1}).$$
(11)

After following all the processes outlined in Eq. 6 for refinement, we update the residuals as described in Eq. 8.



Figure 6. Maintaining consistent geometry, ERP can produce multiple visual representations based on θ^{aug} .

As this refinement stage iterates repeatedly, the predicted $\hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l}$ is back-projected into 3D Cartesian coordinates,

$$\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{l} = \boldsymbol{\pi}^{-1}(\hat{\mathbf{u}}_{\mathcal{A}\to\mathcal{B}}^{l}).$$
(12)

4.3. Augmentation

A single omnidirectional image can be transformed into multiple distinct ERP images, as shown in Fig. 6. This transformation is feasible by capturing the full spectrum of rays and ensuring a seamless representation in the spherical input image, which facilitates the generation of diverse ERP images while maintaining consistent geometric properties in the world space. Consequently, we define a horizontal rotation matrix $T_{\mathcal{A}}^{\text{aug}}$ with a randomly selected azimuth angle $\theta_{\mathcal{A}}^{\text{aug}} \in [0, 2\pi]$ during training. Based on $T_{\mathcal{A}}^{\text{aug}}$, we rotate and redefine the ERP image $I_{\mathcal{A}}$, the depth map $D_{\mathcal{A}}$, and the pose $T_{\mathcal{A}}$. Notably, this transformation adjusts $T_{\mathcal{A}}$ and $D_{\mathcal{A}}$ together, ensuring consistent geometry in the world space. The same process is applied to the counterpart frame \mathcal{B} .

4.4. Loss

Utilizing dense ground truth depth maps and aligned camera poses, we can derive ERP depth $D_{\mathcal{A}\to\mathcal{B}}$ and matches $S_{\mathcal{A}\to\mathcal{B}}$ during the warping process from frame \mathcal{A} to \mathcal{B} within the spherical coordinate system. We adopt the certainty estimation method proposed by Edstedt et al. [15], which involves finding consistent matches using relative depth consistency between frames \mathcal{A} and \mathcal{B} ,

$$c_{\mathcal{A}\to\mathcal{B}} = \left| \frac{D_{\mathcal{A}\to\mathcal{B}} - D_{\mathcal{B}}}{D_{\mathcal{B}}} \right| < \alpha, \tag{13}$$

where α is 0.05. The binary mask $c_{\mathcal{A}\to\mathcal{B}}$ represents the ground truth certainty map. Diverging from the approach outlined in Edstedt et al. [15], our method constrains the predicted matches $\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{l}$, composed of 3D Cartesian coordinates, to reside on the surface of the unit sphere. This implies that the predicted matches can be interpreted as the ray directions of the spherical camera. Instead of defining the loss function based on the Euclidean distance between the predicted matches $\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{l}$ and the ground truth matches $\mathbf{S}_{\mathcal{A}\to\mathcal{B}}^{l}$, we use the angular difference between the ray directions. Consequently, this approach ensures that $\hat{\mathbf{S}}_{\mathcal{A}\to\mathcal{B}}^{l}$

is optimized along the surface of the unit sphere. We define our regression loss L_r^l using cosine similarity to measure the angular difference. For the certainty loss L_c^l , we employ the binary cross-entropy function, as utilized in Edstedt et al. [15],

$$\mathcal{L}_{\mathbf{r}}^{l} = \sum_{\text{grid}} c_{\mathcal{A} \to \mathcal{B}}^{l} \odot \left(1 - \frac{\|\mathbf{S}_{\mathcal{A} \to \mathcal{B}}^{l} \cdot \mathbf{\hat{S}}_{\mathcal{A} \to \mathcal{B}}^{l}\|}{\|\mathbf{S}_{\mathcal{A} \to \mathcal{B}}^{l}\|\|\mathbf{\hat{S}}_{\mathcal{A} \to \mathcal{B}}^{l}\|\|\mathbf{\hat{S}}_{\mathcal{A} \to \mathcal{B}}^{l}\|}\right), \quad (14)$$

$$\mathcal{L}_{\mathbf{c}}^{l} = \sum_{\text{grid}} c_{\mathcal{A} \to \mathcal{B}}^{l} \log \hat{c}_{\mathcal{A} \to \mathcal{B}}^{l} + (1 - c_{\mathcal{A} \to \mathcal{B}}^{l}) \log(1 - \hat{c}_{\mathcal{A} \to \mathcal{B}}^{l}). \quad (15)$$

