

EDFFDNet: Towards Accurate and Efficient Unsupervised Multi-Grid Image Registration

Haokai Zhu^{1,2,*} Bo Qu^{2,*} Si-Yuan Cao^{1,3,†} Runmin Zhang²
Shujie Chen⁴ Bailin Yang⁴ Hui-Liang Shen²

¹Ningbo Global Innovation Center, Zhejiang University

²College of Information Science and Electronic Engineering, Zhejiang University

³NingboTech University ⁴Zhejiang Key Laboratory of Big Data and Future E-Commerce Technology, Hangzhou, China

hkzhu.zju@gmail.com {22431157, cao-siyuan, runmin_zhang}@zju.edu.cn

{chenshujie, ybl}@zjgsu.edu.cn shenhl@zju.edu.cn

Abstract

Previous deep image registration methods that employ single homography, multi-grid homography, or thin-plate spline often struggle with real scenes containing depth disparities due to their inherent limitations. To address this, we propose an Exponential-Decay Free-Form Deformation Network (EDFFDNet), which employs free-form deformation with an exponential-decay basis function. This design achieves higher efficiency and performs well in scenes with depth disparities, benefiting from its inherent locality. We also introduce an Adaptive Sparse Motion Aggregator (ASMA), which replaces the MLP motion aggregator used in previous methods. By transforming dense interactions into sparse ones, ASMA reduces parameters and improves accuracy. Additionally, we propose a progressive correlation refinement strategy that leverages global-local correlation patterns for coarse-to-fine motion estimation, further enhancing efficiency and accuracy. Experiments demonstrate that EDFFDNet reduces parameters, memory, and total runtime by 70.5%, 32.6%, and 33.7%, respectively, while achieving a 0.5 dB PSNR gain over the state-of-the-art method. With an additional local refinement stage, EDFFDNet-2 further improves PSNR by 1.06 dB while maintaining lower computational costs. Our method also demonstrates strong generalization ability across datasets, outperforming previous deep learning methods.

1. Introduction

Image registration is a fundamental task in computer vision, establishing spatial correspondence between images captured under varying conditions. It has been widely ap-

* Equal Contributions. † Corresponding author.

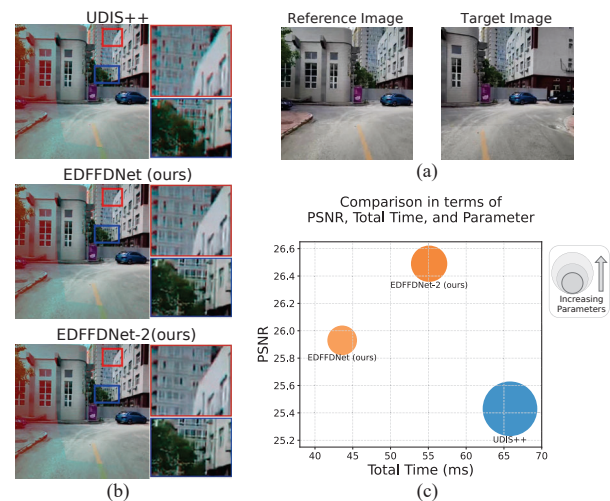


Figure 1. Performance comparison on the UDIS-D dataset [30]. (a) The reference and target images from a sample in the UDIS-D dataset. (b) Warp results of UDIS++ [31], our EDFFDNet, and EDFFDNet-2 on the sample. We combine the green and blue channels of the reference image with the red channel of the warped target image for visualization. Our method achieves better alignment in challenging regions with depth disparities (zoomed-in areas). (c) Performance comparison plot. Our method exhibits lower computational overhead while maintaining competitive performance.

plied in various domains, including image/video stitching [14, 32, 40], video stabilization [26, 43], camera calibration [42], HDR imaging [13], and SLAM [46]. Traditional methods typically rely on the extraction of geometric features [22, 27, 33], followed by robust estimation strategies [2, 3, 11] with outlier rejection to achieve registration. Although these methods are effective in various scenarios, they often exhibit limited robustness in low-texture scenes.

Deep learning-based methods are then proposed. Pioneer work [9] proposes a VGG-style network to estimate

homography. Subsequent studies refine the network architecture [10, 20, 44] or integrate [7] the IC-LK [1] to enhance performance. Cao *et al.* [5] introduce an end-to-end homography estimation network, which is further improved through attention mechanism [6] and multi-scale correlation [45]. However, these methods are trained in a supervised manner using synthetic datasets, which limits the generalization to real-world images. Nguyen *et al.* [28] first propose an unsupervised approach. Given that real-world images often violate the planar assumption of homography, some studies utilize predicted outlier masks [20, 41] or use coplanarity constraint [16] to constrain the estimation. Despite these advancements, single homography still struggles to align scenes with significant depth disparities due to its inherent limitations.

To address this challenge, more adaptive warping methods are proposed. For traditional methods, Gao *et al.* [12] propose to estimate dual-homography for two dominant planes. Zaragoza *et al.* [40] apply global projective warp with local deviations. Lee *et al.* [21] further improves this by applying weighted homographies. Li *et al.* [23] employ the thin plate spline (TPS) model [4] to realize flexible alignment. Learning-based methods have focused on integrating these warping models into deep networks. To stabilize the training of multi-grid homography networks, Wang *et al.* [38] introduce inter-grid consistency and intra-grid regularity terms. Nie *et al.* [29] propose an efficient contextual correlation layer. Since multi-grid homography fails to comprehensively tackle the limitations of homography, Wang *et al.* [39] predicts pseudo plane masks and applies local homographies for each mask. Nie *et al.* [31] achieve more precise local alignment by adopting the TPS model in place of the multi-grid homography scheme. However, TPS constructs a globally smooth deformation field without direct local support, making it less effective in handling significant local deformations.

To address the aforementioned issues, we propose the Exponential-Decay Free-Form Deformation Network (EDFFDNet). The primary motivation is to better handle local deformation while reducing computational costs through more localized processing. Therefore, we adopt the Free-Form Deformation (FFD) model [36], which inherently offers better locality. However, conventional FFD implementations typically employ B-splines as basis functions, which provide smooth deformation but incur substantial computational overhead. To address this, we introduce an Exponential-Decay Free-Form Deformation (EDFFD) model. Experiments demonstrate that it not only improves registration accuracy but also substantially reduces computational overhead. Thanks to the better locality of EDFFD, we can achieve improved alignment performance through additional local refinements.

