# Unifying Correspondence, Pose and NeRF for Generalized Pose-Free Novel View Synthesis

Sunghwan Hong
Korea University

Jaewoo Jung
Korea University

Heeseong Shin
Korea University

Jiaolong Yang
Microsoft Research Asia

Seungryong Kim
Korea University

Chong Luo
Microsoft Research Asia

## Abstract

*This work delves into the task of pose-free novel view synthesis from stereo pairs, a challenging and pioneering task in 3D vision. Our innovative framework, unlike any before, seamlessly integrates 2D correspondence matching, camera pose estimation, and NeRF rendering, fostering a synergistic enhancement of these tasks. We achieve this through designing an architecture that utilizes a shared representation, which serves as a foundation for enhanced 3D geometry understanding. Capitalizing on the inherent interplay between the tasks, our unified framework is trained end-to-end with the proposed training strategy to improve overall model accuracy. Through extensive evaluations across diverse indoor and outdoor scenes from two real-world datasets, we demonstrate that our approach achieves substantial improvement over previous methodologies, especially in scenarios characterized by extreme viewpoint changes and the absence of accurate camera poses. The project page and code will be made available at: https://ku-cvlab.github.io/CoPoNeRF/.*

## 1. Introduction

In real-world scenarios aimed at rendering novel views from unposed images, the initial step often involves employing an off-the-shelf camera pose estimation [46, 50, 56, 67]. These estimated poses are then typically integrated with a pre-trained generalized NeRF model [20, 66] to facilitate view synthesis. However, this approach is not without its drawbacks. The primary limitation stems from the inherent disparities or misalignments that may arise when combining models dedicated to different tasks. This method risks potential inconsistencies, as it treats pose estimation and NeRF rendering as distinct, separate processes, potentially
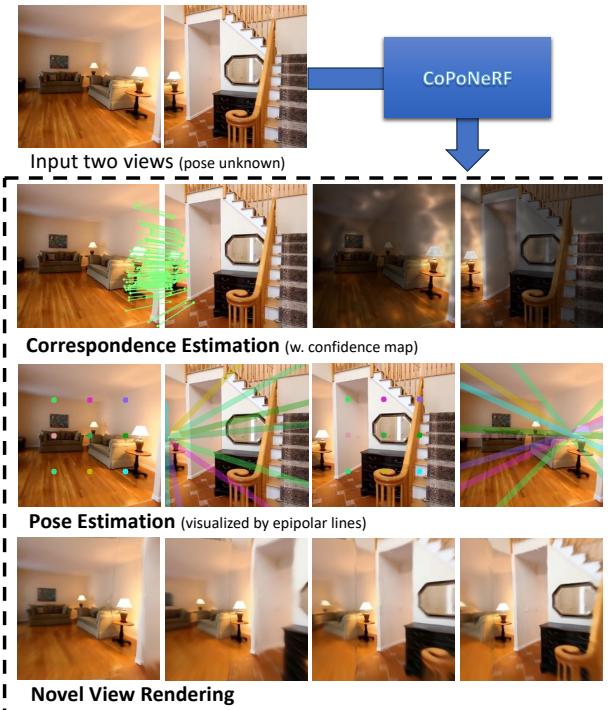


Figure 1. **Overview.** Given an unposed pair of images, possibly under extreme viewpoint changes and with minimal overlapping, our framework synergistically performs and effectively fosters mutual enhancement among three tasks – 2D correspondence estimation, camera pose estimation, and NeRF rendering – to enable high-quality novel view synthesis.

leading to suboptimal results in the synthesized views.

Recent developments in alternative approaches have trended towards the integration of pose estimation with NeRF rendering, thus leading to the advent of pose-free, generalized NeRF approaches [12, 52]. This has been primarily realized through the meticulous assembly of developed modules in a multi-task framework. For example, [52]

exploited correspondence information sourced from a well-established RAFT optical flow model [55] and depth information from a single-view generalizable NeRF model [66] for pose estimation. [12] combined a generalized NeRF module with a RAFT-like recurrent GRU module, responsible for camera pose and depth estimation, and implemented a three-stage training scheme for these two modules. Despite the promising results shown by these seminal approaches, the intrinsic complementarity of the three key tasks, correspondence estimation, pose estimation, and NeRF rendering, are not fully recognized and utilized. This resulted in solutions that were suboptimal, particularly in scenarios characterized by extreme viewpoint changes or minimal overlapping regions.

Acknowledging the shared core objectives between the three tasks, which is the precise interpretation and reconstruction of three-dimensional geometry from two dimensional image data, we emphasize the critical importance of cultivating a *shared representation* among them. To this end, we propose a unified framework, namely CoPoNeRF, designed to estimate three distinct outputs, correspondence, camera pose, and NeRF rendering from this common representation. By adopting joint training, we maximize the synergy between these components, ensuring that each task not only contributes to but also benefits from this shared medium. This integrated approach effectively pushes the boundaries beyond what is achievable when treating each task as an independent and disjoint problem.

We evaluate the effectiveness of our framework using large-scale real-world indoor and outdoor datasets [36, 68]. Our results demonstrate that this framework successfully synthesizes high-quality novel views while simultaneously achieving precise relative camera pose estimation. We also provide extensive ablation studies to validate our choices.
**Our contributions** are summarized as follows:
- We tackle the challenging task of pose-free generalizable novel view synthesis, addressing the minimal view overlap scenarios that are not considered by prior methods. This aspect of our approach illustrates its applicability in handling complex, real-world conditions.
- We propose a unified framework that enhances the processes of pose estimation, correspondence estimation, and NeRF rendering. This framework is designed to exploit the interdependencies of these components with a shared representation learning.
- Leveraging the advanced representations learned by our framework, we achieve state-of-the-art performance not only in pose-free scenarios but also in generalized novel view synthesis with poses.

## 2. Related Work

**Generalized Neural Radiance Fields.** Classical NeRF methodologies rely on numerous multi-view image datasets [4, 42], while recent efforts aim to learn reliable radiance fields from sparse imagery with a single feed-forward pass [10, 13, 20, 33, 58, 66]. These, however, depend heavily on precise camera poses and significant view overlap. To lessen this dependency, various frameworks optimize NeRF by integrating geometry and camera pose refinement, offering a degree of pose flexibility [7, 32, 35, 57, 71].

We focus on generalized frameworks for pose-free view synthesis; DBARF [12], for example, proposes a pose-agnostic solution by combining camera pose estimation with novel view synthesis. However, the network is trained in a staged manner with a local cost volume to encode multi-view information, struggling with minimal overlapping pairs and failing to fully harness the potential synergy between pose estimation and NeRF. FlowCAM [52], on the other hand, leverages a weighted Procrustes analysis [17] and an established optical flow network for point correspondences [55]. Despite its attempt to formulate a multi-task framework, the reliance on the flow model inevitably risks failures in both view synthesis and pose estimation, especially for images with extreme viewpoint changes.

**Establishing Correspondences.** Correspondence estimation, pivotal for various applications such as SLAM [21], SfM [51], and camera pose estimation [43], traditionally entails a sequence involving keypoint detection [6, 18, 39, 48], feature description [22, 45, 64], tentative matching, and outlier filtering [2, 3, 8, 60, 65]. While outlier filtering stage also holds significant importance in relative pose estimation [5], the intrinsic quality of feature descriptors and the validity of matching scores markedly influence the pose prediction outcomes [9, 46]. In this research, we harness the power of meticulously established correspondences, leveraging them to bolster both pose estimation and neural rendering processes, optimizing the overall task efficacy.

**Camera Pose Estimation.** Classic camera pose estimation methods primarily utilize hand-crafted algorithms to solve pose estimations using a set of correspondences, focusing on improving descriptor quality, cost volume, or outlier filtering to enhance correspondence quality [25, 37, 43]. More recent works have shifted towards learning direct mappings from images to poses. Notable advancements include the use of CNN-based networks to solve pose regression, such as the work by [41] and subsequent developments [23, 34]. Our work aligns more closely with methodologies tackling wide-baseline image pairs as inputs, an aspect relatively lesser explored. Some examples include leveraging a 4D correlation map for relative pose regression [9], predicting discrete camera position distributions [11], and modifying the ViT [19] to emulate the 8-point algorithm [46]. Unique in approach, our method pioneers addressing the wide-baseline setting in generalized pose-free novel view synthesis tasks.
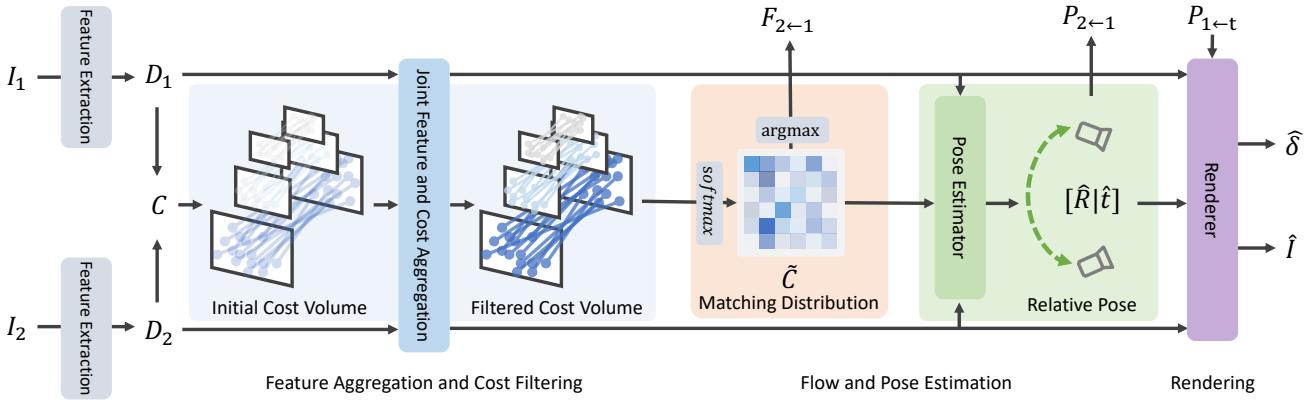
Figure 2. **Overall architecture of the proposed method.** For a pair of images, we extract multi-level feature maps and construct 4D correlation maps at each level, encoding pixel pair similarities. These maps are refined for flow and pose estimation, and the renderer then uses the estimated pose and refined feature maps for color and depth computation.

# 3. Unified Framework for Generalized Pose-Free Novel View Synthesis

## 3.1. Problem Formulation

Assuming an unposed pair of images $I_1, I_2 \in \mathbb{R}^{H \times W \times 3}$ taken from different viewpoints as the input, our goal is to synthesize an image $\hat{I}_t$ from a novel view. In this work, we assume camera intrinsics are given, as it is generally available from modern devices [1]. Different from classical generalized NeRF tasks [10, 58, 66], our task is additionally challenged by the absence of camera pose between the input images. To this end, we estimate the relative camera pose between $I_1, I_2$ as $P_{2 \leftarrow 1} \in \mathbb{R}^{4 \times 4}$, consisting of rotation $R \in \mathbb{R}^{3 \times 3}$ and translation $T \in \mathbb{R}^{3 \times 1}$, and deduce $P_{2 \leftarrow t} = P_{2 \leftarrow 1} P_{1 \leftarrow t}$ with $P_{1 \leftarrow t}$ as the desired rendering viewpoint, which are then used in conjunction with the extracted feature maps to compute the pixel color at the novel view by the renderer.

