Revisit Self-supervised Depth Estimation with Local Structure-from-Motion

Shengjie Zhu and Xiaoming Liu

Department of Computer Science and Engineering, Michigan State University, East Lansing, MI, 48824 zhusheng@msu.edu, liuxm@cse.msu.edu

Abstract. Both self-supervised depth estimation and Structure-from-Motion (SfM) recover scene depth from RGB videos. Despite sharing a similar objective, the two approaches are disconnected. Prior works of self-supervision backpropagate losses defined within immediate neighboring frames. Instead of learning-through-loss, this work proposes an alternative scheme by performing local SfM. First, with calibrated RGB or RGB-D images, we employ a depth and correspondence estimator to infer depthmaps and pair-wise correspondence maps. Then, a novel bundle-RANSAC-adjustment algorithm jointly optimizes camera poses and one depth adjustment for each depthmap. Finally, we fix camera poses and employ a NeRF, however, without a neural network, for dense triangulation and geometric verification. Poses, depth adjustments, and triangulated sparse depths are our outputs. For the first time, we show selfsupervision within 5 frames already benefits SoTA supervised depth and correspondence models. Despite self-supervision, our pose algorithm has certified global optimality, outperforming optimization-based, learningbased, and NeRF-based prior arts. The project page is held in the link.

Keywords: Self-supervision · Depth · Pose · Structure-from-Motion

1 Introduction

Monocular depth estimation [17,29] infers depthmap from a single image. It is an essential vision task with applications in AR/VR [41], autonomous driving [19], and 3D reconstruction [8]. Most methods [5,43,45,69] supervise the model with groundtruth collected from stereo cameras [71] or LiDAR [19]. Recently, self-supervised depth [20,21,76] has drawn significant attention due to its potential to scale up depth learning from massive unlabeled RGB videos.

Classic SfM methods [1,12,47,48,52,65] also reconstruct scene depth from unlabled RGB videos. Despite its relevance, SfM is rarely applied to self-supervised depth learning. We outline two potential reasons. First, SfM is an off-the-shelf algorithm unrelated to the depth estimator. Scale ambiguity renders SfM poses and depths at different scales compared to depth models. Second, self-supervision has a well-defined training scheme to work with universal unlabeled videos. It backpropagates through photometric loss computed within immediate neighboring frames, *e.g.*, red trajectory in Fig. 1. In contrast, SfM is more selective to



Fig. 1: Revisit Self-supervision with Local SfM. The work proposes alternating the learning-through-loss with a local SfM pipeline for self-supervised depth estimation. We summarize our differences. On self-supervision: (1) Instead of using naive two-view camera poses, we propose a Bundle-RANSAC-Adjustment pose optimization algorithm with multi-view constraints. (2) Instead of backpropagating through a loss, we produce a sparse point cloud with explicit triangulation and geometric verification. The point cloud serves as either output or pseudo-groundtruth for self-supervision. On SfM: (1) Our local SfM is adapted to use estimated monocular depthmaps and automatically resolve their scale inconsistency between pairs of images. (2) We maintain accuracy under significant sparse view variations, *e.g.*, red trajectories. We generalize SfM to as few as 5 frames, similar to the number of images used to define self-supervision loss.

input videos. It requires images of diverse view variations (green trajectory in Fig. 1), being inaccurate and unstable when applied to a small frame window.

This work connects self-supervision with SfM. We replace the self-supervision loss with a complete SfM pipeline that maintains robustness to a local window. Shown in Fig. 2, with N frames as input, our algorithm outputs N - 1 camera poses, N - 1 depth adjustments, and the sparse triangulated point cloud. In initialization, N monocular depthmaps and $N \times (N - 1)$ pairwise correspondence maps are inferred. Next, we propose a Bundle-RANSAC-Adjustment pose estimation algorithm that retains accuracy for second-long videos. The algorithm utilizes the 3D priors from mono-depthmap to compensate for the deficient camera views. Correspondingly, we optimize N - 1 depth adjustments to alleviate the depth scale ambiguity by temporally aligning to the root frame depth.

The Bundle-RANSAC-Adjustment extends two-view RANSAC with multiview bundle-adjustment (BA). The algorithm has quadratic complexity, designed for parallel GPU computation. We <u>RAN</u>domly <u>SA</u>mple and hypothesize a set of normalized poses. In <u>C</u>onsensus checking, we apply BA to evaluate a robust inlier-counting scoring function over multi-view images. Camera scales and depth adjustments are determined during BA to maximize the scoring function.

Next, we freeze the optimized poses and employ a Radiance Field (RF), *i.e.*, a NeRFF [39] without a neural network, for triangulation. We optimize RF to achieve multi-view depthmap and correspondence consistency within a shared 3D frustum volume. For outputs, we apply geometric verification to extract multi-view consistent point cloud, *i.e.*, a sparse root depthmap.



Fig. 2: Local Structure-from-Motion. With N neighboring frames, we extract monocular depthmaps and pairwise dense correspondence maps with methods, *e.g.*, ZoeDepth [5] and PDC-Net [57]. Next, skipping the root frame, we optimize the rest N-1 camera poses and depth adjustments. The depth adjustments render input depthmaps **temporally consistent**. Fixing poses and adjustments, we use the Radiance Field (RF) for triangulation and output a geometrically verified sparse root depthmap. Our local SfM applies **self-supervision** with only 5 RGB frames. Yet, our sparse output already outperforms the input supervised depth with SoTA performance.

Fig. 1 contrasts our method with prior self-supervised depth and SfM methods. To our best knowledge, there has not been prior work showing geometrybased self-supervised depth benefits supervised models. However, self-supervision is supposed to augment supervised models with unlabeled data. In Fig. 2, our unique pipeline gives the **first** evident results, that self-supervision with **as few as** 5 frames already benefits supervised models.

On top of depths, our multi-view RANSAC pose has certified global optimality under a robust scoring function. It outperforms prior arts in optimizationbased [48, 77], learning-based [55, 61], and NeRF-based [58] pose algorithms.

Beyond pose and depth, our method has diverse applications. The depth adjustments from our method provide empirically consistent depthmaps, important for AR image compositing. When given RGB-D inputs, our method enables self-supervised correspondence estimation. Our accurate pose estimation gives improved projective correspondence than the SoTA supervised correspondence input. An example is in Fig. 9. We summarize our contributions as:

- We propose a novel local SfM algorithm with Bundle-RANSAC-Adjustment.
- We show the first evident result that self-supervised depth with as few as 5 frames already benefit SoTA supervised models.
- We achieve SoTA sparse-view pose estimation performance.
- We enable self-supervised temporally consistent depthmaps.
- We enable self-supervised correspondence estimation with 5 RGB-D frames.

2 Related Works

Structure-from-Motion. SfM is a comprehensive task [48,65]. A typical pipeline is, correspondence extraction [7,36,59], two-view initialization [3,30,32], triangulation [31,42], and local & global bundle-adjustment [48,65]. Classic methods

require diverse view variations for accurate reconstruction. Our method compensates SfM on scarse camera views via introducing deep depth estimator. Further, we suggest SfM itself is a self-supervised learning pipeline, as in Fig. 1. Finally, our SfM is not up-to-scale and shares the metric space as the input depthmap.

