# MVSTER: Epipolar Transformer for Efficient Multi-View Stereo

Xiaofeng Wang[1], Zheng Zhu[2], Fangbo Qin[1], Yun Ye[2], Guan Huang[2], Xu Chi[2], Yijia He[3], and Xingang Wang[1]

[1] Institute of Automation, Chinese Academy of Sciences
{wangxiaofeng2020,qinfangbo2013,xingang.wang}@ia.ac.cGn
[2] PhiGent Robotics    zhengzhu@ieee.org
{yun.ye,guang.huan,xu.chi}@phigent.ai
[3] Kwai Inc.   heyijia2016@gmail.com

**Abstract.** Learning-based Multi-View Stereo (MVS) methods warp source images into the reference camera frustum to form 3D volumes, which are fused as a cost volume to be regularized by subsequent networks. The fusing step plays a vital role in bridging 2D semantics and 3D spatial associations. However, previous methods utilize extra networks to learn 2D information as fusing cues, underusing 3D spatial correlations and bringing additional computation costs. Therefore, we present MVSTER, which leverages the proposed epipolar Transformer to learn both 2D semantics and 3D spatial associations efficiently. Specifically, the epipolar Transformer utilizes a detachable monocular depth estimator to enhance 2D semantics and uses cross-attention to construct data-dependent 3D associations along epipolar line. Additionally, MVSTER is built in a cascade structure, where entropy-regularized optimal transport is leveraged to propagate finer depth estimations in each stage. Extensive experiments show MVSTER achieves state-of-the-art reconstruction performance with significantly higher efficiency: Compared with MVSNet and CasMVSNet, our MVSTER achieves 34% and 14% relative improvements on the DTU benchmark, with 80% and 51% relative reductions in running time. MVSTER also ranks first on Tanks&Temples-Advanced among all published works. Code is available at https://github.com/JeffWang987/MVSTER.

**Keywords:** Multi-view Stereo, Transformer, Depth Estimation, Optimal Transport

## 1 Introduction

Given multiple 2D RGB observations and camera parameters, Multi-View Stereo (MVS) aims to reconstruct the dense geometry of the scene. MVS is a fundamental task in 3D computer vision, with applications ranging from autonomous navigation to virtual/augmented reality. Despite being extensively studied by traditional geometric methods [21,45,52,59] for years, MVS is still challenged

by unsatisfactory reconstructions under conditions of illumination changes, non-Lambertian surfaces and textureless areas [30,46].

Recent researches [66,67] have relieved the aforementioned problems via learning-based methods. Typically, they extract image features through 2D Convolutional Neural Networks (CNN). Then, source features are warped into reference camera frustum to form source volumes, which are fused as a cost volume to produce depth estimations. Fusing source volumes is an essential step in the whole pipeline and many MVS approaches [66,71,58,55,69] put efforts into it. The core of the fusing step is to explore correlations between multi-view images. MVSNet [66] follows the philosophy that various images contribute equally to the 3D cost volume, and utilizes variance operation to fuse different source volumes. However, such fusing method ignores various illumination and visibility conditions of different views. To alleviate this problem, [55,16,22] enrich 2D feature semnatics via Deformable Convolution Network (DCN) [14], and [69,71] leverage extra networks to learn per-pixel weights as a guidance for fusing multi-view features. However, these methods introduce onerous network parameters and restrict efficiency. Besides, they only concentrate on 2D local similarities as a criteria for correlating multiple views, neglecting depth-wise 3D associations, which could lead to inconsistency in 3D space [27].

Therefore, in this paper, we explore an efficient approach to model 3D spatial associations for fusing source volumes. Our intuition is to learn 3D relations from data itself, without introducing extra learning parameters. Recent success in attention mechanism prompts that Transformer [53] is appropriate for modeling 3D associations. The key advantage of Transformer is it leverages cross-attention to build data-dependent correlations, introducing minimal learnable parameters. Besides, compared with CNN, Transformer has expanded receptive field, which is more adept at constructing long-range 3D relations. Therefore, we propose the epipolar Transformer, which efficiently builds multi-view 3D correlations along the epipolar line. Specifically, we firstly leverage an auxiliary monocular depth estimator to enhance the 2D semantics of the *query* feature. The auxiliary branch guides our network to learn depth-discriminative features, and it can be detached after training, which brings no extra computation cost. Subsequently, cross-attention is utilized to model 3D associations explicitly from features on epipolar lines, without introducing sophisticated networks. Additionally, we formulate the depth estimation as a depth-aware classification problem and solve it with entropy-regularized optimal transport [40], which propagates finer depth maps in a cascade structure.

Owing to the epipolar Transformer, MVSTER obtains enhanced reconstruction results with fewer depth hypotheses. Compared with MVSNet [66] and CasMVSNet [25], our method reduces 88% and 73% relative depth hypotheses, making 80% and 51% relative reduction in running time, yet obtaining 34% and 14% relative improvements on the DTU benchmark, respectively. Besides, our method ranks first among all published works on Tanks&Temples-Advanced. The main technique contributions are four-fold as follows:

- We propose a novel end-to-end Transformer-based method for multi-view stereo, named MVSTER. It leverages the proposed epipolar Transformer to efficiently learn 3D associations along epipolar line.

- An auxiliary monocular depth estimator is utilized to guide the *query* feature to learn depth-discriminative information during training, which enhances feature semantics yet brings no efficiency compromises.

- We formulate depth estimation as a depth-aware classification problem and solve it with the entropy-regularized optimal transport, which produces finer depth estimations propagated in the cascade structure.

- Extensive experiments on DTU, Tanks&Temples, BlendedMVS, and ETH3D show our method achieves superior performance with significantly higher efficiency than existing methods.

## 2    Related Work

**Learning-based MVS**  With the rapid progress of deep learning in 3D perception [41,74,42,48,28,37,18,35], the MVS community is gradually dominated by learning-based methods [66,67,62,55,58,36,25,69]. They achieve better reconstruction results than traditional methods [21,45,5,20,52]. Learning-based MVS approaches project source images into reference camera frustum to form multiple 3D volumes, which are fused through variance operation [66,67,62,25,64,9]. Such a fusing method follows the philosophy that the feature volumes from various source images contribute equally [66], neglecting heterogeneous illumination and scene content variability [69]. To remedy the aforementioned problem, PVA-MVSNet [69] proposes a self-adaptive view aggregation module to learn the different significance in source volumes. Vis-MVSNet [71] computes pixel-visibility to represent matching quality, which serves as a volume fusing weight. AA-RMVSNet [58] leverages expensive DCNs [14] to enhance intra-view semantics, and it aggregates inter-view with pixel-wise weight. However, these methods use CNN-based module aggregating local features as fusing guidance, which lacks long-range 3D associations and thus restricts their performance under challenging conditions. Besides, such aggregation modules bring extra computation cost burdening the network. In contrast, the proposed epipolar Transformer learns both 2D semantics and 3D spatial relations from data itself, without bringing onerous network parameters.

**Efficient MVS**  To construct an efficient MVS pipeline, cascade-structured methods [25,10,64,36] are proposed. They address MVS problem in a coarse to fine manner, assuming decreasing depth hypotheses along reference camera frustum at each stage. PatchmatchNet [55] and IterMVS [54] further decrease hypothesized depth number and discard expensive 3D CNN regularization in the cascade structure. However, they achieve high efficiency with significant performance compromises. Additionally, cascade methods have difficulty to recover from errors introduced at coarse resolutions [25]. In this paper, cascade structure

is leveraged to boost efficiency, and optimal transport is utilized to produce finer depth estimations at each stage of the cascade structure.