The total loss function comprises a weighted sum of the regression loss and the certainty loss, as detailed in Edstedt et al. [15], Melekhov et al. [39], Tan et al. [56], Zhou et al. [70], with λ set at 0.01,

$$L_{\text{total}} = \sum_{l=1}^{L} L_{\text{r}}^{l} + \lambda L_{\text{c}}^{l}.$$
 (16)

5. Experiments 5.1. Experiments Settings

Matterport3D Dataset Training our method requires ERP input images, ground truth depth maps, and aligned poses. The Matterport3D dataset [5] encompasses 90 indoor scenes represented by 10,800 panoramas reconstructed as textured meshes. However, the dataset lacks pose and depth information for *skybox* images, which are essential for creating ERP images. Previous works have addressed this limitation by rendering both images and depth maps from the textured mesh [71] or by employing 360° SfM to estimate poses [45]. In our approach, we generate the poses for *skybox* images directly from the originally proposed camera poses in Matterport3D. Through experimentation, we found that treating the 12th camera pose, out of the 18 viewpoints (comprising 6 rotations and 3 tilt angles) in each panorama, identically to the second skybox image did not result in any issues. We define the remaining poses for the *skybox* images by rotating 90° in each direction from the second pose. We adhere to the official benchmark split, utilizing 61 scenes for training, 11 for validation, and 18 for testing. For two-view pose estimation, it is necessary to create pairs of overlapped images. We achieve this by transforming ERP depth maps between frames within the spherical coordinate system. Pixels where the depth difference is below a specified threshold, e.g. 0.1, are classified as inliers. Subsequently, we compare the ratio of these inliers to the total number of pixels. We organize both the training and testing datasets based on the overlap ratio of image pairs and the benchmark split. Specifically, images with the overlap ratio exceeding 30% are distributed into respective training and testing splits. As a result, the training set contains 44,700 pairs, while the test set comprises 4,575 pairs. We resize the resolution of ERP images and depth maps to 640×320 .

Stanford2D3D Dataset Stanford2D3D [2] consists of data scanned from six large-scale indoor spaces collected from three distinct buildings. This dataset contains a relatively small number of 1,413 panorama images and, therefore, is utilized exclusively for testing purposes. We assess the overlap ratio between frames and include them in the test split if their ratio exceeds 50%. A total of 3,460 pairs are incorporated into the test set. During testing, we resize the resolution to 640×320 .

EgoNeRF and OmniPhotos Dataset EgoNeRF [7] introduces 11 synthetic scenes created with Blender [10] and 11 real scenes captured with a RICOH THETA V camera. OmniPhotos [4] provides a dataset captured with an Insta360 ONE X camera. Both datasets contain egocentric scenes captured with a casually rotating camera stick. Consequently, their rotation axes, pole regions, or camera height change, resulting in different distortions compared to Matterport3D or Stanford2D3D. We present additional qualitative results from these datasets to validate our method.

Implementation Details We employ the AdamW [35] optimizer with a weight-decay factor of 10^{-2} , a learning rate of $5 \cdot 10^{-6}$ for multiscale feature extractor, and 10^{-4} for the SSAM and the Geodesic Flow Refiner. EDM is trained for 300,000 steps with a batch size of 4 in a single RTX 3090 GPU, which takes approximately two days to complete. During evaluation, the balanced sampling approach using kernel density estimation [15] tends to establish correspondences primarily in concentrated areas with high probability distributions, making it unsuitable for omnidirectional images. Thus, we randomly sample up to 5,000 matches after certainty filtering with a threshold of 0.8 to ensure correspondences cover the entire area.

5.2. Experimental Results

We compare our proposed method EDM with four different methods: 1) SPHORB [69] is a hand-crafted keypointbased feature matching algorithm. 2) SphereGlue [19] is a learning-based keypoint matching method. Both SPHORB [69] and SphereGlue [19] are specifically designed for spherical images. 3) DKM [15] and 4) RoMa [16] are state-of-the-art dense matching algorithms for perspective images. To estimate the essential matrix and the relative pose for spherical cameras, Solarte et al. [52] proposed a normalization strategy and non-linear optimization within the classic 8-point algorithm. We adopt this for two-view pose estimation in all quantitative comparisons.