In the previous state-of-the-art (SOTA) method [31],

MLPs are used for motion aggregation. While linear layers effectively achieve global motion aggregation, their dense interactions lead to large computation cost, restricting deployment. Motivated by the sparse processing in depthwise separable convolution [17], we propose the Adaptive Sparse Motion Aggregator (ASMA). ASMA transforms dense interactions into sparse ones through our proposed Group Linear Layers (GLL), followed by a simple linear layer to adaptively fuse the sparse aggregation results. Experiments demonstrate that ASMA significantly reduces parameters by about 66.6% while preserving the accuracy. Furthermore, previous deep registration methods [29, 31] consistently use global correlation across all stages. While this offers a large receptive field for low-overlap cases, it introduces disturbance from distant regions in local refinement stages, reducing accuracy. We thus propose a progressive correlation strategy that applies global correlation for coarse motion estimation and local correlation for fine local motion estimation, realizing higher accuracy and efficiency.

Based on the above improvements, our method achieves significant performance gains compared to the SOTA method [31]. As shown in Fig. 1, our method demonstrates robust performance in regions with significant depth disparities with lower computational costs. Specifically, with only one local refinement, our EDFFDNet reduces the number of parameters, inference memory, and total time by 70.5%, 32.6%, and 33.7%, respectively, while achieving a 0.5 dB PSNR improvement. Furthermore, with an additional local refinement, our EDFFDNet achieves a 1.06 dB PSNR gain, while still maintaining lower computational costs.

In summary, the contributions of our work are as follows:

- We propose EDFFDNet, an unsupervised registration framework that outperforms state-of-the-art methods in accuracy, computational efficiency, and generalization.
- We introduce an exponential-decay free-form deformation model that outperforms conventional TPS and B-spline FFD models in handling local deformation and enhancing computational efficiency.
- We design an adaptive sparse motion aggregator that replaces dense interactions with sparse ones for motion aggregation, reducing parameters and improving accuracy.
- We introduce a progressive correlation strategy that applies global-local correlation patterns for coarse-to-fine motion estimation, achieving better accuracy and efficiency.

2. Related Work

2.1. Single Homography Methods

Homography estimation methods include feature-based and learning-based approaches. Feature-based methods detect feature points [22, 27, 33] and use robust estimators [2, 3, 11] to solve homography, but struggle in low-texture

regions. Learning-based methods learn robust representations effectively. DeTone *et al.* [9] first proposed a VGG-style [35] deep homography network. Recent supervised advances include end-to-end estimators [5], attention mechanisms [6], and multi-scale correlation searching [45] for higher accuracy. However, supervised methods trained on synthetic data generalize poorly to real images. Nguyen *et al.* [28] first trained a homography network on real data unsupervised, while Jiang *et al.* [19] generated realistic datasets from real images. To improve accuracy, Zhang *et al.* [41] and Le *et al.* [20] used predicted masks to eliminate moving objects, and Hong *et al.* [16] introduced coplanarity constraints for dominant plane registration. However, single-homography methods fundamentally lack representational capacity for scenes with depth disparities.

2.2. Adaptive Warping Methods

Various approaches have been proposed to address challenges caused by depth disparities. Gao *et al.* [12] introduced a dual-homography warping model for two dominant planes. Zaragoza *et al.* [40] proposed APAP, which combines global homography with local homographies. Lee *et al.* [21] improved this by partitioning images into superpixels and estimating weighted homographies. Li *et al.* [23] employed the thin plate spline (TPS) model to handle parallax in image warping. However, these feature-based methods often fail in low-texture scenarios. For learning-based approaches, Wang *et al.* [38] introduced inter-grid consistency and intra-grid regularity terms to stabilize the training of multi-grid homography networks. Nie *et al.* [29] proposed a contextual correlation layer to improve efficiency. However, multi-grid homography does not fully address the restriction of homography. Wang *et al.* [39] then proposed multi-plane homography estimation, which applies local homographies for each predicted pseudo-plane. Nie *et al.* [31] replaced the multi-grid homography scheme with TPS to achieve more precise and efficient local alignment. However, TPS constructs a globally smooth deformation field without local support, which limits its performance in scenarios requiring more localized deformations.

3. Method

3.1. Network Architecture

As illustrated in Fig. 2, the proposed Exponential-Decay Free-Form Deformation Network (EDFFDNet) consists of three main modules: multi-scale feature extractor (MFE), global homography estimation module, and local refinement module. Initially, multi-scale features are extracted from the input target image \mathbf{I}_t and reference image \mathbf{I}_r . The global homography module then estimates a global homography \mathbf{H} , which serves as an initial alignment for subsequent local refinements. Taking the number of local re-

finement stages $N_s = 2$ as an example, in each stage i , multi-scale features are utilized to compute local correlations and estimate control point motions. These motions are then used to generate residual displacements $\Delta\mathbf{D}_i$ based on the proposed Exponential-Decay Free-Form Deformation (EDFFD) model. Finally, the residual displacements $\Delta\mathbf{D}_2$ from the last refinement stage are combined with the global homography \mathbf{H} and the residual displacements $\Delta\mathbf{D}_1$ from the previous local refinement stage. This combination is used to align the target image \mathbf{I}_t with the reference image \mathbf{I}_r , resulting in the aligned image $\mathbf{I}_{t \rightarrow r}$.

Feature Extraction. As illustrated in Fig. 2a, the target image \mathbf{I}_t and the reference image \mathbf{I}_r with size $H \times W$ are fed into MFE to extract multi-scale feature maps $\mathbf{F}_t^{(d)}, \mathbf{F}_r^{(d)}$, where $d \in \{4, 8, 16\}$ denotes the downsample scale. MFE mainly consists of ResNet50 blocks [15], and a 1×1 convolution is applied to the feature output at each resolution.

Motion Estimation. As shown in Fig. 2b, the correlation result from the correlation computation is fed into the motion estimator, which estimates motion parameters. Specifically, it estimates 4-point motion for homography and control point motion for free-form deformation (FFD). The motion estimator consists of N_c convolutional blocks and an Adaptive Sparse Motion Aggregator (ASMA). Each convolutional block consists of 3×3 convolutions, ReLU, and max-pooling. Convolutional blocks are used to extract the latent motion feature, while the ASMA transforms the latent motion feature into motion parameters. The details of the ASMA will be discussed in the corresponding section.