Sparse Multi-view Pose Estimation. Estimating poses from sparse frames is crucial for self-supervision [11,20,46,70,74], video depth estimation [22,55,60,77], and sparse-view NeRF [14,27,34,40,58]. Camera poses are estimated either by learning [11,20,22,55], optimization [73,77] or together with NeRF [34,58]. We propose an additional multi-view RANSAC pipeline with improved accuracy.

Self-supervised Depth and Correspondence Estimation. Multiple works improve self-supervised depth in different ways, including learning loss [20,44,63], architecture [23,75], camera pose [6,10,38,73], joint with semantics segmentation [76], and using large-scale data [53,66]. Recently, [53] shows self-supervision only performs on-par with supervised models under substantially more data. [66] shows the benefit of self-supervision via exploiting non-geometry monocular semantic consistency. Our method shows the first evident results where self-supervision benefits supervised models with only 5 consecutive frames.

Consistent Depth Estimation. AR applications necessitate temporally consistent depthmaps, *i.e.*, depthmaps from different temporal frames reside in the same 3D space. Recent works [37,72] align depthmap according to the poses and points from the off-the-shelf COLMAP algorithm. Our method seamlessly integrates SfM with monocular depthmaps, outputting consistent depth and poses. **Test Time Refinement (TTR).** TTR aims to improve self-supervised / supervised depth estimators in testing time with RGB video [9, 10, 28, 50, 64]. Methods [25, 56] rely on off-the-shelf algorithms for pseudo depth and pose labels. Recently, [25] first shows TTR improves supervised models. TTR is our downstream application, which details strategies for utilizing noisy pseudo-labels.

3 Methodology

Our method runs sequentially. From N calibrated images \mathcal{I} , we extract N monocular depthmaps \mathcal{D} and $N \times (N-1)$ pair-wise dense correspondence \mathcal{C} . We split the N images into one root frame \mathbf{I}_o in the center of the N-frame window where $o = \lfloor \frac{N+1}{2} \rfloor$, and N-1 support frames \mathbf{I}_i , where $i \in \mathbb{N}^+ = [1, N] \setminus \{o\}$. In Sec. 3.1, after setting the root frame as identity pose, we use Bundle-RANSAC-Adjustment to optimize N-1 poses \mathcal{P} and N-1 depth adjustments \mathcal{R} . Next, in Sec. 3.2, we apply triangulation by optimizing a frustum Radiance Field (RF) \mathbf{V} , *i.e.*, a NeRF without network. Finally, in Sec. 3.3, we apply geometric verification by rendering multi-view consistent 3D points from RF. An overview is in Fig. 3.

3.1 Bundle-RANSAC-Adjustment Pose Estimation

We generalize two-view RANSAC with multi-view constraints through Bundle-Adjustment. Sec. 3.1.1 describes our pipeline. In Sec. 3.1.2, we propose Hough transform to accelerate computation. We discuss the time complexity in Sec. 3.1.3.



Fig. 3: Algorithm Overview. After extracting monodepths and correspondence maps from inputs: (a) We apply Bundle-RANSAC-Adjustment to optimize N-1 camera poses \mathcal{P} and N-1 depth adjustments \mathcal{R} . (b) We fix poses and depth adjustments and optimize a frustum Radiance Field (RF) for triangulation. (c) We apply geometric verification to extract multi-view consistent 3D points via rendering with RF. We further detail step (a) in Figs 4, 5, and 6, and steps (b) and (c) in Fig. 7.

3.1.1 Optimization Pipeline

<u>RANdom SAmple.</u> We use five-point algorithm [32] as the minimal solver. We execute it between root and each support frame, extracting a pool of $(N-1) \times K$ normalized poses (*i.e.*, pose of unit translation), $\overline{\mathcal{Q}} = \{\overline{\mathbf{P}}_i^k \mid i \in \mathbb{N}^+, k \in [1, K]\}$, where $\overline{\mathbf{P}}_i^k \in \mathbb{R}^{3 \times 4}$. The K is the number of normalized poses extracted per frame. We term a set of N-1 normalized poses as a group $\overline{\mathcal{P}} \in \mathbb{R}^{(N-1) \times 3 \times 4}$. Two-view RANSAC enumerates over single normalized pose $\overline{\mathbf{P}}$. Our multi-view algorithm hence enumerates over normalized pose group $\overline{\mathcal{P}}$. We initialize the optimal group $\overline{\mathcal{P}}^*$ as the top candidate from K poses of $\overline{\mathcal{Q}}$ for each frame. See examples in Fig. 4. **Bundle-Adjustment <u>C</u>onsensus.** While computing consensus counts, the camera scales S and depth adjustments \mathcal{R} are automatically determined with bundle-adjustment to maximize a robust scoring function:

$$\rho_i = \phi(\overline{\mathcal{P}}) = \max_{\mathcal{S},\mathcal{R}} f(\mathcal{S},\mathcal{R} \mid \overline{\mathcal{P}},\mathcal{D},\mathcal{C}).$$
(1)

Search for Optimal Group. Our multi-view RANSAC has a significantly larger solution space than two-view RANSAC. With N view inputs, we determine the optimal group out of K^{N-1} combinations. Hence, we iteratively search for the optimal group with a greedy strategy. For each epoch, we ablate (N-1)(K-1) additional pose groups:

$$\overline{\mathcal{P}}_{i}^{k} = \overline{\mathcal{P}}_{i}^{*} \setminus \{\overline{\mathbf{P}}_{i}^{*}\} \cup \{\overline{\mathbf{P}}_{i}^{k}\}, \qquad (2)$$

where $i \in \mathbb{N}^+$ and $k \in [1, K]$. Combine Eq. (2) and Fig. 4, taking frame *i* as an example, we replace the optimal pose $\overline{\mathbf{P}}_i^*$ by its K - 1 other candidates $\overline{\mathbf{P}}_i^k$, generating K - 1 groups. For *N* frames, we have (N - 1)(K - 1) + 1 groups. We apply bundle-adjustment to each group to evaluate Eq. (1). As shown in Fig. 3 and Fig. 4, we select the normalized pose together with its optimized scales and depth adjustments that maximize the scores as the output,

$$\mathcal{P}_i^* = b(\overline{\mathcal{P}}_i^*, \mathcal{S}_i^*), \ \mathcal{R}_i^* = \mathcal{R}_i^k, \ \text{where} \ k = \arg\max\{\rho_i^k\}, \ \overline{\mathcal{P}}_i^* = \overline{\mathcal{P}}_i^k, \ \mathcal{S}_i^* = \mathcal{S}_i^k, \ (3)$$



Fig. 4: Pose Optimization Pipeline. We show a sample execution when N=3 and K=3. We initialize normalized pose candidates pool $\overline{\mathcal{Q}}$. Optimal group $\overline{\mathcal{P}}^*$ is set to top candidates within $\overline{\mathcal{Q}}$. In each epoch, Eq. (2) ablates pose group $\overline{\mathcal{P}}^k_i$. Each group is scored with Eq. (1) via BA with Hough Transform, detailed in Sec. 3.1.2. The optimal group with the highest score is updated with Eq. (3). Termination occurs when the maximum score stabilizes. We maintain quadratic complexity by avoiding repetitive computation after the first epoch, shown with the Comp. Graph, detailed in Sec. 3.1.3.

where $b(\cdot)$ combines normalized poses with scales. Fig. 2 third column plots an adjusted temporal consistent depthmap after applying \mathcal{R}^* . In Fig. 4, the algorithm terminates when the maximum score stops increasing.

