

Deep PatchMatch MVS with Learned Patch Coplanarity, Geometric Consistency and Adaptive Pixel Sampling

Jae Yong Lee^{1*}, Chuhan Zou² and Derek Hoiem¹

^{1*} Department of Computer Science, University of Illinois at Urbana-Champaign, 201 N Goodwin Ave, Urbana, 61801, Illinois, United States.

²Amazon Inc., 410 Terry Ave N, Seattle, 98109, Washington, United States.

*Corresponding author(s). E-mail(s): lee896@illinois.edu;

Contributing authors: zouchuha@amazon.com; dhoiem@illinois.edu;

Abstract

Recent work in multi-view stereo (MVS) combines learnable photometric scores and regularization with PatchMatch-based optimization to achieve robust pixelwise estimates of depth, normals, and visibility. However, non-learning based methods still outperform for large scenes with sparse views, in part due to use of geometric consistency constraints and ability to optimize over many views at high resolution. In this paper, we build on learning-based approaches to improve photometric scores by learning patch coplanarity and encourage geometric consistency by learning a scaled photometric cost that can be combined with reprojection error. We also propose an adaptive pixel sampling strategy for candidate propagation that reduces memory to enable training on larger resolution with more views and a larger encoder. These modifications lead to 6-15% gains in accuracy and completeness on the challenging ETH3D benchmark, resulting in higher F_1 performance than the widely used state-of-the-art non-learning approaches ACMM and ACPM.

Keywords: Multi-View Stereo, Learning-based Stereo, 3D Reconstruction

1 Introduction

Multi-view stereo (MVS) aims to infer accurate and complete 3D geometry from a set of calibrated images, with many applications such as robotics [7, 20], mixed reality [19, 30], and vision-based inspection [9]. The first deep learning-based approaches work well for densely sampled views with small-scaled scenes, such as the DTU dataset [1] and the Tanks-and-Temples intermediate benchmark [12]. However, in more challenging benchmarks with sparse views, wide baselines, and large depth ranges, such as the ETH3D High-res benchmark [24], non-learning based methods

achieve better performance by jointly optimizing over many views for the depths, normals, and visibilities of each pixel that maximize photometric scores while satisfying geometric consistency. Recent methods incorporate PatchMatch optimization into an end-to-end training framework [14, 25], partially closing the performance gap. The latest work PatchMatch-RL [14] proposes to jointly learn from pixel-wise depth, normal, and visibility, by using reinforcement learning to overcome the non-differentiable PatchMatch optimization. However, methods such as ACMM [27] still outperform, in part due to the

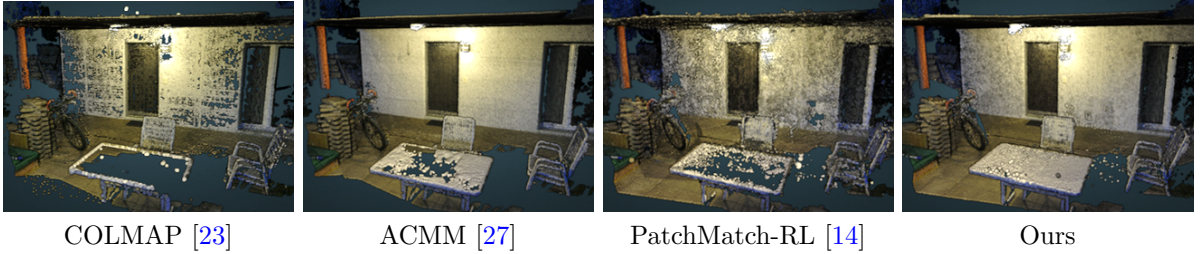


Fig. 1 We propose a learning-based PatchMatch MVS approach with learned patch coplanarity, geometric consistency and adaptive pixel sampling. Our method achieves more complete reconstruction compared to existing MVS methods.

precision of the bilaterally weighted NCC photometric cost, the ability to optimize over many views at high resolution, and the ability to incorporate geometric consistency constraints.

In this paper, we propose three designs to address the above-mentioned issues of learning-based PatchMatch MVS. First, we improve the photometric cost function by directly learning whether nearby pixels are co-planar. Patch-based MVS methods compute photometric score by comparing intensities or features over some neighborhood after rectifying for depth and normal, under the assumption that the surface is locally planar. Existing approaches weight the neighborhood based on color and pixel distance [23] or dot-product attention between feature vectors [14]. We modify the affinity neighborhood network CSPN [5] to directly learn which neighboring points are likely coplanar, which provides weights for more robust photometric costs in both textureless and textured regions.

Second, we incorporate geometric consistency regularization, replacing the recurrent cost regularization of PatchMatch-RL [14]. Although geometric consistency is commonly used in non-learning approaches (e.g., [23, 27, 28]), integration in deep learning approaches is challenging, since the neighborhood geometry is not differentiable with respect to the reference view geometry. We side-step this problem by learning to map the photometric score into a log-disbelief (i.e., negative log likelihood) and model reprojection error with a Gaussian distribution, such that the joint likelihood is the sum of photometric log-disbelief and the L2-norm reprojection error.

Finally, we investigate how to train learning-based PatchMatch MVS to perform inference with high-resolution images, more views, and deeper features. The challenge is that higher resolutions

require larger receptive fields to include sufficiently discriminate texture. Coarse-to-fine architectures [14, 25, 27] help, but more complex features with larger receptive fields are also required, which in turn requires prohibitive quantities of GPU memory. We propose adaptive sampling of the supporting pixels and gradient checkpointing to reduce the memory usage in the training. This enables more PatchMatch iterations and larger sets of view selections and source views during training, as well as use of a U-Net [22]-shaped feature extractor that covers a much larger receptive field compared to the shallow VGG-style backbone used in [14].

In summary, our main **contributions** are:

- **Learned coplanarity** which weights the photometric score support based on estimated coplanarity.
- **Modeling geometry consistency guided inference** with learned photometric score functions which improves completion of the reconstruction.
- **Efficient training using adaptively sampled subset of supporting pixels** that reduces memory usage in the training time to allow: a) more PatchMatch training iterations, b) a larger set of source images, c) a larger number of pixel-wise view selections and d) a more complex feature extraction backbone.

Combined, our method outperforms state-of-the-art learning-based MVS methods on the ETH3D high-res benchmark by a large margin with the same image resolution. We also provide an efficient and scalable inference architecture that enables high-resolution inference on up to 10 source images of 3200×2112 , outperforming the classical counterpart methods ACMM [27] and ACMP [28] at the same image resolution.

2 Related Works

The goal of multi-view stereo (MVS) is to estimate the scene geometry using the images with the known poses. We refer readers to Furukawa and Hernández [6] for a broad review and focus on the most related work here.

Stereo based depth estimation methods hypothesize depth values and score the hypotheses based on photometric similarity (or matching cost). Initial works [2, 11, 15] compute the matching cost as the sum of square differences or absolute differences between corresponding pixels in a local patch. This simple score is unreliable for oblique surfaces, near depth discontinuities, and for low-texture or non-Lambertian surfaces. MVS approaches evolved to address these concerns by jointly estimating depth and normal for each pixel [4], assigning weights to supporting pixels [23], and using robust matching functions, e.g. learned or NCC, as well as regularizing the cost based on neighborhood estimates or scores [10, 16, 31, 32]. The joint depth and normal search has been implemented by using PatchMatch [3] based optimization by Bleyer et al. [4], and extended to multi-view stereo with visibility selection [8, 23, 35]. Incorporation of normal estimates enables a more accurate patch-wise homography to match views of oblique surfaces, but depth discontinuities remain a problem. To address depth discontinuities, bilaterally weighted normalized cross correlation (NCC) is proposed [23], with weights modeling the likelihood of the coplanarity of each supporting pixel as proportional to a product of color and position difference terms. Bilaterally weighted NCC is widely used in recent non-learning PatchMatch-based methods [13, 21, 26–29].