3.2. Exponential-Decay Free-Form Deformation

Restricted by the planar assumption, homography transformation often exhibits inaccurate alignment in non-planar scenes due to its inherent limitation. The multi-grid homography scheme is proposed [40] but lacks efficient parallel acceleration capabilities for deep learning [29]. The thin plate spline (TPS) transformation model [4] is then employed to achieve efficient and flexible deformation [31], but it inherently struggles to handle substantial local deformation. In contrast, B-spline free-form deformation (B-spline FFD) [36] provides better locality essentially compared to TPS. However, the cubic B-spline basis computation incurs significant computational overhead. To address this, we propose the Exponential-Decay Free-Form Deformation (EDFFD) method, which achieves more efficient local deformation through an improved basis function.

Taking the local refinement stage i as an example, for B-spline FFD, let $P = \{\mathbf{p}_{m,n} \mid 0 \leq m \leq M_i, 0 \leq n \leq N_i\}$ denote the set of control points uniformly distributed on the image, forming a mesh grid of size $M_i \times N_i$. The motion of the control point $\mathbf{p}_{m,n}$ estimated by the network is denoted as $\Delta\mathbf{p}_{m,n}$. The deformation of a point $\mathbf{x} = (x_1, x_2)$ in the

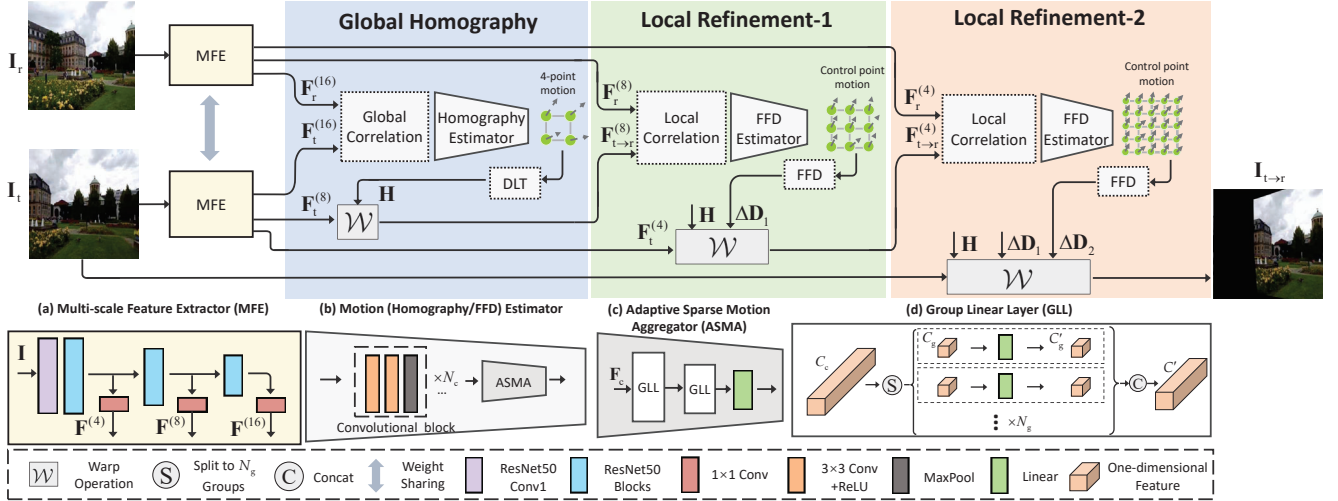


Figure 2. An overview of the proposed Exponential-Decay Free-Form Deformation Network (EDFFDNet). (a) Multi-scale Feature Extractor (MFE). (b) Motion (Homography/FFD) Estimator. (c) Adaptive Sparse Motion Aggregator (ASMA). (d) Group Linear Layer (GLL).

plane is directly given by:

$$\mathbf{x}' = \mathbf{x} + \sum_{m=0}^{M_i} \sum_{n=0}^{N_i} \Delta \mathbf{p}_{m,n} \Phi((\mathbf{x} - \mathbf{p}_{m,n})/\eta), \quad (1)$$

where η represents the grid spacing, \mathbf{x}' is the deformed position, and $\Phi(\cdot)$ denotes the basis product, defined as:

$$\Phi(\mathbf{u}) = \beta^3(u_1)\beta^3(u_2), \quad (2)$$

where $\mathbf{u} = (u_1, u_2)$ denotes the input vector, $\beta^3(\cdot)$ is the cubic B-spline function obtained by three times convolution of the zeroth-order B-spline function [37] (For more details, please refer to Section A.1 of supplementary material.), formulated as:

$$\beta^3(u) = \begin{cases} \frac{2}{3} - |u|^2 + \frac{|u|^3}{2}, & 0 \leq |u| \leq 1 \\ \frac{(2-|u|)^3}{6}, & 1 \leq |u| < 2, \\ 0, & 2 \leq |u|. \end{cases} \quad (3)$$

Cubic B-spline has two principal characteristics that make them effective for deformation tasks [36]: 1) Locality, where each control point influences a limited neighborhood, with its influence diminishing over distance, enabling localized shape manipulation, and 2) C^2 continuity, ensuring smooth deformation fields for natural image alignment. However, the computation presents three issues: 1) High polynomial computation cost, 2) Basis product that requires separate basis computation in two dimensions, and 3) Piecewise computation that impedes GPU parallelism. To overcome these limitations while preserving beneficial properties, we propose EDFFD as an alternative to B-spline FFD, defined as:

$$\mathbf{x}' = \mathbf{x} + \sum_{m=0}^{M_i} \sum_{n=0}^{N_i} \Delta \mathbf{p}_{m,n} \exp(-r_{m,n}/(\theta\eta)), \quad (4)$$

where θ is a factor controlling the decay rate of influence, $r_{m,n}$ represents the Euclidean distance, defined as:

$$r_{m,n} = \|\mathbf{x} - \mathbf{p}_{m,n}\|. \quad (5)$$

This model is motivated by four key principles:

- **Simplified Influence Metric:** Instead of combining control point influences in two dimensions via basis products, we directly use Euclidean distance, eliminating expensive separate basis computation.
- **Computational Efficiency:** The exponential function provides C^∞ smoothness with lower computational overhead than cubic polynomials. Modern GPU architectures further accelerate its implementation through hardware-optimized transcendental units.
- **Parallel Compatibility:** The non-piecewise nature of the exponential function enables fully parallel evaluation across spatial domains, contrasting with B-spline's conditional branching that hinders GPU utilization.
- **Locality Preservation:** The exponential function naturally decays significantly with distance, ensuring that the influence of each control point remains localized.