Scoring Function. Similar to other RANSAC methods, we adopt robust inliercounting based scoring functions. Expand Eq. (1) for a specific group $\overline{\mathcal{P}}$:

$$\phi(\overline{\mathcal{P}}) = \sum_{i,i\neq j} \sum_{j} f_{i,j}(s_i, s_j, r_i, r_j \mid \overline{\mathbf{P}}_i, \overline{\mathbf{P}}_j, \mathbf{D}_i, \mathbf{D}_j, \mathbf{C}_{i,j}),$$
(4)

where i, j are frame index. We set per-frame camera scale, depth, depth adjustment, and correspondence as $s \in S$, $\mathbf{D} \in \mathcal{D}$, $r \in \mathcal{R}$, and $\mathbf{C} \in \mathcal{C}$. The scoring function $f_{i,j}(\cdot)$ has various forms. First, we describe a 2D scoring function:

$$f_{i,j}^{2\mathrm{D}}(\cdot) = \sum_{m} \mathbf{1} \left(\| \pi(s_i, s_j, r_i \mid \overline{\mathbf{P}}_i, \overline{\mathbf{P}}_j, d_i^m) - \mathbf{c}_{i,j}^m \|_2 < \lambda^{2\mathrm{D}} \right),$$
(5)

where $m \in [1, M]$ indexes sampled pixels per frame pair. $f_{i,j}^{2D}(\cdot)$ measures the inlier count between depth projected correspondence and input correspondence. $\pi(\cdot)$ is projection process. Intrinsic is skipped. d and \mathbf{c} are depth and correspondence sampled from \mathbf{D} and \mathbf{C} . An example is in Fig. 5. The $\mathbf{1}(\cdot)$ is the indicator function. The projected pixel is an inlier if it resides within the circle of radius λ^{2D} and center at correspondence $\mathbf{c}_{i,j}^m$ (denoted as \mathbf{p}_j in Fig. 5). $\mathbf{c}_{i,j}^m$ is sampled from correspondence map $\mathbf{C}_{i,j}$. Second, we introduce a 3D scoring function:

$$f_{i,j}^{3D}(\cdot) = \sum_{m} \mathbf{1} \left(\| \pi^{-1}(s_i \mid \overline{\mathbf{P}}_i, r_i, d_i^m) - \pi^{-1}(s_j \mid \overline{\mathbf{P}}_j, r_j, d_j^m) \|_2 < \lambda^{3D} \right).$$
(6)

Depth pair d_i and d_j is determined by correspondence. Unlike the 2D one, the 3D function fixes depth adjustment r. Function $\pi^{-1}(\cdot)$ back-projects 3D point.



Fig. 5: Hough Transform between Two Normalized Poses. With fixed normalized poses, there exists three variables, scales $s_i \& s_j$ of $\overline{\mathbf{P}}_i \& \overline{\mathbf{P}}_j$ and adjustment r_i . Pixel \mathbf{p}_i and \mathbf{p}_j are corresponded. Ablating pose scales maps pixel \mathbf{p}_i to a set of epipolar lines $\{\mathbf{l}_i\}$, however, bounded by Red and Green at infinite scales. We have three observations. First, with fixed normalized poses, epipolar lines \mathbf{l}_i have limited possibilities. Second, scale s and depth adjustment d are equivalent, both adjusting projection on epipolar line. Third, per epipolar line, to be an inlier, the projection has to reside within the line-circle intersection, between $\mathbf{p}_{\pi}^{\text{st}}$ and $\mathbf{p}_{\pi}^{\text{ed}}$. The observations motivate us to discretize the solution space to a 2D matrix, *i.e.*, Hough Transform. Right figure plots an example transformation $\mathbf{H}_{i,j}^m$ from frame i to j on the *m*th pixel \mathbf{p}_i .

3.1.2 Hough Transform Acceleration

Maximizing Eq. (1) for each pose group is computationally prohibitive, as shown in Fig. 4. We propose Hough Transform for acceleration. We use Eq. (5), the 2D function $f^{2D}(\cdot)$ as an example for illustration. See our motivation in Fig. 5. **Hough Transform.** The relative pose between $\overline{\mathbf{P}}_i$ and $\overline{\mathbf{P}}_i$ is defined as:

$$\mathbf{P}_{i,j} = \mathbf{P}_j \mathbf{P}_i^{-1} = \left[\mathbf{R}_{i,j} \ s_{i,j} \overline{\mathbf{t}}_{i,j} \right] = \left[\mathbf{R}_j \mathbf{R}_i^{-1} - s_i \mathbf{R}_j \mathbf{R}_i^{-1} \overline{\mathbf{t}}_i + s_j \overline{\mathbf{t}}_j \right], \tag{7}$$

where $\mathbf{R}, \overline{\mathbf{t}}$, and s are rotation, normalized translation and pose scale. From Eq. (7) and Fig. 5, $\overline{\mathbf{t}}_{i,j}$ is controlled by the scale s_i and s_j , and thus we have:

$$\lim_{s_i \to +\inf} \overline{\mathbf{t}}_{i,j} = -\mathbf{R}_j \mathbf{R}_i^{-1} \overline{\mathbf{t}}_i, \quad \lim_{s_j \to +\inf} \overline{\mathbf{t}}_{i,j} = \overline{\mathbf{t}}_j.$$
(8)

For a pixel \mathbf{p}_i on frame *i*, its corresponding epipolar line \mathbf{l}_i on frame *j* is:

$$\mathbf{l}_{i} = \mathbf{K}^{-\mathsf{T}} [\bar{\mathbf{t}}_{i,j}]_{\mathsf{X}} \mathbf{R}_{i,j} \mathbf{K}^{-1} \mathbf{p}_{i}.$$

$$\tag{9}$$

Eq. (8) and Eq. (9) suggest the epipolar line has limited possibilities. Operation $[\cdot]_{\times}$ is the cross product in matrix form. Further, as the depth re-projected pixel \mathbf{p}_{π} of \mathbf{p}_i always locate on the epipolar line \mathbf{l}_i [24], we have:

$$\mathbf{l}_{i}^{\mathsf{T}}\mathbf{p}_{\pi} = 0, \ \mathbf{p}_{\pi} = \pi(s_{i}, s_{j}, r_{i} \mid \overline{\mathbf{P}}_{i}, \overline{\mathbf{P}}_{j}, d_{i}).$$
(10)

To be an inlier of the scoring function $f^{2D}(\cdot)$, we have:

$$\|\mathbf{p}_{\pi} - \mathbf{p}_{j}\|_{2} \leqslant \lambda^{2\mathrm{D}}.$$
(11)

Combining Eq. (10), Eq. (11) and Fig. 5, to be an inlier, the projected pixel \mathbf{p}_{π} has to reside within the line segment, with two end-points computed by the linecircle intersection. The circle centers at corresponded pixel \mathbf{p}_j on frame j with a radius λ^{2D} . We denote the two end-points $\mathbf{p}_{\pi}^{\text{st}}$ and $\mathbf{p}_{\pi}^{\text{ed}}$. We add their calculation in Supp. Function $J(\cdot)$ follows [77] Supp. Eq. (4), which maps a projected pixel \mathbf{p}_{π} and adjusted depth $r_i d_i$ to camera scale $s_{i,j}$ as: $s_{i,j} = J(\overline{\mathbf{P}}_{i,j}, r_i d_i, \mathbf{p}_{\pi})$.

