

# Attention Aware Cost Volume Pyramid Based Multi-view Stereo Network for 3D Reconstruction

Anzhu Yu \*\*, Wenyue Guo\*, Bing Liu\*, Xin Chen, Xin Wang, Xuefeng Cao, Bingchuan Jiang  
*PLA Strategic Support Force Information Engineering University, Zhengzhou, 450000, China*

---

## Abstract

We present an efficient multi-view stereo (MVS) network for 3D reconstruction from multi-view images. While previous learning based reconstruction approaches performed quite well, most of them estimate depth maps at a fixed resolution using plane sweep volumes with a fixed depth hypothesis at each plane, which requires densely sampled planes for desired accuracy and therefore is difficult to achieve high resolution depth maps. In this paper we introduce a coarse-to-fine depth inference strategy to achieve high resolution depth. This strategy estimates the depth map at coarsest level, while the depth maps at finer levels are considered as the upsampled depth map from previous level with pixel-wise depth residual. Thus, we narrow the depth searching range with priori information from previous level and construct new cost volumes from the pixel-wise depth residual to perform depth map refinement. Then the final depth map could be achieved iteratively since all the parameters are shared between different levels. At each level, the self-attention layer is introduced to the feature extraction block for capturing the long range dependencies for depth inference task, and the cost volume is generated using similarity measurement instead of the variance based methods used in previous work. Experiments were conducted on both the DTU benchmark dataset and recently released BlendedMVS dataset. The results demonstrated that our model could outperform most state-of-the-arts (SOTA) methods. The codebase of this project is at <https://github.com/ArthasMil/AACVP-MVSNet>.

*Keywords:*

Multi-view stereo, 3D Reconstruction, Cost Volume, Coarse-to-fine, Deep Learning.

---

## 1. Introduction

Multi-view stereo targets at reconstructing the observed 3D scene with dense representation from multi-view images and corresponding camera parameters, which has been extensively studied for decades and covers a wide range of applications such as photogrammetry (Malihi et al., 2016; Masiero et al., 2019; Rottensteiner et al., 2014), cartography (Bitelli et al., 2018; Buyukdemircioglu and Kocaman, 2020) and augmented reality (Harazono et al., 2019; Sing and Xie, 2016; Yang et al., 2013). MVS also remains a core problem of photogrammetry and computer vision in many aspects (Hirschmüller, 2007; Tola et al., 2012; Yao et al., 2018; Zbontar and LeCun, 2015; Zhu et al., 2017).

---

\*Prof. A.Yu, Dr. B. Liu and Dr. W. Guo have equal contribution to this work and are co-first authors.

\*\*Corresponding author: Anzhu Yu.

*Email addresses:* anzhu\_yu@126.com (Anzhu Yu ), guowyer@163.com (Wenyue Guo ), liubing220524@126.com (Bing Liu ), xincheng\_cosmos@126.com (Xin Chen), rainbowxdo@163.com (Xin Wang), CAO\_Xue\_Feng@163.com (Xuefeng Cao), jbc021@163.com (Bingchuan Jiang)

The traditional MVS methods introduce the hand-crafted similarity metrics for image association followed by some optimization steps for dense point cloud generation (e.g., normalized cross correlation as similarity metric and semi-global matching for optimization (Hirschmüller, 2007)). Though these methods perform quite well under ideal Lambertian scenarios with high textured regions, they suffer from incomplete reconstruction in some low-textured, reflective regions where the accuracy and robustness of dense matching decrease (Yao et al., 2018). In the meantime, the traditional methods are usually conducted sequentially, which is a time and memory consuming procedure and limits its application in situation where efficiency is needed. The Simultaneous Localization And Mapping (SLAM) is also a popular way for sensing the 3D structure of surroundings. Though SLAM could extract a lot of 2D/3D landmarks (A.K.A the object points in photogrammetry application) along with the transformation matrix of each frame, the density of generated point cloud is usually not enough for 3D reconstruction for the reason that 3D reconstruction usually needs pixel-wise matching and geo-positioning instead of 3D point cloud of sparse landmarks. Overall, accurate and efficient 3D reconstruction is still a hot and challenging topic.

Owing to rapid developments of the Deep Learning technique (LeCun et al., 2015), some learning based MVS models arise in recent years (Hartmann et al., 2017; Ji et al., 2017; Kar et al., 2017; Xiang et al., 2020). The convolutional neural networks (CNNs) that could extract hierarchical features with stronger representation ability make it possible to speed up MVS processing by using the computation power of GPUs. Some end-to-end learning based MVS methods were presented by modeling the regression relationship between multi-view images and corresponding depth maps (Chen et al., 2019a; Xiang et al., 2020; Yao et al., 2018), after which the estimated depth maps could be fused and processed to generate the dense 3D point clouds of target region (e.g., using the Fusibile toolbox (Galliani et al., 2015a)). Though the aforementioned literatures could achieve better results in terms of accuracy and completeness compared to traditional MVS methods, they require huge GPU memory caused by high dimension cost volume (Esteban and Schmitt, 2004; Seitz et al., 2006) for depth inference. A normal way to release memory burden is to inference depth map with downsampled images (Yao et al., 2018, 2019). It comes at the cost of reconstruction accuracy and completeness, however.

To this end, we endeavor to propose an Attention Aware Cost Volume Pyramid Multi-view Stereo Network (AACVP-MVSNet) for 3D reconstruction. In summary, our main contributions are listed below:

- We introduce the self-attention layers for hierarchical features extraction, which may capture important information for the subsequent depth inference task.
- We introduce the similarity measurement for cost volume generation instead of the variance based methods used by most MVS networks.
- We use a coarse-to-fine strategy for depth inference which is applicable for large scaled images, and extensive experiments validate that the proposed approach achieves an overall performance improvement than most SOTA algorithms on DTU dataset.

## 2. Related work

### 2.1. Traditional MVS methods

In 3D reconstruction workflow, the traditional MVS methods are usually implemented following the sparse point cloud generation procedure (Ma and Liu, 2018) (e.g. the Structure from Motion computation (SfM) (Schonberger and Frahm, 2016; Yang et al., 2013)). To reconstruct

the dense 3D point cloud, the recovered intrinsic and extrinsic parameters of cameras for each image and the sparse point cloud obtained from SfM or SLAM are set as inputs. The Clustering Views for Multi-view Stereo (CMVS) (Furukawa et al., 2010) and the Patch-based Multi-view Stereo (PMVS) (Koch et al., 2014) are very popular methods for dense 3D reconstruction. The CMVS introduces SfM filter to merge extracted feature points and decomposes the input images into a set of image clusters of manageable size, after which the MVS software could be used for 3D reconstruction. The PMVS uses clustered images from CMVS as input and generates dense 3D point clouds through matching, expansion and filtering.

The Semi-Global Matching (SGM) is also widely way for 3D reconstruction, which is proposed for estimation of a dense disparity map from rectified stereo image pairs with introducing the penalty of inconsistency (Hirschmüller, 2005). As the SGM algorithm has trade-off between computing time and the quality of results, it's faster than PMVS and has encountered wide adoption in real-time stereo vision applications (Hirschmüller, 2007).

Though the aforementioned works yield impressive results and perform well on the accuracy, they require photometric consistency and would achieve unsatisfactory matching results with hand-crafted features and similarity when dealing with non-Lambertian surfaces, low textured and texture-less regions (Xiang et al., 2020). Therefore, the traditional MVS methods still need to be improved to achieve more robust and complete reconstruction results (Galliani et al., 2015a; Vu et al., 2011).

## 2.2. Learned stereo algorithms

In contrast to traditional stereo matching algorithms used in traditional MVS methods that introduce hand-crafted image features and matching metrics (Hirschmüller and Scharstein, 2007), the learned stereo algorithms introduce deep learning networks to achieve better matching results.

Most of learned stereo algorithms focus on image association procedures, in which the Convolutional Neural Networks (CNNs) are used for hierarchical features extraction and matching. Luo et al. (2016) and Žbontar and LeCun (2016) use CNNs to extract learned features for matching, followed by the SGM for dense reconstruction. The SGMNet (Seki and Pollefeys, 2017) introduces CNNs for penalty-parameters estimation and outperformed state-of-the-art accuracy on KITTI benchmark datasets. Zhang and Lee (2019) introduces Graph Neural Network (GNN) model for feature matching, which transforms coordinates of feature points into local features to replace NP-hard assignment problem in some traditional methods. The CNN-CRF introduces conditional random field into networks and forms an end-to-end matching algorithm (Knobelreiter et al., 2017). The RF-Net (Shen et al., 2019) is an end-to-end trainable matching network based on receptive field to compute sparse correspondence between images. In the work of Kendall et al. (2017), a deep learning architecture is proposed for regressing disparity from a rectified pair of stereo images and the cost volume is used for feature aggregation.