## 3.2. Cost Volume Construction

The first stage of the pipeline is feature extraction, which will be shared across all three tasks. Because our method must be robust for scale differences and extreme viewpoint changes, we use multi-level feature maps to capture both geometric and semantic cues from different levels of features. Given a pair of images $I_1$ and $I_2$, we first extract $l$-levels of deep local features $D_1^l, D_2^l \in \mathbb{R}^{h^l \times w^l \times c^l}$ with ResNet [27]. Subsequent to feature extraction, the extracted features undergo cost volume construction.

Unlike the previous methods that only consider local receptive fields within their cost volumes [10, 12, 33, 63], we consider *all pairs of similarities between features* to handle both small and large displacements [14, 15, 28–30]. Specifically, we compute and store the pairwise cosine similarity between features, obtaining a 4D cost volume $\{C^l\}_{l=1}^L \in \mathbb{R}^{h^l \times w^l \times h^l \times w^l}$, where $L$ is the number of levels.

## 3.3. Feature Aggregation and Cost Filtering

**Joint Feature and Cost Aggregation.** Recent progress in image correspondence has demonstrated the value of self- and cross-attention mechanisms in capturing global context within images and enhancing inter-image feature extraction, vital for understanding multi-view geometry [20, 54, 62]. Studies have also emphasized the importance of cost aggregation for reducing noise in cost volumes and embedding geometric priors [13, 16, 31].

Building upon these developments, we introduce a self-attention-based aggregation block that processes the feature maps and cost volume, *i.e.*, $D_1, D_2$, and $C$ (level indicator $l$ omitted for brevity). Specifically, two augmented cost volumes are first constructed by feature and cost volume concatenation: $[C, D_1]$ and $[C^T, D_2] \in \mathbb{R}^{h \times w \times (hw+c)}$. Then, treating each 2D location in the augmented cost volume as a token, our aggregation block performs self-attention operation $\phi$ using feature maps and cost volumes as *values*. The resulting cost volumes are obtained as $\phi([C, D_1]) + \phi([C^T, D_2])^T$ that ensures consistent matching scores.

**Leveraging Cost Volume as Matching Distribution.** Our method leverages enhanced feature maps and a refined cost volume from the aggregation block to inter-condition the feature maps. Unlike the standard practice of using a cross-attention map from two feature maps [13, 54, 62], we introduce a simple and more effective adaptation by employing the refined cost volume from the aggregation block, rather than computing a separate cross-attention map. Specifically, we apply softmax to this volume to create a cross-attention map, which then guides the alignment of feature maps with matching probabilities. This layer is integrated with the aggregation block in an interleaved manner, crucial for refining and assimilating multi-view information. More details can be found in the *supp. material*.

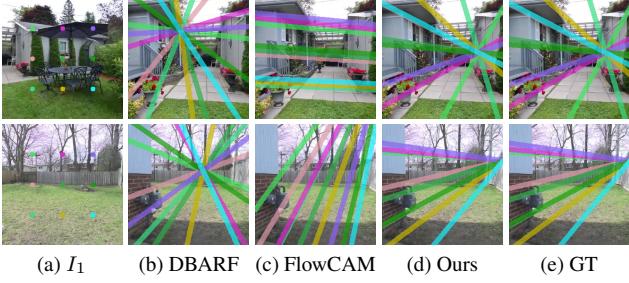(a) $I_1$    (b) DBARF    (c) FlowCAM    (d) Ours    (e) GT

Figure 3. **Visualization of epipolar lines.** We use the relative camera pose to draw epipolar lines based on the points in (a). Our predictions can well follow the ground truth even under large viewpoint changes.
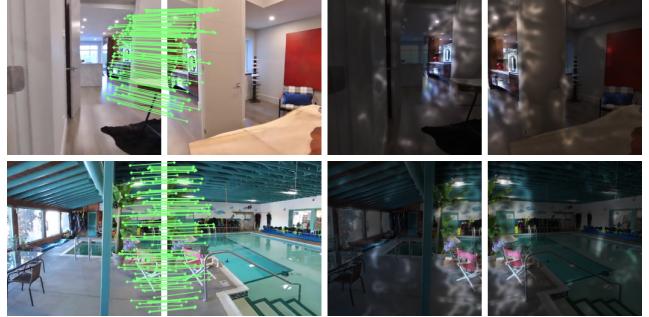


Figure 4. **Visualization of correspondences and confidence.** We show top 100 confident matches between input images and the covisible regions are highlighted based on confidence scores.

The final cost volume, $\frac{1}{L}\sum_l^L \tilde{C}^l$, calculated from each $l$-th level, is then used for relative pose and flow estimation. This cost volume plays a pivotal role in consolidating multi-level feature correspondences, directly impacting the accuracy of our pose and flow estimations.

### 3.4. Flow and Relative Pose Estimation

Making use of the cost volume from previous steps, we define a dense flow field, $F(i)$ that warps all pixels $i$ in image $I_1$ towards $I_2$. We also estimate relative camera pose $P_{2\leftarrow1}$ from this cost volume as it sufficiently embodies confidence scores and spatial correspondences [9, 46]. To estimate the dense flow map $F_{2\leftarrow1}$, we can simply apply argmax to find the highest scoring correspondences. While this may be sufficient for image pairs with large overlapping regions, we account for the potential occlusions by computing a confidence map. Specifically, following [40], we obtain a cyclic consistency map $M_{2\leftarrow1}(i)$ using the reverse field $F_{1\leftarrow2}$ as an additional input, and check if the following condition is met for consistency: $||F_{2\leftarrow1}(i) + F_{1\leftarrow2}(i + F_{2\leftarrow1}(i))||_2 < \tau$, where $||\cdot||_2$ denotes frobenius norm and $\tau$ is a threshold hyperparameter. The reverse cyclic consistency map $M_{1\leftarrow2}$ is computed with similar procedure.

To estimate camera parameters, we use the knowledge taken from previous study [46] and adopt an essential matrix module to output rotation $R$ and translation $t$. The essential matrix module is a mapping module that exploits each transformer token from images to a feature that is used to predict $R$ and $t$. This module contains positional encoding, bilinear attention, and dual softmax over attention map $A$. Following the design in Sec. 3.3, we make a modification to replace $A$ with our cost volume, since it acts as a key for emulating 8-point algorithm, such that the better spatial correspondences encoded in $A$ can aid more accurate camera pose estimation. Subsequent to the essential matrix module, we finally regress 6D rotation representations [69] and 3 translation parameters with scales using MLPs.

### 3.5. Attention-based Renderer

Within our method, the rendering module is tasked with synthesizing novel views, guided by the estimated camera poses and a pair of aggregated features from previous steps. Borrowing from recent advancements, we adopt a strategy of sampling pixel-aligned features along the epipolar lines of each image and augment the features by a corresponding feature in the other image, as suggested by Du et al. [20]. This technique also enables us to ascertain depth $\delta$ by triangulating these features, thereby streamlining the typically resource-intensive 3D sampling process. Given the set of sampled features from epipolar lines, we adopt an attention-based rendering procedure to compute the pixel color, as done similarly in previous methods [20, 33, 53].

### 3.6. Training Objectives

The outputs of our models are colors, depths, an estimated relative camera pose, and a dense flow map. Our model is trained with an objective function consisting of four losses: image reconstruction loss $\mathcal{L}_{\text{img}}$, matching loss $\mathcal{L}_{\text{match}}$, camera pose loss $\mathcal{L}_{\text{pose}}$, and the triplet consistency loss $\mathcal{L}_{\text{tri}}$. For rendering, we use the photometric loss between the rendered color and the target color defined as $L1$ loss.

**Matching Loss.** To provide training signals for correspondence estimation, we adopt self-supervised SSIM loss as a valuable alternative since obtaining ground-truth correspondences between image pairs is often challenging and impractical, since the depth information is required. The SSIM [59] loss computes the structural similarity between the warped image and the target image, offering a data-driven approach to assess the quality of image registration without relying on explicit depth measurements or ground-truth correspondences. The matching loss is defined as:

$$\mathcal{L}_{\text{match}} = \sum_i M_{1\leftarrow2}(i)(1 - \text{SSIM}(F_{1\leftarrow2}(I_2(i)), I_1(i)))$$
$$+ M_{2\leftarrow1}(i)(1 - \text{SSIM}(F_{2\leftarrow1}(I_1(i)), I_2(i))),$$
$$(1)$$

where $\phi(\cdot, \cdot)$ yields the SSIM score between two comparative inputs.

**Pose Loss.** Although we only take a pair of unposed images as input for the inference phase, for the training phase, we incorporate readily available and ubiquitous camera poses, thanks to the extensive availability of video data and the deployment of conventional pose estimation methodologies prevalent in the field, including SfM and SLAM. Our pose loss is a combination of geodesic loss [49] for rotation and $L_2$ distance loss for translation[1]. Specifically, they are combined with addition and defined as:

$$\mathcal{L}_{\text{rot}} = \arccos\left(\frac{\text{trace}(\hat{R}^T R) - 1}{2}\right)$$
$$\mathcal{L}_{\text{trans}}(\hat{t}, t) = \|\hat{t} - t\|_2^2, \quad (2)$$

where $\hat{R}$ and $\hat{t}$ indicates the estimated rotation and translation.

Empirical results indicate that including gradient feedback from other losses alongside the pose loss contributes to unstable training, typically leading to suboptimal model performance. Aligning with the literature [12, 32, 35], our findings also accentuate that the expansive search space and the intrinsic complexities associated with pose optimization increase the difficulty of the learning process.

To address this challenge, we directly incorporate ground-truth pose data into key modules, like rendering or feature projection, during training. This approach restricts gradient flow to the pose loss, proving highly effective in our experiments. Conceptually, this resembles the teacher forcing strategy in RNNs [61], where ground truth, rather than previous network outputs, guides training. This method encourages network parameters are optimized in direct alignment with pose estimation objectives, similar to using ground truth inputs for more direct supervision in teacher forcing.