**Transformers in 3D Vision** Transformers [53,2,51,15,43] find their initial applications in natural language processing and have drawn attention from computer vision community [17,34,6,63,8,56]. In tasks for 3D vision, PYVA [65] and NEAT [11] use cross-attention to build correlations between bird's eye view and front view. STTR [32] formulates stereo depth estimation as a sequence-to-sequence correspondence problem that is optimized by self-attention and cross-attention. Recently, Transformer extends its application to MVS. LANet [72], TransMVSNet [16] and MVSTR [75] introduce an attention mechanism extracting dense features with global contexts, which expands the network receptive field. However, these methods densely correlate each pixel within 2D feature maps, which makes significant efficiency compromises. On the contrary, our epipolar Transformer leverages geometric knowledge, restricting attention associations within the epipolar line, which significantly reduces dispensable feature correlations and makes our pipeline more efficient. Besides, MVSTER only leverages the essential cross-attention of Transformer [53], without introducing sophisticated architecture (*i.e.*, position encoding, Feedforward Neural Network (FNN) and self-attention), which further boosts efficiency.

**Auxiliary Task Learning** Auxiliary branch learning is demonstrated effective in multiple vision tasks [26,73,38]. In general, the auxiliary tasks are selected to be positively related to the main task, thus taking effect during training. In addition, the branch can be discarded after training, bringing no burden during inference. ManyDepth [57] is a self-supervised monocular depth estimator, utilizing MVS cost volume as an auxiliary branch, which enhances estimation reliability. This inspires us that MVS assists monocular depth estimation, and vice versa. Therefore, an auxiliary monocular depth estimation branch is leveraged in MVSTER to learn depth-discriminative features.

## 3   Method

In this section, we give a detailed description of MVSTER. The network architecture is illustrated in Fig. 1. Given a reference image and its corresponding source images, we firstly extract 2D multi-scale features using Feature Pyramid Network (FPN) [33]. Source image features are then warped into reference camera frustum to construct source volumes via differentiable homography (Sec. 3.1). Subsequently, we leverage the epipolar Transformer to aggregate source volumes and produce the cost volume, which is regularized by lightweight 3D CNNs to make depth estimations (Sec. 3.2). Our pipeline is further built in a cascade structure, propagating depth map in a coarse to fine manner (Sec. 3.3). To reduce erroneous depth hypotheses during depth propagating, we formulate depth estimation as a depth-aware classification problem and optimize it with optimal transport. Finally, the network losses are given (Sec. 3.4).
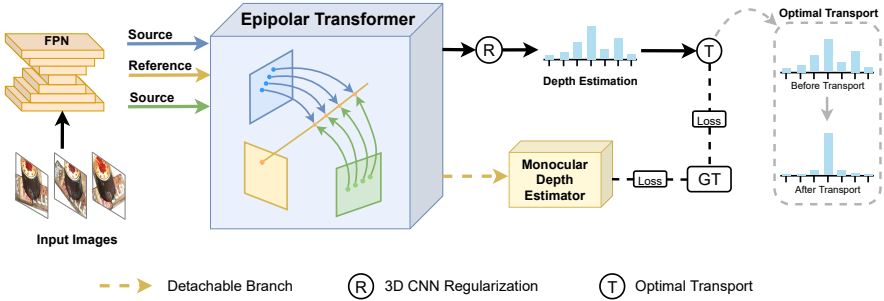
**Fig. 1.** MVSTER architecture. MVSTER firstly extracts features via FPN, then the multi-view features are aggregated by the epipolar Transformer, where the auxiliary branch makes monocular depth estimation to enhance context. Subsequently, the aggregated feature volume is regularized by 3D CNNs, producing depth estimations. Finally, optimal transport is utilized to optimize the predicted depth.

### 3.1   2D Encoder and 3D Homography

Given a reference image $\mathbf{I}_{i=0} \in \mathbb{R}^{H \times W \times 3}$ and its neighboring source images $\mathbf{I}_{i=1,\dots,N-1} \in \mathbb{R}^{H \times W \times 3}$, the first step is to extract the multi-scale 2D features of these inputs. A FPN-like network is applied, where the images are downscaled $M$ times to build deep features $\mathbf{F}_{i=0,\dots,N-1}^{k=0,\dots,M-1} \in \mathbb{R}^{H_k \times W_k \times C_k}$. The scale $k = 0$ denotes the original size of images. The subsequent formulations can be generalized to a specific scale $k$, so $k$ is omitted for simplicity.

Following previous learning-based methods [66,67,55,16], we utilize plane sweep stereo [12] that establishes multiple front-to-parallel planes in the reference view. Specifically, equipped with camera intrinsic parameters $\{\mathbf{K}_i\}_{i=0}^{N-1}$ and transformations parameters $\{[\mathbf{R}_{0,i} \mid \mathbf{t}_{0,i}]\}_{i=1}^{N-1}$ from source views to reference view, source features can be warped into the reference camera frustum:

$$\mathbf{p}_{s_i,j} = \mathbf{K}_i \cdot \left(\mathbf{R}_{0,i} \cdot \left(\mathbf{K}_0^{-1} \cdot \mathbf{p}_r \cdot d_j\right) + \mathbf{t}_{0,i}\right), \tag{1}$$

where $d_j$ denotes $j$-th hypothesized depth of pixel $\mathbf{p}_r$ in the reference feature, and $\mathbf{p}_{s_i,j}$ is the corresponding pixel in the $i$-th source features. After the warping operation, $N-1$ source volumes $\{\mathbf{V}_i\}_{i=1}^{N-1} \in \mathbb{R}^{H \times W \times C \times D}$ are constructed, where $D$ is the total number of hypothesized depths.

### 3.2   Epipolar Transformer

Next, we introduce the epipolar Transformer to aggregate source volumes from different views. The original attention function in Transformer [53] can be described as mapping a *query* and a set of *key-value* pairs to an output. Similarly, in the proposed epipolar Transformer, the reference feature is leveraged as the user *query* to match source features (*keys*) along the epipolar line, thus enhancing the corresponding depth *value*. Specifically, we enrich the reference *query* via

an auxiliary task of monocular depth estimation. Subsequently, cross-attention computes associations between *query* and source volumes under epipolar constraint, generating attention guidance to aggregate the feature volumes from different views. The aggregated features are then regularized by lightweight 3D CNNs. In the following, we firstly give details about the *query* construction, then elaborate on the epipolar Transformer guided feature aggregation. Finally, the lightweight regularization strategy is given.

**Query Construction** As aforementioned, we deem the reference feature as a *query* for the epipolar Transformer. However, features extracted by shallow 2D CNNs become less discriminative at non-Lambertian and low-texture regions. To remedy this problem, [55,16,58,22] utilize expensive DCNs [14] or ASPP [49] to enrich features. In contrast, we propose a more efficient way to enhance our *query*: building an auxiliary monocular depth estimation branch to regularize the *query* and learn depth-discriminative features.

A common decoder [24] used in the monocular depth estimation task is applied in our auxiliary branch. Given multi-scale reference features $\{\mathbf{F}_0^k\}_{k=0}^{M-1}$ that are extracted via FPN, we expand a low resolution feature map through interpolation, and concatenate it with the subsequent scale feature. The aggregated feature maps are fed into regression head [24,23] to make monocular depth estimations:
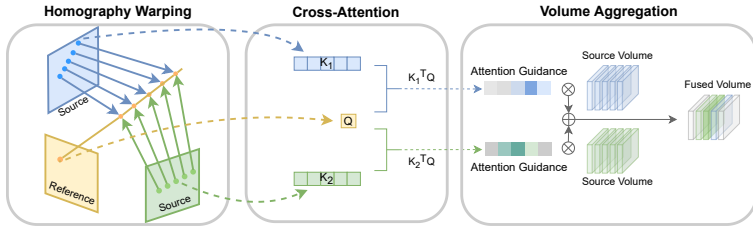
$$\mathbf{M}_k = \mathbf{\Phi}([\mathbf{I}(\mathbf{F}_0^k), \mathbf{F}_0^{k+1}]), \tag{2}$$