## 2.3. Learned MVS algorithms

There are few deep learning algorithms for MVS problem before the work of Hartmann et al. (2017). The SurfaceNet (Ji et al., 2017) is the pioneer to build a learning based pipeline for MVS problem, which builds the cost volume with sophisticated voxel-wise view selection and use 3D regularization for inferencing surface voxels. Yao et al. (2018) proposes MVSNet for the MVS problem, which introduces differentiable homography to build the cost volume for features aggregation and use 3D regularization for depth inference. To reduce the memory burden, Yao et al. (2019) proposes R-MVSNet, which sequentially regularized the 2D cost maps along the depth direction via the gated recurrent unit (GRU). Based on these works,

many learned MVS networks are proposed to improve the accuracy and completeness of the 3D reconstruction results, such as the PointMVSNet (Chen et al., 2019b), the CasMVSNet (Gu et al., 2020), the PVA-MVSNet (Yi et al., 2020) and the CVP-MVSNet (Yang et al., 2020). Some of them use a coarse-to-fine strategy to improve the accuracy and completeness of reconstruction results compared to the of original MVSNet, as the depth map could be estimated with higher resolution input images whereby better 3D reconstruction would be achieved.

Feature extraction is a key issue for learned MVS algorithms. Regarding another key issue, cost volume generation, though the aforementioned literatures introduce CNN blocks for feature extraction, they hardly capture long-range dependencies during coarse-to-fine strategy and fail to capture the important information for depth inference tasks. Though the variance based cost metric is widely used for cost volume generation, it is pointed out by Tulyakov et al. (2018) that the number of channels of the cost volume could be reduced and similar accuracy can still be achieved, which implied that the variance based cost volume with a lot of channels may be redundant and the memory consumption as well as the computational requirement could be reduced. To overcome and address these issues, we introduce the self-attention mechanism into depth inference procedure to improve the overall accuracy of reconstruction results, and implement a similarity based method for cost volume generation to simultaneously reduce the computational requirement and the memory consumption.

### 3. Methodology

This section describes the detailed architecture of the proposed network. The design of AACVP-MVSNet strongly borrows the insights from previous MVS approaches and novel feature extraction methods.

The overall system is depicted in Fig.(1). The multi-view images are first downsampled to form image pyramid, after which a weights-shared feature extraction block is built for feature extraction at each level. The depth inference begins at coarse level (level  $L$ ) by building the cost volume  $\mathbf{C}^L$  using the similarity measurement, namely the cost volume correlation which uses the similarity metrics rather than variance based metrics at bottom in Fig.(1), and the initial depth map is generated by cost volume regularization using the 3D convolution block and the softmax operation. The estimated depth map  $\mathbf{D}^L$  is upsampled to the image size of next level (level  $L - 1$ ) and the cost volume  $\mathbf{C}^{(L-1)}$  is built with depth hypothesis planes estimation followed by cost volume correlation. The depth residual map  $\mathbf{R}^{(L-1)}$  is also estimated with 3D convolution block and the softmax operation, and the depth map  $\mathbf{D}^{(L-1)}$  is upsampled to the image size of the level  $L - 2$  for depth inference at  $(L - 2)$ th level. Thus, an iterative depth map estimation procedure is formed along with a cost volume pyramid  $\{\mathbf{C}^i\}(i = L, L - 1, \dots, 0)$ .

As existing works, we assume the reference image is denoted as  $\mathbf{I}_0 \in \mathbb{R}^{H \times W}$ , where  $H$  and  $W$  are the height and width of the input image respectively. Let  $\{\mathbf{I}_i\}_{i=1}^N$  be input the  $N$  source images for reconstruction. For MVS problem, the corresponding camera intrinsic matrix, rotation matrix, and translation vector denoted as  $\{\mathbf{K}_i, \mathbf{R}_i, \mathbf{t}_i\}_{i=0}^N$  are known for all input views. Our goal is to estimate the depth map  $\mathbf{D}^0$  from  $\{\mathbf{I}_i\}_{i=0}^N$  with given  $\{\mathbf{K}_i, \mathbf{R}_i, \mathbf{t}_i\}_{i=0}^N$ .

#### 3.1. Self-attention based hierarchical feature extraction

##### 3.1.1. Self-attention based feature extraction block

Our feature extraction block consists of 8 convolutional layers and a self-attention layers with 16 output channels, each of which is followed by a leaky rectified linear unit (Leaky ReLU), as shown in Fig.(2). In contrast to feature extraction block of previous MVS networks, the self-attention mechanism is introduced for learning to focus on important information for depth inference.

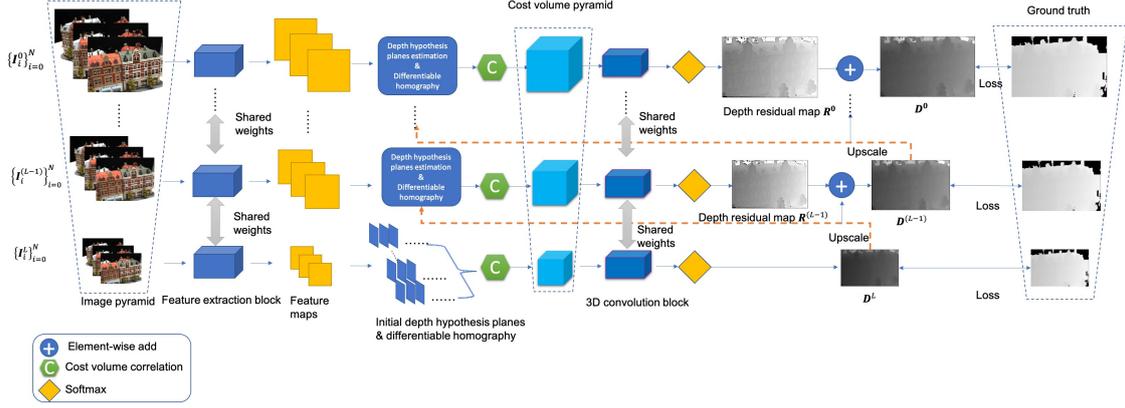


Figure 1: The network structure of AACVP-MVSNet. The feature extraction block and the 3D convolution block are both weights-shared between all levels. The image pyramid is built firstly, and the iterative depth estimation starts at the coarsest level. The depth map estimated at each level is taken as input at next level for depth residual estimation.

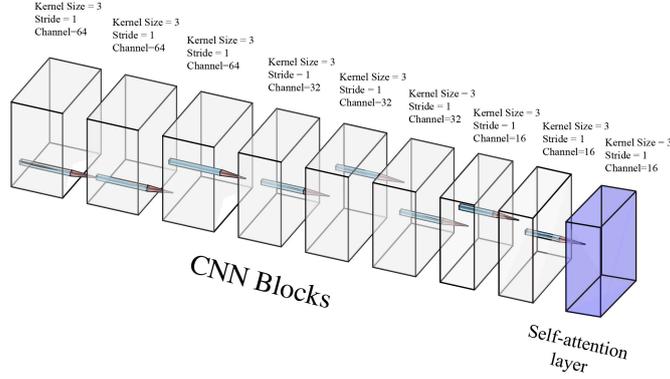


Figure 2: The self-attention based feature extraction block. This block consists of 8 CNN layers and one self-attention layer. The kernel size is set to 3 and the stride is set to 1 for all layers. The out channels decrease every 3 layers and the final output channel is 16.