Table 1 shows the quantitative results of the pose estimation in Matterport3D. Despite SPHORB and SphereGlue

Method	Image	Feature	@5°	AUC @10°	@20°
SPHORB [69]	ERP	sparse	0.38	1.41	3.99
SphereGlue [19]	ERP	sparse	11.29	19.95	31.10
DKM [15]	persepctive perspective	dense	18.43	28.50	38.44
RoMa [16]		dense	12.45	22.37	34.24
EDM (ours)	ERP	dense	45.15	60.99	73.60

Table 1. Quantitative comparison on Matterport3D with recent algorithms. EDM improves AUC@5° by 26.72.

Method	Image	Feature	@5°	AUC @10°	@20°
SPHORB [69]	ERP	sparse	0.14	1.01	4.08
SphereGlue [19]	ERP	sparse	11.25	22.41	36.57
DKM [15]	perspective perspective	dense	12.46	22.18	34.13
RoMa [16]		dense	11.48	22.52	37.07
EDM (ours)	ERP	dense	55.08	71.65	82.72

Table 2. Quantitative comparison on Stanford2D3D with recent algorithms. EDM improve AUC@5° by 42.62.

being designed for the ERP images, the presence of textureless or repetitive regions, which are common in indoor environments of Matterport3D, leads to performance degradation in the keypoint-based methods. SPHORB fails to estimate the essential matrix correctly due to the limited number of matching points. EDM demonstrates significantly higher performance than all the other methods.

Figure 7 illustrates the qualitative results in Matterport3D. The previous methods designed for perspective images, such as DKM and RoMa, exhibit good matching ability but encounter challenges when confronted with the distortions of ERP. While SphereGlue and SPHORB perform well in discriminative regions, their performance deteriorates as the overlap ratio decreases, resulting in numerous false positive matches. In contrast, EDM can estimate dense correspondences regardless of occlusion and textureless areas. Due to the similarity in results between DKM and RoMa, we have only included the former to maintain a concise visualization. Experimental results in Fig.8 depict the relationship between image overlap ratio and AUC@20° performance. As expected, a decrease in the overlap ratio leads to severe performance degradation in the previous works. On the other hand, our proposed method demonstrates robustness in more challenging scenes, maintaining similar performance levels until the overlap decreases to 60%, compared to other methods.

For a fair comparison, we use another benchmark dataset, Stanford2D3D. We validate EDM using a model trained on Matterport3D without additional training on Stanford2D3D. In Table 2, EDM outperforms the previous



Figure 7. Qualitative results on Matterport3D. (a) The blue lines represent the results of matching points from SPHORB [69]; the green lines correspond to SphereGlue [19]. Both (b) DKM [15] and (c) EDM depict the outcomes of multiplying the warped image with the certainty map. EDM can estimate dense and accurate matches even in the presence of distortions and severe occlusions. The numbers beside the images represent the overlap ratio, reflecting the difficulty of matching. Smaller numbers indicate more challenging scenes.



Figure 8. Performance with respect to the overlap ratio. This highlights the robustness of EDM in scenarios with varying levels of overlap, particularly in challenging conditions where the overlap ratio is limited.

works by a significant margin, especially in scenes with severe occlusion. The certainty map demonstrates EDM's robustness, particularly in handling occluded scenes. Additionally, although the panorama images in Stanford2D3D



Figure 9. Qualitative results on Stanford2D3D. The blue and green lines correspond to SPHORB and SphereGlue.

contain missing regions in the upper and lower parts of the sphere, the proposed spherical positional embedding enables the network to predict matching correspondences accurately, as shown in Fig. 9.

5.3. Additional Qualitative Results

To demonstrate the robust performance of our method across diverse environments, we qualitatively validate EDM using additional datasets such as EgoNeRF and OmniPhotos. As it is primarily trained on indoor environments [5] where the camera is oriented parallel to gravity, severely slanted image pairs of rotational scenes or outdoor envi-



Figure 10. Qualitative results on EgoNeRF [7] and OmniPhotos [4]. Despite being primarily trained on indoor scenes, EDM effectively estimates dense matching on these datasets, demonstrating its generalization capability across diverse environments.

ronments may cause EDM to fail in accurately estimating correspondences. However, despite these differences in settings, EDM demonstrates the ability to conduct dense feature matching robustly, as shown in Fig. 10.