In learning-based PatchMatch MVS, Lee et al. [14] use scaled dot-product based attention to model supporting pixel weights. Both bilaterally weighted NCC and scaled dot-product attention rely on center-pixel similarity to measure the local co-planarity of the supporting pixels. However, when pixels across depth boundaries have similar colors or flat surfaces have varying texture, center-pixel similarity does not correctly model coplanarity. Instead, we propose to directly regress the supporting pixel-wise co-planarity using ground-truth supervision.

Local photometric scores are unpredictable or unreliable for textureless or non-Lambertian surfaces, requiring consistency checks across neighborhoods or other views. COLMAP [23] checks the consistency of predicted depths and normals with neighboring views to augment the photometric score and filter inconsistent points, improving accuracy but decreasing recall. ACMM [27] propose a coarse-to-fine geometric consistency regularized inference. Reconstruction at coarse resolutions enables photometric scores with wide receptive fields with more texture, yielding more accurate but imprecise depth estimates, that are then refined at finer scales using a combination of geometric consistency and photometric cost terms. Our work combines the coarse-to-fine geometric consistency driven inference with the learning based PatchMatch MVS by mapping the learned photometric score into a log-disbelief scale, which allows the optional aggregation of additional log-disbelief terms, such as a geometric consistency term.

DeepC-MVS [13] achieves state-of-the-art results on ETH3D by adding a post-process refinement step to a slightly modified version of ACMM, filtering or refining depth values using learned confidences based on the initial estimates of depth, normal, and visibility. Similar refinements would likely benefit learning-based approaches also, but we leave this to future work.

The state-of-the-art learning-based methods, Vis-MVSNet [34] and EPP-MVSNet [17] operate on high-resolution images, which allows more definite reconstruction as pixel-level precision improves as frustums corresponding to each projection pixels become narrower. Both Vis-MVSNet [34] and EPP-MVSNet [17] use U-Net [22] as a shared CNN based encoder, which is known to preserve high-frequency details and provide wider receptive fields that make the encoder more suitable for high-resolution reconstruction.

We additionally integrate the learned patch coplanarity to sample the subset of supporting pixels in training time. This allows time and memory efficient training using policy gradient [14], which lets our method to be trained with more PatchMatch iterations and source view selections per pixel, as well as using 2D U-Net [22] based feature extraction.

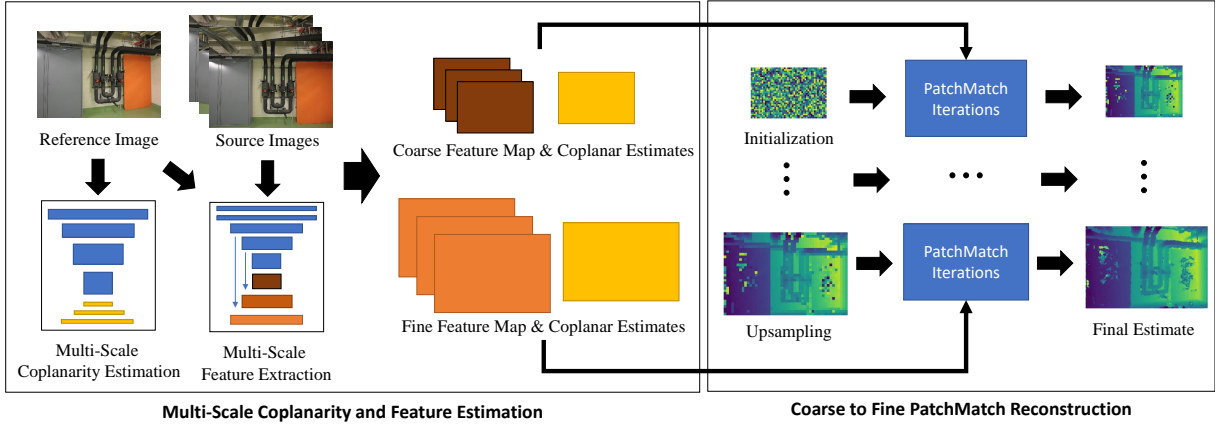


Fig. 2 Training Architecture overview: We extract multiscale features for each image and multiscale coplanarity weights for the reference image. We use features and coplanar maps of corresponding scales to perform coarse-to-fine estimation of depth, normal and visibility maps. At the coarsest stage, we initialize pixelwise depth and normal maps. Then, a series of PatchMatch iterations update the estimates. The PatchMatch iteration consists of four stages: (1) pixelwise view selection; (2) candidate propagation; (3) candidate scoring; and (4) candidate selection. The current solution is then upsampled to a finer level of input through nearest neighbor interpolation, and this procedure continues until the oriented point estimates at the finest level are fused from all images.

3 Method

Figure 2 shows an overview of our architecture during training. We build our approach on PatchMatch-RL (Sec. 3.1), a recent learning-based PatchMatch MVS framework, and propose three designs that lead to large performance gain in both accuracy and completeness: (1) improve the photometric cost function by directly modeling the coplanarity of neighboring pixels (Sec. 3.2); (2) replace the recurrent cost regularization with geometric consistency by learning a scaled photometric cost (Sec. 3.3); (3) design an adaptive pixel sampling strategy for candidate propagation to reduce memory, so that high-res images, more views, and deeper features can be used for training (Sec. 3.4).

3.1 Review of PatchMatch-RL

The input of PatchMatch-RL is a set of reference and source images $\{I_r, I_i\}_{i=1}^N \subset \mathcal{I}$, and their corresponding cameras poses $\{C_r, C_i = (K_i, R_i, t_i)\}_{i=1}^N \subset \mathcal{C}$ including intrinsic and extrinsic parameters. The method uses a coarse-to-fine PatchMatch MVS inference to jointly estimate the depths, normals, and visibilities of the reference image $D_r, N_r, \{V_{r,i}\}_{i=1}^N$, and uses policy-gradient for end-to-end training.

PatchMatch-RL starts by extracting CNN based features $\{f_\theta^c(I_r) = F_r^s, f_\theta^c(I_i) = F_i^s\}_{i=1}^N$ at

multiple scales $s \in \{3, 2, 1\}$, where f_θ^c represents trainable CNN backbone with parameters θ , and s represents the \log_2 scale resolution of the image (e.g. $F_i^s \rightarrow \mathbb{R}^{\frac{H}{2^s} \times \frac{W}{2^s} \times C}$). Then, the method initializes the depth and normal map $d_r^{s,t}, n_r^{s,t}$ at the coarsest scale $s = 3$ and the initial iteration $t = 0$. For the ease of notation, we drop the scale s and let $d_r^t, n_r^t, v_{r,i}^t$ denote the depth, normal and visibility map at the t -th iteration on the same scale in the coarse-to-fine stage from onward.

After initialization, a series of PatchMatch iterations are run at each scale, and each iteration jointly updates the depth, normal, and visibility maps. A single PatchMatch iteration consists of 4 steps:

1. **Pixel-wise view selection** updates the visibility map $v_{r,i}^t \rightarrow v_{r,i}^{t+1}$ using the current set of depth d_r^t and normal n_r^t scored by an MLP f_θ^v , which defines the parameterized view-sampling policy π_θ^v .
2. **Candidate Propagation** generates potential candidates of depth and normal for the next iteration by using PatchMatch-based propagation.
3. **Candidate Scoring** scores each candidate depth / normal pairs using the current set of visibility map $v_{r,i}^{t+1}$ through an RNN-based scorer regularizer f_θ^s , which defines the parameterized candidate-sampling policy π_θ^s .