Given a learned weight matrix  $\mathbf{W} \in \mathbb{R}^{k \times k \times d_{in} \times d_{out}}$  the convolutional output  $y_{ij} \in \mathbb{R}^{d_{out}}$  at pixel  $(i, j)$  is defined by the summing the linear product of input image  $\mathbf{I}$  and weight matrix:

$$y_{ij} = \sum_{a,b \in \mathbf{B}} W_{i-a,j-b} \cdot \mathbf{I}(a, b) \quad (1)$$

where  $k$  denotes kernel size,  $d_{in}$  is the input channels quantity while  $d_{out}$  is that of output channels, and  $\mathbf{B}$  is the image block for convolution computation with the same size of the kernel. Compared to traditional attention mechanism (Bartunov et al., 2018; Chorowski et al., 2015; Xu et al., 2015), the self-attention mechanism is defined as attention applied to a single context instead of across multiple context (Ramachandran et al., 2019), which directly models long-distance interactions and leads to state-of-the-art models for various tasks (Devlin et al., 2018; Shaw et al., 2018; Shazeer et al., 2018). The self-attention could be formulated as (Ramachandran et al., 2019)

$$y_{ij} = \sum_{a,b \in \mathbf{B}} \text{Softmax}_{ab}(\mathbf{q}_{ij}^T \mathbf{k}_{ab}) \mathbf{v}_{ab} \quad (2)$$

where  $\mathbf{q}_{ij} = \mathbf{W}_Q x_{ij}$ ,  $\mathbf{k}_{ab} = \mathbf{W}_K x_{ab}$  and  $\mathbf{v}_{ab} = \mathbf{W}_V x_{ab}$  denote queries, keys and values respectively, and the matrix  $\mathbf{W}_l (l = Q, K, V)$  consists of the learned parameters. Compared to traditional convolution computation (Eq.(1)), Eq.(2) could be decomposed into 3 steps:

- **Step 1.** Compute the queries ( $\mathbf{q}_{ij}$ ), keys ( $\mathbf{k}_{ab}$ ) and values ( $\mathbf{v}_{ab}$ ).
- **Step 2.** Measure the similarity of the queries and keys by calculating their inner product  $\mathbf{q}_{ij}^T \mathbf{k}_{ab}$ , followed by the softmax operation that maps the similarity between 0.0 and 1.0.
- **Step 3.** Weight the values with the similarity in **Step 2** and repeat the all steps for every pixel in  $\mathbf{B}$ .

It is not difficult to figure out that the output  $y_{ij}$  is achieved by linear transformations of the pixel in position  $ij$  and the neighborhood pixels, and this operation aggregates spatial information with a convex combination of value vectors with mixing weights parametrized by content interactions. The difference between the 2D convolution layer and the 2D self-attention layer is illustrated in Fig.(3) when  $k = 3$ .

However, Eq.(2) does not contain the positional information for queries  $\mathbf{q}_{ij}$ , which makes it permutation equivariant, limiting expressivity for vision tasks (Ramachandran et al., 2019). Therefore, the positional information embedding procedure should be introduced to achieve better results. We introduce the relative position embedding (Shaw et al., 2018) rather than the absolute position embedding (Vaswani et al., 2017) methods to bring in the row and column offset (denoted as  $\mathbf{r}_{a-i, b-j}$ ) into Eq.(2), and the self-attention computation could be formulated as

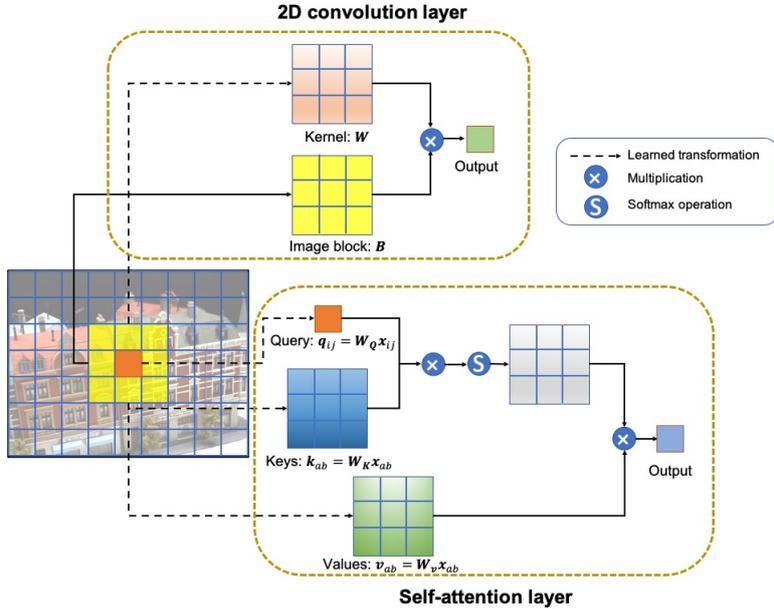


Figure 3: The difference between convolution layer and self-attention layer when the kernel size equals 3. The convolution layer could be considered as the combination of pixels and learned weights  $\mathbf{W}$  while the self-attention layer has three learned weights matrices, namely  $\mathbf{W}_Q$ ,  $\mathbf{W}_K$  and  $\mathbf{W}_V$ . The final output of the self-attention layer is computed through 3 steps.

$$y_{ij} = \sum_{a,b \in \mathbf{B}} \text{Softmax}_{ab}(\mathbf{q}_{ij}^T \mathbf{k}_{ab} + \mathbf{q}_{ij}^T \mathbf{r}_{a-i, b-j}) \mathbf{v}_{ab} \quad (3)$$

According to Eq.(3), each  $y_{ij}$  measures the similarity between the query and an element in  $\mathbf{B}$ , and enjoys translation equivariance (Cordonnier et al., 2019).

### 3.1.2. Hierarchical feature extraction

In contrast to previous work such as MVSNet (Yao et al., 2018) that extract features at a fixed scale to generate depth map with a fixed resolution, our network estimates the depth map by using a coarse-to-fine strategy. Therefore, our feature extraction pipeline consists of two steps.

We denote  $l \in \{0, 1, \dots, L\}$  for image levels from original scales to the coarsest resolution. The first step is to build a total of  $(L + 1)$  levels image pyramid  $\{\mathbf{I}_i^L\}_{i=1}^N$  for all the input source images and reference image. The next step is obtaining hierarchical representations at  $L$ th level using the feature extraction block presented in Fig.(2). Since the extraction block could be implemented on images captured from different views and different scales, the learned weights of the block should be shared as shown in Fig.(1). We define the extracted feature maps at  $l$ th level by  $\{\mathbf{f}_i^L\} \in \mathbb{R}^{H/2^l \times W/2^l \times Ch}$  for following sections, where  $Ch$  denotes the quantity of output channels.

### 3.2. Coarse-to-fine depth estimation

The cost volume is introduced for stereo matching by building the disparity searching space of each pixel with regular grids (Scharstein and Szeliski, 2002; Xu et al., 2017). Since the proposed network introduces pyramid structure in image spaces, the Cost Volume Pyramid (CVP) is intuitively formed. In our proposed network, the CVP is used for depth map inference at coarsest resolution and depth residual estimation at finer scales.

#### 3.2.1. Depth inference at the coarsest resolution

We construct the cost volume at coarsest scales as an initial depth map inference. In most learned MVS approaches, the cost volume is generated by transforming all extracted feature maps to the one generated from the reference image (Gu et al., 2020; Yao et al., 2018, 2019). We follow these research for cost volume generation at  $L$ th level while using different feature aggregation method.

Given the depth range  $(d_{min}, d_{max})$  of the reference image  $\mathbf{I}_0^L$ , the cost volume is constructed by sampling  $M$  fronto-parallel planes uniformly, which could be formulated as

$$d_m = d_{min} + m(d_{max} - d_{min})/M \quad (4)$$

where  $m = 0, 1, \dots, M - 1$  denotes the hypothesized depth planes. Similar to Yao et al. (2019), we introduce the differentiable homography matrix  $\mathbf{H}_i^L(d)$  for cost volume transformation from the  $i$ th source views to reference image at  $L$ th level, that is,

$$\mathbf{H}_i^L(d) = \mathbf{K}_i^L \mathbf{R}_i (\mathbf{E} - \frac{\mathbf{t}_0 - \mathbf{t}_i}{d} \mathbf{n}_0^T) \mathbf{R}_0^T (\mathbf{K}_0^L)^T \quad (5)$$

where the upper-case  $L$  indicates the level of images and  $\mathbf{E}$  denotes the identity matrix.