where $\mathbf{\Phi}(\cdot)$ is monocular depth decoder, $\mathbf{I}(\cdot)$ is the interpolation function and $[\cdot, \cdot]$ denotes concatenation operation. Subsequently, the monocular depth estimation is repeated for queries with different scales. Notably, such auxiliary branch is only used in the training phase, guiding our network to learn depth-aware features.
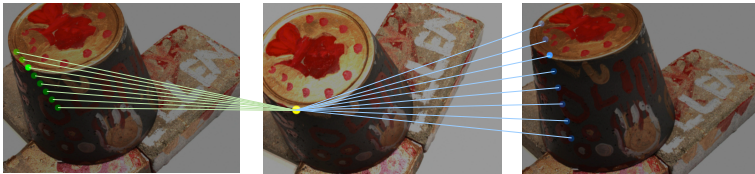
**Epipolar Transformer Guided Aggregation** The aggregation pipeline is depicted in Fig. 2(a), which aims at building 3D associations of the *query* feature. However, depth-wise 3D spatial information is not explicitly delivered by the 2D query feature map, so we firstly restore the depth information via homography warping. According to the warping operation in Equation (1), the hypothesized depth locations of *query* feature $\mathbf{p}_r$ are projected onto the source image epipolar line, resulting in the source volume features $\{\mathbf{p}_{s_i,j}\}_{j=0}^{D-1}$, which are regarded as the *keys* for the epipolar Transformer. Consequently, the *key* features along the epipolar line are leveraged to construct depth-wise 3D associations of the *query* feature, which is implemented with the cross-attention operation:

$$\mathbf{w}_i = \mathrm{softmax}(\frac{\mathbf{v_i}^\mathsf{T}\mathbf{p}_r}{t_e\sqrt{C}}), \tag{3}$$

where $\mathbf{v_i} \in \mathbb{R}^{C \times D}$ is calculated by stacking $\{\mathbf{p}_{s_i,j}\}_{j=0}^{D-1}$ along depth dimension, $t_e$ is the temperature parameter, and $\mathbf{w}_i$ is the attention correlating *query* and *keys*. We visualize an example of real images in Fig. 2(b), where the attention focuses on the most matched location on the epipolar line.

(a) Epipolar Transformer guided aggregation



(b) Visualization of attention on the DTU dataset

**Fig. 2.** (a) The epipolar Transformer aggregation. Homography warping is leveraged to restore depth-wise information of the reference feature, then cross-attention computes 3D associations between query and source volumes under epipolar constraint, generating attention guidance to aggregate the feature volumes from different views. (b) Visualization of the cross-attention score on scan 1 of the DTU dataset, where the opacity of dots on the epipolar line represents the attention score.

The calculated attention $\mathbf{w}_i$ between *query* and *keys* is utilized to aggregate *values*. As for the Transformer *value* design, we follow [55,60,71] to use group-wise correlation, which measures the visual similarity between reference feature and source volumes in an efficient manner:

$$\mathbf{s}_i^g = \frac{1}{G} \left\langle \mathbf{v}_i^g, \mathbf{p}_r^g \right\rangle, \tag{4}$$

where $g = 0, ..., G-1$, $\mathbf{v}_i^g \in \mathbb{R}^{\frac{C}{G} \times D}$ is the $g$-th group feature of $\mathbf{v}_i$, $\mathbf{p}_r^g \in \mathbb{R}^{\frac{C}{G} \times 1}$ is the $g$-th group feature of $\hat{\mathbf{p}}_r$, and $\langle \cdot, \cdot \rangle$ is the inner product. $\{\mathbf{s}_i^g\}_{g=0}^{G-1}$ are then stacked along channel dimension to get $\mathbf{s}_i \in \mathbb{R}^{G \times D}$, which is the *value* for our Transformer. Finally, *values* are aggregated by epipolar attention score $\mathbf{w}_i$ to determine the final cost volume:

$$\mathbf{c} = \frac{\sum_{i=1}^{N-1} \mathbf{w}_i \mathbf{s}_i}{\sum_{i=1}^{N-1} \mathbf{w}_i}. \tag{5}$$

In summary, for the proposed epipolar Transformer, a detachable monocular depth estimation branch is firstly leveraged to enhance depth-discriminative 2D semantics, then the cross-attention between *query* and *keys* is utilized to construct depth-wise 3D associations. Finally, the combined 2D and 3D information serves as guidance for aggregating different views. As shown in Equation (3)-(5),

the epipolar Transformer is designed as an efficient aggregation module, where no learnable parameter is introduced, and the epipolar Transformer only learns data-dependent associations.

**Lightweight Regularization** Due to non-Lambertian surfaces or object occlusions, the raw cost volume is noise-contaminated [66]. To smoothen the final depth map, 3D CNNs are utilized to regularize the cost volume. Considering we have embedded 3D associations into the cost volume, depth-wise feature encoding is omitted in our 3D CNNs, which makes it more efficient. Specifically, we reduce convolution kernel size from $3 \times 3 \times 3$ to $3 \times 3 \times 1$, only aggregating cost volume along feature width and height. The regularized probability volume $\mathbf{P} \in \mathbb{R}^{H \times W \times D}$ is highly desirable in per-pixel depth confidence prediction, which is leveraged to make depth estimations in the cascade structure.

### 3.3   Cascade Depth Map Propagation

Cascade structure is proven effective in stereo depth estimation [47,19,50], monocular reconstruction [4] and MVS [25,10,55], which brings efficiency and enhanced performance. Following [55], a four-stage searching pipeline is set for MVSTER, where the resolutions of inputs for the four stages are $H \times W \times 64$, $\frac{H}{2} \times \frac{W}{2} \times 32$, $\frac{H}{4} \times \frac{W}{4} \times 16$ and $\frac{H}{8} \times \frac{W}{8} \times 8$ respectively. Following [60,55], the inverse depth sampling is utilized to initialize depth hypotheses in the first stage, which is equivalent to equidistant sampling in pixel space. To propagate depth map in a coarse to fine manner, the depth hypotheses of each stage are centered at the previous stage's depth prediction, and $D_k$ hypotheses are uniformly generated within the hypothesized depth range.

### 3.4   Loss

Although cascade structure benefits from coarse to fine pipeline, it has difficulty recovering from errors introduced at previous stages [25]. To alleviate this problem, a straightforward way is to generate a finer depth map at each stage, especially avoiding predicting depth far away from the ground truth. However, previous methods [67,16,58] simply regard depth estimation as a multi-class classification problem, which treats each hypothesized depth equally without considering the distance relationship between them. For example in Fig. 3, given a ground truth depth probability distribution, the cross-entropy losses of case 1 and case 2 are the same. However, the depth prediction of case 1 is out of the valid range and can not be properly propagated to the next stage.

In this paper, the depth prediction is formulated as a depth-aware classification problem, which emphasizes the penalty of the predicted depth that is distant from the ground truth. Specifically, we measure the distance between the predicted distribution $\mathbf{P}_i \in \mathbb{R}^D$ and the ground truth distribution $\mathbf{P}_{\theta,i} \in \mathbb{R}^D$ with the off-the-shelf Wasserstein distance [3]:

$$d_w(\mathbf{P}_i, \mathbf{P}_{\theta,i}) = \inf_{\gamma \in \Pi(\mathbf{P}_i, \mathbf{P}_{\theta,i})} \sum_{x,y} |x - y| \gamma(x, y), \qquad (6)$$
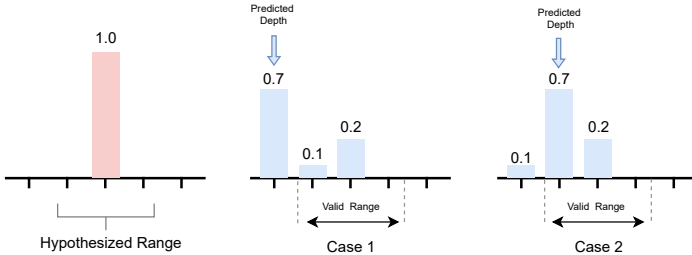
**Fig. 3.** Example illustrating that cross-entropy loss is not aware of the relative distance between each hypothesized depth. The left-most subfigure is the ground truth. Case 1, Case 2 are two predicted depth distributions.

where inf stands for infimum, and $\Pi(\mathbf{P}_i, \mathbf{P}_{\theta,i})$ is the set of all possible distributions whose marginal distributions are $\mathbf{P}_i$ and $\mathbf{P}_{\theta,i}$, which satisfies $\sum_x \gamma(x,y) = \mathbf{P}_i(y)$ and $\sum_y \gamma(x,y) = \mathbf{P}_{\theta,i}(x)$. Such formulation is inspired by the optimal transport problem [40] that calculates the minimum work transporting $\mathbf{P}_i$ to $\mathbf{P}_{\theta,i}$, which can be differentially solved via the sinkhorn algorithm [13].