Corollary 1. A pixel is an inlier iff:

$$J(\overline{\mathbf{P}}_{i,j}, r_i d_i, \mathbf{p}_{\pi}^{st}) \leqslant s_{i,j} \leqslant J(\overline{\mathbf{P}}_{i,j}, r_i d_i, \mathbf{p}_{\pi}^{ed}).$$
(12)

Corollary 2. Scale and depth are equivalent as;

$$s_{i,j} = J(\overline{\mathbf{P}}_{i,j}, r_i d_i, \mathbf{p}_{\pi}) = r_i \cdot J(\overline{\mathbf{P}}_{i,j}, d_i, \mathbf{p}_{\pi}).$$
(13)

See Fig. 5 and proof in Supp. Combine Eqs. (12) and (13),

$$J(\overline{\mathbf{P}}_{i,j}, d_i, \mathbf{p}_{\pi}^{\mathrm{st}}) \leqslant \frac{s_{i,j}}{r_i} \leqslant J(\overline{\mathbf{P}}_{i,j}, d_i, \mathbf{p}_{\pi}^{\mathrm{ed}}).$$
(14)

Set $g(\cdot)$ maps the variables under optimization to intermediate term $\frac{s_{i,j}}{r_i}$.

$$J(\overline{\mathbf{P}}_{i,j}, d_i, \mathbf{p}_{\pi}^{\mathrm{st}}) \leq g(r_i, s_i, s_j \mid \overline{\mathbf{P}}_i, \overline{\mathbf{P}}_j) \leq J(\overline{\mathbf{P}}_{i,j}, d_i, \mathbf{p}_{\pi}^{\mathrm{ed}}).$$
(15)

The *i*th pixel is an inlier if and only if its projection satisfies Eq. (15). Note, the value space of function $g(\cdot)$ is mapped to a 2D space **H** after Hough Transform:

$$x = g(r_i, s_i, s_j \mid \overline{\mathbf{P}}_i, \overline{\mathbf{P}}_j), \ y = \arccos(\overline{\mathbf{t}}_{i,j}^{\mathsf{T}} \overline{\mathbf{t}}_j),$$
(16)

where x and y are transformed coordinates. From Eq. (16), x is a synthesized translation magnitude and y is angular variable. We then set $x \in [0, x_{\max}]$, and $y \in [0, \theta_{\max}]$, where $\theta_{\max} = \arccos(-\overline{\mathbf{t}}_i^{\mathsf{T}} \mathbf{R}_{i,j} \overline{\mathbf{t}}_i)$. Finally, the value of **H** is:

$$\forall y \in [0, \theta_{\max}], \ \mathbf{H}(x \mid y) = 1, \text{if } x \in [J_{\min}, J_{\max}],$$
(17)

where J_{\min} and J_{\max} are the two bounds from Eq. (15). The transformation over the scoring function $f_{i,j}^{2D}$ with all M sampled pixels between frame \mathbf{I}_i and \mathbf{I}_j :

$$\mathbf{H}_{i,j} = \sum_{m} \mathbf{H}_{i,j}^{m}, \ f_{i,j}^{2\mathrm{D}}(s_i, s_j, r_i \mid \overline{\mathbf{P}}_i, \overline{\mathbf{P}}_j) = \mathbf{H}_{i,j}(x, y),$$
(18)

where x and y are functions of s_i, s_j, r_i . Eq. (1) becomes:

$$\phi(\overline{\mathcal{P}}) = \max_{\mathcal{S},\mathcal{R}} \sum_{i} \sum_{j,j \neq i} \mathbf{H}_{i,j}(x(\mathcal{S},\mathcal{R}), y(\mathcal{S},\mathcal{R})).$$
(19)

In our implementation, we discretize $\mathbf{H}_{i,j}$ to a 2D matrix.

Accelerate Bundle-Adjustment Consensus. The BA determines N-1 camera scales and N-1 depth adjustments to maximize the scoring function $\phi(\cdot)$ in Eq. (19). With Hough transform, BA maximizes the summarized intensity via



Fig. 6: Visualize Hough Transform Matrix \mathbf{H}_{i}^{j} from Eq. (18). Area with higher intensity suggests more inlier counts. Given normalized pose group, for N views, there exists $N \times (N-1)$ matrices \mathbf{H}_{i}^{j} , constraining N-1 scale and N-1 adjustments. We plot the start and end points after optimizing Eq. (19) in the figure.

indexing $N \times (N-1)$ Hough transform matrices **H**. It avoids BA repetitively enumerating all sampled pixels. Fig. 6 shows an example optimization process. Certified Global Optimality of robust inlier-counts scoring function Eq. (5)

and Eq. (6) are achieved after optimization. See Fig. 8 for more analysis.

Optimization with RGB-D. With GT depthmap, the algorithm switches to the 3D scoring function $f_{i,j}^{3D}(\cdot)$. The depth adjustment is fixed to 1 and the 2D line-circle intersection becomes 3D line-sphere intersection. See Supp. for details.

3.1.3 Computational Complexity

Naive Time Complexity. From Eq. (2) and Fig. 4, in each epoch, we evaluate (N-1)(K-1) pose groups with Hough Transform Acceleration. Suppose each group takes T iterations to optimize Eq. (19), the time complexity is:

$$\mathcal{O}((N-1)(K-1) \cdot N(N-1) \cdot (M+T)),$$
 (20)

where each group computes N(N-1) Hough matrices **H**. Each matrix enumerates M sampled pixels, see Eq. (18). Maximizing Eq. (19) becomes indexing **H**, hence has a constant time complexity T, where $T \ll M$.

Counting Unique Hough Matrices. Most computation is spent on Hough matrices. In Fig. 4, each connection in the computation graph suggests two unique Hough matrices. We minimize time complexity by only computing **unique** Hough matrices. In Fig. 4 first epoch, the initial optimal group $\overline{\mathcal{P}}^*$ has N(N-1) matrices. Each ablated group only differs by one pose, hence introducing 2(N-1)(N-1)(K-1) matrices. The first-epoch complexity is then:

$$\mathcal{O}^{H}(N(N-1)M + 2(N-1)^{2}(K-1)M) + \mathcal{O}^{BA}(N(N-1)(K-1)T).$$
(21)

Only the Hough transform is accelerated. As $T \ll M$, the complexity of BA is neglectable. After the first epoch, $\overline{\mathcal{P}}^*$ only updates one pose per epoch, hence introducing 2(N-2)(K-1) matrices. The complexity for the rest epochs is,

$$\mathcal{O}^{H}(2(N-2)(K-1)M) + \mathcal{O}^{BA}(N(N-1)(K-1)T).$$
(22)

While Eq. (22) has linear complexity, our method only updates one pose per epoch. Updating poses in all frames like other SfM methods is still quadratic.