**Triplet Consistency Loss.** Finally, we propose a triplet consistency loss $\mathcal{L}_{\text{tri}}$ seamlessly incorporating all the outputs of our model, flow, pose, and depth. Our loss extends the cycle consistency [70] loss, which has been prevalent in the field of computer vision, by incorporating the three outputs from our network. Specifically, at each iteration, given a set of randomly selected coordinates in the target frame that are projected to $I_1$ and $I_2$ using depth $\delta$, camera pose $P$, and camera intrinsic $K$, we define the p: $i'_1 = K_1 P_{1 \leftarrow t} \delta(i) K_1^{-1} i$ and $i'_2$ is defined similarly. Using the set of projected points $i'_1$, we employ $F_{2 \leftarrow 1}$ to warp them towards $I_2$ to obtain warped coordinates $\hat{i}'_2$. We finally

---

[1] Although two-view geometry inherently lacks the capability to discern translation scales, we let the model learn to align all predictions to true scales via recognition, as done in [46].

apply Huber loss function [26] $\psi$ such that $M(i'_1)\psi(\hat{i}'_2, i'_2)$, where $M(\cdot)$ eliminates the unconfident matches. This loss effectively measures the consistency between the estimated depth and optical flow across the views. If the estimated depth and flow are accurate, the projected and warped points should coincide, resulting in a small loss value.

In summary, our total training loss is defined as

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{img}} + \mathcal{L}_{\text{match}} + \mathcal{L}_{\text{pose}} + \lambda_{\text{tri}} \mathcal{L}_{\text{tri}}, \quad (3)$$

where $\lambda_{\text{tri}}$ is a scaling factor.

# 4. Experiments

## 4.1. Implementation Details

Our encoder uses ResNet-34, taking 256×256 image as an input, and extracts 3 levels of feature maps with a spatial resolution of 16, 32, and 64. We set $\lambda_{\text{tri}} = 0.01$. Our network is implemented using PyTorch [44] and trained with the AdamW [38] optimizer. We set the base learning rate as 2e−4 and use an effective batch size of 64. The network is trained for 50K iterations, taking approximately 2 days. We exponentially decay the learning rate with $\gamma = 0.95$ after every epoch. We train and evaluate all other baselines on the same datasets for fair comparisons. We provide training and evaluation details of ours and our competitors in the *supp. material* for completeness.

## 4.2. Experimental Settings

**Datasets.** We train and evaluate our method on RealEstate10K [68], a large-scale dataset of both indoor and outdoor scenes, and ACID [36], a large-scale dataset of outdoor coastline scenes. For RealEstate10K, we employ a subset of the complete dataset, resulting in a training set comprising 21,618 scenes and a testing set consisting of 7,200 scenes, while for ACID, we use 10,935 scenes for training and 1,893 scenes for testing.

**Tasks and Baselines.** We assess the performance of our method on two tasks: novel-view synthesis and relative pose estimation. The latter task serves a dual purpose as it also assesses the quality of our correspondences, a methodological approach that aligns with prior image matching studies. We first compare with established generalized NeRF variants, specifically PixelNeRF [66] and [20], highlighting the complexities of wide-baseline inputs. Subsequently, we compare with existing pose-free generalized NeRF methods, namely DBARF [12] and FlowCAM [52]. For relative pose estimation, the evaluation includes several comparisons: matching methods [18, 50, 56] followed by 5-point algorithm [43] and RANSAC [24] and end-to-end pose estimation frameworks [46, 67]. Finally, we conduct a comparative analysis with pose-free generalized NeRF methods [12, 52].

|  (a) $I_1$ | (b) $I_2$ | (c) DBARF | (d) FlowCAM | (e) Ours | (f) Ground Truth |

Figure 5. **Qualitative comparison on RealEstate10K.**

**Evaluation Metrics.** We use the standard image quality metrics (PSNR, SSIM, LPIPS, MSE) for novel view synthesis evaluation. For relative pose estimation, we use geodesic rotation error and angular difference for translation[2] as done in classical methods [41, 43]. Our statistical analysis includes average, median, and standard deviation of errors, with the median offering robustness against outliers and standard deviation indicating error variability.

**Evaluation Protocol.** As our approach is the first work to this extremely challenging task, we introduce a new evaluation protocol. By default, we assume $I_1, I_2$ and $P_{1 \leftarrow t}$ are provided to the model, while $I_t$ *is used solely for computing the metrics*[3]. First, when presented with a video sequence of a scene, we employ a frame-skipping strategy. The value of $N$ frames skipped between each frame is dynamically determined based on the total number of frames. For sequences with fewer than 100 frames, $N$ is calculated as one-third of the total frame count; otherwise, we set $N = 50$. This gives us three images $I_1, I_t, I_2$, which are taken from $N = 0, 50, 100$, respectively. Next, for evaluation, while a common practice in relative pose estimation tasks is to leverage rotation variance [9, 46], this approach disregards translations, often leading to controversial classification of the images into one of the distributions in some cases, *e.g.,* image pairs with zoomed-in and out cameras. We thus partition the test set into three subsets named *Small*, *Medium*, and *Large* to denote the *degree of overlap in the scenes*. With such a splitting scheme, for the RealEstate10K dataset, we obtain subsets containing 3593, 1264, and 2343 scenes

---

[2]Since translation scale is theoretically indeterminable in two-view camera pose estimation, evaluating it could potentially lead to inconclusive or erroneous interpretations (see the *supp. material* for more results).

[3]Existing pose-free generalized NeRF methods use target frames for additional geometric cues during evaluation [12, 52]. For practicality, we assume target frames are *only available for metric calculation*, not for method operation, applying this uniformly across all methods. This aligns with real-world scenarios where target views are not accessible.

respectively, whereas for the ACID dataset, they encompass 559, 429, and 905 scenes each. We show the visualization of the splits and their distributions in the *supp. material*.

To quantitatively compute the overlapping regions, we employ a pre-trained state-of-the-art dense image matching method [56] to find the overlapping ratio $o_{12}$ within the image, defined as the ratio of pixels in $I_1$ whose correspondence with pixels in $I_2$ is found with high confidence. We define the overlap between two images to be the intersection over union of two images as $overlap = \frac{1}{o_{12}^{-1} + o_{21}^{-1} - 1}$, and consider images with overlap greater than 0.75 as *Large*, less than 0.5 as *Small*, and the in-between as *Medium*. Finally, the evaluation metrics are computed using synthesized novel view $\hat{I}_t$ and the estimated relative pose between $I_1$ and $I_2$, $\hat{P}_{1 \leftarrow 2}$ and those of the ground-truths.

### 4.3. Experimental Results

**Relative Pose Estimation.** We report quantitative results in Tab. 1 and the visualization of epipolar lines from the estimated camera poses are shown in Fig. 4. From the results, we observe that our framework significantly outperforms the existing pose-free NeRFs, where they fail to estimate reliable camera pose which can lead to poor view-synthesis quality. Moreover, we observe that compared to pose estimation methods [46, 67], our framework achieve significantly better accuracy, demonstrating the effectiveness of the captured synergy between pose estimation, rendering and image matching. However, it is also notable that PDC-Net+ [56] achieves better performance. This is because PDC-Net was learned with GT correspondences that are obtained using depth information, which indicates further improvements in all our three tasks can be promoted if depth information is incorporated in our framework.

**Novel View Synthesis.** Tab. 2 shows quantitative comparisons, whereas Fig. 5 show qualitative comparisons. From the results, compared to previous pose-free approaches [12,

| Overlap | Task | Method | RealEstate-10K | | | | | | ACID | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Rotation | | | Translation | | | Rotation | | | Translation | | |
| | | | Avg(°↓) | Med(°↓) | STD(°↓) | Avg(°↓) | Med(°↓) | STD(°↓) | Avg(°↓) | Med(°↓) | STD(°↓) | Avg(°↓) | Med(°↓) | STD(°↓) |
| Small | COLMAP (Matching) | SP+SG [18, 24, 50] | 9.793 | 2.270 | 22.084 | 12.549 | 4.638 | 23.048 | 10.920 | 2.797 | 22.761 | 22.214 | 7.526 | 33.719 |
| | | PDC-Net+ [24, 56] | 3.460 | 1.128 | 7.717 | 6.913 | 2.752 | 15.558 | 2.520 | 0.579 | 6.372 | 15.664 | 4.215 | 29.640 |
| | Pose Estimation | Rockwell *et al.* [46] | 12.604 | 6.860 | 14.502 | 91.455 | 91.499 | 56.872 | 8.466 | 3.151 | 13.380 | 88.421 | 88.958 | 36.212 |
| | | RelPose [67] | 12.102 | 4.803 | 21.686 | - | - | - | 10.081 | 4.753 | 13.343 | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 17.520 | 13.218 | 15.946 | 126.282 | 140.358 | 43.691 | 8.721 | 3.205 | 12.916 | 95.149 | 99.490 | 47.576 |
| | | FlowCAM [52] | 11.883 | 6.778 | 15.676 | 87.119 | 58.245 | 26.895 | 8.663 | 6.675 | 7.930 | 92.130 | 85.846 | 40.821 |
| | | Ours | **5.471** | **2.551** | **11.733** | **11.862** | **5.344** | **21.080** | **3.548** | **1.129** | **8.619** | 23.689 | 11.289 | 30.391 |
| Medium | COLMAP (Matching) | SP+SG [18, 24, 50] | 1.789 | 0.969 | 3.502 | 9.295 | 3.279 | 20.456 | 3.275 | 1.306 | 6.474 | 16.455 | 5.426 | 29.035 |
| | | PDC-Net+ [24, 56] | 1.038 | 0.607 | 1.841 | 6.667 | 2.262 | 18.247 | 2.378 | 0.688 | 5.841 | 14.940 | 4.301 | 27.379 |
| | Pose Estimation | Rockwell *et al.* [46] | 12.168 | 6.552 | 14.385 | 82.478 | 82.920 | 55.094 | 4.325 | 1.564 | 6.177 | 90.555 | 90.799 | 51.469 |
| | | RelPose [67] | 4.942 | 3.476 | 6.206 | - | - | - | 5.801 | 2.803 | 6.574 | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 7.254 | 4.379 | 7.009 | 79.402 | 75.408 | 54.485 | 4.424 | 1.685 | 6.164 | 77.324 | 77.291 | 49.735 |
| | | FlowCAM [52] | 4.154 | 3.346 | 3.466 | 42.287 | 41.594 | 24.862 | 8.778 | 6.589 | 7.489 | 95.444 | 87.308 | 43.198 |
| | | Ours | **2.183** | **1.485** | **2.419** | **10.187** | **5.749** | **15.801** | **2.573** | **1.169** | **3.741** | 21.401 | 10.656 | 28.243 |
| Large | COLMAP (Matching) | SP+SG [18, 24, 50] | 1.416 | 0.847 | 1.984 | 21.415 | 7.190 | 34.044 | 1.851 | 0.745 | 3.346 | 22.018 | 7.309 | 33.775 |
| | | PDC-Net+ [24, 56] | 0.981 | 0.533 | 1.938 | 16.567 | 5.447 | 29.883 | 1.953 | 0.636 | 4.133 | 18.447 | 4.357 | 35.564 |
| | Pose Estimation | Rockwell *et al.* [46] | 12.771 | 7.214 | 14.863 | 91.851 | 88.923 | 57.444 | 2.280 | 0.699 | 3.512 | 86.580 | 87.559 | 50.369 |
| | | RelPose [67] | 4.217 | 2.447 | 5.621 | - | - | - | 4.309 | 2.011 | 5.288 | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 3.455 | 1.937 | 3.862 | 50.094 | 33.959 | 43.659 | 2.303 | 0.859 | 3.409 | 54.523 | 38.829 | 45.453 |
| | | FlowCAM [52] | 2.349 | 1.524 | 2.641 | 34.472 | 27.791 | 31.615 | 9.305 | 6.898 | 9.929 | 97.392 | 89.359 | 43.777 |
| | | Ours | **1.529** | **0.991** | **1.822** | **15.544** | **7.907** | **24.626** | 3.455 | 1.129 | 7.265 | 22.935 | 10.588 | 30.974 |
| *Avg* | COLMAP (Matching) | SP+SG [18, 24, 50] | 5.605 | 1.301 | 16.129 | 14.887 | 5.058 | 27.238 | 4.819 | 1.203 | 13.473 | 20.802 | 6.878 | 32.834 |
| | | PDC-Net+ [24, 56] | 2.189 | 0.751 | 5.678 | 10.100 | 3.243 | 22.317 | 2.315 | 0.619 | 5.655 | 16.461 | 4.292 | 31.391 |
| | Pose Estimation | Rockwell *et al.* [46] | 12.585 | 6.881 | 14.587 | 90.115 | 88.648 | 40.948 | 4.568 | 1.312 | 8.358 | 88.433 | 88.961 | 36.197 |
| | | RelPose [67] | 8.285 | 3.845 | 16.329 | - | - | - | 6.348 | 2.567 | 9.047 | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 11.144 | 5.385 | 13.516 | 93.300 | 102.467 | 57.290 | 4.681 | 1.421 | 8.417 | 71.711 | 68.892 | 50.277 |
| | | FlowCAM [52] | 7.426 | 4.051 | 12.135 | 50.659 | 46.281 | 52.321 | 9.001 | 6.749 | 8.864 | 95.405 | 88.133 | 42.849 |
| | | Ours | **3.610** | **1.759** | **8.617** | **12.766** | **7.534** | **15.510** | **3.283** | **1.134** | **7.093** | 22.809 | 14.502 | 21.572 |