In summary, the loss function consists of two parts: Wasserstein loss measuring the distance between predicted depth distribution and ground truth, and $L_1$ loss optimizing monocular depth estimation:

$$Loss = \sum_{k=0}^{M-1} \sum_{i \in \mathbf{p}_{\text{valid}}} d_w(\mathbf{P}_i^k, \mathbf{P}_{\theta,i}^k) + \lambda L_1(\mathbf{M}_i^k, \mathbf{P}_{\theta,i}^k), \tag{7}$$

where $\mathbf{p}_{\text{valid}}$ refers to the set of valid ground truth pixels, and $\lambda$ is the loss weight. The total loss is calculated for $M$ stages.

## 4 Experiments

### 4.1 Datasets

MVSTER is evaluated on DTU [1], Tanks&Temples [30], BlendedMVS [68] and ETH3D [46] to verify its effectiveness. Among the four datasets, DTU is an indoor dataset under laboratory conditions, which contains 124 scenes with 49 views and 7 illumination conditions. Following MVSNet [66], DTU is split into `training`, `validation` and `test` set. Tanks&Temples is a public benchmark providing realistic video sequences, which is divided into the intermediate set and a more challenging advanced set. BlendedMVS is a large-scale synthetic dataset that contains 106 `training` scans and 7 `validation` scans. ETH3D benchmark introduces high-resolution images with strong view-point variations, which is split into `training` and `test` sets. As for the evaluation metrics, DTU, Tanks&Temples, and ETH3D evaluate point cloud reconstructions using overall metrics [1] and $F_1$ score [30,46]. BlendedMVS evaluates depth map estimations

using depth-wise metric [68]: EPE stands for $L_1$ distance between predicted depth map and ground truth, $e_1$ and $e_3$ represent the proportion of pixels with depth error larger than 1 and 3.

## 4.2   Implementation Details

Following the common practice [39,54,16], MVSTER is firstly trained on DTU `training` set and evaluated on DTU `test` set, then it is finetuned on BlendedMVS before being tested on Tanks&Temples and ETH3D benchmark. For DTU training, we use ground truth provided by MVSNet [66], whose depth range is sampled from 425mm to 935mm. The input view selection and data pre-processing are the same as [55]. For BlendedMVS, we use the original image resolution and the number of input images is set as 7.

The hypothesized depth numbers $\{D_k\}_{k=0,...,3}$ for each stage are set as 8, 8, 4, 4. Following [39,58], the hypothesized number of the 1st stage is doubled when MVSTER is tested on Tanks&Temples and ETH3D. The group correlation $\{G_k\}_{k=0,...,3}$ are set as 8, 8, 4, 4. For inverse depth sampling, the inverse depth range $R_k$ satisfies $\frac{1}{R_k} = \frac{1}{D_{k-1}-1}\frac{1}{R_{k-1}}$. For the epipolar Transformer, the temperature parameter $t_e$ is set as 2. And the loss weight $\lambda$ is set as 0.0003 in the experiments. We train MVSTER for 10 epochs and optimize it with Adam [29] ($\beta_1 = 0.9, \beta_2 = 0.999$). MVSTER is trained on four NVIDIA RTX 3090 GPUs with batch size 2 on each GPU. The learning rate is initially set as 0.001, which decays by a factor of 2 after 6, 8 and 9 epochs.

For point cloud reconstruction, we follow previous methods [67,66,25] to use both geometric and photometric constraints for depth filtering. We set the view consistency number and the photometric probability threshold as 4 and 0.5, respectively. The final depth fusion steps also follow previous methods [67,66,25].

## 4.3   Benchmark Performance

**Evaluation on DTU** We compare MVSTER with traditional methods [21,45,52], published learning-based methods [66,67,70,69,25,10,55,58] and approaches from recent technical reports [16,75,39,54]. The input images are set as different resolutions (MVSTER*: 1600×1200 and MVSTER: 864×1152) to compare with previous methods, and the number of views is set as 5. The quantitative results are shown in Table 1, where MVSTER* achieves a state-of-the-art overall score and completeness score among all the competitors. Significantly, the inference time of MVSTER* is 0.17s, which is faster than the previous fastest method [55]. Additionally, MVSTER with lower resolution (864×1152) still outperforms all published works, and it runs at 0.09s per image with 2764 MB GPU memory consumption, which sets a new state of the art for efficient learning-based MVS. Qualitative comparisons are shown in Fig. 4, where MVSTER reconstructions provide denser results with finer details. Especially, compared with CasMVSNet [25] and PatchmatchNet [55], MVSTER recovers more details at object boundaries and textureless areas.

**Table 1.** Quantitative results of different methods on the DTU `evaluation` set. Methods with * denote their input resolution is $1600 \times 1200$. The last four methods with gray font come from technical reports.

| Method | Acc.↓ | Comp.↓ | Overall↓ | Runtime (s)↓ |
|---|---|---|---|---|
| Gipuma [21] | **0.283** | 0.873 | 0.578 | - |
| COLMAP [45] | 0.400 | 0.664 | 0.532 | - |
| Tola [52] | 0.342 | 1.190 | 0.766 | - |
| MVSNet [66] | 0.396 | 0.527 | 0.462 | 0.85 |
| R-MVSNet [67] | 0.383 | 0.452 | 0.417 | 0.89 |
| Fast-MVSNet [70] | 0.336 | 0.403 | 0.370 | 0.37 |
| CVP-MVSNet* [69] | 0.296 | 0.406 | 0.351 | 1.12 |
| CasMVSNet [25] | 0.325 | 0.385 | 0.355 | 0.35 |
| UCS-Net* [10] | 0.338 | 0.349 | 0.344 | 0.32 |
| PatchmatchNet* [55] | 0.427 | 0.277 | 0.352 | 0.18 |
| AA-RMVSNet [58] | 0.376 | 0.339 | 0.357 | - |
| MVSTER | 0.350 | 0.276 | 0.313 | **0.09** |
| MVSTER* | 0.340 | **0.266** | **0.303** | 0.17 |
| TransMVSNet [16] | 0.321 | 0.289 | 0.305 | 0.99 |
| MVSTR [75] | 0.356 | 0.295 | 0.326 | 0.81 |
| UniMVSNet [39] | 0.352 | 0.278 | 0.315 | - |
| IterMVS* [54] | 0.373 | 0.354 | 0.363 | 0.18 |



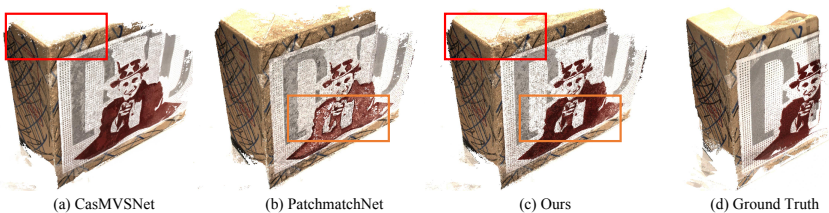(a) CasMVSNet        (b) PatchmatchNet        (c) Ours        (d) Ground Truth

**Fig. 4.** Reconstruction results on DTU (scan 13). Our method delivers more accurate boundaries in the red bounding box than CasMVSNet [25], and preserves more detials of textureless area in the orange bounding box than PatchmatchNet [55].