3.3. Adaptive Sparse Motion Aggregation

Previous work [31] utilizes convolutional blocks and MLP to construct the motion estimator. Although linear layers are computationally efficient and perform well in motion aggregation, they suffer from high parameter counts, limiting their practical deployment. To address this limitation, we propose an Adaptive Sparse Motion Aggregator (ASMA) inspired by the sparse feature processing in depth-wise separable convolutions [17], which reduces computational overhead while preserving accuracy. ASMA comprises two Group Linear Layers (GLL) and a single linear layer. The group linear layers transform dense interactions into sparse ones, while the linear layer adaptively fuses the

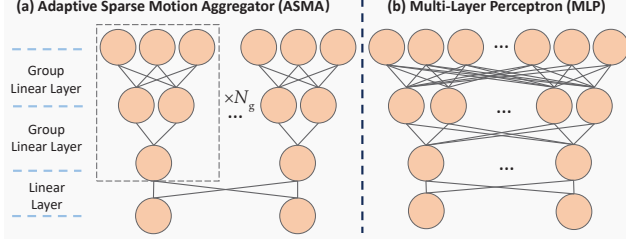


Figure 3. Comparison of neuron interactions between our Adaptive Sparse Motion Aggregator (ASMA) and a Multi-Layer Perceptron (MLP). ASMA achieves efficient motion aggregation by sparsely processing grouped features.

sparse motion aggregation results. Fig. 2c and Fig. 2d illustrate the structure of ASMA, and Fig. 3 highlights the differences between ASMA and traditional MLP.

Let $\mathbf{F}_c \in \mathbb{R}^{C_c}$ denote the flattened latent motion feature extracted from N_c convolutional blocks. This feature is first fed into GLL, which partitions \mathbf{F}_c into N_g groups to form the group feature $\mathbf{F}_{g,k} \in \mathbb{R}^{C_g}$, $k = 1, \dots, N_g$, where $C_g = C_c/N_g$ assuming C_c is divisible by N_g . Each group feature $\mathbf{F}_{g,k}$ is then passed through an independent linear layer to produce output group feature $\mathbf{F}'_{g,k} \in \mathbb{R}^{C'_g}$ as follows:

$$\mathbf{F}'_{g,k} = \mathbf{W}_k(\mathbf{F}_{g,k}) + \mathbf{b}_k, \quad (6)$$

where \mathbf{W}_k and \mathbf{b}_k are the weight and bias of the k -th linear layer. All the output group features are concatenated and then passed through the ReLU activation σ to obtain the output feature $\mathbf{F}' \in \mathbb{R}^{C'}$ of the current GLL as follows:

$$\mathbf{F}' = \sigma(\text{Concat}(\mathbf{F}'_{g,1}, \dots, \mathbf{F}'_{g,N_g})). \quad (7)$$

GLL effectively conducts feature interactions within each group by transforming the dense connections of a traditional linear layer into sparse ones essentially. After passing through two GLLs, the sparse aggregation motion results are adaptively fused by a simple linear layer. This fusion serves as a global interaction mechanism and is responsible for outputting the motion parameters.

3.4. Progressive Correlation Strategy

Low-overlap cases in real-world images require global correlation for wide search ranges, but this is computationally expensive and disturbs local refinement. While [29] proposes efficient global correlation, successive computations still incur significant overhead [31]. Inspired by the observation that required search ranges decrease as estimation accuracy improves, we propose a progressive correlation strategy: using global correlation in the global homography stage for broad search, then transitioning to local correlation in refinement stages for improved efficiency and accuracy.

For global correlation computation, we adopt the patch-to-patch correlation approach [29]. For the target feature $\mathbf{F}_t^{(d)}$, dense patches of size $K \times K$ with stride 1 are extracted as convolutional filters. These filters are applied to

the reference feature $\mathbf{F}_r^{(d)}$ to form a global correlation volume \mathbf{C}^g with shape $H^{(d)} \times W^{(d)} \times H^{(d)}W^{(d)}$. Each value in the column can be formulated as:

$$\mathbf{C}_{(x_r, y_r, x_t, y_t)}^g = \sum_{i, j = -\lfloor \frac{K}{2} \rfloor}^{\lfloor \frac{K}{2} \rfloor} \frac{\langle \mathbf{F}_{r, (x_r+i, y_r+j)}^{(d)}, \mathbf{F}_{t, (x_t+i, y_t+j)}^{(d)} \rangle}{\|\mathbf{F}_{r, (x_r+i, y_r+j)}^{(d)}\| \|\mathbf{F}_{t, (x_t+i, y_t+j)}^{(d)}\|}. \quad (8)$$

Each position in \mathbf{C}^g can be regarded as a vector of length $H^{(d)}W^{(d)}$. The vectors are first scaled by a constant α to increase the in-class distance and then passed through a softmax to obtain matching probabilities. The positions with the highest probabilities are thereby identified, from which we derive the feature flow \mathbf{V} with shape $H^{(d)} \times W^{(d)} \times 2$ as the global correlation result.

For local correlation computation, we adopt the local correlation method utilized in [45] that directly computing the correlation between the feature at position \mathbf{p} in the reference feature $\mathbf{F}_r^{(d)}$ and features within the local area in the target feature $\mathbf{F}_t^{(d)}$. The local correlation is formulated as:

$$\mathbf{C}^l(\mathbf{p}, \mathbf{p}') = \mathbf{F}_r^{(d)}(\mathbf{p})^\top \mathbf{F}_t^{(d)}(\mathcal{A}(\mathbf{p}', r)), \quad (9)$$

where $\mathcal{A}(\mathbf{p}', r)$ represents the sampling local area centered at \mathbf{p}' with radius r , \mathbf{C}^l is the local correlation result with shape $H^{(d)} \times W^{(d)} \times (2r+1)^2$.

3.5. Optimization

To achieve both precise content alignment and natural shape preservation in the image registration results, we optimize the network using two terms. The content alignment term aligns the input images according to their content, while the shape preservation term prevents unnatural distortions.

Content Alignment. Consider a reference image \mathbf{I}_r and a target image \mathbf{I}_t . Let $\mathcal{W}(\cdot, \cdot)$ represent the warping operation, and \mathbf{J} denote an all-one matrix with the same resolution as the image. We define the content alignment loss as follows:

$$\begin{aligned} \mathcal{L}_{\text{content}} = & \lambda_0 \|\mathbf{I}_r \cdot \mathcal{W}(\mathbf{J}, \mathcal{H}) - \mathcal{W}(\mathbf{I}_t, \mathcal{H})\|_1 \quad (10) \\ & + \lambda_0 \|\mathbf{I}_t \cdot \mathcal{W}(\mathbf{J}, \mathcal{H}^{-1}) - \mathcal{W}(\mathbf{I}_r, \mathcal{H}^{-1})\|_1 \\ & + \sum_{i=1}^{N_s} \lambda_i \|\mathbf{I}_r \cdot \mathcal{W}(\mathbf{J}, \mathcal{FFD}_i) - \mathcal{W}(\mathbf{I}_t, \mathcal{FFD}_i)\|_1, \end{aligned}$$

where λ_0 and λ_i are weights, \mathcal{H} and \mathcal{FFD}_i denote the warping parameters of global homography and free-form deformation at the local refinement stage i , respectively.