Table 1. **Pose estimation performance on RealEstate-10K and ACID.** Gray color indicates methods not directly comparable as they supervise correspondence with ground-truth depth; they are included for reference only. We also specify the targeted task for each method.

| Overlap | GT Pose | Method | RealEstate-10K | | | | ACID | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | PSNR↑ | LPIPS↓ | SSIM↑ | MSE↓ | PSNR↑ | LPIPS↓ | SSIM↑ | MSE↓ |
| Small | ✓ | PixelNeRF [66] | 13.126 | 0.639 | 0.466 | 0.058 | 16.996 | 0.528 | 0.487 | 0.030 |
| | | Du *et al.* [20] | 18.733 | 0.378 | 0.661 | 0.018 | 25.553 | 0.301 | 0.773 | 0.005 |
| | ✗ | DBARF [12] | 13.453 | 0.563 | 0.522 | 0.045 | 14.306 | 0.503 | 0.541 | 0.037 |
| | | FlowCAM [52] | 15.435 | 0.528 | 0.570 | 0.034 | 20.153 | 0.475 | 0.594 | 0.016 |
| | | Ours | **17.153** | **0.459** | **0.577** | **0.025** | **22.322** | **0.358** | **0.649** | **0.010** |
| Medium | ✓ | PixelNeRF [66] | 13.999 | 0.582 | 0.462 | 0.042 | 17.228 | 0.534 | 0.501 | 0.029 |
| | | Du *et al.* [20] | 22.552 | 0.263 | 0.764 | 0.008 | 25.694 | 0.303 | 0.769 | 0.005 |
| | ✗ | DBARF [12] | 15.201 | 0.487 | 0.560 | 0.030 | 14.253 | 0.457 | 0.538 | 0.038 |
| | | FlowCAM [52] | 18.481 | 0.592 | 0.441 | 0.018 | 20.158 | 0.476 | 0.585 | 0.015 |
| | | Ours | **19.965** | **0.343** | **0.645** | **0.013** | **22.407** | **0.352** | **0.648** | **0.009** |
| Large | ✓ | PixelNeRF [66] | 15.448 | 0.479 | 0.470 | 0.031 | 17.229 | 0.522 | 0.500 | 0.028 |
| | | Du *et al.* [20] | 26.199 | 0.182 | 0.836 | 0.004 | 25.338 | 0.307 | 0.763 | 0.005 |
| | ✗ | DBARF [12] | 16.615 | 0.380 | 0.648 | 0.022 | 14.086 | 0.419 | 0.534 | 0.039 |
| | | FlowCAM [52] | 22.418 | 0.707 | 0.287 | 0.009 | 20.073 | 0.478 | 0.580 | 0.016 |
| | | Ours | **22.542** | **0.250** | **0.724** | **0.008** | **22.529** | **0.351** | **0.649** | **0.009** |
| *Avg* | ✓ | PixelNeRF [66] | 14.438 | 0.577 | 0.467 | 0.047 | 17.160 | 0.527 | 0.496 | 0.029 |
| | | Du *et al.* [20] | 21.833 | 0.294 | 0.736 | 0.011 | 25.482 | 0.304 | 0.769 | 0.005 |
| | ✗ | DBARF [12] | 14.789 | 0.490 | 0.570 | 0.033 | 14.189 | 0.452 | 0.537 | 0.038 |
| | | FlowCAM [52] | 18.242 | 0.597 | 0.455 | 0.023 | 20.116 | 0.477 | 0.585 | 0.016 |
| | | Ours | **19.536** | **0.398** | **0.638** | **0.016** | **22.440** | **0.323** | **0.649** | **0.010** |

Table 2. **Novel view rendering performance on RealEstate-10K and ACID.** Gray text indicates methods not directly comparable for their use of ground-truth pose at evaluation.

52], our approach outperforms them all. Note that we also include results from generalized NeRFs [20, 66] to highlight the complexity and challenge of this task. It's important to note that while our method may not surpass the state-of-the-art [20], the proximity of our results to it underlines the potential and effectiveness of our approach in the absence of camera pose information.

## 4.4. Ablation Study and Analysis

**Component ablation.** In Tab. 3, we validate the effectiveness of each component within our framework. The base-line in the first row represents a variant equipped with only the feature backbone, renderer, pose head and image reconstruction loss. We then progressively add each component. From the results, we observe clear improvements on performance for every component, demonstrating that they all contribute to the final performance. A particularly illustrative comparison are (I) vs (II) and (II) vs (III), where simply adding each loss leads to apparent improvements, indicating that the process of finding correspondences, learning 3D geometry through rendering and estimating camera all contribute largely to the performance. However, it is also notable that PDC-Net+. + [20] reports lower rotation and translation angular errors. This can be attributed to the use of classical solvers [24, 43] that are known to output more precise transformations given a sufficient numbe of correspondences [47].

**Will our method be more effective than the combination of readily available models from separate tasks?** Unless our framework achieves more competitive rendering quality, the practicality of our method will be rather limited. In this analysis, we compare our framework with the variants that adopt two separate methods for camera pose estimation and rendering. The results are reported in Table 4a. Specifically, for the first row, we combine pretrained generalized NeRF [20] with an off-the-shelf matching network for pose estimation. The second row shows the outcomes obtained with [46]. Note that RelPose [67] does not predict translations, and we thus leverage [46]. Summarizing the results, we observe that our approach outperforms the

| | Components | Avg | | | | | Large | | | | | Medium | | | | | Small | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR | SSIM | LPIPS | Mean Rot.(°) | Mean Trans.(°) | PSNR | SSIM | LPIPS | Mean Rot.(°) | Mean Trans.(°) | PSNR | SSIM | LPIPS | Mean Rot.(°) | Mean Trans.(°) | PSNR | SSIM | LPIPS | Mean Rot.(°) | Mean Trans.(°) |
| (I) | Baseline | 14.646 | 0.553 | 0.513 | 29.123 | 52.237 | 16.430 | 0.626 | 0.406 | 27.622 | 71.047 | 14.246 | 0.532 | 0.524 | 28.636 | 51.675 | 13.624 | 0.522 | 0.578 | 30.275 | 40.164 |
| (II) | + pose loss | 16.03 | 0.547 | 0.485 | 8.755 | 62.246 | 18.872 | 0.638 | 0.367 | 4.303 | 74.553 | 16.311 | 0.548 | 0.476 | 8.193 | 62.510 | 14.087 | 0.488 | 0.565 | 11.855 | 54.127 |
| (III) | + flow head (SSIM loss) | 17.393 | 0.578 | 0.440 | 6.737 | 43.104 | 17.964 | 0.588 | 0.419 | 4.257 | 34.584 | 20.083 | 0.658 | 0.329 | 2.578 | 54.666 | 15.439 | 0.522 | 0.520 | 10.325 | 38.554 |
| (IV) | + cycle loss | 17.899 | 0.593 | 0.432 | 6.675 | 43.132 | 20.555 | 0.662 | 0.321 | 2.624 | 50.541 | 18.882 | 0.590 | 0.412 | 4.137 | 34.882 | 15.821 | 0.549 | 0.512 | 10.210 | 38.107 |
| (V) | + aggregation module | 18.629 | 0.611 | 0.406 | 5.008 | 24.769 | 21.402 | 0.689 | 0.294 | 1.977 | 31.295 | 19.219 | 0.621 | 0.383 | 2.951 | 19.509 | 16.614 | 0.556 | 0.487 | 7.706 | 22.363 |
| (VI) | + matching distribution | **19.536** | **0.638** | **0.398** | **3.610** | **12.766** | **22.542** | **0.724** | **0.250** | **1.530** | **10.187** | **19.965** | **0.645** | **0.343** | **2.183** | **11.860** | **17.153** | **0.577** | **0.459** | **3.610** | **12.766** |

Table 3. **Component ablations on RealEstate10K.**

| | PSNR | SSIM | LPIPS | R (°) | t (°) | t (m) |
|---|---|---|---|---|---|---|
| PDC-Net+. + [20] | 18.140 | 0.606 | 0.366 | **2.091** | **8.817** | 0.696 |
| Rockwell et al. [46] + [20] | 14.892 | 0.562 | 0.500 | 12.588 | 90.189 | 0.524 |
| **Ours** | **19.526** | **0.641** | **0.312** | 2.739 | 11.362 | **0.290** |

(a) **Fixed Pose and Generalized NeRF**

| | PSNR | SSIM | LPIPS | R (°) | t (°) |
|---|---|---|---|---|---|
| DBARF [12] + pose loss ($\mathcal{L}_{pose}$) | 12.998 | 0.468 | 0.566 | 11.82 | 80.66 |
| FlowCAM [52] + pose loss ($\mathcal{L}_{pose}$) | 18.646 | 0.589 | 0.433 | 7.505 | 44.347 |
| **Ours** | **19.536** | **0.638** | **0.398** | **3.610** | **12.766** |

(b) **Pose Supervision**

| | PSNR | SSIM | LPIPS | R (°) | t (°) |
|---|---|---|---|---|---|
| Baseline + w/o Teacher Forcing | 13.856 | 0.502 | 0.577 | 31.112 | 104.490 |
| Baseline | 14.646 | 0.553 | 0.513 | 29.123 | 52.237 |
| Ours + w/o Teacher Forcing | 18.785 | 0.635 | 0.415 | 5.254 | 40.571 |
| **Ours** | **19.536** | **0.638** | **0.398** | **3.610** | **12.766** |

(c) **Training Strategy**

| | PSNR | SSIM | LPIPS | R (°) | t (°) |
|---|---|---|---|---|---|
| Du et al. [20] + Noisy Pose ($\sigma = 0.025$) | 18.850 | 0.618 | 0.363 | 2.292 | 6.171 |
| Du et al. [20] + GT Pose | 21.833 | 0.736 | **0.294** | - | - |
| Ours + Noisy Pose ($\sigma = 0.025$) | 19.500 | 0.633 | 0.353 | 2.292 | 6.171 |
| Ours + GT Pose | **22.781** | **0.758** | 0.314 | - | - |

(d) **The learned representation**

Table 4. **More Ablations and insights. See text for details.**

other variants by large margin, highlighting the importance of capturing the underlying synergy between the tasks and the practicality of the proposed approach.