Eq.(5) suggests possible pixel corresponding between feature maps of source views and the reference image. For  $N$  input images,  $N$  4D tensors  $\mathbf{f}_i^L$  with the size of  $W/2^L \times H/2^L \times M \times Ch$  are generated, which is a memory consuming procedure. Thus, the feature aggregation process is usually implemented. In contrast to previous learned MVS networks that use the variance based feature aggregation (Chen et al., 2019a; Gu et al., 2020; Yang et al., 2020; Yao et al., 2019), we introduce the average group-wise correlation (Guo et al., 2019) that build the cost volumes by the similarity measurement for image matching tasks, whose basic idea is splitting

the features into groups and computing correlation maps group by group. Thus, the first step is to divide the feature channels of feature maps into  $G$  groups and compute the similarity between the  $i$ th group feature maps between reference image and the  $j$ th ( $j \in \{1, 2, \dots, N-1\}$ ) wrapped image at hypothesized depth plane  $d_m$  as:

$$\mathbf{S}_{j,d_m}^{i,L} = \frac{1}{Ch/G} \left\langle \mathbf{f}_{ref}^{i,L}(d_m), \mathbf{f}_j^{i,L}(d_m) \right\rangle \quad (6)$$

where  $i \in \{0, 1, \dots, G-1\}$ ,  $\langle \cdot, \cdot \rangle$  is the inner product,  $\mathbf{f}_j^L(d_m)$  indicates the interpolated feature map after wrapping  $\mathbf{f}_j^L$  to reference images by using Eq.(5), and all operations above are element-wise. When the feature similarities of all the  $G$  groups are computed by Eq.(6) for  $\mathbf{f}_{ref}^L(d_m)$  and  $\mathbf{f}_j^L(d_m)$ , the original feature maps  $\{\mathbf{f}_i^L\}$  could be compressed into a  $G$ -channel similarity tensor  $\mathbf{S}_{j,d_m}^L$  of size  $G \times H/2^L \times W/2^L$ . After calculating the similarity measurement at all the  $M$  hypothesized depth plane, the cost volume of  $j$ th source image  $\mathbf{C}_j^L = \text{concat}(\mathbf{S}_{j,0}^L, \mathbf{S}_{j,1}^L, \dots, \mathbf{S}_{j,M-1}^L)$  would be a  $G \times H/2^L \times W/2^L \times M$  sized tensor. Thus, the final aggregated cost volume could be computed as the average similarity of all the views, that is:

$$\mathbf{C}^L = \frac{1}{N-1} \sum_{j=1}^{N-1} \mathbf{C}_j^L \quad (7)$$

When the cost volume at  $L$ th level is estimated, the probability volume  $\mathbf{P}^L$  could be generated by a 3D convolution block (shown in Fig.(1) and Fig.(5)) similar to the previous learned MVS networks (Gu et al., 2020; Yang et al., 2020) and the depth map of each pixel  $\mathbf{p}$  at the coarsest level could be estimated as

$$\mathbf{D}^L(\mathbf{p}) = \sum_{m=0}^{M-1} d_m \mathbf{P}^L(\mathbf{p}, d_m) \quad (8)$$

where  $d_m$  is calculated by using Eq.(4).

### 3.2.2. Depth residual estimation at finer scales

Since the depth map  $\mathbf{D}^L(\mathbf{p})$  at  $L$ th level is obtained at the lowest resolution, the quality of 3D reconstruction would be limited. Thus, we try to refine  $\mathbf{D}^L(\mathbf{p})$  at a finer level, and the residual map estimation is intuitively implemented. We take residual map estimation at  $(L-1)$ th level ( $\mathbf{R}^{(L-1)}$ ) as an example, and the mathematical model could be summarized as

$$\begin{cases} \mathbf{R}^{(L-1)} = \sum_{m=-M/2}^{M/2} r_{\mathbf{p}}(m) \mathbf{P}_{\mathbf{p}}^{(L-1)}(r_{\mathbf{p}}) \\ \mathbf{D}^{(L-1)}(\mathbf{p}) = \mathbf{R}^{(L-1)} + \mathbf{D}_{upscale}^{(L)}(\mathbf{p}) \end{cases} \quad (9)$$

where  $M$  is the number of hypothesized depth planes,  $r_{\mathbf{p}} = m\Delta d_{\mathbf{p}}$  represents the depth residual,  $\Delta d_{\mathbf{p}} = l_{\mathbf{p}}/M$  is the depth interval,  $\mathbf{D}_{upscale}^{(L)}$  is the upscaled depth map from  $L$ th level and  $l_{\mathbf{p}}$  denotes the depth searching range at  $\mathbf{p}(u, v)$ . As shown in Fig.(4),  $\Delta d_{\mathbf{p}}$  and  $l_{\mathbf{p}}$  determine the depth estimation result at each pixel  $\mathbf{p}$ , and are the key parameters for depth residual estimation. There exist several methods for depth searching range and depth interval determination, such as iteratively range narrowing (Gu et al., 2020) and uncertainty based searching boundary (Cheng et al., 2020). Here we use the determination method used by Yang et al. that projects  $\mathbf{p}(u, v)$  in reference image and the corresponding points in source images into object space and

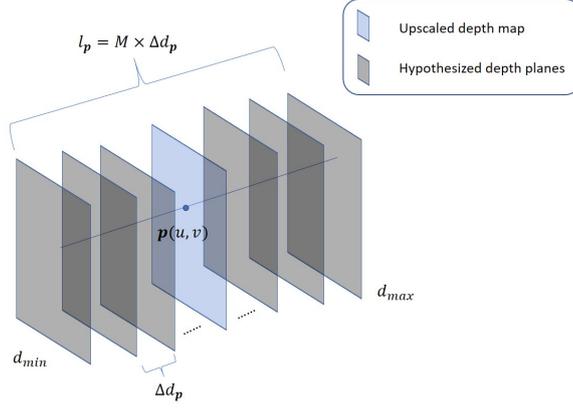


Figure 4: The depth searching range. The upscaled plane is in the middle of searching space, while there are  $M/2$  hypothesized planes at each side.  $\Delta d_p$  is the depth interval and the total searching distance is  $l_p$

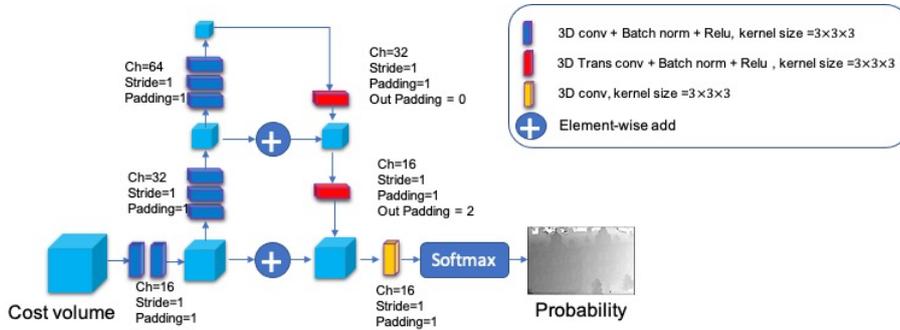


Figure 5: The structure of 3D convolution block. Similar to Yang et al. (2020), a multi-scale 3D convolutional network is applied to estimate the probabilities of different depth or residual hypothesis for each pixel.

determine depth interval  $\Delta d_p$  as the distance of the projection of two neighbor pixels along the epipolar line, because it is pointed by Yang et al. that it is not necessary to sample depth planes densely as the projections of those sampled 3D points in the image are too close to provide extra information for depth inference. Thus, the  $l_p$  could be calculated as the product of  $\Delta d_p$  and the given parameter  $M$  at level  $(L - 1)$ .

The aggregated cost volume  $\mathbf{C}^{(L-1)}$  could be built by using Eq.(6) and Eq.(7), and depth map  $\mathbf{D}^{(L-1)}$  could be achieved by using Eq.(9) after 3D convolution block and softmax operation for  $\mathbf{P}^{(L-1)}$ .

### 3.2.3. Iterative depth map estimation

Recall that our ultimate goal is to estimate  $\mathbf{D}^0$  from images at finest level. Because our network structure (Fig.(1)) only uses 2D and 3D convolutional layers that capture local features and all the weights are shared between different levels, we could estimate the depth map iteratively by input  $\{\mathbf{I}_i^l\}_{i=0}^N$  ( $0 \leq l < L-1$ ) to the feature extraction block for hierarchical feature maps extraction, followed by depth hypothesis estimation with upscaled depth map  $\mathbf{D}^{l+1}$  and cost volume generation. The residual depth  $\mathbf{R}^l$  could be generated with 3D convolution block and the softmax operation along with the depth map  $\mathbf{D}^l$ . Thus, we put  $\mathbf{D}^l$  as input at  $(l - 1)$ th level and an iterative depth map estimation process is formed. The final depth map  $\mathbf{D}^0$  will be achieved when we reach the top level in Fig.(1).