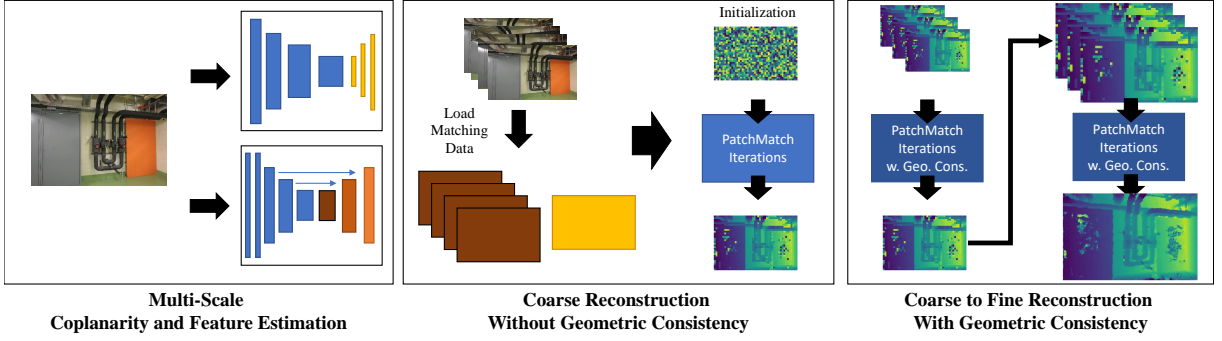


Fig. 3 Inference Architecture: First, features and coplanar maps for each image are generated and stored in files. Then, coarse-level reconstruction without geometric consistency constraints is performed for each reference image. Specifically, the relevant coarse coplanar and feature maps are loaded and, after initializing depths and normals, a PatchMatch iteration (Figure 2) is performed to infer geometry. Next, coarse-to-fine reconstruction is performed with geometric consistency. For each stage, results from the previous stage are used to initialize estimates and compute geometric consistency scores. All images in the scene are processed stagewise, and the inference is complete once the finest scale is processed.

4. **Candidate Selection** updates the current depth and normal map to d_r^{t+1} and n_r^{t+1} ,

The parameterized policies represents the sampling procedure defined by the distribution formed by the scores obtained by the MLP f_θ^v and the RNN f_θ^s . The view selection score distribution takes the supporting pixel-aggregated group-wise correlation scores for each source images as input. The supporting pixel-aggregated group-wise correlation scores use supporting pixel weights, obtained by the dot-product attention between the features at the center pixel and the supporting pixels, to compute weighted mean of the group-wise correlation scores between the pixels of reference feature map and the corresponding pixels of the source feature map. Given supporting pixels $q \in W(p)$ around the pixel p and its corresponding weights $w(q)$, the supporting pixel-aggregated group-wise correlation of the pixel p can be defined as:

$$\mathcal{G}(p) = \sum_{s=1}^N \sum_{q \in W(p)} w(q)(F^r(q) \otimes F^s(q'))/N \quad (1)$$

where q' indicates the corresponding pixel of q in the source image warped by the depth $d_r(p)$ and normal $n_r(p)$ value at p .

Similarly, the candidate scoring distribution obtained through f_θ^s , takes view-weighted, supporting pixel aggregated group-wise correlation

scores for selected views as input:

$$\mathcal{G}_\phi^v(p) = \sum_{s \in \mathcal{V}} \sum_{q \in W(p)} v_{r,s}(p)w(q)(F^r(q) \otimes F^s(q^\phi)) \quad (2)$$

where \mathcal{V} is the set of selected views, and $v_{r,s}(p)$ indicates visibility estimate of the pixel p in the reference image r from the source image s , and q^ϕ denotes the corresponding pixel of q in the source image, warped by the candidate depth $d_\phi(p)$ and normal $n_\phi(p)$. The output value of the RNN f_θ^s measures the photometric score of each candidate as $S_{pho}^\phi(p) = f_\theta^s(\mathcal{G}_\phi^v(p))$.

In inference, the policy performs arg-max selection based on the set of computed scores. In training, the policy performs selection based on the decaying ϵ -greedy method, where it samples from the distribution with the probability of ϵ and samples using arg-max with the probability of $1 - \epsilon$.

After a sufficient number of iterations for each scale, the depth and normal maps are upsampled using nearest neighbor upsampling. Then, the next stage of the coarse-to-fine architecture starts, and the process is repeated until the most finest-level iteration is completed.

To train the model in an end-to-end differentiable manner, PatchMatch-RL uses policy-gradient algorithm REINFORCE to jointly train $f_\theta^c, f_\theta^v, f_\theta^s$. The reward for each iteration is defined as $r^t = \mathcal{N}(d_r^t; D_r^*, \sigma_d) \cdot \mathcal{N}(n_r^t; N_r^*, \sigma_n)$, where D_r^* and N_r^* represent the ground-truth depth

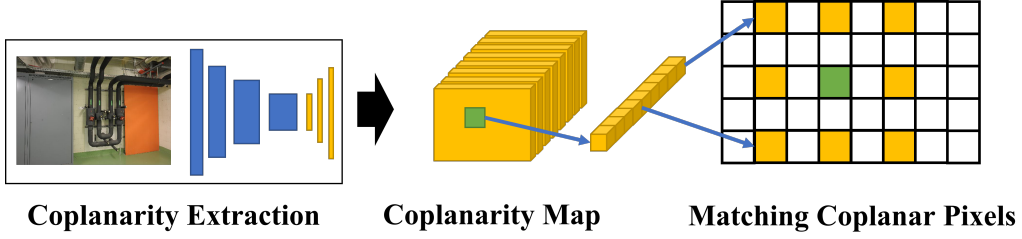


Fig. 4 Coplanarity Extraction and Representation: We use a coplanar extractor to estimate the coplanar values of the supporting pixels of each pixel in the image. The output coplanarity map is an array of shape $H \times W \times 9$, where each position represents the corresponding 3×3 supporting pixels of the image.

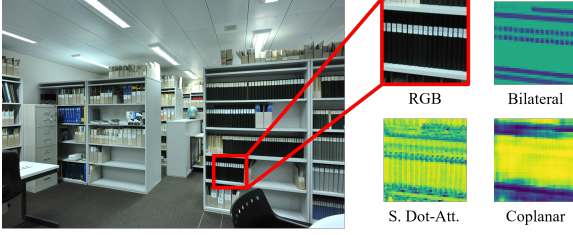


Fig. 5 Comparison of Supporting Weights: We visualize different types of the supporting pixel weights that are 3 pixel above from the center. From the left top, RGB denotes the extracted patch for the reference, Bilateral denotes the bilateral weighting [23], S. Dot-Att. denotes scaled dot-product attention [14] and Coplanar denotes our learned coplanar weights. Our method more correctly identifies which pixels correspond to coplanar surfaces, e.g. tags on binders are coplanar with dark part of binders, while shelves are not.

and normal maps, and $\mathcal{N}(x; \mu, \sigma)$ represents the probability of observing x from the normal distribution with mean μ and standard deviation σ . The update for each iteration is:

$$\nabla_{\theta} J = \sum_t \sum_{t' \geq t} \nabla_{\theta} \gamma^{t'-t} r^t \log \pi_{\theta}. \quad (3)$$

We let γ represent the discount coefficient of future reward. PatchMatch-RL additionally modifies the candidate selection loss to use cross-entropy between the ground-truth distribution of observing each candidate (e.g. $\forall(d, n), \mathcal{N}(d; D_r^*, \sigma_d) \cdot \mathcal{N}(n; N_r^*, \sigma_n)$) and the distribution formed by the scores of each candidates using the candidate scoring policy for robustness in training.