**Direct pose supervision.** To assess the impact of direct pose supervision on the rendering and pose estimation performance of existing methods, we explore the potential enhancement of [12, 52] through the integration of direct supervision. For this experiment, we modify them by incorporating the same loss signals as our approach, specifically the geodesic rotation loss and L2 distance loss for translation. The results are reported in Table 4b.

Although we have found that the use of direct pose supervision aiming to harness benefits from synergistic relationships among different tasks is crucial, when applied to existing frameworks, we have observed only marginal improvements in image quality and a decline in the performance of pose estimation for FlowCAM, while overall decline is obeserved in the performance of DBARF. This outcome is primarily attributed to the architectural design of FlowCAM, where each module operates in a relatively isolated manner without a focus on seamless integration. Conversely, the performance reduction observed in DBARF is multifaceted, with specific causative factors being challenging to pinpoint. These findings are expounded in the *supp. material* for further discussion. In contrast, our proposed framework demonstrates an inherent advantage in harnessing the benefits of pose supervision without requiring further considerations.

**Ablation study on our training strategy.** As explained in Section 3.6, we adopt a special training strategy that conceptually bears similarity to the teacher forcing strategy. In Table 4c, we validate whether this strategy actually helps. For this experiments, we evaluate two variants that builds upon either *Baseline* or *Ours*: the variant that uses the estimated camera poses at training phase and the other that uses the ground-truth. Comparing the results, we observe clear performance differences between the variants, demonstrating the effectiveness of the proposed strategy.

**The learned representation.** In Table 4d, we compare four variants: the first two rows show results from the state-of-the-art generalized NeRF method with ground-truth and noisy poses, respectively, while the next two rows detail results using our framework under the same conditions. The noisy poses are synthetically perturbed from the ground-truth poses by adding Gaussian noise with $\sigma = 0.025$. From the results, we can observe that that our unified framework's learned representation markedly improves rendering performance, outperforming the current state-of-the-art generalized NeRF method when using ground-truth poses. These results further highlight the importance of jointly learning the three tasks in improving the capabilities of the shared representation.

## 5. Conclusion

In this work, we have presented a novel unified framework that integrates camera pose estimation, NeRF rendering and correspondence estimation. This approach effectively overcomes the limitations of existing approaches, particularly in scenarios with limited data and complex geometries. Our experimental results, encompassing both indoor and outdoor scenes with only a pair of wide-baseline images, demonstrate the framework's robustness and adaptability in achieving high-quality novel view synthesis and precise camera pose estimation. Extensive ablation studies further validates our choices and highlight the potential of our method to set a new standard in this task.

## Acknowledgement

Figure 6. **Overview of aggregation module.**



Figure 7. **Illustration of our training losses.**
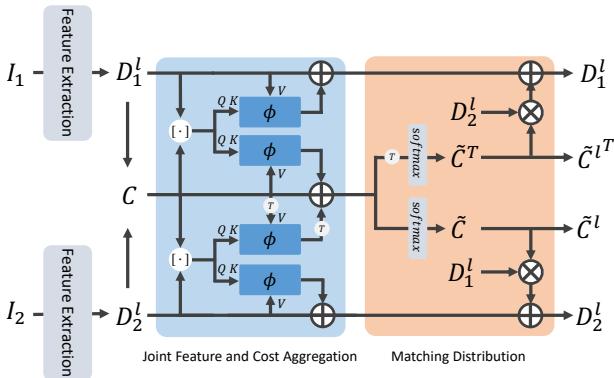
This document includes the following contents: 1) more architectural details of our method, 2) more training and evaluation details of our method and others, 3) distribution of our overlap-based data splitting, 4) more discussions about the experimental results, and 5) additional quantitative and qualitative results for the comparison with other methods and our ablation study.

## A. Architectural Details

### A.1. Feature Aggregation and Cost Filtering

Analogous to existing methods, we use attention-based operations for refining both feature and cost volume. We present an overview of the adopted aggregation module in Fig. 6.

### A.2. Loss Signals

In Fig. 7, we show an illustration of our training losses. As shown in the figure, the rendering loss is computed between $\hat{I}$ and $I$, pose loss is computed using the estimated camera pose $[\hat{R}|\hat{t}]$, flow loss is computed using the estimated flow $F_{1\rightarrow 2}$ and the triplet consistency loss is computed using $\hat{\delta}, [\hat{R}|\hat{t}]$, and $F_{1\rightarrow 2}$.

## B. More Training and Evaluation Details

### B.1. Training Details

#### B.1.1. Our Method

**Training Strategy.** Our training procedure closely resembles that of Du et al. [20], with distinctions in data usage and augmentation. Instead of applying random cropping and flipping, as done by Du et al., we used a subset of datasets without data augmentation. We train the model with 4 A6000 GPUs for $1 \sim 2$ days, iterating for 50K iterations, with 192 rays and 64 sampled points on epipolar lines. With this configuration, our rendering speed is approximately 0.4 FPS.
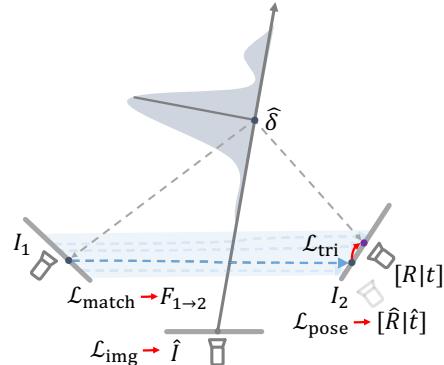
#### B.1.2. FlowCAM [52]

**Architecture.** The inference process of FlowCAM is divided into three distinct steps. First, each frame within a video sequence undergoes feature extraction, yielding deep features and backward flows. Subsequently, using the single-view pixelNeRF algorithm to a surface point cloud, representing the anticipated 3D termination point for each pixel. Next, within each frame, confidence weights are computed among the points. This is achieved by utilizing the RAFT [55] predictions. Finally, these computed confidence weights are fed into a sequence of linear layers to derive the final confidence weights required for solving the Weighted Procrustes formulation [17], and then the desired views are re-rendered.

**Training Strategy.** For training, we take $256 \times 256$ input images. While RealEstate10K was already trained by the authors and the pretrained weights were available, we verify that the training scheme authors provide can reliably transfer to ACID, we attempted reproducing the results on RealEstate first, for which we were able to reproduce the results close to those reported in the paper. Given this success, we followed the same training procedure used and released by the authors for RealEstate to train on ACID. We train for 50K iterations with a single A6000 GPU, which takes approximately 1.5 days, with leaving other hyperparameters unchanged.

#### B.1.3. Rockwell *et al.* [46]

**Architecture.** In their architecture, there are three main components: Image encoder, ViT layer and Essential Matrix Module followed by MLPs. Taking $256 \times 256$ input image pairs as inputs, the image is resized to $224 \times 224$, and then the model first extracts deep features via vanilla resnet-18. Then the feature maps from the coarsest layer are fed to ViT-Tiny for self-attention operations. Subsequently, these feature maps are fed to the Essential Matrix Module, which performs the cross-attention that emulates the 8-point algo-

rithm, and finally, the output is reshaped and fed to the pose regression MLPs.

**Training Strategy.** For training, we follow the same procedure and adopt the default hyperparameters used in the training scripts, as provided in the official github repository that the authors provide for both RealEstate10K and ACID. We use the same data sampling strategy as the one we used to train our model. Specifically, for each scene consisting of a video sequence, we use the first and the last frame as the input images, and the ground-truth relative pose for supervision is computed between them. We trained the network for a total of 120K iterations with batch size set to 32 using a single A6000 GPU.

### B.1.3. RelPose [67]

**Architecture.** Relpose inference is divided into two steps. A pairwise pose prediction step is followed by a joint reasoning step of multiple pairwise estimated relative poses. By taking a set of images as input, they first group all possible pairs of images to estimate all the pairwise relative poses between images. Leveraging an energy-based model, the estimated pairwise relative poses recover a probability distribution over conditional relative rotations where the condition is given as the uniformly sampled relative pose $R \in \mathrm{SO}(3)$. Estimated poses are further refined in the joint reasoning step by inducing a joint likelihood over the camera transformations across multiple images and iteratively improving an initial estimate by maximizing this likelihood.

**Training Strategy.** For training, we follow the same training strategy as [20] and ours, since we aim to compare the performance of relative pose estimation given stereo pairs. However, when using only stereo pairs as input, the joint reasoning step cannot be done as there is only one estimated pose. To make a fair comparison, we increased the number of uniformly sampled relative pose $R \in \mathrm{SO}(3)$ from $N = 36864$ to $N = 250000$, which is the number of queries used in the second stage of the framework. The training was done for 400K iterations of batch size set to 64, using four A6000 GPUs.

### B.1.4. DBARF [12]

**Architecture.** The architecture of DBARF consists of three components: an image encoder, a Pose and Depth Estimation Module, and a Renderer for novel view synthesis. By selecting a target image and nearby images from a scene graph, the ResNet-like [27] image encoder first extracts a feature map used for the subsequent steps. The feature maps of the nearby images are then warped to the target view using the currently estimated camera poses and depths to construct a local cost map for pose and depth estimation done with training a recurrent GRU. The estimated pose is

then used as an input of the Renderer, where they use the IBRNet [58] to render novel views. To enable robust optimization of both the Pose and Depth Estimation Module and the Renderer, they adopt a staged training strategy of dividing the overall training process into three steps: training only the Pose and Depth Estimation Module, training only the Renderer, and jointly training the two components.