Furthermore, we demonstrate the applicability of our method to various omnidirectional downstream tasks. As shown in Fig. 11, our approach successfully performs triangulation from pairs of omnidirectional images. By leveraging EDM's capability to predict dense correspondences, the triangulated points yield a dense 3D reconstruction. For a more comprehensive discussion, please refer to the supplementary materials.

5.4. Ablation Study

DKM's dependence on the pinhole camera model makes it inherently unsuitable for learning with ERP images. To ensure the fair comparison, we modified the warping process in the loss function of DKM to support spherical cameras, resulting in DKM*. As shown in Table 3, this demonstrates the structural effectiveness of our proposed bidirectional coordinate transformation. The proposed positional embeddings result in improvements based on the coordinate system of the spherical camera model. We observe that utilizing a 3D grid input of Cartesian coordinates yields better performance than 2D spherical ones. Additionally, in our method, positional embedding with a linear layer slightly outperforms spherical positional encoding with sinusoidal [31]. Table 3 also confirms the advantage of our rotational augmentation. Through this augmentation technique, we can effectively address the challenge of a limited number of datasets for omnidirectional images in dense matching tasks.

6. Conclusion, Limitations, and Future Work

In this paper, we present, for the first time, a novel dense feature matching method tailored for omnidirectional im-

Method	Positional Embedding	Bidirectional Transformation	Rotational Augmentation	@5°	AUC @10°	@20°
DKM*	2D linear	-	-	19.83	33.06	46.24
Ours	2D linear	\checkmark	-	29.67	45.90	60.82
Ours	2D linear	\checkmark	\checkmark	35.03	51.14	65.07
Ours	3D linear	\checkmark	-	34.64	50.82	65.16
Ours	3D linear	\checkmark	\checkmark	45.15	60.99	73.60
Ours	3D sinusoidal	\checkmark	\checkmark	42.39	58.27	70.98

Table 3. Ablation study for the proposed method. DKM* indicates the DKM model trained on Matterport3D with a modified loss function for ERP images. Compared to DKM*, our method enhances performance through the proposed spherical positional embedding in SSAM, bidirectional transformation via Geodesic Flow Refinement, and rotational augmentation.



Figure 11. Triangulation results on Matterport3D and Stanford2D3D. These point clouds are generated through spherical triangulation using the estimated poses between two omnidirectional images. Our method can reconstruct dense point clouds in textureless regions, which are particularly challenging in indoor environments.

ages. Leveraging the foundational principles of DKM, we integrate the inherent characteristics of the spherical camera model into our dense matching process using geodesic flow fields. This integration instills distortion awareness within the network, thereby enhancing its performance specifically for ERP images. However, it is important to note that our method is predominantly trained on indoor datasets where the camera is vertically oriented, rendering it somewhat vulnerable to extreme rotations or outdoor environments. To address this limitation, future endeavors will focus on diversifying the training data and data augmentation to encompass a wider range of environments, fortifying the robustness of our network. Furthermore, we aim to extend our method into downstream tasks, particularly for visual localization and mapping applications for omnidirectional images.