**Table 2.** Quantitative results on Tanks&Temples-advanced. The evaluation metric is the mean F-score and the last four methods with gray font come from technical reports.

| Method | Mean F-score | Aud. | Bal. | Cou. | Mus. | Pal. | Tem. |
|---|---|---|---|---|---|---|---|
| COLMAP [45] | 27.24 | 16.02 | 25.23 | 34.70 | 41.51 | 18.05 | 27.94 |
| ACMH [59] | 34.02 | 23.41 | 32.91 | **41.17** | 48.13 | 23.87 | 34.60 |
| R-MVSNet [67] | 29.55 | 19.49 | 31.45 | 29.99 | 42.31 | 22.94 | 31.10 |
| CasMVSNet [25] | 31.12 | 19.81 | 38.46 | 29.10 | 43.87 | 27.36 | 28.11 |
| PatchmatchNet [55] | 32.31 | 23.69 | 37.73 | 30.04 | 41.80 | 28.31 | 32.29 |
| EPP-MVSNet [36] | 35.72 | 21.28 | 39.74 | 35.34 | 49.21 | 30.00 | 38.75 |
| AA R-MVSNet [58] | 33.53 | 20.96 | 40.15 | 32.05 | 46.01 | 29.28 | 32.71 |
| MVSTER | **37.53** | 26.68 | 42.14 | 35.65 | 49.37 | 32.16 | 39.19 |
| TransMVSNet [16] | 37.00 | 24.84 | **44.69** | 34.77 | 46.49 | **34.69** | 36.62 |
| MVSTR [75] | 32.85 | 22.83 | 39.04 | 33.87 | 45.46 | 27.95 | 27.97 |
| UniMVSNet [39] | **38.96** | **28.33** | 44.36 | 39.74 | **52.89** | 33.80 | 34.63 |
| IterMVS [54] | 34.17 | 25.90 | 38.41 | 31.16 | 44.83 | 29.59 | 35.15 |

**Evaluation on Tanks&Temples** MVSTER is tested on Tanks&Temples to demonstrate the generalization ability. We use the original image resolution and set the number of views as 7. The depth range, camera parameters, and view selection strategies are aligned with PatchmatchNet [55]. And we follow the dynamic consistency checking method in depth filtering [62]. We compare MVSTER with traditional methods [45,59], published learning-based methods [67,25,55,36,58], and approaches from recent technique reports [16,75,39,54]. Advanced set quantitative results are shown in Table 2, where MVSTER achieves the highest mean F-score among all published works, and the inference time per image is 0.26s. Although our performance is 1.4% lower than the recent UniMVS-Net [39], the inference speed of MVSTER is 3× faster than UniMVSNet[4]. For Tanks&Temples-Intermediate, MVSTER achieves a 60.92 mean F-score, which is 7.8% better than the previous most efficient method [55]. Overall, MVSTER shows strong generalization ability with great efficiency.

**Evaluation on ETH3D** For evaluation on the ETD3D dataset, the input images are resized to $1920 \times 1280$ and the number of inputs is set as 7. The depth range, camera parameters, and view selection strategies are aligned with Patch-matchNet [55]. We compare MVSTER with traditional methods [45,20,45,59], published learning-based methods [55,31] and approaches from technique reports [61,54]. The running time per image is 0.30s and the quantitative results are shown in Table 3. On the `training` set, MVSTER achieves better $F_1$-score than the most competitive traditional method ACMH [59] and the recent IterMVS [54]. On the `test` set, our method obtains 8.9% improvement over the previous most efficient method [55], which is comparable to the recent IterMVS [54]. This demonstrates MVSTER can be well generalized to high-resolution images.

**Table 3.** Quantitative results on the ETH3D benchmark, which is split into a `training` set and a `test` set. The last two methods with gray font come from technical reports.

| Methods | Training set | | | Test set | | |
|---|---|---|---|---|---|---|
| | Acc. | Comp. | $F_1$ -score | Acc. | Comp. | $F_1$ -score |
| Gipuma [45] | 84.44 | 34.91 | 36.38 | 86.47 | 24.91 | 45.18 |
| PMVS [20] | 90.23 | 32.08 | 46.06 | 90.08 | 31.84 | 44.16 |
| COLMAP [45] | **91.85** | 55.13 | 67.66 | 91.97 | 62.98 | 73.01 |
| ACMH [59] | 88.94 | 61.59 | 70.71 | **89.34** | 68.62 | 75.89 |
| PatchmatchNet [55] | 64.81 | 65.43 | 64.21 | 69.71 | 77.46 | 73.12 |
| PatchMatch-RL [31] | 76.05 | 62.22 | 67.78 | 74.48 | 72.89 | 72.38 |
| MVSTER | 76.92 | **68.08** | **72.06** | 77.09 | **82.47** | **79.01** |
| PVSNet [61] | 67.84 | 69.66 | 67.48 | 66.41 | 80.05 | 72.08 |
| IterMVS [54] | 79.79 | 66.08 | 71.69 | 84.73 | 76.49 | **80.06** |

## 4.4   Ablation Study

The ablation study is conducted to analyze the effectiveness of each component, which is measured with DTU's point cloud reconstruction metric [1] and Blend-

---

[4] The inference time of UniMVSNet [39] is not reported, but its baseline is CasMVSNet [25], whose inference speed is 3× slower than MVSTER.

edMVS's depth estimation metric [68]. Unless specified, the image resolutions for DTU and BlendedMVS are $864 \times 1152$ and $576 \times 768$.

**Epipolar Transformer (ET)** Existing methods for aggregating different views in learning-based MVS can be categorized as two types: (*i*) variance fusing [66,25,10,64,67], (*ii*) CNN-based fusing [58,69,71,55]. In this experiment, the aforementioned two methods are compared with the ET module under three conditions[5]. The quantitative results are concluded in Table 4. We observe that reducing hypothesized depth number can significantly decrease inference time. Compared with the hypothesized number used by MVSNet [66] and CasMVS-Net [25], our method relatively reduces 70% and 53% running time. However, the variance fusing strategy shows restricted improvement in the third condition with fewest hypothesized number. CNN-based fusing alleviates the problem by enhancing local visual similarity, but it relatively brings 45%, 46%, 89% computation cost in three cases. In contrast, the ET module shows consistent performance improvement under different hypothesized cases, which demonstrates that 3D spatial associations are beneficial for aggregating multi-view features. Significantly, the ET module learns data-dependent fusing guidance, introducing minimal network parameters and bringing no extra computation cost.

**Table 4.** Quantitative results of different fusing methods under conditions with different hypothesized depth numbers.

| Method | Hypo. Num. | Overall↓ | EPE ↓ | $e_1$↓ | $e_3$↓ | Runtime (s)↓ | Param (M)↓ |
|---|---|---|---|---|---|---|---|
| Variance Fusion | 192 | 0.460 | 1.62 | 19.34 | 9.84 | 0.40 | 0.34 |
| CNN Fusion | 192 | 0.442 | 1.58 | **17.89** | 9.47 | 0.58 | 0.35 |
| ET (Ours) | 192 | **0.435** | **1.54** | 17.93 | **9.32** | **0.40** | **0.34** |
| Variance Fusion | 48,32,8 | 0.335 | 1.28 | 14.82 | 7.55 | 0.28 | 0.93 |
| CNN Fusion | 48,32,8 | 0.327 | **1.07** | 14.33 | 7.03 | 0.41 | 0.94 |
| ET (Ours) | 48,32,8 | **0.323** | 1.09 | **14.17** | **6.89** | **0.28** | **0.93** |
| Variance Fusion | 8,8,4,4 | 0.334 | 1.39 | 15.32 | 7.92 | 0.09 | 0.98 |
| CNN Fusion | 8,8,4,4 | 0.320 | 1.33 | **14.80** | 7.32 | 0.17 | 1.01 |
| ET (Ours) | 8,8,4,4 | **0.313** | **1.31** | 14.98 | **7.27** | **0.09** | **0.98** |

**Monocular Depth Estimator (MDE)** In this experiment, the proposed MDE module is compared with DCN used by [55,16,58] and ASPP used by [22]. The results are shown in Table 5. We observe that the ASPP restricts the reconstruction performance, and DCN enhances reconstruction results with reduced depth error $e_1$. However, DCN brings high computation costs and introduces onerous learning parameters. In contrast, the MDE module shows comparable performance improvement with DCN but introduces no extra computation burden. We also provide an example in Fig. 5, where the MDE module enhances features details at object boundaries, which could reduce ambiguity for depth estimations within boundary areas.