**Training Strategy.** For training, we first take 256×256 input images and then resize them to 224×224. As there are no provided weights for DBARF on RealEstate10K and ACID, we trained the network from scratch following the process provided by the authors. For both datasets, we selected six nearby views of the target view during training by selecting three frames before the target frame with a 10, 20, and 30 frame difference each and three frames after the target frame with a 10, 20, and 30 frame difference. We trained the network for a total of 200K iterations, where the three stages of their proposed staged training were repeated every 10K steps. The training was done with a single A6000 GPU.

## B.2. Evaluation Details

**Differences of the evaluation strategy in original work of DBARF and FlowCAM and our adopted evaluation strategy.** In the original evaluation strategy of DBARF, given a sequence of frames, DBARF picks an abitrary view, treating it as a target view, and selects nearby views, considering them as context images. Subsequently, pairwise pose estimation and depth estimation are performed between the target image and each of the context images. The estimated values are then fed to GeNeRF for rendering and evaluation. In our evaluation approach, we assume that only two context images and a relative camera pose to the target view are provided. This means that in contrast to the original evaluation setting, where target view is accessible for depth and pose estimation, our evaluation setting does not employ the target view for the model, but only accessible for metric computation.

In the evaluation setting of FlowCAM, FlowCAM first feeds all the frames to single-view PixelNeRF, and they are used for warping by an off-the-shelf optical flow network to output matching confidence that will be used for Weighted Procrustes formulation. Similar to DBARF, the target views are accessed for pose estimation in this evaluation setting. In our evaluation setting, we substitute the estimated poses with ground-truth poses except for the relative pose between $I_1$ and $I_2$, which is their prediction and remains unchanged.

**Fixed Pose and Generalized NeRF.** In Table 4 of the main paper, we conducted additional ablation study and

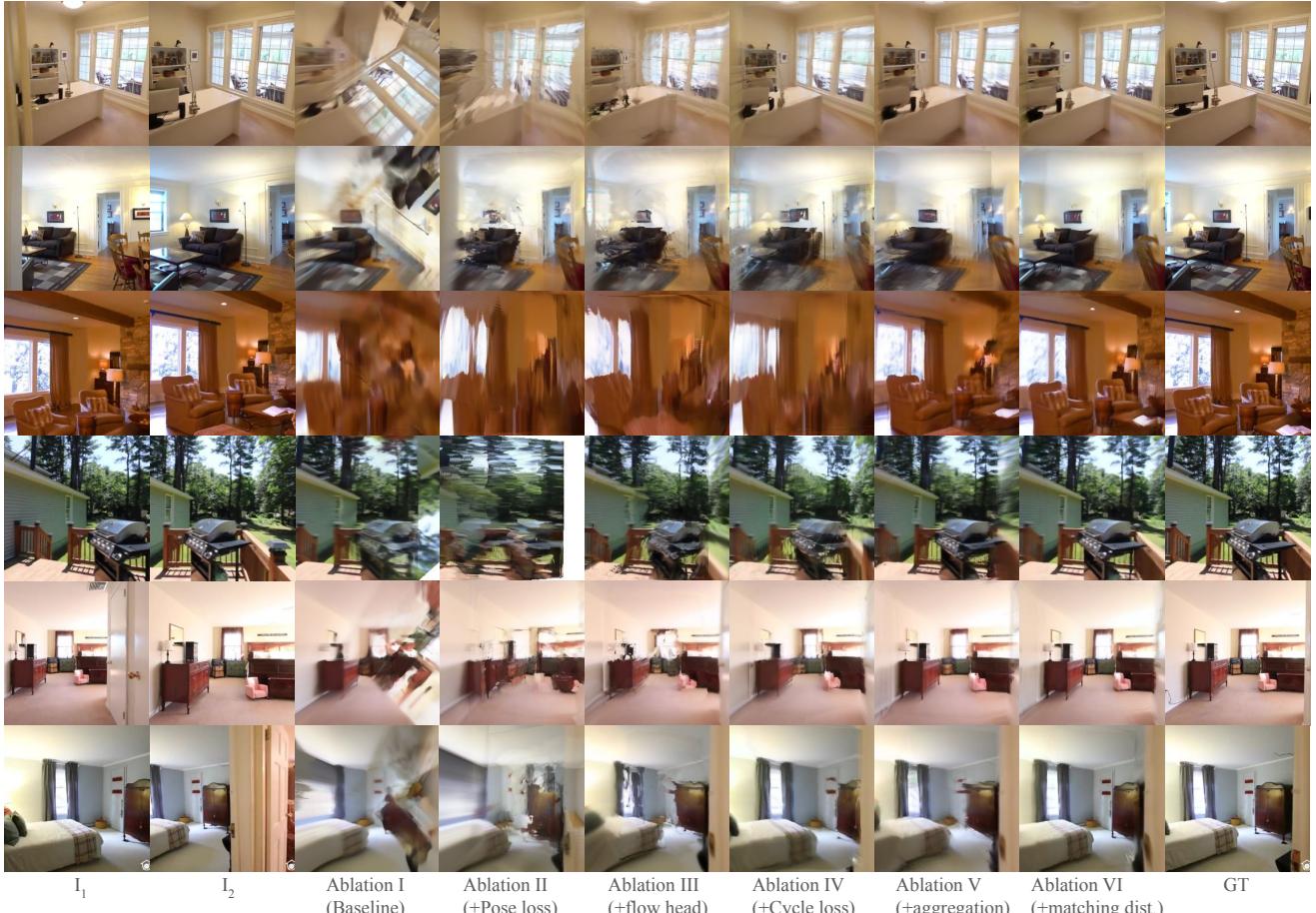| $I_1$ | $I_2$ | Ablation I (Baseline) | Ablation II (+Pose loss) | Ablation III (+flow head) | Ablation IV (+Cycle loss) | Ablation V (+aggregation) | Ablation VI (+matching dist.) | GT |

Figure 8. **Qualitative results for component ablation study.** Consistent with the quantitative results (Table 3 of the main paper), each variant exhibits apparent differences in qualitative comparisons and shows the efficacy of our designed components.
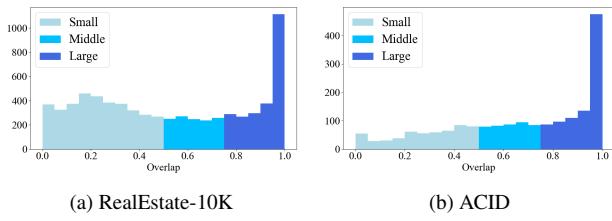


(a) RealEstate-10K  (b) ACID

Figure 9. **Distribution of the data splits.**

analysis. While all other experiments follow the same evaluation protocol, some subtle changes are made for (a). Specifically, as PDC-Net encountered some cases where RANSAC failed to converge due to extremely small overlapping regions with barely any correspondence, we exclude such cases during metric calculations. For a fair comparison, for the other variants, *i.e.*, Rockwell *et al.*+[20] and ours, we also disregard those scenes, which leads to different results compared to our main table.

## C. Data Split Details

We provide the statistics of each overlap-based data split of RealEstate10K and ACID in Fig. 9.

## D. More Discussions

**Pose Supervision.** In this section, we discuss the results we obtained from the experiment that combines direct pose supervision to FlowCAM and DBARF. For FlowCAM, we hypothesize that the disjoint design among each modules, *i.e.,* pose and renderer, and the lack of coherence among the modules results in limited overall benefits. Because off-the-shelf flow model is left frozen for inference, it will inevitably risk finding inaccurate poses when there are only minimal overlapping regions between image pairs, which underscores the effectiveness of our proposed approach that emphasizes the meticulous unification of various tasks to effectively harness their synergistic potential. For DBARF, contrary to initial expectations, adding direct pose supervision results in a decline in their pose estimation perfor-

mance despite our meticulous efforts to optimize the hyperparameters and conduct multiple trials. We hypothesize that the performance degradation could be attributed to their original training strategy, which only assumes image sequences with small viewpoint changes and could have influenced overall performance largely. Another potential reason is that their incorporation of pose estimation and the rendering requires non-trivial implementation considerations to obtain a boosted synergy, or the confidence score produced by the off-the-shelf model might've made a disparity between the updated renderer and the optical flow prediction. Despite the challenges we faced and the hypothesis we made, we leave this exploration as future work since such investigation is beyond the scope of this paper.

## E. More Results

### E.1. Absolute Translation Error with Scales

Table 5 presents the results of translation estimation evaluated with both absolute error in meters and angular difference in degrees. Note that, as mentioned in the main paper, evaluating absolute error requires assessing the models' ability to gauge scales via, *e.g.*, object recognition, since translation scale is theoretically indeterminable in two-view geometry. This could potentially result in erroneous interpretations regarding the models' proficiency in estimating relative camera pose from two views.

### E.2. More Qualitative Results

Fig. 10 shows the correspondences built by our method and the overlapping regions characterized by high confidence scores. As we can see, our method can detect matching points robustly across different scenarios.

Fig. 11 visualizes the epipolar lines with the relative poses estimated by different methods. Visually inspected, our method yields more accurate results especially on challenging cases with small overlap.

Fig. 12 and Fig. 13 present more novel view rendering results of different methods. On both datasets, our method yields outcomes that are sharper and more geometrically accurate.

### E.3. Qualitative Results for Ablation Study

In Fig. 8, we provide qualitative comparisons for each variant introduced for the component ablation study. Consistent with the quantitative results, each variant exhibits apparent differences in qualitative comparisons as well.