References

- Friska Amalia and Ahmad Fitriyansah. Case study of 360 image viewer software utilization in interior design presentation to improve product immersion. In *ICCED*. IEEE, 2023.
- [2] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. arXiv preprint arXiv:1702.01105, 2017. 2, 6
- [3] Matheus A Bergmann, Paulo GL Pinto, Thiago LT da Silveira, and Cláudio R Jung. Gravity alignment for single panorama depth inference. In *SIBGRAPI*. IEEE, 2021. 2
- [4] Tobias Bertel, Mingze Yuan, Reuben Lindroos, and Christian Richardt. Omniphotos: casual 360 vr photography. ACM Transactions on Graphics (TOG), 39(6):1–12, 2020. 2, 6, 8
- [5] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. arXiv preprint arXiv:1709.06158, 2017. 1, 2, 5, 7
- [6] Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. Text2light: Zero-shot text-driven hdr panorama generation. ACM Transactions on Graphics (TOG), 41(6):1–16, 2022. 4
- [7] Changwoon Choi, Sang Min Kim, and Young Min Kim. Balanced spherical grid for egocentric view synthesis. In *CVPR*, 2023. 2, 6, 8
- [8] Dongyoung Choi, Hyeonjoong Jang, and Min H Kim. Omnilocalrf: Omnidirectional local radiance fields from dynamic videos. arXiv preprint arXiv:2404.00676, 2024. 2
- [9] Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. arXiv preprint arXiv:1801.10130, 2018. 2
- [10] Blender Online Community. Blender a 3D modelling and rendering package. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 6
- [11] Benjamin Coors, Alexandru Paul Condurache, and Andreas Geiger. Spherenet: Learning spherical representations for detection and classification in omnidirectional images. In ECCV, 2018. 2
- [12] Thiago LT da Silveira, Paulo GL Pinto, Jeffri Murrugarra-Llerena, and Cláudio R Jung. 3d scene geometry estimation from 360 imagery: A survey. ACM Computing Surveys, 55 (4):1–39, 2022. 1
- [13] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *CVPR Workshops*, 2018. 2
- [14] Marc Eder, Mykhailo Shvets, John Lim, and Jan-Michael Frahm. Tangent images for mitigating spherical distortion. In CVPR, 2020. 3
- [15] Johan Edstedt, Ioannis Athanasiadis, Mårten Wadenbäck, and Michael Felsberg. DKM: Dense kernelized feature matching for geometry estimation. In *CVPR*, 2023. 1, 2, 3, 4, 5, 6, 7
- [16] Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. Roma: Revisiting robust losses for dense feature matching. arXiv preprint arXiv:2305.15404, 2023. 2, 6

- [17] Carlos Esteves, Christine Allen-Blanchette, Ameesh Makadia, and Kostas Daniilidis. Learning so (3) equivariant representations with spherical cnns. In *ECCV*, 2018. 2
- [18] Cornelia Fermüller and Yiannis Aloimonos. Geometry of eye design: Biology and technology. In Multi-Image Analysis: 10th International Workshop on Theoretical Foundations of Computer Vision Dagstuhl Castle, Germany, March 12–17, 2000 Revised Papers, pages 22–38. Springer, 2001. 1
- [19] Christiano Gava, Vishal Mukunda, Tewodros Habtegebrial, Federico Raue, Sebastian Palacio, and Andreas Dengel. Sphereglue: Learning keypoint matching on high resolution spherical images. In CVPR Workshops, 2023. 2, 3, 6, 7
- [20] Christiano Gava, Yunmin Cho, Federico Raue, Sebastian Palacio, Alain Pagani, and Andreas Dengel. Spherecraft: A dataset for spherical keypoint detection, matching and camera pose estimation. In WACV, 2024. 3
- [21] Julia Guerrero-Viu, Clara Fernandez-Labrador, Cédric Demonceaux, and Jose J Guerrero. What's in my room? object recognition on indoor panoramic images. In *ICRA*. IEEE, 2020. 1
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016. 3
- [23] Will Hutchcroft, Yuguang Li, Ivaylo Boyadzhiev, Zhiqiang Wan, Haiyan Wang, and Sing Bing Kang. Covispose: Covisibility pose transformer for wide-baseline relative pose estimation in 360° indoor panoramas. In ECCV. Springer, 2022. 3
- [24] Chiyu Jiang, Jingwei Huang, Karthik Kashinath, Philip Marcus, Matthias Niessner, et al. Spherical cnns on unstructured grids. arXiv preprint arXiv:1901.02039, 2019. 2
- [25] Hualie Jiang, Zhe Sheng, Siyu Zhu, Zilong Dong, and Rui Huang. Unifuse: Unidirectional fusion for 360 panorama depth estimation. *IEEE Robotics and Automation Letters*, 6 (2):1519–1526, 2021. 2
- [26] Hakyeong Kim, Andreas Meuleman, Hyeonjoong Jang, James Tompkin, and Min H Kim. Omnisdf: Scene reconstruction using omnidirectional signed distance functions and adaptive binoctrees. arXiv preprint arXiv:2404.00678, 2024. 2
- [27] Jong Weon Lee, Suya You, and Ulrich Neumann. Large motion estimation for omnidirectional vision. In Proceedings IEEE Workshop on Omnidirectional Vision (Cat. No. PR00704), pages 161–168. IEEE, 2000. 1
- [28] Kunhong Li, Longguang Wang, Li Liu, Qing Ran, Kai Xu, and Yulan Guo. Decoupling makes weakly supervised local feature better. In *CVPR*, 2022. 2
- [29] Longwei Li, Huajian Huang, Sai-Kit Yeung, and Hui Cheng. Omnigs: Omnidirectional gaussian splatting for fast radiance field reconstruction using omnidirectional images. arXiv preprint arXiv:2404.03202, 2024. 2
- [30] Meng Li, Senbo Wang, Weihao Yuan, Weichao Shen, Zhe Sheng, and Zilong Dong. S2Net: Accurate panorama depth estimation on spherical surface. *IEEE Robotics and Automation Letters*, 8(2):1053–1060, 2023. 4
- [31] Xiang Li, Haoyuan Cao, Shijie Zhao, Junlin Li, Li Zhang, and Bhiksha Raj. Panoramic video salient object detection with ambisonic audio guidance. In *AAAI*, 2023. 4, 8