(b) Geometric Verification

Fig. 7: Triangulation optimizes frustum RF for multiview consistency w.r.t. depth and correspondence. Geometric Verification inferences RF for sparse multiview consistent 3D points. For simplicity, in (a), we only plot L_c defined from the root frame.

3.2Frustum Radiance Field Triangulation

Frustum Radiance Field. Now, we fix the optimized pose \mathcal{P}^* . Then we employ a frustum radiance field **V** of size $H \times W \times D$ for dense triangulation. Field \mathbf{V} is defined over the root frame \mathbf{I}_o and shares similarity with the categorical depthmap [4,17]. We follow [58,62] in rendering the depth d. The RGB estimation is skipped as unrelated. A 3D ray originated from pixel \mathbf{p}_i at frame *i* is discretized into a set of 3D points and depth labels. With slight abuse of notation, we denote $\{\hat{\mathbf{p}}_{i,t} = \mathbf{o} + d_t \mathbf{r} \mid t \in [1,T]\}$, where $\hat{\mathbf{p}}$ is a 3D point, d_t is depth label and \mathbf{r} is ray direction. Set integration interval $\delta_t = d_{t+1} - d_t$, depth d is:

$$d(\mathbf{p}_i) = \sum_t \alpha_t d_t, \ \alpha_t = T_t (1 - \exp\left(-\sigma_t \delta_t\right)), \ T_t = \exp\left(-\sum_{t' \in [1,t]} \sigma_{t'} \delta_{t'}\right).$$
(23)

We set the camera origin of frame i as **o**. Instead of regressing occupancy δ with MLP [58, 62], we directly interpolate the radiance field V:

$$\delta_t = \mathbf{V}(u, v, w), \text{ where } \begin{bmatrix} u \ v \ w \end{bmatrix}^{\mathsf{T}} = \pi(\mathbf{E}, \hat{\mathbf{p}}_{i,t}).$$
(24)

Matrix **E** is the identity matrix. Function $\pi(\cdot)$ is projection function. Compared to using the MLP, frustum radiance field \mathbf{V} is more computationally efficient [16]. Triangulation. Classic triangulation method [48] operates on a single 3D point. The RF provides additional constraints where all optimized points share a canonical 3D volume. In Fig. 7, we supervise \mathbf{V} for multi-view consistency between dense depthmap \mathcal{D} and correspondence map \mathcal{C} . On depth:

$$L_D = \frac{1}{NM} \sum_{i} \sum_{m} \|\pi(\mathbf{P}_i, \hat{\mathbf{p}}^m) - d_i^m\|_1.$$
(25)

Here, $\hat{\mathbf{p}}^m$ is rendered from the root frame, following depth computed with Eq. (23). To apply correspondence consistency, we have:

$$L_{C} = \frac{1}{N(N-1)M} \sum_{i} \sum_{j,j\neq i} \sum_{m} \|\pi(\mathbf{P}_{j}, \hat{\mathbf{p}}_{i}^{m}) - \mathbf{q}_{i,j}^{m}\|_{1},$$
(26)

where $\hat{\mathbf{p}}_{i}^{m} = \pi^{-1}(\mathbf{P}_{i}, \mathbf{p}_{i}^{m}, d(\mathbf{p}_{i}^{m})), \mathbf{p}_{i}^{m} = \pi(\mathbf{P}_{i}, \hat{\mathbf{p}}^{m})$. With slight abuse of notation, function $\pi(\cdot)$ returns depth for L_D , and location for L_C . We always first render from the root frame and subsequently project to N frames. From there, we project to other supported frames again, forming N(N-1) pairs.

Table 1: Self-Supervised Depth Estimation. We apply self-supervision with 5 frames via executing the local SfM. We output improved sparse depthmaps over SoTA supervised inputs. The evaluation is conducted over the root frame.

Dataset	Method	Density	δ0.5	δ_1	SI_{log}	A.Rel	S.Rel	\mathbf{RMS}	${\rm RMS}_{\rm log}$
ScanNet [13]	ZoeDepth [5]	9.1%	0.877	0.963	6.655	0.056	0.016	0.154	0.075
	L Ours		0.902	0.976	5.901	0.050	0.014	0.149	0.070
	ZeroDepth [35]	5.6%	0.641	0.834	12.860	0.124	$0.\overline{086}$	$\overline{0}.\overline{3}\overline{3}\overline{7}$	0.152
	L Ours		0.686	0.877	9.463	0.106	0.067	0.295	0.133
	Metric3D [68]	2.6%	0.804	0.946	6.708	0.067	0.020	0.150	0.084
	L Ours		0.854	0.968	4.170	0.055	0.014	0.125	0.068
KITTI360 [33]	ZoeDepth [5]	4.0%	0.677	0.899	14.154	0.103	0.490	3.521	0.153
	L Ours		0.719	0.910	13.220	0.094	0.474	3.499	0.145
	ZeroDepth [35]	4.5%	0.584	0.844	16.468	0.132	0.819	$\bar{3}.\bar{4}8\bar{6}$	0.183
	L Ours		0.654	0.877	13.881	0.115	0.772	3.395	0.164
	Metric3D [68]	3.0%	0.846	0.958	9.226	0.072	0.508	$\overline{2.194}$	0.104
	ightarrow Ours	3.270	0.860	0.963	8.896	0.068	0.487	2.139	0.101

Table 2: Consistent Depth Estimation. We measure the numerical improvement by aligning the support frame depthmaps to the root frame with our depth adjustment scalars. The evaluation is conducted on support frames on ScanNet [13].

Method	$\delta_{0.5}$	δ_1	SI_{log}	A.Rel	S.Rel	\mathbf{RMS}	$\mathbf{RMS}_{\mathbf{log}}$
ZoeDepth [5]	0.658	0.894	9.242	0.104	0.039	0.255	0.128
\sqcup Ours	0.793	0.942	9.242	0.079	0.024	0.203	0.105
ZeroDepth [35]	0.351	0.589	20.145	0.254	0.223	$0.\overline{5}6\overline{5}$	0.287
\sqcup Ours	0.490	0.725	20.145	0.199	0.156	0.457	0.237
Metric3D [68]	0.533	0.753	$\overline{12.425}$	0.216	0.339	$0.\overline{4}9\overline{5}$	0.228
\sqcup Ours	0.664	0.838	12.425	0.137	0.126	0.345	0.175

3.3 Geometric Verification

With the RF optimized, we apply geometric verification to acquire sparse multiview consistent 3D points, as in Fig. 7:

$$\mathcal{C} = \{ \sum_{i,i\neq o} c_i^m \ge n^c \}, \ c_i^m = 1 \text{ if } \sum_{i,i\neq o} \|\hat{\mathbf{p}}_i^m - \hat{\mathbf{p}}^m\|_2 \le \lambda^c.$$
(27)

We follow the same rendering process as training, where $\hat{\mathbf{p}}_i^m$ is computed with Eq. (26). First, we render 3D points from the root frame, project them to other views, and render 3D points from there again. A point is valid if a minimum of n^c views are consistent with the root.