---

[5] Three hypothesized depth number: (*i*) $D : 192$ used by one-stage methods [66,67,58], (*ii*) $D : 48, 32, 8$ used by three-stage methods [25,16], and (*iii*) $D : 8, 8, 4, 4$ used by MVSTER. All of these conditions follow implementation details described in Sec. 4.2.

**Optimal Transport in Depth Propagation (OT)** Learning-based MVS methods usually use $L_1$ loss to regress the depth map [66,25] or use cross-entropy loss to classify depth [67,16]. In this experiment, the two losses are compared with the Wasserstein loss computed by OT. Apart from the aforementioned evaluation metrics, we introduce $S_3$ EPE and $S_4$ EPE, which stand for EPE of stage 3 and stage 4 depth estimations on DTU. As shown in Table 6, OT improves point cloud reconstruction performance and reduces depth error. Especially, it greatly reduces depth error of the last two stages, which demonstrates OT is effective in propagating finer depth maps to later stages.

**Table 5.** Comparison of our MDE module with DCN and ASPP.

| Method | Overall↓ | EPE ↓ | $e_1$↓ | $e_3$↓ | Runtime (s)↓ | Param (M)↓ |
|---|---|---|---|---|---|---|
| Raw feature | 0.317 | 1.35 | 15.00 | 7.37 | 0.09 | 0.98 |
| DCN | **0.313** | 1.33 | **14.82** | 7.55 | 0.23 | 1.51 |
| ASPP | 0.327 | 1.34 | 15.30 | 7.53 | 0.15 | 1.16 |
| MDE (Ours) | **0.313** | **1.31** | 14.98 | **7.27** | **0.09** | **0.98** |



(a) Raw image        (b) Feature without MDE guidance        (c) Feature with MDE guidance
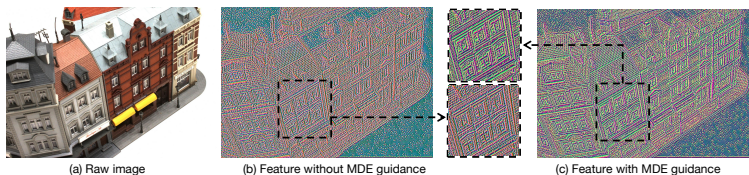
**Fig. 5.** An example shows that MDE guides the FPN feature to deliver more details at object boundaries. For better visualization, we leverage PCA to reduce the number of feature channels to 3 and color the channels with RGB.

**Table 6.** Comparison of optimal transport with $L_1$ loss and cross-entropy loss.

| Method | Overall↓ | EPE ↓ | $e_1$ ↓ | $e_3$ ↓ | $S_3$ EPE ↓ | $S_4$ EPE↓ |
|---|---|---|---|---|---|---|
| $L_1$ Loss | 0.321 | 1.47 | 15.32 | 7.53 | 7.02 | 6.32 |
| CE Loss | 0.314 | 1.34 | **14.96** | 7.40 | 7.12 | 6.64 |
| OT (Ours) | **0.313** | **1.31** | 14.98 | **7.27** | **6.41** | **5.90** |

## 5   Conclusions

In this paper, we present the epipolar Transformer for efficient MVS, termed as MVSTER. The proposed epipolar Transformer leverages both 2D semantics and 3D spatial associations to efficiently aggregate multi-view features. Specifically, MVSTER enriches 2D depth-discriminative semantics via an auxiliary monocular depth estimator. And the cross-attention on the epipolar line constructs 3D associations without learnable parameters. The combined 2D and 3D information serves as guidance to aggregate different views. Moreover, we formulate depth estimation as a depth-aware classification problem, which produces finer depth estimations propagated in the cascade structure. Extensive experiments on DTU, Tanks&Temple, BlendedMVS, and ETH3D show our method achieves

stage-of-the-art performance with significantly higher efficiency. We hope that MVSTER can serve as an efficient baseline for learning-based MVS, and further work may focus on simplifying 2D extractors and 3D CNNs.

## 6   Additional Implementation Details

**Network Architecture of Feature Extractor** We use a four-stage Feature Pyramid Network (FPN) [33] to extract image features, and the detailed parameters with layer descriptions are summarized in Table 7. For Deformable Convolutional Networks (DCN) [14] and Atrous Spatial Pyramid Pooling (ASPP) [7] that are used in the ablation study of our main text, the network parameters are listed in Table 8.

**Table 7.** The detailed parameters of FPN, where S denotes stride, and if not specified with *, each convolution layer is followed by a Batch Normalization layer (BN) and a Rectified Linear Unit (ReLU).

| Stage Description | Layer Description | Output Size |
|---|---|---|
| - | Input Images | $H \times W \times 3$ |
| FPN Stage 1 | Conv2D, $3 \times 3$, S1, 8 | $H \times W \times 8$ |
| FPN Stage 1 | Conv2D, $3 \times 3$, S1, 8 | $H \times W \times 8$ |
| FPN Stage 1 Inner Layer* | Conv2D, $1 \times 1$, S1, 64 | $H \times W \times 64$ |
| FPN Stage 1 Output Layer* | Conv2D, $1 \times 1$, S1, 8 | $H \times W \times 8$ |
| FPN Stage 2 | Conv2D, $5 \times 5$, S2, 16 | $H/2 \times W/2 \times 16$ |
| FPN Stage 2 | Conv2D, $3 \times 3$, S1, 16 | $H/2 \times W/2 \times 16$ |
| FPN Stage 2 | Conv2D, $3 \times 3$, S1, 16 | $H/2 \times W/2 \times 16$ |
| FPN Stage 2 Inner Layer* | Conv2D, $1 \times 1$, S1, 64 | $H/2 \times W/2 \times 64$ |
| FPN Stage 2 Output Layer* | Conv2D, $1 \times 1$, S1, 16 | $H/2 \times W/2 \times 16$ |
| FPN Stage 3 | Conv2D, $5 \times 5$, S2, 32 | $H/4 \times W/4 \times 32$ |
| FPN Stage 3 | Conv2D, $3 \times 3$, S1, 32 | $H/4 \times W/4 \times 32$ |
| FPN Stage 3 | Conv2D, $3 \times 3$, S1, 32 | $H/4 \times W/4 \times 32$ |
| FPN Stage 3 Inner Layer* | Conv2D, $1 \times 1$, S1, 64 | $H/4 \times W/4 \times 64$ |
| FPN Stage 3 Output Layer* | Conv2D, $1 \times 1$, S1, 32 | $H/4 \times W/4 \times 32$ |
| FPN Stage 4 | Conv2D, $5 \times 5$, S2, 64 | $H/8 \times W/8 \times 64$ |
| FPN Stage 4 | Conv2D, $3 \times 3$, S1, 64 | $H/8 \times W/8 \times 64$ |
| FPN Stage 4 | Conv2D, $3 \times 3$, S1, 64 | $H/8 \times W/8 \times 64$ |
| FPN Stage 4 Inner Layer* | Conv2D, $1 \times 1$, S1, 64 | $H/8 \times W/8 \times 64$ |
| FPN Stage 4 Output Layer* | Conv2D, $1 \times 1$, S1, 64 | $H/8 \times W/8 \times 64$ |

**Network Architecture of Light-Weight 3D CNN** An UNet [44] structured 3D CNN is applied for cost volume regularization at each stage, where the kernel size $3 \times 3 \times 3$ is partially replaced with $3 \times 3 \times 1$ in MVSTER for a more efficient

**Table 8.** The detailed parameters of DCN and ASPP, where S denotes stride, and D denotes dilation parameter for ASPP.

| Stage Description | Layer Description | Output Size |
|---|---|---|
| DCN Stage 1 | DCN2D, $3 \times 3$, S1, 8 | $H \times W \times 8$ |
| DCN Stage 2 | DCN2D, $3 \times 3$, S1, 16 | $H/2 \times W/2 \times 16$ |
| DCN Stage 3 | DCN2D, $3 \times 3$, S1, 32 | $H/4 \times W/2 \times 32$ |
| DCN Stage 4 | DCN2D, $3 \times 3$, S1, 64 | $H/8 \times W/8 \times 64$ |
| ASPP Stage 1 | Conv2D, $3 \times 3$, S1, D{1,6,12}, 8 | $H \times W \times 8$ |
| ASPP Stage 2 | Conv2D, $3 \times 3$, S1, D{1,6,12}, 16 | $H/2 \times W/2 \times 16$ |
| ASPP Stage 3 | Conv2D, $3 \times 3$, S1, D{1,6,12}, 32 | $H/4 \times W/2 \times 32$ |
| ASPP Stage 4 | Conv2D, $3 \times 3$, S1, D{1,6,12}, 64 | $H/8 \times W/8 \times 64$ |

pipeline. Apart from the input cost volume size, the network architectures are the same for each stage in the cascade structure, so we only report the detailed parameters of the 4th stage in Table 9.