- [32] Yuyan Li, Zhixin Yan, Ye Duan, and Liu Ren. Panodepth: A two-stage approach for monocular omnidirectional depth estimation. In *3DV*. IEEE, 2021. 2
- [33] Yuyan Li, Yuliang Guo, Zhixin Yan, Xinyu Huang, Ye Duan, and Liu Ren. Omnifusion: 360 monocular depth estimation via geometry-aware fusion. In CVPR, 2022. 2
- [34] 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. 2
- [35] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [36] David G Lowe. Distinctive image features from scaleinvariant keypoints. *International journal of computer vision*, 60:91–110, 2004. 2
- [37] Yikun Ma, Dandan Zhan, and Zhi Jin. Fastscene: Text-driven fast 3d indoor scene generation via panoramic gaussian splatting. *arXiv preprint arXiv:2405.05768*, 2024. 2
- [38] Kevin Matzen, Michael F Cohen, Bryce Evans, Johannes Kopf, and Richard Szeliski. Low-cost 360 stereo photography and video capture. ACM Transactions on Graphics (TOG), 36(4):1–12, 2017. 1
- [39] Iaroslav Melekhov, Aleksei Tiulpin, Torsten Sattler, Marc Pollefeys, Esa Rahtu, and Juho Kannala. Dgc-net: Dense geometric correspondence network. In WACV. IEEE, 2019. 2, 5
- [40] Emanuele Menegatti, Takeshi Maeda, and Hiroshi Ishiguro. Image-based memory for robot navigation using properties of omnidirectional images. *Robotics and Autonomous Systems*, 47(4):251–267, 2004. 1
- [41] Negar Nejatishahidin, Will Hutchcroft, Manjunath Narayana, Ivaylo Boyadzhiev, Yuguang Li, Naji Khosravan, Jana Košecká, and Sing Bing Kang. Graph-covis: Gnn-based multi-view panorama global pose estimation. In *CVPR*, 2023. 3
- [42] Randal C Nelson and John Aloimonos. Finding motion parameters from spherical motion fields (or the advantages of having eyes in the back of your head). *Biological cybernetics*, 58(4):261–273, 1988. 1
- [43] Gaurav Pandey, James R McBride, and Ryan M Eustice. Ford campus vision and lidar data set. *The International Journal of Robotics Research*, 30(13):1543–1552, 2011.
- [44] Jerome Revaud, Cesar De Souza, Martin Humenberger, and Philippe Weinzaepfel. R2d2: Reliable and repeatable detector and descriptor. Advances in neural information processing systems, 32, 2019. 2
- [45] Manuel Rey-Area, Mingze Yuan, and Christian Richardt. 360monodepth: High-resolution 360deg monocular depth estimation. In *CVPR*, 2022. 2, 5
- [46] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *ICCV*. Ieee, 2011. 2, 3
- [47] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In CVPR, pages 4938– 4947, 2020. 2