3.2 Learning Patch Coplanarity

Estimating the surface normal for each pixel requires a patch-based photometric score, computed as a weighted sum of similarities over

supporting interpolated feature positions. Rather than computing weights using scaled dot-product attention like PatchMatch-RL, we propose to set the weights to the likelihood of coplanarity with supporting pixels, which is estimated using a separate CNN architecture. (i.e., we directly estimate $w(q)$ for $q \in W(p)$ in Eqs. 1 and 2) We modify the CSPN architecture [5] used to learn the affinity of neighboring points, and change its last layer to a 3×3 convolution with a dilation of 3, so that the receptive field matches the supporting pixel weight locations. Figure 4 depicts an overview of our coplanarity extraction and output feature representation. In addition to the policy gradient loss, we use mean squared error term from the ground-truth coplanarity points obtained from the ground-truth depth to supervise the the training. We show in Figure 5 that compared with existing supporting pixel weighting schemes, our approach can better identify coplanarity within local patches with different textures.

3.3 Encouraging Geometric Consistency in Inference

We propose geometric consistency based regularization to replace the recurrent cost regularization used in PatchMatch-RL. The key idea of geometric consistency is to use the inferred geometry of the source images to validate the reference image geometry. We adopt the geometric consistency score used in ACMM [27], which defines an additional score term based on the reprojection error of the visible source view pixels.

Although geometric consistency is commonly used for non-learning based approaches [23, 27], it is challenging to integrate in learning-based approaches, since the source view geometry is not differentiable with respect to the reference view

geometry. In order to allow the optional adding of geometric log-disbelief scores in candidate selection stage without having to train additional module to combine the two scores, we remap the photometric score and the geometric consistency score range to the log-disbelief range. In detail, we let the photometric score ranges logarithmically from 0 to S_{max} , where 0 implies perfect photometric match and S_{max} implies complete disagreement. We use MLP as the photometric consistency scorer to replace the RNN-based scorer f_θ^s , since the history of candidate selection cannot handle the additional geometric consistency term in inference.

We let the pixel-wise reprojection error distance be the log scale disbelief score of the geometric consistency term. Given the candidate depth and normal map $d_\phi, n_\phi = \phi$, we use the homography between image j and image k at pixel p defined as $H_{j \rightarrow k}^\phi(p) = K_k \cdot (R_{j \rightarrow k} - \frac{t_{j \rightarrow k} \cdot n_\phi^T(p)}{d_\phi(p)}) \cdot K_j^{-1}$, and compute the reprojection error:

$$E_s^\phi(p) = H_{s \rightarrow r}^\phi(H_{r \rightarrow s}^\phi(p) \cdot p) - p_2 \quad (4)$$

We define G_{max} as the maximum log-disbelief term that can be added to geometric consistency, and combine the score weighted by the visibility mapping weights:

$$S_{geo}^\phi(p) = \sum_{i=1}^N v_{r,i}^t \cdot (S_{pho}^\phi(p) + \min(E_i^\phi(p), G_{max})) \quad (5)$$

where $S_{pho}^\phi(p)$ denotes the photometric score obtained from the MLP. We then use $S_{geo}^\phi(p)$ to select candidate for the pixel p , updating the depth $d_r(p)$ and normal $n_r(p)$ at the corresponding iteration.

Since we need to have the source view geometry available to use geometric consistency based regularization, we use a modified version of the coarse-to-fine architecture, as illustrated in Figure 3. In inference, the source view and its inferred geometry in the previous stage are used to calculate S_{geo}^ϕ for the candidate selection step of each PatchMatch iteration. The inferred depth and normal maps at each stage are stored in a file, and loaded in the next stage of reconstruction with geometric consistency.

3.4 Adaptive Supporting Pixel Sampling

Training the learning-based PatchMatch MVS model using policy gradient requires memory linear to the number of depth, normal and visibility map evaluations (*e.g.* number of PatchMatch iterations), since the history of gradients need to be stored. Hence, PatchMatch-RL [14] uses a shallow VGG-style backbone, with a small number of PatchMatch iterations per scale. We resolve this problem by using a sampled subset of the supporting pixels in training, which effectively reduces the training computation by half and the memory by roughly 30%, so that a larger feature extraction backbone can be used with faster training. We use the distribution of the coplanar values of the supporting pixels as a sampling distribution. We sample 3 points out of 8 supporting pixels, (excluding the center pixel of the 3×3 local supporting pixels), which is the minimum number of points that supports estimates of surface normal orientation. The adaptive supporting pixel sampling can be used in the inference time for faster runtime, at the cost of a slight drop in reconstruction quality. Table 3 and Table 4 in Section 4 show the detailed analysis of the effect of adaptive supporting pixel sampling in inference.

4 Experiments

We evaluate our method on two standard MVS benchmarks, ETH3D High-res Multi-View benchmark [24] and Tanks-and-Temples(TnT) [12] intermediate and advanced benchmark.

4.1 Implementation Details

For all benchmarks, we use BlendedMVS [33] dataset for training. BlendedMVS is a large-scale MVS dataset that contains 113 indoor, outdoor and object-level scenes. We use 7 images (*i.e.*, 1 reference image and 6 source images) with an input resolution of 768×576 , and an output resolution of 384×288 for training. Among the 6 source images, we use 3 images sampled from the 10 best matching source images given by the global view selection provided by the dataset, and 3 images randomly selected from the remaining images in the scene. In training local view selection, we consider the three best views to be

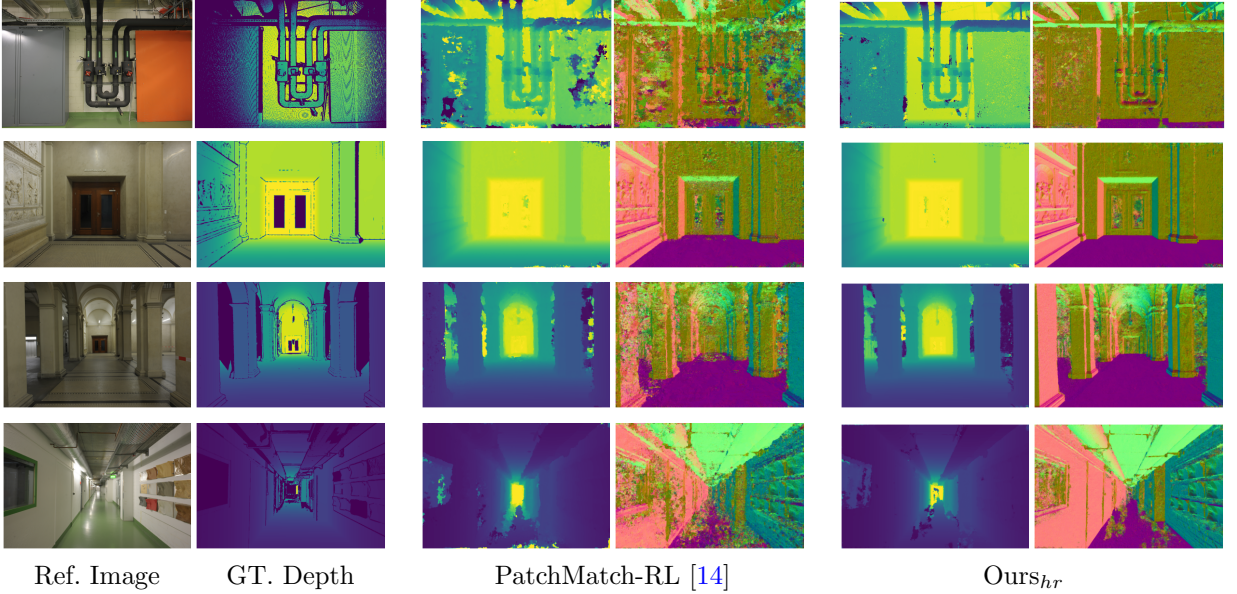


Fig. 6 Qualitative comparison against PatchMatch-RL on the ETH3D high-res multi-view benchmark training scenes. We show from the left the reference image, the ground-truth depth, the depth and normal estimated by PatchMatch-RL, the depth and normal estimated by Ours_{hr}. All depth maps share the same color scale based on the ground-truth depth ranges. Our method can reconstruct more accurate points in thin structure regions, and have less noise in texture-less surfaces.

visible and the two worst views to be not visible. In inference, we select the three best source views for reconstruction. Regarding the number of PatchMatch iterations for each scale, we use $N_{iter}^3, N_{iter}^2, N_{iter}^1 = 8, 2, 1$ in training, and 8, 3, 3 in inference respectively. We finally fuse the estimated depth and normal maps using the fusion method of Galliani et al. [8], where we filter depths by the number of minimum consistent views containing relative depth error below 1%, reprojection error less than 2px and normal angle difference less than 10° . We apply nearest neighbor upsampling with 5×5 median filtering to upscale the depth map to the same resolution of the original image during fusion. We implemented our method in PyTorch and use RTX 3090 for training and evaluation. We use a customized CUDA kernel to accelerate the computation of supporting pixel-aggregated group-wise correlation in inference. These design choices are evaluated in our ablation study (Sec. 4.4, Table 5).