For network backward propagation, we build the loss function as the sum of  $L_1$ -norm be-

tween ground truth and estimated depth map, that is,

$$\mathbb{L} = \sum_{l=0}^L \sum_{\mathbf{p} \in \Omega} \|\mathbf{D}^l(\mathbf{p}) - \mathbf{D}_{GT}^l(\mathbf{p})\|_1 \quad (10)$$

where  $GT$  represents ground truth of depth map and  $\Omega$  is the set of valid pixels. Therefore, the weights could be trained by optimizing Eq.(10) and the estimated depth map could be achieved by forward propagation using the trained model.

## 4. Experiments

In this section, we present the datasets used in our experiments and training configurations. Afterwards the analysis of the experiments and ablation studies are presented.

### 4.1. Datasets

We use the DTU dataset, which is now a widely used benchmark for 3D reconstruction, and the recent released BlendedMVS dataset for experiments.

- **DTU dataset.** The DTU dataset consists of 124 different indoor scenes including a variety of objects scanned by fixed camera trajectories in 7 different lighting conditions (Aanaes et al., 2016). The original size of image is  $1600 \times 1200$  pixels. Following the common configurations, we use the same training and evaluation dataset split used by Cheng et al. (2020); Gu et al. (2020); Yao et al. (2019) and the ground truth is provided by Yao et al. (2018). This dataset could be used for quantitative analysis.
- **BlendedMVS dataset.** The BlendedMVS dataset is a large-scale MVS dataset for generalized multi-view stereo networks (Yao et al., 2020). This dataset contains more than 17k MVS training samples covering a variety of 113 scenes, including outdoor buildings, architectures, sculptures and small objects. However, there are no official ground truth provided, nor does the evaluation toolbox. Thus, we could only compare results qualitatively.

### 4.2. Training and evaluation

The training of learned MVS methods are usually memory consuming and slow. Since our network is built for depth map estimation iteratively and all the weights are shared between different layers, we could train our network use the downsampled images and ground truth to boost the training procedure.

The training on DTU dataset was implemented with the image size of  $160 \times 128$  pixels along with the corresponding camera parameters including camera intrinsic matrices, rotation matrices and translation vectors provided by Yang et al. (2020) while the trained weights were evaluated on full sized images. It is notable that the width and height should be dividable by 16 to make the input size suitable for the 3D convolution block (shown in Fig.(5)). We adopt 3 views for training as is set in Chen et al. (2019a); Yao et al. (2018), and 3 views are used for evaluation. More results with different number of views will be shown in Section 4.6.2. Both the image pyramid and ground truth pyramid have 2 levels for training and the coarsest resolution of images is  $80 \times 64$  pixels. We evaluate our model with a similar size with the training dataset at the coarsest level, so we first crop the images to the resolution  $1600 \times 1184$  pixels and build image pyramid with 5 levels with the resolution  $100 \times 74$  pixels at the coarsest level.

The training on BlendedMVS dataset was similar to DTU training procedure. Here we use the images with officially provided low resolution ( $768 \times 576$  pixels) for following experiments.

The images and the ground truth were downsampled to  $384 \times 288$  pixels and the parameters in intrinsic matrices were adjusted correspondingly. We build image pyramid with 3 levels for training and 4 levels for evaluation, and the coarsest resolution is  $96 \times 72$  pixels for each experiment.

In all experiments, we set the number of hypothesized depth plane  $M = 48$  at coarsest level and  $M = 8$  at other levels for both training and evaluation. The networks are trained on 4 Nvidia GeForce RTX 2080Ti graphics cards for 40 epochs with batch size set to 36 (mini-batch size of 9 per GPU). We used Adam (Kingma and Ba, 2014) to optimize the proposed network and the initial learning rate is set to  $1 \times 10^{-3}$  which is multiplied by 0.5 at 10th, 25th, 32nd epoch.

#### 4.3. Post processing and accuracy evaluation

We fuse all the depth maps into a complete one and generate the dense point cloud by the Fusibile toolbox (Galliani et al., 2015b) as used in Cheng et al. (2020); Gu et al. (2020); Yang et al. (2020), which consists of three steps for point cloud generation: photometric filtering, geometric consistency filtering, and depth fusion.

To quantitatively evaluate the 3D-reconstruction performance on DTU dataset, we calculate both the mean accuracy (referred as Acc.), the mean completeness (referred as Comp.) and the overall accuracy (referred as OA) which is defined as:

$$OA = \frac{Acc. + Comp.}{2} \quad (11)$$

The accuracy and completeness could be calculated using the official MATLAB script provided by the DTU dataset (Aanæs et al., 2016).

#### 4.4. Results on DTU dataset

We first compare our results to those reported by traditional geometric-based methods and other learning-based baseline methods. As summarized in Tab.(1), our method outperforms all methods in terms of completeness and overall accuracy with the number of group is set to  $G = 4$ , while Gipuma (Galliani et al., 2015b) performs the best in terms of accuracy. Here we found an interesting phenomenon that the number of group is a quarter of original channel quantity which is the same as channel deduction implemented in the MVSNet that reduces the 32-channel cost volume to an 8-channel one, before taken into 3D regularization block. Meanwhile, it is also demonstrated by Tulyakov et al. (2018) that a compressed 8-channel cost volume could achieve similar accuracy compared with original 32-channel cost volume in image matching task. This makes us believe that the raw cost volume representation may be redundant, and the cost volume could be compressed to the one with fewer channel quantity. Whether a quarter of original channel quantity is the best choice should be proved theoretically or validated with more experiments.

We compare our 3D reconstruction results with MVSNet and CasMVSNet in Fig.(6), Fig.(7) and Fig.(8). Both the AACVP-MVSNet and CasMVSNet achieve comparable completeness in these examples while the point cloud generated by MVSNet is sparser as they were extracted with low resolution depth maps. As we can see in Fig.(6) and Fig.(8), our results are smoother on the surfaces. Moreover, the letters in Fig.(7) is more complete and could be more easily recognized from 3D reconstruction results.

We also show the memory usage by AACVP-MVSNet in Tab.(2), and the baseline is the network whose cost volumes were generated by variance based method used in Chen et al. (2019a); Gu et al. (2020); Xiang et al. (2020); Yao et al. (2018) instead of the similarity measurement based method used in our network while the other parts stay the same with Fig.(1).

As indicated in Tab.(2), the memory usage is around 10% less than the baseline, and the memory usage decreases slightly with as the parameter  $G$  decreases, which could be easily illustrated by Eq.(6) and Eq.(7).

Table 1: Quantitative results of reconstruction quality on the DTU evaluation dataset (lower is better).

Methods	Acc.(mm)	Comp. (mm)	OA (mm)
Camp (Campbell et al., 2008)	0.835	0.554	0.695
Gipuma (Galliani et al., 2015b)	<b>0.283</b>	0.873	0.578
Furu (Furukawa and Ponce, 2009)	0.613	0.941	0.777
SurfaceNet (Ji et al., 2017)	0.450	1.040	0.745
MVSNet* (Yao et al., 2018)	0.396	0.527	0.462
R-MVSNet* (Yao et al., 2019)	0.383	0.452	0.418
PruMVSNet (Xiang et al., 2020)	0.495	0.433	0.464
PointMVSNet (Chen et al., 2019a)	0.361	0.421	0.391
CasMVSNet (Gu et al., 2020)	0.325	0.385	0.355
CVP-MVSNet (Yang et al., 2020)	0.296	0.406	0.351
Ours( $G = 8$ )	0.363	0.332	0.347
Ours( $G = 2$ )	0.360	0.341	0.351
Ours( $G = 4$ )	0.357	<b>0.326</b>	<b>0.341</b>

\*Official MVSNet and R-MVSNet implementation used the Altizure internal library for post-processing and could achieve higher accuracy compared to the Fusibile toolbox.