Shape Preservation. As in [31], we preserve the shape by leveraging the inter-grid constraint ℓ_{inter} and intra-grid constraint ℓ_{intra} based on the grid edge \vec{e} as follows:

$$\mathcal{L}_{\text{shape}} = \sum_{i=1}^{N_s} (\ell_{\text{intra}, i} + \ell_{\text{inter}, i}). \quad (11)$$

Let $\{\vec{e}_{h,i}\}$ and $\{\vec{e}_{v,i}\}$ represent the sets of horizontal and vertical edges of local refinement stage i , respectively. The intra-grid constraint is defined as follows:

$$\ell_{\text{intra},i} = \frac{1}{(M_i + 1) \times N_i} \left(\sum_{\{\vec{e}_{h,i}\}} \sigma(\langle \vec{e}, \vec{u} \rangle - \frac{2W}{N_i}) + \right. \quad (12)$$

$$\left. \frac{1}{M_i \times (N_i + 1)} \sum_{\{\vec{e}_{v,i}\}} \sigma(\langle \vec{e}, \vec{v} \rangle - \frac{2H}{M_i}) \right),$$

where σ denotes the ReLU activation function, \vec{u} and \vec{v} are unit vectors along the x and y directions, respectively. For the non-overlapping region, the consecutive edges are encouraged to be colinear to preserve the structures. The inter-grid constraint is defined as follows:

$$\ell_{\text{inter},i} = \frac{1}{E} \sum_{\{\vec{e}_{c1,i}, \vec{e}_{c2,i}\}} Q_{c1,c2} \cdot \left(1 - \frac{\langle \vec{e}_{c1}, \vec{e}_{c2} \rangle}{\|\vec{e}_{c1}\| \cdot \|\vec{e}_{c2}\|} \right), \quad (13)$$

where $\{\vec{e}_{c1,i}, \vec{e}_{c2,i}\}$ represents the set of edge pairs that are consecutive in either the horizontal or vertical direction of the local refinement stage i , E denotes the number of edge pairs, and $Q_{c1,c2}$ is 1 for edge pairs in non-overlapping regions and 0 otherwise. The entire loss function is defined as follows:

$$\mathcal{L} = \mathcal{L}_{\text{content}} + \omega \mathcal{L}_{\text{shape}}, \quad (14)$$

where ω is the weight to balance the loss.

4. Experiments

4.1. Implementation Details and Datasets

Network training proceeds in two phases: first, the global homography module is trained for 10 epochs, then training extends to 100 epochs with joint optimization of both global and local modules. For global correlation, patch size $K = 3$ and $\alpha = 10$; for local correlation, radius $r = 4$. Loss weights are $\lambda_0 = 1$, $\lambda_1 = 1.3$, $\lambda_2 = 1.7$, and $\omega = 10$. Training uses batch size 4 and Adam optimizer with learning rate 10^{-4} on PyTorch. Experiments run on NVIDIA RTX 4090 GPU, Intel Xeon Platinum 8352V CPU @ 2.10GHz, and 512GB RAM.

We evaluate our network on the UDIS-D dataset [30], which contains 10,440 training and 1,106 testing image pairs with diverse overlap ratios and scenes. Following prior work [29–31], we assess performance using PSNR and SSIM within overlapping regions. We also evaluate zero-shot performance on other datasets [8, 12, 34, 40].

4.2. Ablation Studies

We conduct ablation studies on the deformation model, motion aggregation, and correlation computation. Unless specified, the default configuration uses EDFFD with $N_s = 1$, $M_1 = 12$, $N_1 = 12$, $\theta = 0.75$, ASMA with $N_g = 8$,

and progressive correlation. **Inference time** is the time for the network to estimate motion parameters from input image pairs, **warp time** is the time for the deformation model to output aligned images from motion parameters, and **total time** is their sum.

4.2.1. Deformation Model

As shown in Table 1, we compare the Thin Plate Spline (TPS), B-spline Free-Form Deformation (B-spline FFD), and our proposed EDFFD under identical network settings. The results demonstrate that B-Spline FFD achieves more precise registration than the TPS model due to its better locality. However, B-spline FFD incurs a significantly larger computational overhead. Our proposed EDFFD achieves nearly the same registration performance as B-spline FFD, while significantly reducing the warp time and inference memory by 47.3% and 34.0%, respectively. This highlights the effectiveness of our improvements in the basis function. Additionally, our EDFFD further reduces warp time by 32.4% compared with the TPS model, due to its more simplified computation. This further indicates the efficiency and accuracy of our proposed EDFFD.

Table 1. Ablation on deformation models.

Model	PSNR	SSIM	Warp time (ms)	Memory (GB)
TPS	25.49	0.838	30.5	3.3
B-spline FFD	25.95	0.850	39.1	4.7
EDFFD	25.93	0.852	20.6	3.1

Table 2. Ablation on influence factor θ settings.

θ	0.25	0.50	0.75	1.00	1.25	1.50
PSNR	25.30	25.93	25.93	25.91	25.85	25.81
SSIM	0.834	0.852	0.852	0.851	0.849	0.848

Table 3. Ablation on local refinement settings.

Stage1 Grid size	Stage2 Grid size	PSNR	SSIM	Warp time (ms)
12×12	N/A	25.93	0.852	20.6
12×12	12×12	26.43	0.866	22.1
12×12	18×18	26.49	0.868	23.8
12×12	24×24	26.55	0.869	25.2

Impact of locality. By design, smaller θ increases influence range, enhancing smoothness but weakening locality, and vice versa. Table 2 shows a noticeable drop at $\theta = 0.25$, while larger θ limits influence range, causing poor smoothness and reduced accuracy. These results indicate that our model’s improved performance primarily benefits from well-balanced locality scope.

Additional Local Refinement. As shown in Table 3, we achieve significant accuracy improvements through an additional local refinement stage with denser grid settings, while incurring minor increases in warp time. This highlights the efficiency and potential performance of our EDFFD.