4.2 ETH3D High-res Multi-View Benchmark

ETH3D High-res Multi-View Benchmark is one of the more challenging benchmark in MVS, covering

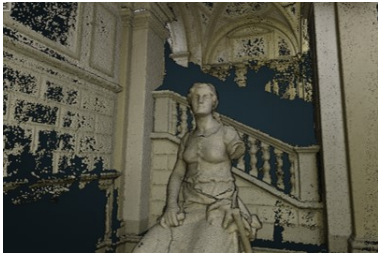
large scale scenes with sparsely sampled, high-resolution images of 6048×4032 . The benchmark contains training and testing scenes, where training scenes contain 7 indoor and 6 outdoor scenes, and testing scenes contain 9 indoor scenes and 3 outdoor scenes. We evaluate our method for both sets of scenes at two different resolution; higher resolution of 3200×2112 (Ours_{hr}) which matches the resolution of ACMM [27], and lower resolution of 1920×1280 (Ours_{lr}) which matches the resolution of Vis-MVSNet [34] and PatchMatch-RL [14]. We use 15 source views and 10 source views for Ours_{lr} and Ours_{hr} respectively.

Table 1 shows the quantitative comparison of our method. Our method achieves the state-of-the-art result among all the learning-based methods even with using the **lower-resolution** images, outperforming EPP-MVSNet [17] by 4.1% and 3.5% in the F_1 score for the train set in the 2cm and 5cm threshold, and differing by 0.8% for the test set in the 5cm threshold. Moreover, Ours_{hr} achieves better F_1 score than ACMP [28] and ACMM [27] for all settings. Specifically, our F_1 scores are higher by 1.2% and 2.9% at 2cm and 5cm thresholds on the training sets and by 4.3% and 4.4% at 2cm and 5cm thresholds on the test sets, compared to ACMM [27]. In Figure 6,

Table 1 Results on the ETH3D high-res Multi-View benchmark train and test sets. All methods do not train on any ETH3D data. Bold and underline denote the method with the highest and the second highest score for each setting, respectively. Ours_{hr} method outperforms all other methods in both 2cm and 5cm criteria in train and test scenes except DeepC-MVS [13], which was trained using the ETH3D train set. Ours_{lr} outperforms all existing learning-based methods at all settings except for EPP-MVSNet [17] at the 2cm benchmark on the test set.

Method	Resolution	Test 2cm: Accuracy / Completeness / F1			Test 5cm: Accuracy / Completeness / F1		
		Indoor	Outdoor	Combined	Indoor	Outdoor	Combined
DeepC-MVS [13]	-	89.1 / <u>85.2</u> / 86.9	89.4 / <u>86.4</u> / 87.7	89.2 / <u>85.5</u> / 87.1	95.3 / 91.1 / <u>93.0</u>	95.7 / <u>92.9</u> / 94.3	95.4 / 91.5 / <u>93.3</u>
ACMP [28]	3200x2130	90.6 / 74.2 / 80.6	<u>90.4</u> / 79.6 / 84.4	90.5 / 75.6 / 81.5	95.4 / 83.2 / 88.4	96.5 / 86.3 / 91.0	95.7 / 84.0 / 89.0
ACMM [27]	3200x2130	91.0 / 72.7 / 79.8	89.6 / 79.2 / 83.6	<u>90.7</u> / 74.3 / 80.8	96.1 / 82.9 / 88.5	<u>97.0</u> / 86.3 / 91.1	96.3 / 83.7 / 89.1
ACMH [27]	3200x2130	<u>91.1</u> / 64.8 / 73.9	84.0 / 80.0 / 81.8	89.3 / 68.6 / 75.9	97.4 / 75.0 / 83.7	94.1 / 87.1 / 90.4	<u>96.6</u> / 78.0 / 85.4
COLMAP [23]	3200x2130	92.0 / 59.7 / 70.4	92.0 / 73.0 / 80.8	92.0 / 63.0 / 73.0	<u>96.6</u> / 73.0 / 82.0	97.1 / 83.9 / 89.7	96.8 / 75.7 / 84.0
PatchmatchNet [25]	2688x1792	68.8 / 74.6 / 71.3	72.3 / 86.0 / 78.5	69.7 / 77.5 / 73.1	84.6 / 85.1 / 84.7	87.0 / 92.0 / 89.3	85.2 / 86.8 / 85.9
EPP-MVSNet [17]	3072x2048	85.0 / 81.5 / 83.0	86.8 / 82.6 / 84.6	85.5 / 81.8 / 83.4	93.7 / 89.7 / 91.6	94.6 / 89.9 / 92.1	93.9 / 89.8 / 91.7
PatchMatch-RL [14]	1920x1280	73.2 / 70.0 / 70.9	78.3 / 78.3 / 76.8	74.5 / 72.1 / 72.4	88.0 / 83.7 / 85.5	92.6 / 89.0 / 90.5	89.2 / 85.0 / 86.8
Ours _{lr}	1920x1280	83.8 / 81.5 / 82.1	85.8 / 83.5 / 84.1	84.3 / 82.0 / 82.6	93.3 / <u>91.6</u> / 92.3	94.8 / 91.5 / 93.0	93.7 / <u>91.6</u> / 92.5
Ours _{hr}	3200x2112	84.3 / 85.8 / <u>84.8</u>	85.5 / 87.4 / <u>86.3</u>	84.6 / 86.2 / <u>85.1</u>	93.4 / 93.4 / 93.3	94.3 / 93.5 / <u>93.8</u>	93.6 / 93.4 / 93.5