Table 2: Memory used per batch by AACVP-MVSNet with different parameter  $G$ . The training images were resized to  $160 \times 128$  pixels.

	Baseline	$G = 4$	$G = 8$	$G = 16$
Memory usage (Mb)	1171.44	<b>1048.44</b>	1065.89	1087.93
Compared to baseline (%)	-0	<b>-10.50</b>	-9.01	-7.13

#### 4.5. Results on BlendedMVS dataset

Since the BlendedMVS dataset does not contain officially provided 3D reconstruction results, nor does any SOTA algorithms provide officially BlendedMVS dataset training implementation except MVSNet and R-MVSNet that only provide pre-trained weights without fusion results, we only show some qualitative results. We first pick up all the large scale outdoor scenes from all samples, and split the picked scenes with the ratio 4 : 1 for training and evaluation.

We show some 3D reconstruction results in Fig.(9), and it is easily known that our generated point clouds are smooth and complete. The comparison of depth map generation results between AACVP-MVSNet and MVSNet is shown in Fig.(10), and it’s obvious that our depth map has a higher resolution and while captures more high-frequency details in edgy areas.

#### 4.6. Ablation Studies

In this section, we provide ablation experiments and quantitative analysis to evaluate the strengths and limitations of the key components in our framework. For all the following studies, experiments are performed and evaluated on the DTU dataset, and both accuracy and completeness are used to measure the reconstruction quality. We set the number of groups  $G = 4$ , and all settings are the same as used in Section 4.2.

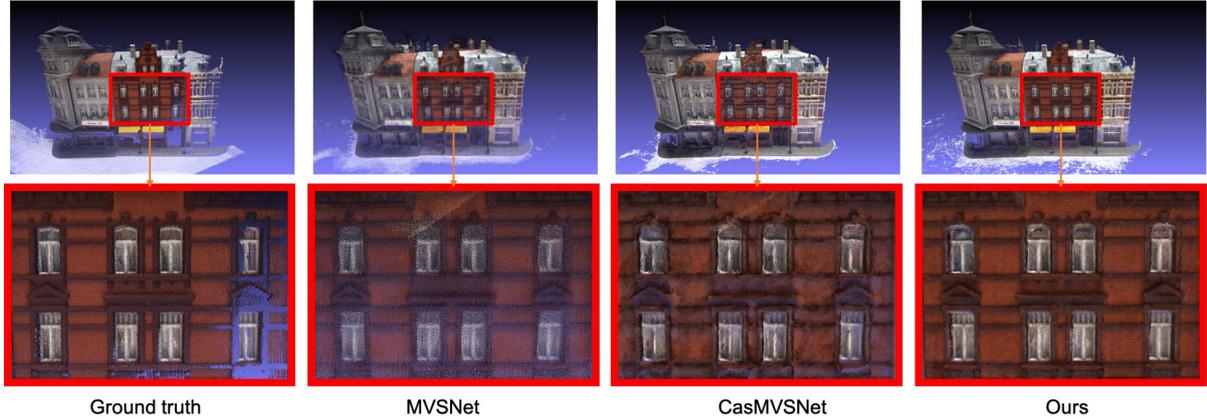


Figure 6: 3D reconstruction result of 9th scene in DTU dataset.

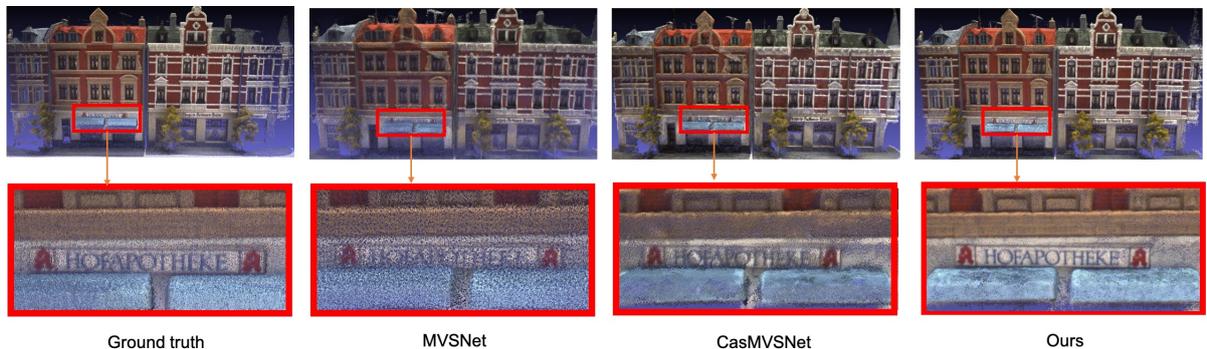


Figure 7: 3D reconstruction result of 15th scene in DTU dataset.

#### 4.6.1. Multi-head self-attention layers

In practice, multiple attention heads are used to learn multiple distinct representations of the input in many applications. It could be implemented by partitioning the input  $x_{ij}$  in Eq.(2) and Eq.(3) into depth-wise into  $H$  groups  $x_{ij}^h \in \mathbb{R}^{d_{in}/H}$  ( $h = 0, 1, \dots, H - 1$ ) (the head), and compute the learned parameters  $\mathbf{W}_l^h$  ( $l = Q, K, V$ ) for all the heads, followed by concatenating the output representations into  $y_{ij}$  as output (Ramachandran et al., 2019).

We set  $H = 2, 4, 8$  respectively, and the reconstruction quality stays almost the same when  $H = 2, 4$ , as shown in Tab.(3) where the single-headed denotes the self-attention layer used in Fig.(2). When  $H$  is set to 8, the overall accuracy of reconstruction decreases slightly, which illustrates that the number of channels may not be too few in each group (2 channels in a group when  $H = 8$ ).

Table 3: Reconstruction quality on DTU dataset by AACVP-MVSNet with different parameter  $H$ .

	Singe-headed	$H = 2$	$H = 4$	$H = 8$
Acc. (mm)	<b>0.357</b>	0.362	0.359	0.375
Comp. (mm)	0.326	<b>0.325</b>	0.332	0.339
OA (mm)	<b>0.341</b>	0.343	0.345	0.357

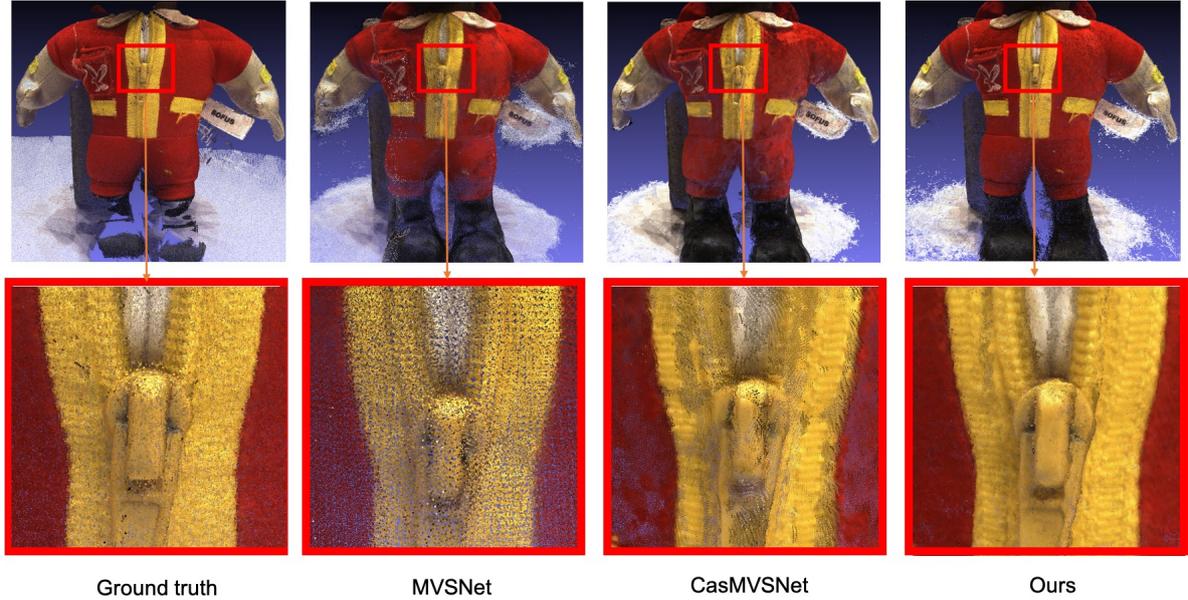


Figure 8: 3D reconstruction result of 49th scene in DTU dataset.