- [48] Olivier Saurer, Friedrich Fraundorfer, and Marc Pollefeys. Omnitour: Semi-automatic generation of interactive virtual tours from omnidirectional video. In *3DPVT*, 2010. 1
- [49] Johannes L Schonberger and Jan-Michael Frahm. Structurefrom-motion revisited. In CVPR, 2016. 1
- [50] Zhijie Shen, Chunyu Lin, Kang Liao, Lang Nie, Zishuo Zheng, and Yao Zhao. Panoformer: Panorama transformer for indoor 360° depth estimation. In *ECCV*. Springer, 2022.
 2
- [51] Herman P Snippe and Jan J Koenderink. Discrimination thresholds for channel-coded systems. *Biological cybernetics*, 66(6):543–551, 1992. 3
- [52] Bolivar Solarte, Chin-Hsuan Wu, Kuan-Wei Lu, Yi-Hsuan Tsai, Wei-Chen Chiu, and Min Sun. Robust 360-8pa: Redesigning the normalized 8-point algorithm for 360-fov images. In *ICRA*. IEEE, 2021. 6
- [53] Yu-Chuan Su and Kristen Grauman. Learning spherical convolution for fast features from 360 imagery. Advances in neural information processing systems, 30, 2017. 2
- [54] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Hohonet: 360 indoor holistic understanding with latent horizontal features. In *CVPR*, 2021. 2
- [55] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In CVPR, 2021. 2
- [56] Dongli Tan, Jiang-Jiang Liu, Xingyu Chen, Chao Chen, Ruixin Zhang, Yunhang Shen, Shouhong Ding, and Rongrong Ji. Eco-tr: Efficient correspondences finding via coarse-to-fine refinement. In ECCV. Springer, 2022. 5
- [57] Prune Truong, Martin Danelljan, and Radu Timofte. Glunet: Global-local universal network for dense flow and correspondences. In CVPR, 2020. 2
- [58] Prune Truong, Martin Danelljan, Luc Van Gool, and Radu Timofte. Learning accurate dense correspondences and when to trust them. In *CVPR*, 2021. 2
- [59] Michał Tyszkiewicz, Pascal Fua, and Eduard Trulls. Disk: Learning local features with policy gradient. Advances in Neural Information Processing Systems, 33:14254–14265, 2020. 2
- [60] Fu-En Wang, Yu-Hsuan Yeh, Min Sun, Wei-Chen Chiu, and Yi-Hsuan Tsai. Bifuse: Monocular 360 depth estimation via bi-projection fusion. In *CVPR*, 2020. 2
- [61] Niall Winters, José Gaspar, Gerard Lacey, and José Santos-Victor. Omni-directional vision for robot navigation. In Proceedings IEEE Workshop on Omnidirectional Vision (Cat. No. PR00704), pages 21–28. IEEE, 2000. 1
- [62] Changhee Won, Hochang Seok, Zhaopeng Cui, Marc Pollefeys, and Jongwoo Lim. Omnislam: Omnidirectional localization and dense mapping for wide-baseline multi-camera systems. In *ICRA*. IEEE, 2020. 2
- [63] Mai Xu, Chen Li, Shanyi Zhang, and Patrick Le Callet. State-of-the-art in 360 video/image processing: Perception, assessment and compression. *IEEE Journal of Selected Topics in Signal Processing*, 14(1):5–26, 2020. 1, 2
- [64] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Learning a discriminative feature network for semantic segmentation. In *CVPR*, 2018. 3

- [65] Ilwi Yun, Hyuk-Jae Lee, and Chae Eun Rhee. Improving 360 monocular depth estimation via non-local dense prediction transformer and joint supervised and self-supervised learning. In AAAI, 2022. 2
- [66] Chao Zhang, Stephan Liwicki, William Smith, and Roberto Cipolla. Orientation-aware semantic segmentation on icosahedron spheres. In *ICCV*, 2019. 2
- [67] Fanglue Zhang, Junhong Zhao, Yun Zhang, and Stefanie Zollmann. A survey on 360° images and videos in mixed reality: Algorithms and applications. *Journal of Computer Science and Technology*, 38(3):473–491, 2023. 1
- [68] Hengzhi Zhang, Hong Yi, Haijing Jia, Wei Wang, and Makoto Odamaki. Panopoint: Self-supervised feature points detection and description for 360deg panorama. In CVPR Workshops, 2023. 3
- [69] Qiang Zhao, Wei Feng, Liang Wan, and Jiawan Zhang. Sphorb: A fast and robust binary feature on the sphere. *International journal of computer vision*, 113:143–159, 2015. 2, 3, 6, 7
- [70] Qunjie Zhou, Torsten Sattler, and Laura Leal-Taixe. Patch2pix: Epipolar-guided pixel-level correspondences. In *CVPR*, 2021. 5
- [71] Nikolaos Zioulis, Antonis Karakottas, Dimitrios Zarpalas, and Petros Daras. Omnidepth: Dense depth estimation for indoors spherical panoramas. In ECCV, 2018. 5