Method	Resolution	Train 2cm: Accuracy / Completeness / F1			Train 5cm: Accuracy / Completeness / F1		
		Indoor	Outdoor	Combined	Indoor	Outdoor	Combined
ACMP [28]	3200x2130	92.3 / 72.3 / <u>80.5</u>	87.6 / 72.0 / <u>78.9</u>	90.1 / 72.2 / <u>79.8</u>	96.3 / 81.8 / 88.1	95.6 / 82.7 / 88.6	96.0 / 82.2 / 88.3
ACMM [27]	3200x2130	92.5 / 68.5 / 78.1	88.6 / <u>72.7</u> / 79.7	<u>90.7</u> / 70.4 / 78.9	96.4 / 78.4 / 86.1	96.2 / 83.9 / <u>89.5</u>	96.3 / 80.9 / 87.7
ACMH [27]	3200x2130	<u>92.6</u> / 59.2 / 70.0	84.7 / 64.4 / 71.5	88.9 / 61.6 / 70.7	<u>97.7</u> / 70.1 / 80.5	95.4 / 75.6 / 83.5	<u>96.6</u> / 72.7 / 81.9
COLMAP [23]	3200x2130	95.0 / 52.9 / 66.8	<u>88.2</u> / 57.7 / 68.7	91.9 / 55.1 / 67.7	98.0 / 66.6 / 78.5	<u>96.1</u> / 73.8 / 82.9	97.1 / 69.9 / 80.5
PatchmatchNet [25]	2688x1792	63.7 / 67.7 / 64.7	66.1 / 62.8 / 63.7	64.8 / 65.4 / 64.2	78.7 / 80.0 / 78.9	86.8 / 73.2 / 78.5	82.4 / 76.9 / 78.7
EPP-MVSNet [17]	3072x2048	86.5 / 71.1 / 77.4	78.5 / 63.5 / 70.0	82.8 / 67.6 / 74.0	94.0 / 82.8 / 87.8	93.3 / 78.5 / 85.2	93.6 / 80.8 / 86.6
PatchMatch-RL [14]	1920x1280	76.6 / 60.7 / 66.7	75.4 / 64.0 / 69.1	76.1 / 62.2 / 67.8	89.6 / 76.5 / 81.4	88.8 / 81.4 / 85.7	90.5 / 78.8 / 83.3
Ours _{lr}	1920x1280	85.3 / <u>76.0</u> / 79.9	80.7 / 72.3 / 76.1	83.2 / <u>74.3</u> / 78.1	93.7 / <u>88.7</u> / <u>90.9</u>	92.6 / <u>86.3</u> / 89.3	93.2 / <u>87.6</u> / <u>90.1</u>
Ours _{hr}	3200x2112	86.1 / 78.4 / 81.5	81.4 / 76.3 / 78.6	84.0 / 77.4 / 80.1	94.3 / 88.8 / 91.3	93.3 / 86.6 / 89.7	93.8 / 87.8 / 90.6



ACMM [27]



EPP-MVS [17]



Ours

Fig. 7 Point cloud reconstruction results on ETH3D Dataset. From left to right, we show the reconstruction results of ACMM [27], EPP-MVSNet [17] and our method in the *Statue* from the ETH3D benchmark [24]. We refer to the https://www.eth3d.net/result_details?id=302 for all reconstruction results on the benchmark scenes.

we compare the depth and normal maps of our method against PatchMatch-RL [14]. In Figure 7 we compare the point cloud reconstruction with ACMM [27], EPP-MVSNet [17] and ours. We show that our method achieves far less noise and is more complete in challenging texture-less surfaces compared to the baseline.

4.3 Tanks and Temples Benchmark

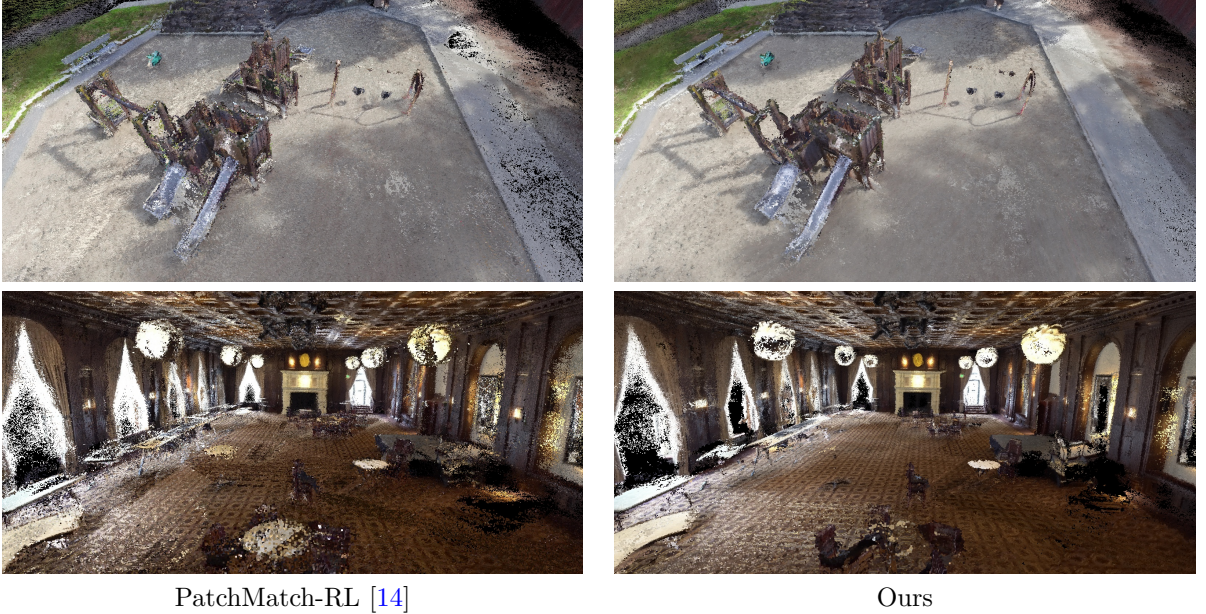
We further evaluate our method on Tanks and Temples [12] intermediate and advanced benchmark. The intermediate benchmark contains 8

scenes captured inward-facing around the center of the scene, and the advanced benchmark contains 6 scenes covering larger area. We use 15 images of resolution 1920×1072 for all benchmarks. Table 2 compares related work.

We improve the F_1 score by 3.5% and 2.0% on Intermediate and Advanced scene compared to PatchMatch-RL [14]. Our method achieves slightly worse results on the intermediate scenes and on-par performance on the advanced scenes compared to ACMM [27], differing by 2.0% and 0.2% in F_1 score respectively. We also note that we

Table 2 Results on the Tanks and Temples benchmark. We show precision, recall and F_1 scores on the intermediate and advanced scenes for each methods. The best performing model is marked as bold.

Method	Precision / Recall / F1	
	Intermediate	Advanced
DeepC-MVS [13]	59.1 / 61.2 / 59.8	40.7 / 31.3 / 34.5
ACMP [28]	49.1 / 73.6 / 58.4	42.5 / 34.6 / 37.4
ACMM [27]	49.2 / 70.9 / 57.3	35.6 / 34.9 / 34.0
COLMAP [23]	43.2 / 44.5 / 42.1	33.7 / 24.0 / 27.2
R-MVSNet [32]	43.7 / 57.6 / 48.4	31.5 / 22.1 / 24.9
CasMVSNet [10]	47.6 / 74.0 / 56.8	29.7 / 35.2 / 31.1
AttMVS [16]	61.9 / 58.9 / 60.1	40.6 / 27.3 / 31.9
PatchmatchNet [25]	43.6 / 69.4 / 53.2	27.3 / 41.7 / 32.3
EPP-MVSNet [17]	53.1 / 75.6 / 61.7	34.6 / 35.6 / 35.7
PatchMatch-RL [14]	45.9 / 62.3 / 51.8	30.6 / 36.8 / 31.8
Ours	50.5 / 63.3 / 55.3	38.8 / 32.4 / 33.8

**Fig. 8 Point Cloud reconstruction results on Tanks and Temples Dataset** We show PatchMatch-RL [14] and Our reconstruction results on the *Playground* and *Ballroom* scene from Tanks and Temples dataset. Our method achieves more detail and less noise.

are the second best learning based method after EPP-MVSNet [17] on the advanced benchmark.

Figure 8 shows the point cloud reconstruction of our method compared to PatchMatch-RL.

4.4 Ablations

Ablations of design choices. We provide an extensive study on how different component affects the overall reconstruction quality by using

ETH3D High-Res multi-view benchmark [24] on the training sets. We show ablation studies at higher resolution 3200×2112 (Ours_{hr}) and lower resolution 1920×1280 (Ours_{lr}) in Table 3 and Table 4 respectively. For a fair comparison, we use the same image resolution and number of views used as our approach in each table. We show that the ablation on Ours_{hr} shows similar results to Ours_{lr} .

Table 3 Ablation Study on ETH3D High-Res Training Set at lower resolution 1920×1280 (Ours_{lr}). We compare how much each module contributes to the overall reconstruction quality. “G. C.” denotes geometric consistency, “U-Net” denotes the feature extractor using shallow VGG-style feature extractor, and “Copl.” denotes dot-product attention instead of learned coplanarity to weight supporting pixels. The memory measures peak GPU reserved memory, and time measures average time used for inference on one reference image. We use the author provided values for PatchMatch-RL [14]. The model with the highest score is marked with bold.