| Overlap | Task | Method | RealEstate-10K | | | | | | ACID | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Translation | | | Translation | | | Translation | | | Translation | | |
| | | | Avg(m)↓ | Med(m)↓ | STD(m)↓ | Avg(°)↓ | Med(°)↓ | STD(°)↓ | Avg(m)↓ | Med(m)↓ | STD(m)↓ | Avg(°)↓ | Med(°)↓ | STD(°)↓ |
| Small | Matching | SP+SG [18, 24, 50] | 0.973 | 0.759 | 0.840 | 12.549 | 4.638 | 23.048 | 0.979 | 0.661 | 1.094 | 22.214 | 7.526 | 33.719 |
| | | PDC-Net+ [24, 56] | 0.696 | 0.597 | 0.591 | 6.913 | 2.752 | 15.558 | 0.667 | 0.573 | 0.714 | 15.664 | 4.215 | 29.640 |
| | Pose Estimation | Rockwell *et al.* [46] | 1.692 | 1.459 | 1.119 | 91.455 | 91.499 | 56.872 | 1.576 | 1.057 | 3.557 | 88.421 | 88.958 | 36.212 |
| | | RelPose [67] | - | - | - | - | - | - | - | - | - | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 2.782 | 2.549 | 1.803 | 126.282 | 140.358 | 43.691 | 2.134 | 1.187 | 6.959 | 95.149 | 99.490 | 47.576 |
| | | FlowCAM* [52] | 1.543 | 1.400 | 0.901 | 87.119 | 58.245 | 26.895 | 0.732 | 0.487 | 0.810 | 92.130 | 85.846 | 40.821 |
| | | Ours | **0.532** | **0.353** | **0.642** | **11.862** | **5.344** | **21.080** | **0.378** | **0.171** | **0.533** | **23.689** | **11.289** | **30.391** |
| Medium | Matching | SP+SG [18, 24, 50] | 0.390 | 0.344 | 0.261 | 9.295 | 3.279 | 20.456 | 0.528 | 0.466 | 0.431 | 16.455 | 5.426 | 29.035 |
| | | PDC-Net+ [24, 56] | 0.360 | 0.322 | 0.253 | 6.667 | 2.262 | 18.247 | 0.612 | 0.563 | 0.482 | 14.940 | 4.301 | 27.379 |
| | Pose Estimation | Rockwell *et al.* [46] | 0.842 | 0.705 | 0.581 | 82.478 | 82.920 | 55.094 | 0.713 | 0.554 | 0.649 | 90.555 | 90.799 | 51.469 |
| | | RelPose [67] | - | - | - | - | - | - | - | - | - | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 0.816 | 0.574 | 0.782 | 79.402 | 75.408 | 54.485 | 0.772 | 0.473 | 0.931 | 77.324 | 77.291 | 49.735 |
| | | FlowCAM* [52] | 0.563 | 0.510 | 0.384 | 42.287 | 41.594 | 24.862 | 0.654 | 0.447 | 0.768 | 95.444 | 87.308 | 43.198 |
| | | Ours | **0.203** | **0.150** | **0.178** | **10.187** | **5.749** | **15.801** | **0.324** | **0.133** | **0.615** | **21.401** | **10.656** | **28.243** |
| Large | Matching | SP+SG [18, 24, 50] | 0.612 | 0.665 | 0.202 | 21.415 | 7.190 | 34.044 | 0.619 | 0.641 | 0.260 | 22.018 | 7.309 | 33.775 |
| | | PDC-Net+ [24, 56] | 0.601 | 0.659 | 0.200 | 16.567 | 5.447 | 29.883 | 0.707 | 0.606 | 0.882 | 18.447 | 4.357 | 35.564 |
| | Pose Estimation | Rockwell *et al.* [46] | 0.468 | 0.363 | 0.377 | 91.851 | 88.923 | 57.444 | 0.431 | 0.304 | 0.457 | 86.580 | 87.559 | 50.369 |
| | | RelPose [67] | - | - | - | - | - | - | - | - | - | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 0.217 | 0.098 | 0.318 | 50.094 | 33.959 | 43.659 | **0.281** | **0.111** | 0.488 | 54.523 | 38.829 | 45.453 |
| | | FlowCAM* [52] | 0.179 | 0.107 | 0.197 | 34.472 | 27.791 | 30.615 | 0.778 | 0.442 | 2.789 | 97.392 | 89.359 | 43.777 |
| | | Ours | **0.095** | **0.067** | **0.102** | **15.544** | **7.907** | **24.626** | 0.456 | 0.146 | **0.276** | **22.935** | **10.588** | **30.974** |
| *Avg* | Matching | SP+SG [18, 24, 50] | 0.749 | 0.629 | 0.654 | 14.887 | 5.058 | 27.238 | 0.703 | 0.610 | 0.676 | 20.802 | 6.878 | 32.834 |
| | | PDC-Net+ [24, 56] | 0.696 | 0.0.597 | 0.591 | 10.100 | 3.243 | 22.317 | 0.671 | 0.587 | 0.744 | 16.461 | 4.292 | 31.391 |
| | Pose Estimation | Rockwell *et al.* [46] | 1.145 | 0.821 | 1.022 | 90.115 | 88.648 | 40.948 | 0.833 | 0.500 | 2.041 | 88.433 | 88.961 | 36.197 |
| | | RelPose [67] | - | - | - | - | - | - | - | - | - | - | - | - |
| | Pose-Free NeRF | DBARF [12] | 1.603 | 0.930 | 1.787 | 93.300 | 102.467 | 57.290 | 0.939 | 0.366 | 3.901 | 71.711 | 68.892 | 50.277 |
| | | FlowCAM* [52] | 0.916 | 0.647 | 0.913 | 50.659 | 46.281 | 52.321 | 1.665 | 1.538 | 0.748 | 95.405 | 88.133 | 42.849 |
| | | Ours | **0.332** | **0.177** | **0.506** | **12.766** | **7.534** | **15.510** | **0.404** | **0.150** | **0.197** | **22.809** | **14.502** | **21.572** |

Table 5. Translation estimation performance evaluated with both absolute error (in meters) and angular error (in degrees). Note that since translation scale is theoretically indeterminable in two-view geometry, evaluating absolute error requires assessing the models' ability to gauge scales via, *e.g.*, object recognition. This could potentially result in erroneous interpretations regarding the models' proficiency in estimating relative pose from two views. *: FlowCAM [52] results have been updated to rectify an error in the numerical values originally presented.
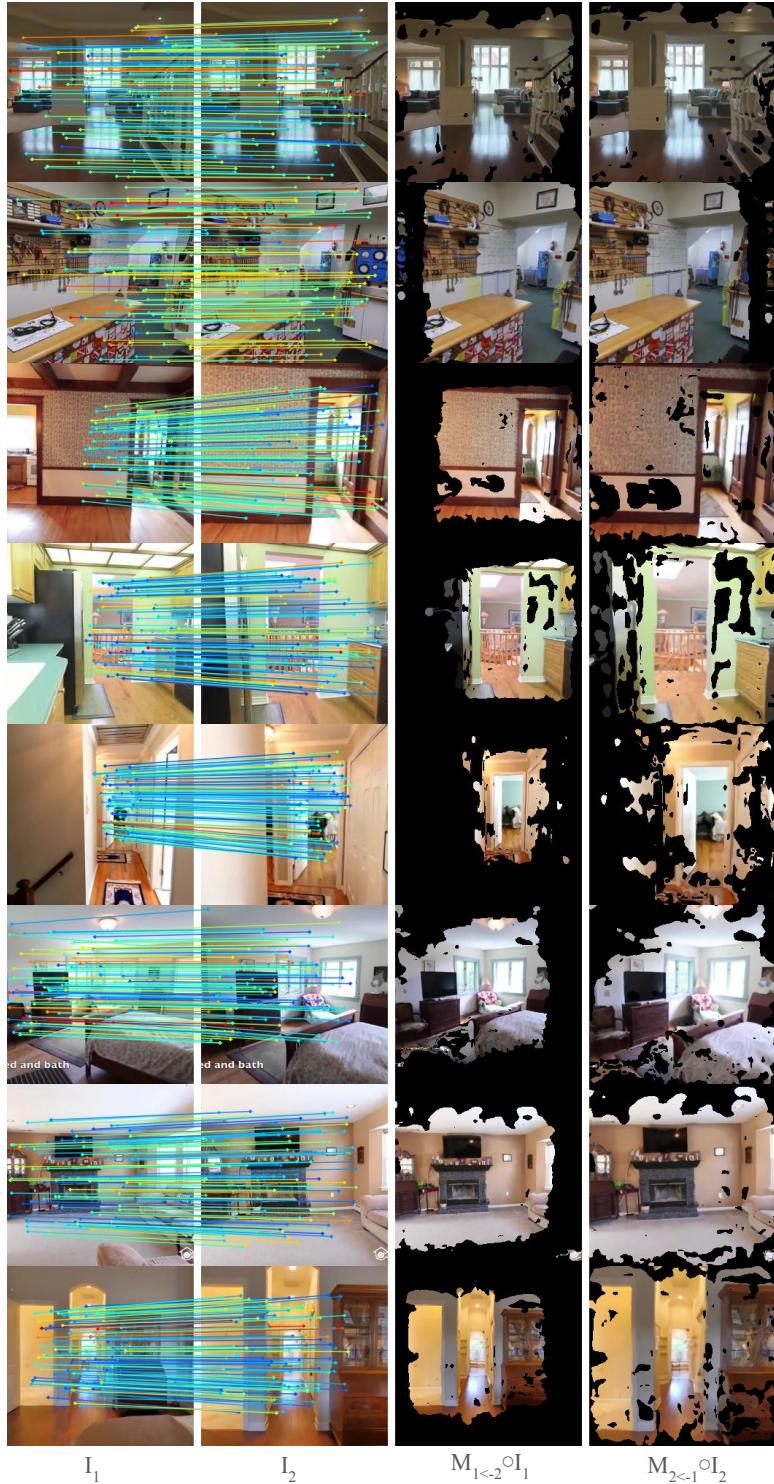
Figure 10. **Randomly selected correspondences and confident regions.** For each pair of images, we visualize a set of randomly selected correspondences (left), and from the complete set of correspondences, and those with confidence score of higher than a threshold $\tau$ are shown as visible (right).
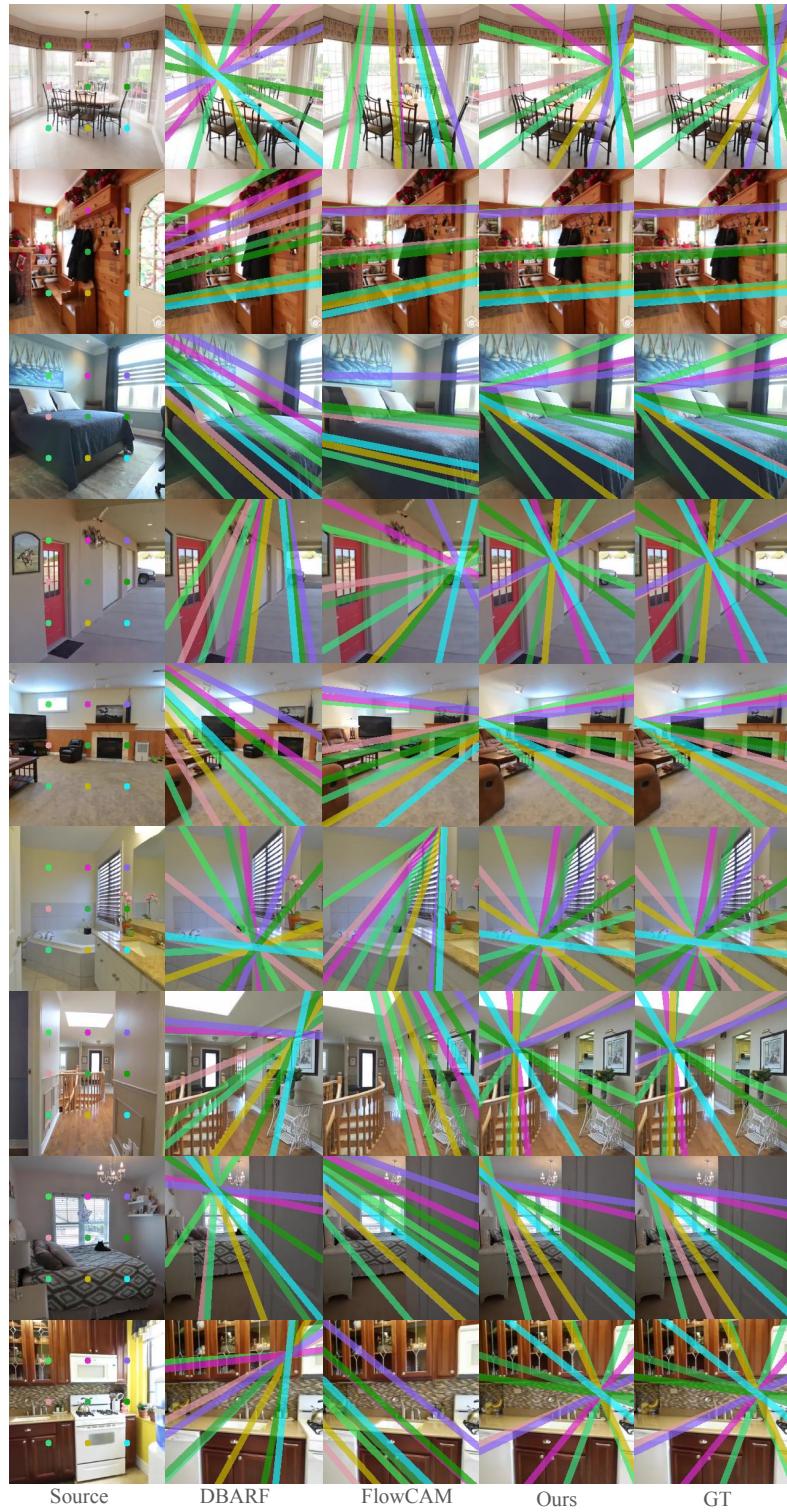
| Source | DBARF | FlowCAM | Ours | GT |