4 Experiments

4.1 Self-supervised Depth Estimation

We benchmark whether self-supervision benefits supervised depth in unseen test data. For the correspondence estimator, we use PDC-Net [57]. For depth estimators, we adopt recently published in-the-wild depth estimator, including

Table 3: Self-Supervised Correspondence Estimation. We improve correspondence with RGB-D inputs, using metrics from [57]. The entry train and test are training and testing datasets of correspondence estimators. [Key: M=MegaDepth, S=ScanNet]



ZoeDepth [5], ZeroDepth [35], and Metric3D [68]. We evaluate with ScanNet [13] and KITTI360 [33] where all models perform zero-shot prediction.

Test Data. In dense correspondence estimation, methods [57, 58, 78] output confidence score per correspondence. We follow [57, 58] to set a minimum threshold of 0.95. We run on ScanNet test split and it returns 92 sequences with sufficient correspondence. We form our test split by sampling 5 neighboring frames per valid sequence. Similarly, we run on KITTI360 data and randomly select 100×5 test split, *i.e.*, 100 sequences with 5 frames each. We consider it a comprehensive experiment. Similar to SPARF [58], our triangulation trains a NeRF-like structure. For reference, SPARF experiment on DTU dataset [26] includes only 15 sequences each with 3 images. In comparison, we include around 100 sequences. **Evaluation Protocols.** We evaluate on **root** frame. We remove the scale ambiguity in the local SfM system to correctly reflect depth improvement. Specifically, we adjust all 5 depthmaps by an identical scalar computed between estimated root and GT depthmap, *i.e.*, the median scaling [21]. This eliminates scale ambiguity in the root frame while preserving it in support frames.

Results. In Tab. 1, our point cloud has a density of 2.6% - 9.1%, which amounts to 10 - 30k points on a 480×640 image. On accuracy, we have **unanimous** improvement over all supervised models of both datasets. Especially, we outperform strong baselines of ZoeDepth on ScanNet and Metric3D on KITTI360.

4.2 Consistent Depth Estimation

We evaluate on ScanNet. We follow Sec. 4.1 data split but evaluate the **support** frames. Temporal consistent depth is essential for AR applications [37]. Tab. 2

13



Fig. 9: Self-supervised Correspondence Estimation enabled by our method with RGB-D inputs. The correspondence error is marked by the radius of the circle.

reflects the performance gain by aligning support frames to root with adjustments, which are jointly estimated with camera poses, see Fig. 2 and Fig. 3.

4.3 Self-supervised Correspondence Estimation

Real-world image correspondence label is expensive, *e.g.* KITTI provides only 200 optical flow labels. Existing datasets, such as MegaDepth and ScanNet, require large-scale 3D reconstruction with manual verification. Hence, correspondence estimators can not fine-tune on general RGB-D datasets like NYUv2 [51] or KITTI [18]. But our method enables self-supervised correspondence estimation on RGB-D data when using 3D scoring function Eq. (6). The camera poses are optimized with the point cloud specified by depthmap and correspondence. The accurate pose in turn improves projective correspondence over inputs. We use the same test split as Sec. 4.1. The evaluation accumulates correspondence of each frame pair. Fig. 8a shows our improvement is **unanimous** over both confident and unconfident estimation. A visual example is in Fig. 9.

4.4 Sparse-view Pose Estimation

Comparison with Optimization-based and Learning-based Poses. Previous studies either evaluate two-view pose [22,55], or SLAM-like odometry [61]. For more comparison, following Sec. 4.1 ScanNet split, we keep root frame and gradually add neighboring frames. In Tab. 4, LightedDepth [77] and ours both use PDC-Net [57] correspondence and ZoeDepth [5] mono-depth. COLMAP [48] uses PDC-Net correspondence. In evaluation, we follow [58] in aligning to GT poses. In Tab. 4, our **zero-shot** pose accuracy significantly outperforms all prior arts, including [22,55,61] with ScanNet [13] or ScanNet++ [67] in their training set. See Supp. Tab. 4 for complete comparison from 3 to 9 frames. In Fig. 8, we attribute our superiority to certified global optimality over robust measurements. **Comparison with NeRF-based Poses.** Sparse view NeRF methods optimize NeRF jointly with camera poses, mandating a sophisticated and time-consuming optimization scheme. *E.g.*, SPARF [58], takes one day to optimize the pose and NeRF. Typically, their poses are initialized with COLMAP. Our method provides an alternative initialization with superior performance. In Tab. 5, our

initialization achieves better or on-par pose performance than SoTA [58] while only taking ~ 3 minutes (Fig. 7). Our lower performance on Replica dataset might

Table 4: Sparse-view Pose Comparison with optimization-based and learningbased methods. We only compare against COLMAP on its success sequences. Please see the complete comparison from 3 to 9 frames in Supp. Tab. 1. Our method performs zero-shot testing on ScanNet while outperforming DeepV2D [55], DRO [22] with ScaNet [13] in training set. DUSt3R [61] trains on a similar dataset ScanNet++ [67].

Frames Method		Zero-shot	Suc. (%)	PCK-3	C3D-3	Rot.	Trans.
	COLMAP [48]	1	36.7	0.584	0.863	0.577	1.296
5	Ours	1	100.0	0.727	0.904	0.422	1.062
	$\overline{\text{DeepV2D}}$ [55] - ScanNet	×		0.526	0.805	0.945	1.496
	DeepV2D $[55]$ - NYUv2	1		0.530	0.771	1.041	1.568
	DeepV2D $[55]$ - KITTI	/ / X		0.125	0.387	4.908	4.231
	LightedDepth [77]			0.651	0.832	0.469	1.550
	DRO [22] - ScanNet		100.0	0.656	0.853	0.385	1.200
	DRO [22] - KITTI	1		0.003	0.211	3.610	5.469
	DUSt3R $[61]$ w.o. Intrinsic	1		0.364	0.705	0.487	2.074
	DUSt3R $[61]$ w.t. Intrinsic	 ✓ 		0.594	0.824	0.570	1.759
	Ours	1		0.799	0.900	0.368	1.120

Table 5: Sparse-view Pose Comparison with NeRF-based methods following [58].

Mothod	Frames	LLF	F [49]	Replica [54]		
Method		Rot.	Trans.	Rot.	Trans.	
BARF [34]		2.04	11.6	3.35	16.96	
RegBARF $[34, 40]$		1.52	5.0	3.66	20.87	
DistBARF $[2, 34]$	9	5.59	26.5	2.36	7.73	
SCNeRF [27]	3	1.93	11.4	0.65	4.12	
SPARF [58]		0.53	<u>2.8</u>	0.15	0.76	
Ours		0.46	1.9	0.52	<u>4.09</u>	

be due to ZoeDepth not being trained on synthetic data. Our work suggests the straightforward "first-pose-then-NeRF" scheme also applies to short videos.