G. C.	Model		Accuracy / Completeness / F1		Mem. (MiB)	Time (s)
	U-Net	Copl.	Train 2CM	Train 5cm		
	PatchMatch-RL		76.1 / 62.2 / 67.8	90.5 / 78.8 / 83.3	7,693	13.54
×	×	×	84.1 / 64.7 / 72.3	94.0 / 78.8 / 85.2	4,520	5.02
✓	×	×	83.9 / 67.0 / 73.8	93.6 / 80.7 / 86.3	4,952	8.78
×	✓	×	84.0 / 70.5 / 76.1	93.9 / 84.6 / 88.8	4,566	5.43
×	×	✓	82.8 / 67.2 / 73.5	92.9 / 81.0 / 86.2	5,058	6.68
×	✓	✓	83.6 / 73.3 / 77.3	93.3 / 86.8 / 89.8	5,092	7.00
✓	×	✓	82.8 / 72.9 / 77.0	92.8 / 86.2 / 89.2	5,320	10.61
✓	✓	×	83.3 / 73.2 / 77.2	93.2 / 86.7 / 89.6	5,012	9.22
✓	✓	✓	83.2 / 74.3 / 78.1	93.2 / 87.6 / 90.1	5,658	10.98

Table 4 Ablation Study on ETH3D High-Res Training Set at higher resolution 3200×2112 (Ours_{hr}). We let “Geo. Cons.” to denote the geometric consistency, “U-Net” to denote the feature extractor using shallow VGG-style feature extractor, and “Copl.” to denote the method using dot-product attention instead of learned coplanarity to weight supporting pixels.

G. C.	Model		Accuracy / Completeness / F1		Mem. (MiB)	Time (s)
	U-Net	Copl.	Train 2CM	Train 5cm		
×	×	×	85.4 / 68.2 / 75.8	94.8 / 80.7 / 87.1	8,502	14.03
✓	×	×	84.8 / 70.8 / 77.4	94.4 / 83.4 / 88.3	12,160	22.12
×	✓	×	85.2 / 73.1 / 78.1	93.8 / 84.8 / 88.9	8,514	16.52
×	×	✓	83.6 / 70.5 / 77.2	93.9 / 83.5 / 88.1	10,050	19.13
×	✓	✓	84.3 / 76.0 / 79.4	94.1 / 86.8 / 90.1	9,523	20.25
✓	×	✓	83.5 / 75.7 / 78.7	93.3 / 86.1 / 89.3	12,982	27.71
✓	✓	×	83.9 / 76.1 / 79.1	93.9 / 86.2 / 89.6	12,230	24.89
✓	✓	✓	84.0 / 77.4 / 80.1	93.8 / 87.8 / 90.6	13,484	28.65

Table 5 Additional Ablation Studies. “Bilinear-JPEG” denotes the image downsampled with bilinear sampling and stored with JPEG format; “Trained with Less-View & Iter.”, the model trained with the same hyperparameters as PatchMatch-RL [14]; “No Fusion Upsampling”, the model fused without nearest-neighbor upsampling followed by median filtering; “Adapt-3” and “Adapt-5”, inference with 3 and 5 adaptively sampled supporting pixels; and “Without CUDA Kernel”, inference without using the custom CUDA kernel. The model with the highest score is marked with bold.

Model	Accuracy / Completeness / F1								Mem. (MiB)	Time (s)		
	Train 2CM				Train 5cm							
Bilinear-JPEG	79.9	/	66.0	/	71.6	91.3	/	81.7	/	85.9	5,658	11.02
Trained with Less-View & Iter.	81.9	/	72.2	/	76.3	92.7	/	86.4	/	89.3	5,658	11.02
No Fusion Upsampling	85.6	/	66.8	/	74.8	94.2	/	85.6	/	89.6	5,658	10.92
Adapt-3	84.3	/	69.6	/	75.7	94.1	/	84.7	/	89.1	5,658	9.45
Adapt-5	84.1	/	73.2	/	77.8	93.9	/	87.0	/	90.1	5,658	10.09
Without CUDA Kernel	83.2	/	74.3	/	78.1	93.2	/	87.6	/	90.1	10,839	24.33
Ours _{lr}	83.2	/	74.3	/	78.1	93.2	/	87.6	/	90.1	5,658	10.98

Importance of geometric consistency: Taking the higher resolution as an example, having

Geometric consistency improves the F_1 score by



Fig. 9 Texture Artifact from JPEG compression. The left image shows the reference image, and the right two shows image downsampled with area interpolation stored with PNG format and image downsampled with bilinear interpolation stored with JPEG format. We show that there is an artifact in texture-light area in JPEG stored files, which causes huge drop in accuracy as shown in Table 5.

1.6% at the 2cm threshold, and 1.2% at the 5cm threshold. Without Geometric consistency, the F_1 score drops by 0.7% at the 2cm threshold and 0.5% at the 5cm threshold. However, the overall runtime of the method increases, taking more than 8 seconds on average per image, since geometric consistency requires multi-stage reconstruction.

Importance of feature extraction backbone:

We compare the CNN-based feature extractor used in PatchMatch-RL [14], which uses shallow VGG-style layers with FPN connections, to our more complex U-Net [22]-based feature extractor. From the baseline, at higher resolution, U-Net [22] based feature extractor improves the F_1 score by 2.3% and 1.8% at the 2cm and the 5cm threshold. Without U-Net [22] based feature extractor, the F_1 score drops by 1.4% at the 2cm threshold and 1.3% at the 5cm threshold.

Dot-product attention vs. Learned coplanarity: We introduce learned coplanarity for computing weights of each supporting pixels, and compare with the dot-product attention based weighting scheme used in PatchMatch-RL [14]. Compared to dot-product attention, at higher resolution, learned coplanarity improves the overall F_1 score by 1.4% at the 2cm threshold and 1.0% at the 5cm threshold. Without the learned coplanarity, the F_1 score drops by 1.0% at both thresholds.

Ablations of other engineering tweaks.

We additionally evaluate other factors that contribute to the improvement of our model over PatchMatch-RL [14]. We experiment on Ours_{lr}. We first show that saving bilinearly downsampled

images with JPEG format reduces the overall F_1 score by 6.5% at the 2cm threshold and 4.2% at the 5cm threshold.

This is due to the artifact in JPEG compression as shown in Figure 9. Next, we show that training with fewer of view selections with fewer iterations (*i.e.*, 1 view as visible from 5 source views with 3, 1, 1 iterations, instead of 3 views as visible from 6 source views with 8, 2, 1 iterations) drops the overall F_1 score by 1.8% and 0.8% at 2cm and 5cm threshold. Using upsampling followed by median filtering boosts the F_1 score at the 2cm threshold by 3.3% and by 0.5% at the 5cm threshold.

We also evaluate with adaptive sampling at inference time. We show that using 3 supporting pixels selected by adaptive sampling (“Adapt-3”) drops F_1 score by 2.4% and 1.0% at the 2cm threshold and the 5cm threshold, but reduces 1.5 seconds for time taken per view. With 5 supporting pixels (“Adapt-5”), the F_1 drops slightly by 0.3% at the 2cm threshold with time taken per view reduced by 0.89 seconds per view.

Finally, we compare the runtime statistics of the model without CUDA kernel code (*i.e.*, pure PyTorch [18] implementation). The CUDA kernel code reduces per-view inference time by 13.35 seconds and memory by 5,181 MiB.