#### 4.6.2. Number of views in training and evaluation

Since multi-view images would provide more information for depth inference task, we choose the number of views for training  $nViews_T = 3, 5, 7$  and the number of views for evaluation  $nViews_E = 3, 4, 5$  for experiments. Tab.(4) shows the comparison results, which illustrates that the result of reconstruction would improve when the number of views for evaluation increases while the accuracy stays almost the same with variant of views quantity for training. The best overall accuracy in our experiments is 0.326mm when  $nViews_T = 7$  and  $nViews_E = 5$ , which is the best result on DTU dataset so far as we know.

Table 4: Reconstruction quality on DTU dataset with different number of views.

	Acc. (mm)	Comp. (mm)	OA (mm)
$nViews_T = 3, nViews_E = 3$	0.357	0.326	0.341
$nViews_T = 3, nViews_E = 4$	0.355	0.320	0.338
$nViews_T = 3, nViews_E = 5$	0.359	0.319	0.339
$nViews_T = 5, nViews_E = 3$	0.361	0.326	0.343
$nViews_T = 5, nViews_E = 4$	0.363	0.309	0.336
$nViews_T = 5, nViews_E = 5$	0.365	0.305	0.335
$nViews_T = 7, nViews_E = 3$	0.361	0.326	0.343
$nViews_T = 7, nViews_E = 4$	<b>0.349</b>	0.306	0.328
$nViews_T = 7, nViews_E = 5$	0.353	<b>0.299</b>	<b>0.326</b>

#### 4.6.3. Convergence

The train-loss curves is plotted in Fig.(11) with  $nViews_T = 3, 5, 7$ , which illustrates that our model could converge and achieve similar loss at 40th epoch.



(a) Outdoor Scene1.



(b) Outdoor Scene2.



(c) Outdoor Scene3.



(d) Outdoor Scene4.

Figure 9: Results on the BlendedMVS dataset.



Figure 10: Comparison of depth inference results between MVSNet and AACVP-MVSNet.

## 5. Discussion

Based on the results shown above, we could confirm that the proposed architecture, which take advantages of self-attention and similarity measurement based cost volume generation in 3D reconstruction, could be trained iteratively with a coarse-to-fine strategy and achieves better performance than some state-of-the-art methods by extensive evaluation on benchmark datasets. In this section, we discuss the limitations of our work and some possible future works.

### 5.1. Depth searching range at finer levels

The depth searching range in AACVP-MVSNet is determined by the distance of the projection of two neighbor pixels along the epipolar lines in images at finer levels in our architecture. Though it is tested in [Yang et al. \(2020\)](#) that evaluated the 3D reconstruction accuracy as a function of the interval setting, this method is still an empirical way for depth search space determination. In the future, we will investigate whether the depth searching space could be learned in the network.

### 5.2. Number of hypothesized depth planes

Both in our work and some previous learned MVS methods such as [Chen et al., 2019a](#); [Yang et al., 2020](#), the number of hypothesized depth planes  $M$  is determined empirically. Thus, the interval of hypothesized planes might not be the most reasonable ones for depth map estimation at each level, which may limits the usage of the proposed architecture in different scenes. We would like to research and propose a self-adaptive method for the determination of the parameter  $M$  in the future.

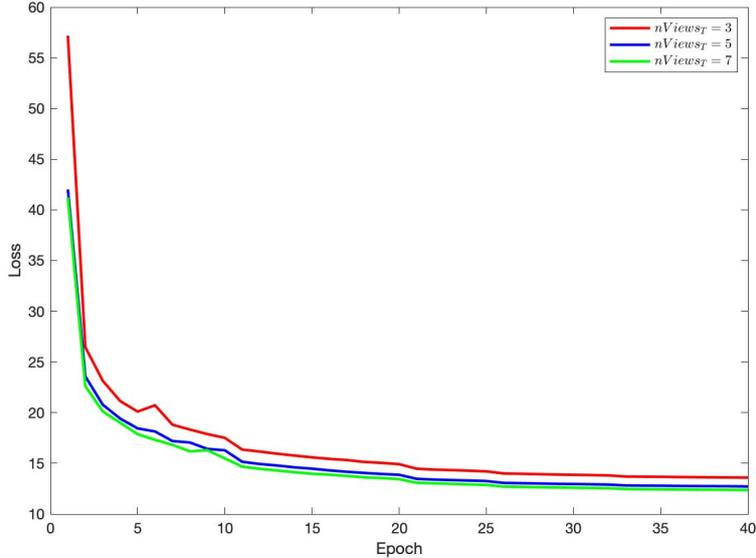


Figure 11: Training loss with  $nViews_T = 3, 5, 7$ .

### 5.3. Learned MVS methods without ground truth

As to our knowledge, most MVS networks are based on supervised learning that requires ground truth which may result in lack of accuracy when the scenes are not similar to the dataset that used for achieving the pre-trained weights. Though some unsupervised MVS methods are proposed recently (Dai et al., 2019; Huang et al., 2020), the accuracy and completeness are still far lower than those of the supervised methods. In the future, we will research on this topic and extend the application of MVS networks to more situations.

## 6. Conclusion

In this paper, we proposed AACVP-MVSNet for MVS problems, which is equipped with the self-attention based feature extraction and similarity measurement based cost volume generation method. The AACVP-MVSNet can estimate the depth map by using a coarse-to-fine strategy. The experimental results show that AACVP-MVSNet outperforms some state-of-the-arts MVS networks after an extensive evaluation on two challenging benchmark datasets. In the future, we want to further improve our network, especially replacing those empirical parameters with the self-adaptive ones.

## 7. Acknowledgment

This work is supported by National Natural Science Foundation of China (No.41801388 and No.41801319).

## References

Aanæs, H., Jensen, R., Vogiatzis, G., Tola, E., Dahl, A., 2016. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision* 120. doi:[10.1007/s11263-016-0902-9](https://doi.org/10.1007/s11263-016-0902-9).

- Bartunov, S., Santoro, A., Richards, B., Marris, L., Hinton, G.E., Lillicrap, T., 2018. Assessing the scalability of biologically-motivated deep learning algorithms and architectures, in: *Advances in Neural Information Processing Systems*, pp. 9368–9378.
- Bitelli, G., Girelli, V.A., Lambertini, A., 2018. Integrated use of remote sensed data and numerical cartography for the generation of 3d city models. *International Archives of the Photogrammetry, Remote Sensing & Spatial Information Sciences* 42.
- Buyukdemircioglu, M., Kocaman, S., 2020. Reconstruction and efficient visualization of heterogeneous 3d city models. *Remote Sensing* 12, 2128.
- Campbell, N.D., Vogiatzis, G., Hernández, C., Cipolla, R., 2008. Using multiple hypotheses to improve depth-maps for multi-view stereo, in: *European Conference on Computer Vision*, Springer. pp. 766–779.
- Chen, R., Han, S., Xu, J., Su, H., 2019a. Point-based multi-view stereo network, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1538–1547.
- Chen, R., Han, S., Xu, J., Su, H., 2019b. Point-based multi-view stereo network, in: *The IEEE International Conference on Computer Vision (ICCV)*.
- Cheng, S., Xu, Z., Zhu, S., Li, Z., Li, L.E., Ramamoorthi, R., Su, H., 2020. Deep stereo using adaptive thin volume representation with uncertainty awareness, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2524–2534.
- Chorowski, J.K., Bahdanau, D., Serdyuk, D., Cho, K., Bengio, Y., 2015. Attention-based models for speech recognition, in: *Advances in neural information processing systems*, pp. 577–585.
- Cordonnier, J.B., Loukas, A., Jaggi, M., 2019. On the relationship between self-attention and convolutional layers. *arXiv* .
- Dai, Y., Zhu, Z., Rao, Z., Li, B., 2019. Mvs2: Deep unsupervised multi-view stereo with multi-view symmetry, in: *2019 International Conference on 3D Vision (3DV)*, IEEE. pp. 1–8.
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* .
- Esteban, C.H., Schmitt, F., 2004. Silhouette and stereo fusion for 3d object modeling. *Computer Vision and Image Understanding* 96, 367–392.
- Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R., 2010. Towards internet-scale multi-view stereo, in: *Computer Vision and Pattern Recognition*.
- Furukawa, Y., Ponce, J., 2009. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence* 32, 1362–1376.
- Galliani, S., Lasinger, K., Schindler, K., 2015a. Massively parallel multiview stereopsis by surface normal diffusion, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 873–881.
- Galliani, S., Lasinger, K., Schindler, K., 2015b. Massively parallel multiview stereopsis by surface normal diffusion, in: *IEEE International Conference on Computer Vision*.