Certified Global Optimality. In Fig. 8b, our Bundle-RANSAC-Adjustment **always** finds more inliers than groundtruth poses. To our best knowledge, we are the **first** work that extends RANSAC to a multi-view system.

Run-time. In Fig. 8c, we run approximately $3 \times$ slower than COLMAP. But both have quadratic complexity. With 3/5/7/9 frames, we take 0.8/2.0/5.3/9.4 minutes on RTX 2080 Ti GPU, while COLMAP uses 0.3/0.9/1.8/3.6 minutes on Intel Xeon 4216 CPU. COLMAP runs sequentially. But our method is highly parallelized. Our core operation Hougn Transform scales up with more GPUs.

5 Conclusion

By revisiting self-supervision with local SfM, we first show self-supervised depth benefits SoTA supervised model with only 5 frames. We have SoTA sparseview pose accuracy, applicable to NeRF rendering. We have diverse applications including self-supervised correspondence and consistent depth estimation.

Limitation. The NeRF-like triangulation constrains our method from applying to large-scale self-supervised learning. Its efficiency requires improvement.

¹⁴ Shengjie Zhu and Xiaoming Liu

References

- Agarwal, S., Furukawa, Y., Snavely, N., Simon, I., Curless, B., Seitz, S.M., Szeliski, R.: Building rome in a day. Communications of the ACM (2011) 1
- Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In: CVPR (2022) 14
- Beder, C., Steffen, R.: Determining an initial image pair for fixing the scale of a 3d reconstruction from an image sequence. In: Joint Pattern Recognition Symposium 3
- 4. Bhat, S.F., Alhashim, I., Wonka, P.: Adabins: Depth estimation using adaptive bins. In: CVPR (2021) 10
- Bhat, S.F., Birkl, R., Wofk, D., Wonka, P., Müller, M.: Zoedepth: Zero-shot transfer by combining relative and metric depth. arXiv preprint arXiv:2302.12288 (2023) 1, 3, 11, 12, 13
- Bian, J., Li, Z., Wang, N., Zhan, H., Shen, C., Cheng, M.M., Reid, I.: Unsupervised scale-consistent depth and ego-motion learning from monocular video. NeurIPS (2019) 4
- 7. Brown, M., Hua, G., Winder, S.: Discriminative learning of local image descriptors. PAMI (2010) 3
- Butime, J., Gutierrez, I., Corzo, L.G., Espronceda, C.F.: 3d reconstruction methods, a survey. In: VISAPP (2006) 1
- Casser, V., Pirk, S., Mahjourian, R., Angelova, A.: Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In: AAAI (2019) 4
- Chen, Y., Schmid, C., Sminchisescu, C.: Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In: ICCV (2019) 4
- 11. Clark, R., Bloesch, M., Czarnowski, J., Leutenegger, S., Davison, A.J.: Ls-net: Learning to solve nonlinear least squares for monocular stereo. ECCV (2018) 4
- 12. Crandall, D., Owens, A., Snavely, N., Huttenlocher, D.: Discrete-continuous optimization for large-scale structure from motion. In: CVPR (2011) 1
- Dai, A., Chang, A., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: ScanNet: Richly-annotated 3D reconstructions of indoor scenes. In: PAMI (2017) 11, 12, 13, 14
- 14. Deng, K., Liu, A., Zhu, J.Y., Ramanan, D.: Depth-supervised nerf: Fewer views and faster training for free. In: CVPR (2022) 4
- Edstedt, J., Sun, Q., Bökman, G., Wadenbäck, M., Felsberg, M.: Roma: Revisiting robust losses for dense feature matching. arXiv preprint arXiv:2305.15404 (2023) 12, 13
- Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR (2022) 10
- 17. Fu, H., Gong, M., Wang, C., Batmanghelich, K., Tao, D.: Deep ordinal regression network for monocular depth estimation. In: CVPR (2018) 1, 10
- Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The KITTI dataset. IJRR (2013) 13
- 19. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: CVPR (2012) 1
- Godard, C., Mac Aodha, O., Firman, M., Brostow, G.J.: Digging into selfsupervised monocular depth estimation. In: ICCV (2019) 1, 4

- 16 Shengjie Zhu and Xiaoming Liu
- Gordon, A., Li, H., Jonschkowski, R., Angelova, A.: Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In: ICCV (2019) 1, 12
- Gu, X., Yuan, W., Dai, Z., Zhu, S., Tang, C., Dong, Z., Tan, P.: Dro: Deep recurrent optimizer for video to depth. IEEE Robotics and Automation Letters (2023) 4, 13, 14
- 23. Guizilini, V., Ambruș, R., Chen, D., Zakharov, S., Gaidon, A.: Multi-frame selfsupervised depth with transformers. In: CVPR (2022) 4
- Hartley, R., Zisserman, A.: Multiple view geometry in computer vision. Cambridge university press (2003) 7
- Izquierdo, S., Civera, J.: Sfm-ttr: Using structure from motion for test-time refinement of single-view depth networks. In: CVPR (2023) 4
- Jensen, R., Dahl, A., Vogiatzis, G., Tola, E., Aanæs, H.: Large scale multi-view stereopsis evaluation. In: CVPR (2014) 12
- Jeong, Y., Ahn, S., Choy, C., Anandkumar, A., Cho, M., Park, J.: Self-calibrating neural radiance fields. In: CVPR (2021) 4, 14
- Kuznietsov, Y., Proesmans, M., Van Gool, L.: Comoda: Continuous monocular depth adaptation using past experiences. In: WACV (2021) 4
- Lee, J.H., Han, M.K., Ko, D.W., Suh, I.H.: From big to small: Multi-scale local planar guidance for monocular depth estimation. arXiv preprint arXiv:1907.10326 (2019) 1
- 30. Lepetit, V., Moreno-Noguer, F., Fua, P.: Ep n p: An accurate o (n) solution to the p n p problem. IJCV (2009) 3
- 31. Li, H.: A practical algorithm for l triangulation with outliers. In: CVPR (2007) 3
- 32. Li, H., Hartley, R.: Five-point motion estimation made easy. In: ICPR (2006) 3, 5
- Liao, Y., Xie, J., Geiger, A.: Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. IEEE Transactions on Pattern Analysis and Machine Intelligence (2022) 11, 12
- Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf: Bundle-adjusting neural radiance fields. In: ICCV (2021) 4, 14
- Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. arXiv preprint arXiv:2303.11328 (2023) 11, 12
- 36. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. IJCV (2004) 3
- Luo, X., Huang, J.B., Szeliski, R., Matzen, K., Kopf, J.: Consistent video depth estimation. ToG (2020) 4, 12
- Mahjourian, R., Wicke, M., Angelova, A.: Unsupervised learning of depth and egomotion from monocular video using 3d geometric constraints. In: CVPR (2018) 4
- Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM (2021) 2
- Niemeyer, M., Barron, J.T., Mildenhall, B., Sajjadi, M.S., Geiger, A., Radwan, N.: Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In: CVPR (2022) 4, 14
- 41. Niu, L., Cong, W., Liu, L., Hong, Y., Zhang, B., Liang, J., Zhang, L.: Making images real again: A comprehensive survey on deep image composition. arXiv preprint arXiv:2106.14490 (2021) 1
- Olsson, C., Eriksson, A., Hartley, R.: Outlier removal using duality. In: CVPR (2010) 3