4.2.2. Motion Aggregation

We investigate different motion aggregation methods in Table 4. ASMA outperforms MLP across all settings, achieving better accuracy while reducing parameters by 66.6% for

$N_g = 8$. More importantly, when MLP’s hidden dimension is reduced by 4 times to match ASMA’s parameter count, ASMA still maintains superior performance. This demonstrates the effectiveness of our ASMA.

Table 4. Ablation on motion aggregation methods.

Motion aggregation	Hidden dimension reduction ratio	N_g	Parameters (M)	PSNR	SSIM
MLP	1	N/A	68.9	25.87	0.850
ASMA	1	4	24.3	25.91	0.851
ASMA	1	8	23.0	25.93	0.852
ASMA	1	16	22.3	25.90	0.851
MLP	4	N/A	27.1	25.76	0.845
ASMA	4	8	17.3	25.89	0.850

Table 5. Ablation on correlation settings.

Homography correlation	FFD correlation	PSNR	SSIM	Inference time (ms)
Global	Global	25.54	0.843	32.8
Local	Local	N/A	N/A	N/A
Global	Local	25.93	0.852	23.0

4.2.3. Correlation Computation

In Table 5, we compare correlation methods for homography and FFD estimation. Global correlation in FFD reduces accuracy due to global field disturbances and increases inference time, while local correlation for homography struggles in low-overlap scenarios from insufficient search range. Our progressive strategy addresses these by providing adequate search range during global homography and finer range in refinement, improving accuracy by 0.39 dB PSNR and reducing inference time by 29.8%.

4.3. Comparative Experiments

We conduct comprehensive comparative experiments for warping accuracy and computational overhead. In the experiments, we set the influence factor $\theta = 0.75$, the number of groups $N_g = 8$, and the number of local refinement stages N_s to either 1 or 2, corresponding to two versions: EDFFDNet and EDFFDNet-2. EDFFDNet uses $M_1 = N_1 = 12$, while EDFFDNet-2 extends EDFFDNet with an additional local refinement stage where $M_2 = N_2 = 18$.

4.3.1. Comparison of Warping Accuracy

We conduct a comprehensive comparison covering both traditional and deep learning-based approaches. The traditional methods included in our comparison are SIFT [27] + RANSAC [11], APAP [40], ELA [23], SPW [24], and LPC [18]. For deep learning methods, we benchmark against UDIS [30], MGDH [29], and UDIS++ [31].

Quantitative Evaluation. Table 6 presents quantitative results on UDIS-D dataset [30], where $I_{3 \times 3}$ denotes identity mapping. Metrics are categorized into three groups following prior works [29–31]. For cases where traditional methods fail, we use identity mapping for evaluation. Our method outperforms previous approaches, especially in challenging scenarios, with significant improvements in

hard and moderate categories, highlighting our model’s enhanced locality.

Qualitative Evaluation. Fig. 4 shows qualitative results on UDIS-D dataset [30], where we combine the green and blue channels of reference image I_r with the red channel of result image $I_{t \rightarrow r}$ for visual assessment. Traditional methods, MGDH, and UDIS++ achieve moderate alignment in hard scenarios with multiple planes and depth disparities requiring localized deformations. EDFFDNet demonstrates better performance through explicit locality, with EDFFDNet-2 achieving well-aligned results in challenging scenarios.

Cross-dataset Evaluation. We further assess the generalization capability of our pre-trained model through cross-dataset validation. We first visualize widely used cases from [40] and [12] in Fig. 5, which demonstrates that our method, despite being trained only on the UDIS-D dataset [30] without cross-dataset fine-tuning, achieves superior generalization performance. We also evaluate zero-shot performance on the ScanNet [8] and ETH3D [34] datasets in Table 7, which demonstrates the strong generalization capability of our method. Moreover, our method maintains competitive performance with traditional methods while requiring substantially less computational time, as shown in Table 8, highlighting its potential for practical deployment.

Table 7. Zero-shot results on the testset of ScanNet and ETH3D. For ScanNet, 10k test image pairs are randomly selected.

	ScanNet		ETH3D	
	PSNR	SSIM	PSNR	SSIM
UDIS [30]	21.82	0.747	19.11	0.615
MGDH [29]	22.08	0.748	19.50	0.641
UDIS++ [31]	21.79	0.729	19.41	0.647
EDFFDNet (ours)	23.37	0.786	20.54	0.693
EDFFDNet-2 (ours)	24.32	0.808	21.47	0.731

Table 8. Comparison of total time (s) required to generate the warping results.

Dataset	Railtrack [40]	Fence [25]	Carpark [12]
Resolution	1500 × 2000	1088 × 816	490 × 653
APAP [40]	159.583	55.560	23.743
ELA [23]	19.251	8.867	5.760
SPW [24]	116.530	9.701	13.711
LPC [18]	2114.944	11.913	68.947
EDFFDNet	0.078	0.055	0.053
EDFFDNet-2	0.097	0.064	0.063

4.3.2. Comparison of Computational Overhead

We compare our method with deep learning approaches UDIS [30], MGDH [29], and UDIS++ [31] in terms of computational costs. Table 9 shows EDFFDNet reduces parameters by 70.5%, inference memory by 32.6%, and total time by 33.7% compared to UDIS++, while improving PSNR by 0.5 dB. With additional local refinement, EDFFDNet-2 achieves 1.06 dB PSNR improvement with low computational costs. These results highlight our method’s accuracy and efficiency.

Table 6. Quantitative comparison of warp on UDIS-D dataset [30] The best is marked in red and the second best is in blue.

	PSNR \uparrow				SSIM \uparrow			
	Easy	Moderate	Hard	Average	Easy	Moderate	Hard	Average
$I_{3\times 3}$	15.87	12.76	10.68	12.86	0.530	0.286	0.146	0.303
SIFT [27]+RANSAC [11]	27.75	24.03	18.46	22.98	0.906	0.828	0.627	0.758
APAP [40]	27.01	23.39	19.54	23.00	0.885	0.802	0.663	0.773
ELA [23]	29.87	25.51	19.68	24.47	0.924	0.865	0.713	0.821
SPW [24]	27.29	22.83	16.94	21.80	0.888	0.764	0.504	0.696
LPC [18]	27.04	22.72	19.34	22.65	0.879	0.768	0.610	0.738
UDIS [30]	27.84	23.95	20.70	23.80	0.902	0.830	0.685	0.793
MGDH [29]	29.52	25.24	21.20	24.89	0.923	0.859	0.708	0.817
UDIS++ [31]	30.19	25.84	21.57	25.43	0.933	0.875	0.739	0.838
EDFFDNet (ours)	30.63	26.31	22.15	25.93	0.938	0.886	0.763	0.852
EDFFDNet-2 (ours)	31.09	26.85	22.79	26.49	0.943	0.898	0.790	0.868