**Table 9.** The detailed parameters of 3D CNN, where S denotes stride, and if not specified with *, each convolution layer is followed by a Batch Normalization layer (BN) and a Rectified Linear Unit (ReLU).

| Stage Description | Layer Description | Output Size |
|---|---|---|
| - | Input Cost Volume | $H \times W \times 4 \times 8$ |
| UNet Stage 1 | Conv3D, $3 \times 3 \times 1$, S1, 8 | $H \times W \times 4 \times 8$ |
| UNet Stage 1 | Conv3D, $3 \times 3 \times 1$, S2, 16 | $H/2 \times W/2 \times 4 \times 16$ |
| UNet Stage 1 | Conv3D, $3 \times 3 \times 3$, S1, 16 | $H/2 \times W/2 \times 4 \times 16$ |
| UNet Stage 1 Inner Layer | TransposeConv3D, $3 \times 3 \times 1$, S2, 8 | $H \times W \times 4 \times 8$ |
| UNet Stage 1 Output Layer* | TransposeConv3D, $3 \times 3 \times 3$, S1, 8 | $H \times W \times 4 \times 8$ |
| UNet Stage 2 | Conv3D, $3 \times 3 \times 1$, S2, 32 | $H/4 \times W/4 \times 4 \times 32$ |
| UNet Stage 2 | Conv3D, $3 \times 3 \times 3$, S1, 32 | $H/4 \times W/4 \times 4 \times 32$ |
| UNet Stage 2 Inner Layer | TransposeConv3D, $3 \times 3 \times 1$, S2, 16 | $H/2 \times W/2 \times 4 \times 16$ |
| UNet Stage 3 | Conv3D, $3 \times 3 \times 1$, S2, 64 | $H/8 \times W/8 \times 4 \times 64$ |
| UNet Stage 3 | Conv3D, $3 \times 3 \times 3$, S1, 64 | $H/8 \times W/8 \times 4 \times 64$ |
| UNet Stage 3 Inner Layer | TransposeConv3D, $3 \times 3 \times 1$, S2, 32 | $H/4 \times W/4 \times 4 \times 32$ |

# 7   Additional Ablation Study

**Ablation Study on Hyperparameters** We conduct an ablation study on loss weight $\lambda$ and temperature parameter $t_e$. As shown in Table 10, $\lambda = 3 \times 10^{-4}$ is a proper loss weight for jointly optimizing monocular depth estimation and multi-view stereo. As shown in Table 11, MVSTER produces a finer depth map when

slowly increasing $t_e$, and the network shows best reconstruction performance on DTU when $t_e = 2$.

**Table 10.** Ablation study on loss weight $\lambda$.

| $\lambda$ | Acc.↓ | Comp.↓ | Overall↓ | EPE↓ | $e_1$ ↓ | $e_3$ ↓ |
|---|---|---|---|---|---|---|
| $1 \times 10^{-2}$ | 0.361 | 0.312 | 0.337 | 1.56 | 16.47 | 8.32 |
| $1 \times 10^{-3}$ | 0.354 | 0.287 | 0.321 | 1.33 | 15.01 | **7.26** |
| $3 \times 10^{-4}$ | **0.350** | 0.276 | **0.313** | **1.31** | 14.98 | 7.27 |
| $1 \times 10^{-4}$ | 0.354 | **0.275** | 0.314 | 1.32 | **14.97** | 7.28 |

**Table 11.** Ablation study on temperature parameter $t_e$.

| $\lambda$ | Acc.↓ | Comp.↓ | Overall↓ | EPE↓ | $e_1$ ↓ | $e_3$ ↓ |
|---|---|---|---|---|---|---|
| 0.5 | 0.354 | 0.279 | 0.317 | 1.33 | 15.09 | 7.52 |
| 1.0 | 0.353 | 0.279 | 0.314 | 1.33 | 15.03 | 7.47 |
| 2.0 | **0.350** | 0.276 | **0.313** | **1.31** | 14.98 | 7.27 |
| 3.0 | 0.353 | **0.274** | 0.314 | **1.31** | **14.89** | **7.08** |

# 8   Point Cloud Visualizations

We visualize point cloud reconstruction results of DTU [1], ETH3D [46] and Tanks&Temples [30] in Fig. 6, Fig. 7 and Fig. 8, respectively. MVSTER shows its robustness on scenes with varying input image resolutions and depth ranges.

**Fig. 6.** Point clouds on DTU [1] reconstructed by MVSTER.

**Fig. 7.** Point clouds on ETH3D [46] reconstructed by MVSTER.

**Fig. 8.** Point clouds on Tanks&Temples [30] reconstructed by MVSTER.

# 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 (2016)
2. Abnar, S., Zuidema, W.H.: Quantifying attention flow in transformers. In: Association for Computational Linguistics (2020)
3. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein GAN. arXiv preprint arXiv:1701.07875 (2017)
4. Bozic, A., Palafox, P., Thies, J., Dai, A., Nießner, M.: Transformerfusion: Monocular rgb scene reconstruction using transformers. In: Advances in Neural Information Processing Systems (2021)
5. Campbell, N.D.F., Vogiatzis, G., Hernández, C., Cipolla, R.: Using multiple hypotheses to improve depth-maps for multi-view stereo. In: European Conference on Computer Vision (2008)
6. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European Conference on Computer Vision (2020)
7. Chen, L., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587 (2017)
8. Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., Sutskever, I.: Generative pretraining from pixels. In: International Conference on Machine Learning (2020)
9. Chen, R., Han, S., Xu, J., Su, H.: Point-based multi-view stereo network. In: IEEE International Conference on Computer Vision (2019)
10. Cheng, S., Xu, Z., Zhu, S., Li, Z., Li, L.E., Ramamoorthi, R., Su, H.: Deep stereo using adaptive thin volume representation with uncertainty awareness. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
11. Chitta, K., Prakash, A., Geiger, A.: Neat: Neural attention fields for end-to-end autonomous driving. In: IEEE International Conference on Computer Vision (2021)
12. Collins, R.T.: A space-sweep approach to true multi-image matching. In: IEEE Conference on Computer Vision and Pattern Recognition (1996)
13. Cuturi, M.: Sinkhorn distances: Lightspeed computation of optimal transport. In: Advances in Neural Information Processing Systems (2013)
14. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: IEEE International Conference on Computer Vision (2017)
15. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. In: Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (2019)
16. Ding, Y., Yuan, W., Zhu, Q., Zhang, H., Liu, X., Wang, Y., Liu, X.: Transmvsnet: Global context-aware multi-view stereo network with transformers. arXiv preprint arXiv:2111.14600 (2021)
17. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021)
18. Dosovitskiy, A., Fischer, P., Ilg, E., Häusser, P., Hazirbas, C., Golkov, V., van der Smagt, P., Cremers, D., Brox, T.: Flownet: Learning optical flow with convolutional networks. In: IEEE International Conference on Computer Vision (2015)