Figure 11. **Comparisons of visualized epipolar lines.**

(a) $I_1$　　(b) $I_2$　　(c) DBARF　　(d) FlowCAM　　(e) Ours　　(f) Ground Truth

Figure 12. **Qualitative comparison on RealEstate10K.**

(a) $I_1$     (b) $I_2$     (c) DBARF     (d) FlowCAM     (e) Ours     (f) Ground Truth

Figure 13. **Qualitative comparison on ACID.**

# References

[1] Eduardo Arnold, Jamie Wynn, Sara Vicente, Guillermo Garcia-Hernando, Aron Monszpart, Victor Prisacariu, Daniyar Turmukhambetov, and Eric Brachmann. Map-free visual relocalization: Metric pose relative to a single image. In *European Conference on Computer Vision*, pages 690–708. Springer, 2022. 3

[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. 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. 2

[4] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 2

[5] Axel Barroso-Laguna, Eric Brachmann, Victor Adrian Prisacariu, Gabriel J Brostow, and Daniyar Turmukhambetov. Two-view geometry scoring without correspondences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8979–8989, 2023. 2

[6] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *Computer Vision–ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7-13, 2006. Proceedings, Part I 9*, pages 404–417. Springer, 2006. 2

[7] Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. Nope-nerf: Optimising neural radiance field with no pose prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4160–4169, 2023. 2

[8] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. Dsac-differentiable ransac for camera localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6684–6692, 2017. 2

[9] Ruojin Cai, Bharath Hariharan, Noah Snavely, and Hadar Averbuch-Elor. Extreme rotation estimation using dense correlation volumes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14566–14575, 2021. 2, 4, 6

[10] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. 2, 3

[11] Kefan Chen, Noah Snavely, and Ameesh Makadia. Wide-baseline relative camera pose estimation with directional learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3258–3268, 2021. 2

[12] Yu Chen and Gim Hee Lee. Dbarf: Deep bundle-adjusting generalizable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24–34, 2023. 1, 2, 3, 5, 6, 7, 8, 10, 13

[13] Yuedong Chen, Haofei Xu, Qianyi Wu, Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Explicit correspondence matching for generalizable neural radiance fields. *arXiv preprint arXiv:2304.12294*, 2023. 2, 3

[14] Seokju Cho, Sunghwan Hong, Sangryul Jeon, Yunsung Lee, Kwanghoon Sohn, and Seungryong Kim. Cats: Cost aggregation transformers for visual correspondence. *Advances in Neural Information Processing Systems*, 34:9011–9023, 2021. 3

[15] Seokju Cho, Sunghwan Hong, and Seungryong Kim. Cats++: Boosting cost aggregation with convolutions and transformers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 3

[16] Seokju Cho, Heeseong Shin, Sunghwan Hong, Seungjun An, Seungjun Lee, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation. *arXiv preprint arXiv:2303.11797*, 2023. 3

[17] Christopher Choy, Wei Dong, and Vladlen Koltun. Deep global registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2514–2523, 2020. 2, 9

[18] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 224–236, 2018. 2, 5, 7, 13

[19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2

[20] Yilun Du, Cameron Smith, Ayush Tewari, and Vincent Sitzmann. Learning to render novel views from wide-baseline stereo pairs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4970–4980, 2023. 1, 2, 3, 4, 5, 7, 8, 9, 10, 11

[21] Hugh Durrant-Whyte and Tim Bailey. Simultaneous localization and mapping: part i. *IEEE robotics & automation magazine*, 13(2):99–110, 2006. 2

[22] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-net: A trainable cnn for joint detection and description of local features. *arXiv preprint arXiv:1905.03561*, 2019. 2

[23] Sovann En, Alexis Lechervy, and Frédéric Jurie. Rpnet: An end-to-end network for relative camera pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 2

[24] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 5, 7, 13

[25] Richard I Hartley. In defense of the eight-point algorithm. *IEEE Transactions on pattern analysis and machine intelligence*, 19(6):580–593, 1997. 2

[26] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2009. 5

[27] 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, 10

[28] Sunghwan Hong and Seungryong Kim. Deep matching prior: Test-time optimization for dense correspondence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9907–9917, 2021. 3

[29] Sunghwan Hong, Seokju Cho, Jisu Nam, Stephen Lin, and Seungryong Kim. Cost aggregation with 4d convolutional swin transformer for few-shot segmentation. In *European Conference on Computer Vision*, pages 108–126. Springer, 2022.

[30] Sunghwan Hong, Jisu Nam, Seokju Cho, Susung Hong, Sangryul Jeon, Dongbo Min, and Seungryong Kim. Neural matching fields: Implicit representation of matching fields for visual correspondence. *Advances in Neural Information Processing Systems*, 35:13512–13526, 2022. 3

[31] Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Flowformer: A transformer architecture for optical flow. *arXiv preprint arXiv:2203.16194*, 2022. 3

[32] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5846–5854, 2021. 2, 5

[33] Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. Geonerf: Generalizing nerf with geometry priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18365–18375, 2022. 2, 3, 4

[34] Zakaria Laskar, Iaroslav Melekhov, Surya Kalia, and Juho Kannala. Camera relocalization by computing pairwise relative poses using convolutional neural network. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 929–938, 2017. 2

[35] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021. 2, 5

[36] Andrew Liu, Richard Tucker, Varun Jampani, Ameesh Makadia, Noah Snavely, and Angjoo Kanazawa. Infinite nature: Perpetual view generation of natural scenes from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14458–14467, 2021. 2, 5

[37] H Christopher Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293 (5828):133–135, 1981. 2

[38] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5

[39] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004. 2

[40] Simon Meister, Junhwa Hur, and Stefan Roth. Unflow: Unsupervised learning of optical flow with a bidirectional census loss. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 4

[41] Iaroslav Melekhov, Juha Ylioinas, Juho Kannala, and Esa Rahtu. Relative camera pose estimation using convolutional neural networks. In *Advanced Concepts for Intelligent Vision Systems: 18th International Conference, ACIVS 2017, Antwerp, Belgium, September 18-21, 2017, Proceedings 18*, pages 675–687. Springer, 2017. 2, 6

[42] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2

[43] David Nistér. An efficient solution to the five-point relative pose problem. *IEEE transactions on pattern analysis and machine intelligence*, 26(6):756–770, 2004. 2, 5, 6, 7

[44] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5

[45] 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

[46] Chris Rockwell, Justin Johnson, and David F Fouhey. The 8-point algorithm as an inductive bias for relative pose prediction by vits. In *2022 International Conference on 3D Vision (3DV)*, pages 1–11. IEEE, 2022. 1, 2, 4, 5, 6, 7, 8, 9, 13

[47] Chris Rockwell, Nilesh Kulkarni, Linyi Jin, Jeong Joon Park, Justin Johnson, and David F. Fouhey. Far: Flexible, accurate and robust 6dof relative camera pose estimation. 2024. 7

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

[49] Seyed Sadegh Mohseni Salehi, Shadab Khan, Deniz Erdogmus, and Ali Gholipour. Real-time deep pose estimation with geodesic loss for image-to-template rigid registration. *IEEE transactions on medical imaging*, 38(2):470–481, 2018. 5

[50] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020. 1, 5, 7, 13

[51] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 2

[52] Cameron Smith, Yilun Du, Ayush Tewari, and Vincent Sitzmann. Flowcam: Training generalizable 3d radiance fields without camera poses via pixel-aligned scene flow. *arXiv preprint arXiv:2306.00180*, 2023. 1, 2, 5, 6, 7, 8, 9, 13

[53] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Generalizable patch-based neural rendering. In *European Conference on Computer Vision*, pages 156–174. Springer, 2022. 4

[54] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8922–8931, 2021. 3

[55] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 2, 9

[56] Prune Truong, Martin Danelljan, Radu Timofte, and Luc Van Gool. Pdc-net+: Enhanced probabilistic dense correspondence network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 1, 5, 6, 7, 13

[57] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. Sparf: Neural radiance fields from sparse and noisy poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4190–4200, 2023. 2

[58] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 2, 3, 10

[59] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 4

[60] Tong Wei, Yash Patel, Alexander Shekhovtsov, Jiri Matas, and Daniel Barath. Generalized differentiable ransac. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17649–17660, 2023. 2

[61] Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989. 5

[62] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8121–8130, 2022. 3

[63] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 3

[64] Kwang Moo Yi, Eduard Trulls, Vincent Lepetit, and Pascal Fua. Lift: Learned invariant feature transform. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VI 14*, pages 467–483. Springer, 2016. 2

[65] Kwang Moo Yi, Eduard Trulls, Yuki Ono, Vincent Lepetit, Mathieu Salzmann, and Pascal Fua. Learning to find good correspondences. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2666–2674, 2018. 2

[66] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 1, 2, 3, 5, 7

[67] Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Relpose: Predicting probabilistic relative rotation for single objects in the wild. In *European Conference on Computer Vision*, pages 592–611. Springer, 2022. 1, 5, 6, 7, 10, 13

[68] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 2, 5

[69] Yi Zhou, Connelly Barnes, Lu Jingwan, Yang Jimei, and Li Hao. On the continuity of rotation representations in neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 4

[70] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 5

[71] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022. 2