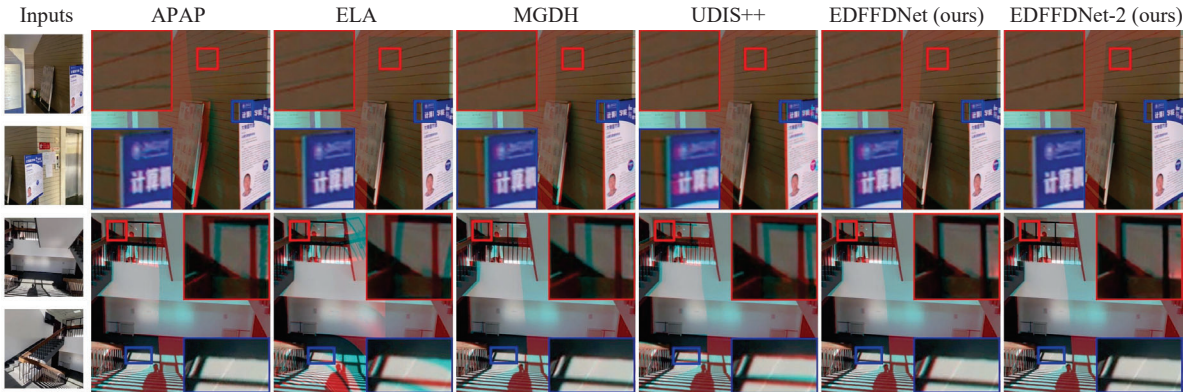


Figure 4. Qualitative results on the UDIS-D dataset [30]. Zoomed-in regions from two distinct planes show multi-plane structures and depth disparities, highlighting the effectiveness of the proposed method in handling local deformations.

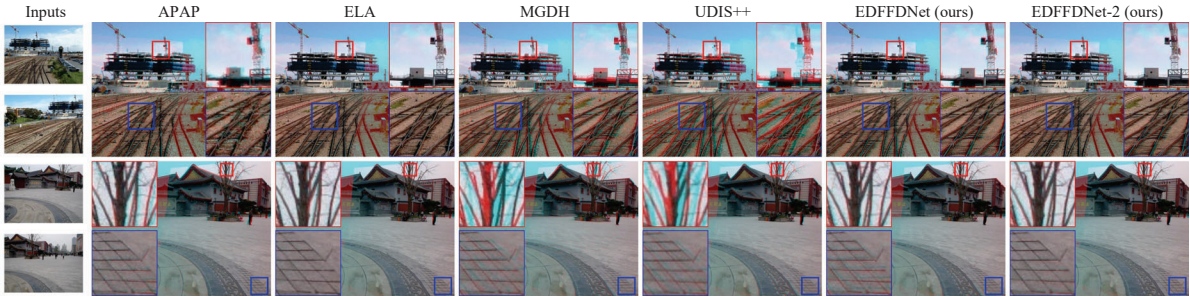


Figure 5. Qualitative results on the widely used cross-dataset cases: “railtrack” (row-1) [40] and “temple” (row-2) [12].

Table 9. Comparison of computational overhead.

Method	PSNR	SSIM	Parameters (M)	Memory (GB)	Total time (ms)
UDIS [30]	23.80	0.793	188.8	7.1	66.7
MGDH [29]	24.89	0.817	16.4	5.3	90.3
UDIS++ [31]	25.43	0.838	78.0	4.6	65.8
EDFFDNet (ours)	25.93	0.852	23.0	3.1	43.6
EDFFDNet-2 (ours)	26.49	0.868	34.5	4.3	55.1

5. Conclusions

We have proposed EDFFDNet, an unsupervised registration framework designed to handle large local deformations. This is achieved through an exponential-decay free-form deformation model for improved locality. Additionally, we

have introduced an adaptive sparse motion aggregator that converts dense interactions into sparse ones for efficiency. We have also developed a progressive correlation strategy for coarse-to-fine estimation. Our approach delivers a significant advancement in both accuracy and efficiency.

Acknowledgements. This work was supported in part by the Zhejiang Provincial Natural Science Foundation of China under grant LD24F020003, in part by the National Natural Science Foundation of China under grant 62301484, in part by the Ningbo Natural Science Foundation of China under grant 2024J454, and in part by the National Key Research and Development Program of China under grant 2023YFB3209800.