19. Duggal, S., Wang, S., Ma, W., Hu, R., Urtasun, R.: Deeppruner: Learning efficient stereo matching via differentiable patchmatch. In: IEEE International Conference on Computer Vision (2019)
20. Furukawa, Y., Ponce, J.: Accurate, dense, and robust multiview stereopsis. IEEE Transactions on Pattern Analysis and Machine Intelligence (2010)
21. Galliani, S., Lasinger, K., Schindler, K.: Massively parallel multiview stereopsis by surface normal diffusion. In: IEEE International Conference on Computer Vision (2015)
22. Giang, K.T., Song, S., Jo, S.: Curvature-guided dynamic scale networks for multi-view stereo. arXiv preprint arXiv:2112.05999 (2021)
23. Godard, C., Aodha, O.M., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
24. Godard, C., Aodha, O.M., Firman, M., Brostow, G.J.: Digging into self-supervised monocular depth estimation. In: IEEE International Conference on Computer Vision (2019)
25. 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: IEEE Conference on Computer Vision and Pattern Recognition (2020)
26. He, C., Zeng, H., Huang, J., Hua, X., Zhang, L.: Structure aware single-stage 3d object detection from point cloud. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
27. He, Y., Yan, R., Fragkiadaki, K., Yu, S.: Epipolar transformer for multi-view human pose estimation. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
28. Ke, Q., Bennamoun, M., An, S., Sohel, F.A., Boussaïd, F.: A new representation of skeleton sequences for 3d action recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
29. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (2015)
30. Knapitsch, A., Park, J., Zhou, Q., Koltun, V.: Tanks and temples: benchmarking large-scale scene reconstruction. ACM Transactions on Graphics (2017)
31. Lee, J.Y., DeGol, J., Zou, C., Hoiem, D.: Patchmatch-rl: Deep mvs with pixelwise depth, normal, and visibility. In: IEEE International Conference on Computer Vision (2021)
32. Li, Z., Liu, X., Drenkow, N., Ding, A., Creighton, F.X., Taylor, R.H., Unberath, M.: Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers. In: IEEE International Conference on Computer Vision (2021)
33. Lin, T., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.: Feature pyramid networks for object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
34. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. IEEE International Conference on Computer Vision (2021)
35. Luo, S., Hu, W.: Diffusion probabilistic models for 3d point cloud generation. In: IEEE Conference on Computer Vision and Pattern Recognition (2021)
36. Ma, X., Gong, Y., Wang, Q., Huang, J., Chen, L., Yu, F.: Epp-mvsnet: Epipolar-assembling based depth prediction for multi-view stereo. In: IEEE International Conference on Computer Vision (2021)

37. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: European Conference on Computer Vision (2020)
38. Mordan, T., Thome, N., Hénaff, G., Cord, M.: Revisiting multi-task learning with ROCK: a deep residual auxiliary block for visual detection. In: Advances in Neural Information Processing Systems (2018)
39. Peng, R., Wang, R., Wang, Z., Lai, Y., Wang, R.: Rethinking depth estimation for multi-view stereo: A unified representation and focal loss. arXiv preprint arXiv:2201.01501 (2022)
40. Peyré, G., Cuturi, M.: Computational optimal transport. Foundations and Trends in Machine Learning (2019)
41. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
42. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in Neural Information Processing Systems (2017)
43. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding by generative pre-training. OpenAI Preprint (2018)
44. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention (2015)
45. Schönberger, J.L., Frahm, J.: Structure-from-motion revisited. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
46. 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: IEEE Conference on Computer Vision and Pattern Recognition (2017)
47. Shen, Z., Dai, Y., Rao, Z.: Cfnet: Cascade and fused cost volume for robust stereo matching. In: IEEE Conference on Computer Vision and Pattern Recognition (2021)
48. Shi, S., Guo, C., Jiang, L., Wang, Z., Shi, J., Wang, X., Li, H.: PV-RCNN: point-voxel feature set abstraction for 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
49. Sinha, A., Murez, Z., Bartolozzi, J., Badrinarayanan, V., Rabinovich, A.: Deltas: Depth estimation by learning triangulation and densification of sparse points. In: European Conference on Computer Vision (2020)
50. Tankovich, V., Hane, C., Zhang, Y., Kowdle, A., Fanello, S.R., Bouaziz, S.: Hitnet: Hierarchical iterative tile refinement network for real-time stereo matching. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 14362–14372 (2021)
51. Tenney, I., Das, D., Pavlick, E.: BERT rediscovers the classical NLP pipeline. In: Association for Computational Linguistics (2019)
52. Tola, E., Strecha, C., Fua, P.: Efficient large-scale multi-view stereo for ultra high-resolution image sets. Machine Vision and Applications (2012)
53. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in Neural Information Processing Systems (2017)
54. Wang, F., Galliani, S., Vogel, C., Pollefeys, M.: Itermvs: Iterative probability estimation for efficient multi-view stereo. arXiv preprint arXiv:2112.05126 (2021)

55. Wang, F., Galliani, S., Vogel, C., Speciale, P., Pollefeys, M.: Patchmatchnet: Learned multi-view patchmatch stereo. In: IEEE Conference on Computer Vision and Pattern Recognition (2021)
56. Wang, H., Zhu, Y., Green, B., Adam, H., Yuille, A.L., Chen, L.: Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In: European Conference on Computer Vision (2020)
57. Watson, J., Aodha, O.M., Prisacariu, V., Brostow, G.J., Firman, M.: The temporal opportunist: Self-supervised multi-frame monocular depth. In: IEEE Conference on Computer Vision and Pattern Recognition (2021)
58. Wei, Z., Zhu, Q., Min, C., Chen, Y., Wang, G.: Aa-rmvsnet: Adaptive aggregation recurrent multi-view stereo network. In: IEEE International Conference on Computer Vision (2021)
59. Xu, Q., Tao, W.: Multi-scale geometric consistency guided multi-view stereo. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
60. Xu, Q., Tao, W.: Learning inverse depth regression for multi-view stereo with correlation cost volume. In: AAAI Conference on Artificial Intelligence (2020)
61. Xu, Q., Tao, W.: Pvsnet: Pixelwise visibility-aware multi-view stereo network. arXiv preprint arXiv:2007.07714 (2020)
62. Yan, J., Wei, Z., Yi, H., Ding, M., Zhang, R., Chen, Y., Wang, G., Tai, Y.: Dense hybrid recurrent multi-view stereo net with dynamic consistency checking. In: European Conference on Computer Vision (2020)
63. Yang, F., Yang, H., Fu, J., Lu, H., Guo, B.: Learning texture transformer network for image super-resolution. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
64. Yang, J., Mao, W., Alvarez, J.M., Liu, M.: Cost volume pyramid based depth inference for multi-view stereo. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
65. Yang, W., Li, Q., Liu, W., Yu, Y., Ma, Y., He, S., Pan, J.: Projecting your view attentively: Monocular road scene layout estimation via cross-view transformation. In: IEEE Conference on Computer Vision and Pattern Recognition (2021)
66. Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: European Conference on Computer Vision (2018)
67. Yao, Y., Luo, Z., Li, S., Shen, T., Fang, T., Quan, L.: Recurrent mvsnet for high-resolution multi-view stereo depth inference. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
68. 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. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
69. Yi, H., Wei, Z., Ding, M., Zhang, R., Chen, Y., Wang, G., Tai, Y.: Pyramid multi-view stereo net with self-adaptive view aggregation. In: European Conference on Computer Vision (2020)
70. Yu, Z., Gao, S.: Fast-mvsnet: Sparse-to-dense multi-view stereo with learned propagation and gauss-newton refinement. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
71. Zhang, J., Yao, Y., Li, S., Luo, Z., Fang, T.: Visibility-aware multi-view stereo network. In: British Machine Vision Conference (2020)
72. Zhang, X., Hu, Y., Wang, H., Cao, X., Zhang, B.: Long-range attention network for multi-view stereo. In: IEEE Winter Conference on Applications of Computer Vision (2021)

73. Zhao, M., Zhang, J., Zhang, C., Zhang, W.: Leveraging heterogeneous auxiliary tasks to assist crowd counting. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
74. Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)
75. Zhu, J., Peng, B., Li, W., Shen, H., Zhang, Z., Lei, J.: Multi-view stereo with transformer. arXiv preprint arXiv:2112.00336 (2021)