Revisit Self-supervised Depth Estimation with Local Structure-from-Motion

- Piccinelli, L., Sakaridis, C., Yu, F.: idisc: Internal discretization for monocular depth estimation. In: CVPR (2023) 1
- 44. Pillai, S., Ambruş, R., Gaidon, A.: Superdepth: Self-supervised, super-resolved monocular depth estimation. In: ICRA (2019) 4
- Ranftl, R., Lasinger, K., Hafner, D., Schindler, K., Koltun, V.: Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. PAMI 1
- Ranjan, A., Jampani, V., Balles, L., Kim, K., Sun, D., Wulff, J., Black, M.J.: Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In: CVPR (2019) 4
- Sarlin, P.E., Lindenberger, P., Larsson, V., Pollefeys, M.: Pixel-perfect structurefrom-motion with featuremetric refinement. PAMI (2023) 1
- Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: CVPR (2016) 1, 3, 10, 13, 14
- Shafiei, M., Bi, S., Li, Z., Liaudanskas, A., Ortiz-Cayon, R., Ramamoorthi, R.: Learning neural transmittance for efficient rendering of reflectance fields (2021) 14
- 50. Shu, C., Yu, K., Duan, Z., Yang, K.: Feature-metric loss for self-supervised learning of depth and egomotion. In: ECCV (2020) 4
- Silberman, N., Hoiem, D., Kohli, P., Fergus, R.: Indoor segmentation and support inference from rgbd images. ECCV (2012) 13
- 52. Snavely, N., Seitz, S.M., Szeliski, R.: Photo tourism: exploring photo collections in 3d. In: Siggraph (2006) 1
- 53. Spencer, J., Russell, C., Hadfield, S., Bowden, R.: Kick back & relax: Learning to reconstruct the world by watching slowtv. In: ICCV (2023) 4
- 54. Straub, J., Whelan, T., Ma, L., Chen, Y., Wijmans, E., Green, S., Engel, J.J., Mur-Artal, R., Ren, C., Verma, S., Clarkson, A., Yan, M., Budge, B., Yan, Y., Pan, X., Yon, J., Zou, Y., Leon, K., Carter, N., Briales, J., Gillingham, T., Mueggler, E., Pesqueira, L., Savva, M., Batra, D., Strasdat, H.M., Nardi, R.D., Goesele, M., Lovegrove, S., Newcombe, R.: The Replica dataset: A digital replica of indoor spaces. arXiv preprint arXiv:1906.05797 (2019) 14
- 55. Teed, Z., Deng, J.: Deepv2d: Video to depth with differentiable structure from motion. In: ICLR (2020) 3, 4, 13, 14
- 56. Tiwari, L., Ji, P., Tran, Q.H., Zhuang, B., Anand, S., Chandraker, M.: Pseudo rgb-d for self-improving monocular slam and depth prediction. In: ECCV (2020) 4
- 57. Truong, P., Danelljan, M., Timofte, R., Van Gool, L.: Pdc-net+: Enhanced probabilistic dense correspondence network. PAMI (2023) 3, 11, 12, 13
- 58. Truong, P., Rakotosaona, M.J., Manhardt, F., Tombari, F.: Sparf: Neural radiance fields from sparse and noisy poses. In: CVPR (2023) 3, 4, 10, 12, 13, 14
- Tuytelaars, T., Mikolajczyk, K., et al.: Local invariant feature detectors: a survey. Foundations and trends[®] in computer graphics and vision (2008) 3
- Ummenhofer, B., Zhou, H., Uhrig, J., Mayer, N., Ilg, E., Dosovitskiy, A., Brox, T.: Demon: Depth and motion network for learning monocular stereo. In: CVPR (2017) 4
- Wang, S., Leroy, V., Cabon, Y., Chidlovskii, B., Revaud, J.: Dust3r: Geometric 3d vision made easy. arXiv preprint arXiv:2312.14132 (2023) 3, 13, 14
- Wang, Z., Wu, S., Xie, W., Chen, M., Prisacariu, V.A.: Nerf-: Neural radiance fields without known camera parameters. arXiv preprint arXiv:2102.07064 (2021) 10
- Watson, J., Firman, M., Brostow, G.J., Turmukhambetov, D.: Self-supervised monocular depth hints. In: ICCV (2019) 4

- 18 Shengjie Zhu and Xiaoming Liu
- 64. Watson, J., Mac Aodha, O., Prisacariu, V., Brostow, G., Firman, M.: The temporal opportunist: Self-supervised multi-frame monocular depth. In: CVPR (2021) 4
- 65. Wu, C.: Towards linear-time incremental structure from motion. In: 3DV (2013) 1, 3
- Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. arXiv preprint arXiv:2401.10891 (2024) 4
- Yeshwanth, C., Liu, Y.C., Nießner, M., Dai, A.: Scannet++: A high-fidelity dataset of 3d indoor scenes. In: ICCV (2023) 13, 14
- Yin, W., Zhang, C., Chen, H., Cai, Z., Yu, G., Wang, K., Chen, X., Shen, C.: Metric3d: Towards zero-shot metric 3d prediction from a single image. In: ICCV (2023) 11, 12
- 69. Yuan, W., Gu, X., Dai, Z., Zhu, S., Tan, P.: New crfs: Neural window fullyconnected crfs for monocular depth estimation. In: CVPR (2022) 1
- Zhan, H., Garg, R., Weerasekera, C.S., Li, K., Agarwal, H., Reid, I.: Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In: CVPR (2018) 4
- Zhang, Z.: Microsoft kinect sensor and its effect. IEEE multimedia 19(2), 4–10 (2012) 1
- Zhang, Z., Cole, F., Tucker, R., Freeman, W.T., Dekel, T.: Consistent depth of moving objects in video. TOG (2021) 4
- 73. Zhao, W., Liu, S., Shu, Y., Liu, Y.J.: Towards better generalization: Joint depthpose learning without posenet. In: CVPR (2020) 4
- 74. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: CVPR (2017) 4
- 75. Zhou, Z., Dong, Q.: Two-in-one depth: Bridging the gap between monocular and binocular self-supervised depth estimation. In: ICCV (2023) 4
- Zhu, S., Brazil, G., Liu, X.: The edge of depth: Explicit constraints between segmentation and depth. In: CVPR (2020) 1, 4
- 77. Zhu, S., Liu, X.: Lighteddepth: Video depth estimation in light of limited inference view angles. In: CVPR (2023) 3, 4, 8, 12, 13, 14
- Zhu, S., Liu, X.: Pmatch: Paired masked image modeling for dense geometric matching. In: CVPR (2023) 12