4.5 Additional Qualitative Results

Figure 10 and 11 show our point cloud reconstruction of the ETH3D [24] High-res train and test set for Ours_{hr} and Ours_{lr} respectively. Figure 12



Fig. 10 Additional Qualitative Results for ETH3D High-Res [24] Train set (top two rows) and Test set (bottom two rows) for Ours_{hr}.

and Figure 13 show our results for the Tanks and Temples [12] dataset.

4.6 Limitations

Our method performs well for high resolution sparse views of large scenes in a balanced accuracy and completeness score. However, non-learning based methods such as ACMM [27] achieve higher accuracy, and cost-volume based deep approaches, such as EPP-MVSNet [17], are better suited to smaller scale scenes with dense views. Also, our inference with geometric-consistency is slower than some other deep MVS algorithms (*e.g.* [25]),

due to file IO constraints in the multiscale inference.

5 Conclusion

We have proposed a learning-based PatchMatch MVS method that benefits from learned patch coplanarity, geometric consistency and adaptive pixel sampling. Our experiments show that each of these individually contribute to higher performance and together lead to much more complete models while retaining high accuracy, comparing well to state-of-the-art learning and non-learning based methods on challenging benchmarks.



Fig. 11 Additional Qualitative Results for ETH3D High-Res [24] Train set (top two rows) and Test set (bottom two rows) for Ours_{lr} .

Data Availability Statement. The data that support the findings of this study are openly available at <https://doi.org/10.1109/CVPR42600.2020.00186> (BlendedMVS) [33], <https://doi.org/10.1145/3072959.3073599> (TanksAndTemples) [12] and <https://doi.org/10.1109/CVPR.2017.272> (ETH3D) [24].

Acknowledgments. This research is supported in part by ONR Award N00014-16-1-2007, NSF award IIS 2020227, and gift from Amazon Go Research.

References

- [1] Aanæs, H., Jensen, R.R., Vogiatzis, G., Tola, E., Dahl, A.B.: Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision* **120**(2), 153–168 (2016)
- [2] Anandan, P.: A computational framework and an algorithm for the measurement of visual motion. *International Journal of Computer Vision* **2**(3), 283–310 (1989)



Fig. 12 Additional Qualitative Results for Tanks and Temples [12] Intermediate set.

- [3] Barnes, C., Shechtman, E., Finkelstein, A., Goldman, D.B.: Patchmatch: A randomized correspondence algorithm for structural image editing
- [4] Bleyer, M., Rhemann, C., Rother, C.: Patch-match stereo-stereo matching with slanted support windows.
- [5] Cheng, X., Wang, P., Yang, R.: Learning depth with convolutional spatial propagation network. *IEEE transactions on pattern analysis and machine intelligence* **42**(10), 2361–2379 (2019)
- [6] Furukawa, Y., Hernández, C.: Multi-view stereo: A tutorial. *Found. Trends Comput. Graph. Vis.* **9**, 1–148 (2015)
- [7] Furukawa, Y., Ponce, J.: Accurate camera calibration from multi-view stereo and bundle adjustment. *International Journal of Computer Vision* **84**(3), 257–268 (2009)
- [8] Galliani, S., Lasinger, K., Schindler, K.: Massively parallel multiview stereopsis by surface normal diffusion. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 873–881 (2015)
- [9] Golparvar-Fard, M., Pena-Mora, F., Savarese, S.: Monitoring changes of 3d building elements from unordered photo collections. In: *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*. pp. 249–256. IEEE (2011)
- [10] Gu, X., Fan, Z., Zhu, S., Dai, Z., Tan, F., Tan, P.: Cascade cost volume for high-resolution multi-view stereo and stereo matching. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2020)
- [11] Hannah, M.J.: *Computer matching of areas in stereo images*. Stanford University (1974)
- [12] Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics* **36**(4) (2017)
- [13] Kuhn, A., Sormann, C., Rossi, M., Erdler, O., Fraundorfer, F.: Deepc-mvs: Deep confidence



Fig. 13 Additional Qualitative Results for Tanks and Temples [12] Advanced set.

prediction for multi-view stereo reconstruction. In: 2020 International Conference on 3D Vision (3DV). pp. 404–413. IEEE (2020)

- [14] Lee, J.Y., DeGol, J., Zou, C., Hoiem, D.: Patchmatch-rl: Deep mvs with pixelwise depth, normal, and visibility. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (October 2021)
- [15] Lucas, B.D., Kanade, T.: An iterative image registration technique with an application to stereo vision. In: Proc. International Joint Conference on Artificial Intelligence (IJCAI). p. 674–679 (1981)
- [16] Luo, K., Guan, T., Ju, L., Wang, Y., Chen, Z., Luo, Y.: Attention-aware multi-view stereo. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1590–1599 (2020)
- [17] Ma, X., Gong, Y., Wang, Q., Huang, J., Chen, L., Yu, F.: Epp-mvsnet: Epipolar-assembling based depth prediction for multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 5732–5740 (October 2021)
- [18] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) Advances in Neural Information Processing Systems 32, pp. 8024–8035. Curran Associates, Inc. (2019)
- [19] Prokopenko, K., Dupont, R.: Towards dense 3d reconstruction for mixed reality in healthcare: Classical multi-view stereo vs deep learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops. pp. 0–0 (2019)
- [20] Rebecq, H., Gallego, G., Mueggler, E., Scaramuzza, D.: Emvs: Event-based multi-view stereo—3d reconstruction with an event camera in real-time. International Journal of Computer Vision **126**(12), 1394–1414 (2018)

- [21] Romanoni, A., Matteucci, M.: Tapa-mvs: Textureless-aware patchmatch multi-view stereo. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 10413–10422 (2019)
- [22] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)
- [23] Schönberger, J.L., Zheng, E., Pollefeys, M., Frahm, J.M.: Pixelwise View Selection for Unstructured Multi-View Stereo. In: European Conference on Computer Vision (ECCV) (2016)
- [24] Schöps, T., Schönberger, J.L., Galliani, S., Sattler, T., Schindler, K., Pollefeys, M., Geiger, A.: A multi-view stereo benchmark with high-resolution images and multi-camera videos. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- [25] Wang, F., Galliani, S., Vogel, C., Speciale, P., Pollefeys, M.: Patchmatchnet: Learned multi-view patchmatch stereo (2020)
- [26] Wang, Y., Guan, T., Chen, Z., Luo, Y., Luo, K., Ju, L.: Mesh-guided multi-view stereo with pyramid architecture. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2039–2048 (2020)
- [27] Xu, Q., Tao, W.: Multi-scale geometric consistency guided multi-view stereo. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5483–5492 (2019)
- [28] Xu, Q., Tao, W.: Planar prior assisted patchmatch multi-view stereo. AAAI Conference on Artificial Intelligence (AAAI) (2020)
- [29] Xu, Z., Liu, Y., Shi, X., Wang, Y., Zheng, Y.: Marmvs: Matching ambiguity reduced multiple view stereo for efficient large scale scene reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020)
- [30] Yang, M.D., Chao, C.F., Huang, K.S., Lu, L.Y., Chen, Y.P.: Image-based 3d scene reconstruction and exploration in augmented reality. *Automation in Construction* **33**, 48–60 (2013)
- [31] Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 767–783 (2018)
- [32] Yao, Y., Luo, Z., Li, S., Shen, T., Fang, T., Quan, L.: Recurrent mvsnet for high-resolution multi-view stereo depth inference. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5525–5534 (2019)
- [33] Yao, Y., Luo, Z., Li, S., Zhang, J., Ren, Y., Zhou, L., Fang, T., Quan, L.: Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. *Computer Vision and Pattern Recognition (CVPR)* (2020)
- [34] Zhang, J., Yao, Y., Li, S., Luo, Z., Fang, T.: Visibility-aware multi-view stereo network. *British Machine Vision Conference (BMVC)* (2020)
- [35] Zheng, E., Dunn, E., Jovic, V., Frahm, J.M.: Patchmatch based joint view selection and depthmap estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1510–1517 (2014)