- Gu, X., Fan, Z., Zhu, S., Dai, Z., Tan, F., Tan, P., 2020. 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, pp. 2495–2504.
- Guo, X., Yang, K., Yang, W., Wang, X., Li, H., 2019. Group-wise correlation stereo network, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- Harazono, Y., Ishii, H., Shimoda, H., Kouda, Y., 2019. Development of a scanning support system using augmented reality for 3d environment model reconstruction, in: International Conference on Intelligent Human Systems Integration, Springer. pp. 460–464.
- Hartmann, W., Galliani, S., Havlena, M., Van Gool, L., Schindler, K., 2017. Learned multi-patch similarity, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 1586–1594.
- Hirschmüller, H., 2005. Accurate and efficient stereo processing by semi-global matching and mutual information, in: CVPR 2005.
- Hirschmüller, H., 2007. Stereo processing by semiglobal matching and mutual information. IEEE Transactions on pattern analysis and machine intelligence 30, 328–341.
- Hirschmüller, H., Scharstein, D., 2007. Evaluation of cost functions for stereo matching, in: 2007 IEEE Conference on Computer Vision and Pattern Recognition, IEEE. pp. 1–8.
- Huang, B., Yi, H., Huang, C., He, Y., Liu, J., Liu, X., 2020. M3vsnet: Unsupervised multi-metric multi-view stereo network. ArXiv abs/2004.09722v2.
- Ji, M., Gall, J., Zheng, H., Liu, Y., Fang, L., 2017. Surfacenet: An end-to-end 3d neural network for multiview stereopsis, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 2307–2315.
- Kar, A., Häne, C., Malik, J., 2017. Learning a multi-view stereo machine, in: Advances in neural information processing systems, pp. 365–376.
- Kendall, A., Martirosyan, H., Dasgupta, S., Henry, P., Kennedy, R., Bachrach, A., Bry, A., 2017. End-to-end learning of geometry and context for deep stereo regression, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 66–75.
- Kingma, D., Ba, J., 2014. Adam: A method for stochastic optimization. Computer Science .
- Knobelreiter, P., Reinbacher, C., Shekhovtsov, A., Pock, T., 2017. End-to-end training of hybrid cnn-crf models for stereo, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2339–2348.
- Koch, C., Paal, S.G., Rashidi, A., Zhu, Z., Koenig, M., Brilakis, I., 2014. Achievements and challenges in machine vision-based inspection of large concrete structures. Advances in Structural Engineering 17, 303–318.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. nature 521, 436–444.
- Luo, W., Schwing, A.G., Urtasun, R., 2016. Efficient deep learning for stereo matching, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5695–5703.

- Ma, Z., Liu, S., 2018. A review of 3d reconstruction techniques in civil engineering and their applications. *Advanced Engineering Informatics* 37, 163–174.
- Malihi, S., Valadan Zoej, M., Hahn, M., Mokhtarzade, M., Arefi, H., 2016. 3d building reconstruction using dense photogrammetric point cloud. *Proceedings of the International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences* 3, 71–74.
- Masiero, A., Chiabrando, F., Lingua, A., Marino, B., Fissore, F., Guarnieri, A., Vettore, A., 2019. 3d modeling of girifalco fortress. *International Archives of the Photogrammetry, Remote Sensing & Spatial Information Sciences* .
- Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., Shlens, J., 2019. Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909* .
- Rottensteiner, F., Sohn, G., Gerke, M., Wegner, J.D., Breitkopf, U., Jung, J., 2014. Results of the isprs benchmark on urban object detection and 3d building reconstruction. *ISPRS journal of photogrammetry and remote sensing* 93, 256–271.
- Scharstein, D., Szeliski, R., 2002. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International journal of computer vision* 47, 7–42.
- Schonberger, J.L., Frahm, J.M., 2016. Structure-from-motion revisited, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4104–4113.
- Seitz, S.M., Curless, B., Diebel, J., Scharstein, D., Szeliski, R., 2006. A comparison and evaluation of multi-view stereo reconstruction algorithms, in: *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, IEEE. pp. 519–528.
- Seki, A., Pollefeys, M., 2017. Sgm-nets: Semi-global matching with neural networks, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 231–240.
- Shaw, P., Uszkoreit, J., Vaswani, A., 2018. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155* .
- Shazeer, N., Cheng, Y., Parmar, N., Tran, D., Vaswani, A., Koanantakool, P., Hawkins, P., Lee, H., Hong, M., Young, C., et al., 2018. Mesh-tensorflow: Deep learning for supercomputers, in: *Advances in Neural Information Processing Systems*, pp. 10414–10423.
- Shen, X., Wang, C., Li, X., Yu, Z., Li, J., Wen, C., Cheng, M., He, Z., 2019. Rf-net: An end-to-end image matching network based on receptive field, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8132–8140.
- Sing, K.H., Xie, W., 2016. Garden: a mixed reality experience combining virtual reality and 3d reconstruction, in: *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in Computing Systems*, pp. 180–183.
- Tola, E., Strecha, C., Fua, P., 2012. Efficient large-scale multi-view stereo for ultra high-resolution image sets. *Machine Vision and Applications* 23, 903–920.
- Tulyakov, S., Ivanov, A., Fleuret, F., 2018. Practical deep stereo (pds): Toward applications-friendly deep stereo matching, in: *Advances in Neural Information Processing Systems*, pp. 5871–5881.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: *Advances in neural information processing systems*, pp. 5998–6008.
- Vu, H.H., Labetut, P., Pons, J.P., Keriven, R., 2011. High accuracy and visibility-consistent dense multiview stereo. *IEEE transactions on pattern analysis and machine intelligence* 34, 889–901.
- Xiang, X., Wang, Z., Lao, S., Zhang, B., 2020. Pruning multi-view stereo net for efficient 3d reconstruction. *ISPRS Journal of Photogrammetry and Remote Sensing* 168, 17–27.
- Xu, J., Ranftl, R., Koltun, V., 2017. Accurate optical flow via direct cost volume processing, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1289–1297.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y., 2015. Show, attend and tell: Neural image caption generation with visual attention, in: *International conference on machine learning*, pp. 2048–2057.
- Yang, J., Mao, W., Alvarez, J.M., Liu, M., 2020. Cost volume pyramid based depth inference for multi-view stereo, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4877–4886.
- Yang, M.D., Chao, C.F., Huang, K.S., Lu, L.Y., Chen, Y.P., 2013. Image-based 3d scene reconstruction and exploration in augmented reality. *Automation in Construction* 33, 48–60.
- Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L., 2018. Mvsnet: Depth inference for unstructured multi-view stereo, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 767–783.
- Yao, Y., Luo, Z., Li, S., Shen, T., Fang, T., Quan, L., 2019. 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.
- Yao, Y., Luo, Z., Li, S., Zhang, J., Ren, Y., Zhou, L., Fang, T., Quan, L., 2020. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. *Computer Vision and Pattern Recognition (CVPR)*.
- Yi, H., Wei, Z., Ding, M., Zhang, R., Chen, Y., Wang, G., Tai, Y.W., 2020. Pyramid multi-view stereo net with self-adaptive view aggregation, in: *ECCV*.
- Zbontar, J., LeCun, Y., 2015. Computing the stereo matching cost with a convolutional neural network, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1592–1599.
- Žbontar, J., LeCun, Y., 2016. Stereo matching by training a convolutional neural network to compare image patches. *The journal of machine learning research* 17, 2287–2318.
- Zhang, Z., Lee, W.S., 2019. Deep graphical feature learning for the feature matching problem, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5087–5096.
- Zhu, X.X., Tuia, D., Mou, L., Xia, G.S., Zhang, L., Xu, F., Fraundorfer, F., 2017. Deep learning in remote sensing: A comprehensive review and list of resources. *IEEE Geoscience and Remote Sensing Magazine* 5, 8–36.