References

- [1] Simon Baker and Iain Matthews. Lucas-Kanade 20 years on: A unifying framework. *International Journal of Computer Vision*, 56(3):221–255, 2004. 2
- [2] Daniel Barath, Jiri Matas, and Jana Noskova. MAGSAC: marginalizing sample consensus. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10197–10205, 2019. 1, 2
- [3] Daniel Barath, Jana Noskova, Maksym Ivashechkin, and Jiri Matas. Magsac++, a fast, reliable and accurate robust estimator. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1304–1312, 2020. 1, 2
- [4] Fred L Bookstein and WDK Green. A thin-plate spline and the decomposition of deformations. *Mathematical Methods in Medical Imaging*, 2(14-28):3, 1993. 2, 3
- [5] Si-Yuan Cao, Jianxin Hu, Zehua Sheng, and Hui-Liang Shen. Iterative deep homography estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1879–1888, 2022. 2, 3
- [6] Si-Yuan Cao, Runmin Zhang, Lun Luo, Beinan Yu, Zehua Sheng, Junwei Li, and Hui-Liang Shen. Recurrent homography estimation using homography-guided image warping and focus transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9833–9842, 2023. 2, 3
- [7] Che-Han Chang, Chun-Nan Chou, and Edward Y Chang. CLKN: Cascaded lucas-kanade networks for image alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2213–2221, 2017. 2
- [8] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. 6, 7
- [9] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Deep image homography estimation. *arXiv preprint arXiv:1606.03798*, 2016. 1, 3
- [10] Farzan Erlik Nowruzi, Robert Laganieri, and Nathalie Japkowicz. Homography estimation from image pairs with hierarchical convolutional networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 913–920, 2017. 2
- [11] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 1, 2, 7, 8
- [12] Junhong Gao, Seon Joo Kim, and Michael S Brown. Constructing image panoramas using dual-homography warping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 49–56. IEEE, 2011. 2, 3, 6, 7, 8
- [13] Natasha Gelfand, Andrew Adams, Sung Hee Park, and Kari Pulli. Multi-exposure imaging on mobile devices. In *Proceedings of the 18th ACM international conference on Multimedia*, pages 823–826, 2010. 1
- [14] Heng Guo, Shuaicheng Liu, Tong He, Shuyuan Zhu, Bing Zeng, and Moncef Gabbouj. Joint video stitching and stabilization from moving cameras. *IEEE Transactions on Image Processing*, 25(11):5491–5503, 2016. 1
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3
- [16] Mingbo Hong, Yuhang Lu, Nianjin Ye, Chunyu Lin, Qijun Zhao, and Shuaicheng Liu. Unsupervised homography estimation with coplanarity-aware GAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17663–17672, 2022. 2, 3
- [17] Andrew G Howard. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 2, 4
- [18] Qi Jia, ZhengJun Li, Xin Fan, Haotian Zhao, Shiyu Teng, Xinchun Ye, and Longin Jan Latecki. Leveraging line-point consistence to preserve structures for wide parallax image stitching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12186–12195, 2021. 7, 8
- [19] Hai Jiang, Haipeng Li, Songchen Han, Haoqiang Fan, Bing Zeng, and Shuaicheng Liu. Supervised homography learning with realistic dataset generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9806–9815, 2023. 3
- [20] Hoang Le, Feng Liu, Shu Zhang, and Aseem Agarwala. Deep homography estimation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7652–7661, 2020. 2, 3
- [21] Kyu-Yul Lee and Jae-Young Sim. Warping residual based image stitching for large parallax. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8198–8206, 2020. 2, 3
- [22] 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. 1, 2
- [23] Jing Li, Zhengming Wang, Shiming Lai, Yongping Zhai, and Maojun Zhang. Parallax-tolerant image stitching based on robust elastic warping. *IEEE Transactions on multimedia*, 20(7):1672–1687, 2017. 2, 3, 7, 8
- [24] Tianli Liao and Nan Li. Single-perspective warps in natural image stitching. *IEEE transactions on image processing*, 29: 724–735, 2019. 7, 8
- [25] Chung-Ching Lin, Sharathchandra U Pankanti, Karthikeyan Natesan Ramamurthy, and Aleksandr Y Aravkin. Adaptive as-natural-as-possible image stitching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1155–1163, 2015. 7
- [26] Shuaicheng Liu, Lu Yuan, Ping Tan, and Jian Sun. Bundled camera paths for video stabilization. *ACM transactions on graphics (TOG)*, 32(4):1–10, 2013. 1
- [27] David G Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91–110, 2004. 1, 2, 7, 8

- [28] Ty Nguyen, Steven W Chen, Shreyas S Shivakumar, Camillo Jose Taylor, and Vijay Kumar. Unsupervised deep homography: A fast and robust homography estimation model. *IEEE Robotics and Automation Letters*, 3(3):2346–2353, 2018. 2, 3
- [29] Lang Nie, Chunyu Lin, Kang Liao, Shuaicheng Liu, and Yao Zhao. Depth-aware multi-grid deep homography estimation with contextual correlation. *IEEE transactions on circuits and systems for video technology*, 32(7):4460–4472, 2021. 2, 3, 5, 6, 7, 8
- [30] Lang Nie, Chunyu Lin, Kang Liao, Shuaicheng Liu, and Yao Zhao. Unsupervised deep image stitching: Reconstructing stitched features to images. *IEEE Transactions on Image Processing*, 30:6184–6197, 2021. 1, 6, 7, 8
- [31] Lang Nie, Chunyu Lin, Kang Liao, Shuaicheng Liu, and Yao Zhao. Parallax-tolerant unsupervised deep image stitching. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7399–7408, 2023. 1, 2, 3, 4, 5, 6, 7, 8
- [32] Lang Nie, Chunyu Lin, Kang Liao, Yun Zhang, Shuaicheng Liu, Rui Ai, and Yao Zhao. Eliminating warping shakes for unsupervised online video stitching. In *European Conference on Computer Vision*, pages 390–407. Springer, 2024. 1
- [33] 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. 1, 2
- [34] Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3260–3269, 2017. 6, 7
- [35] Karen Simonyan. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3
- [36] Nicholas J Tustison, Brian B Avants, and James C Gee. Directly manipulated free-form deformation image registration. *IEEE transactions on image processing*, 18(3):624–635, 2009. 2, 3, 4
- [37] Michael Unser. Splines: A perfect fit for signal and image processing. *IEEE Signal processing magazine*, 16(6):22–38, 1999. 4
- [38] Miao Wang, Guo-Ye Yang, Jin-Kun Lin, Song-Hai Zhang, Ariel Shamir, Shao-Ping Lu, and Shi-Min Hu. Deep online video stabilization with multi-grid warping transformation learning. *IEEE Transactions on Image Processing*, 28(5):2283–2292, 2018. 2, 3
- [39] Yasi Wang, Hong Liu, Chao Zhang, Lu Xu, and Qiang Wang. Mask-homo: Pseudo plane mask-guided unsupervised multi-homography estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5678–5685, 2024. 2, 3
- [40] Julio Zaragoza, Tat-Jun Chin, Michael S Brown, and David Suter. As-projective-as-possible image stitching with moving dlt. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2339–2346, 2013. 1, 2, 3, 6, 7, 8
- [41] Jirong Zhang, Chuan Wang, Shuaicheng Liu, Lanpeng Jia, Nianjin Ye, Jue Wang, Ji Zhou, and Jian Sun. Content-aware unsupervised deep homography estimation. In *Proceedings of the European Conference on Computer Vision*, pages 653–669. Springer, 2020. 2, 3
- [42] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2002. 1
- [43] Zhuofan Zhang, Zhen Liu, Ping Tan, Bing Zeng, and Shuaicheng Liu. Minimum latency deep online video stabilization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23030–23039, 2023. 1
- [44] Yuan Zhou, Anand Rangarajan, and Paul D Gader. An integrated approach to registration and fusion of hyperspectral and multispectral images. *IEEE Transactions on Geoscience and Remote Sensing*, 58(5):3020–3033, 2019. 2
- [45] Haokai Zhu, Si-Yuan Cao, Jianxin Hu, Sitong Zuo, Beinan Yu, Jiacheng Ying, Junwei Li, and Hui-Liang Shen. Mcnet: Rethinking the core ingredients for accurate and efficient homography estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25932–25941, 2024. 2, 3, 5
- [46] Danping Zou and Ping Tan. Coslam: Collaborative visual slam in dynamic environments. *IEEE transactions on pattern analysis and machine intelligence*, 35(2):354–366, 2012. 1