

# MVPNet: Multi-View Point Regression Networks for 3D Object Reconstruction from A Single Image

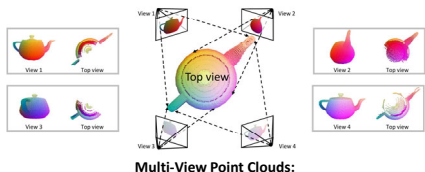


Jinglu Wang (MSRA) Bo Sun (PKU) Yan Lu (MSRA)

## Motivation

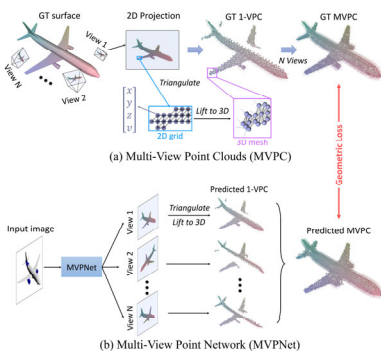
**Why multi-view representation?** Existing representations for 3D reconstruction are mainly of four categories. While 3D volumetric grids suffer high computational complexity, unordered point sets require to solve point-wise mapping, meshes are difficult for CNNs to encode and decoder, multi-view based representation are convolution-favored, ordered and can depict dense and detailed surface.

**Overcome limitations of multi-view representation.** View based representation handing features in 2D projective space can neglect information loss through dimension reduction from 3D to 2D. The proposed MVPC constructs meshes in 3D from 2D grids. MVPC allows us to discretize integrals of surface variations over the constructed triangular mesh and to enforce multi-view consistency with view correspondences.



## Overview

We reconstruct an object's surface from a single image by 1) representing a 3D surface as MVPC, 2) regressing MVPC with MVPNet, and 3) proposing a geometric loss to interpret discrepancy over 3D surfaces.



- Multi-View Point Clouds:** A surface is represented by Multi-View Point Cloud (MVPC). Each pixel in a 1-VPC stores the backprojected surface point  $(x, y, z)$  from this pixel and its visibility  $v$ . The stored 3D points are triangulated according to the 2D grid on the image plane and their normals are shown to indicate surface orientation.
- Multi-View Point Network:** Given an RGB image, the MVPNet generates a set of 1-VPCs and their union forms the predicted MVPC. The geometric loss measures discrepancy between predicted and ground truth MVPC.

## Results

We present qualitative and quantitative results of the reconstruction, comparing the proposed method two collections of state-of-the-art methods according to result representations, i.e., point clouds and volumetric grids.

**Qualitative comparison to point generation method [1]**

**Quantitative comparison with voxel IoU**

Method	plane	branch	object	car	chair	dishbowl	lamp	speaker	teapot	cup	table	chair	sofa	bed
ReconNet at 2017 CVPR	0.019	0.021	0.022	0.024	0.026	0.028	0.030	0.032	0.034	0.036	0.038	0.040	0.042	0.044
ReconNet at 2018 CVPR	0.024	0.027	0.029	0.031	0.033	0.035	0.037	0.039	0.041	0.043	0.045	0.047	0.049	0.051
CVXNet at 2018 CVPR	0.025	0.028	0.030	0.032	0.034	0.036	0.038	0.040	0.042	0.044	0.046	0.048	0.050	0.052
SurfaceNet at 2018 CVPR	0.026	0.029	0.031	0.033	0.035	0.037	0.039	0.041	0.043	0.045	0.047	0.049	0.051	0.053
SurfaceNet and Depth 2017	0.027	0.030	0.032	0.034	0.036	0.038	0.040	0.042	0.044	0.046	0.048	0.050	0.052	0.054
Ours	0.028	0.031	0.033	0.035	0.037	0.039	0.041	0.043	0.045	0.047	0.049	0.051	0.053	0.055

**Quantitative comparison with chamfer distance metric**

Method	plane	branch	object	car	chair	dishbowl	lamp	speaker	teapot	cup	table	chair	sofa	bed
ReconNet at 2017 CVPR	1.200	1.800	2.400	3.000	3.600	4.200	4.800	5.400	6.000	6.600	7.200	7.800	8.400	9.000
ReconNet at 2018 CVPR	1.100	1.700	2.300	2.900	3.500	4.100	4.700	5.300	5.900	6.500	7.100	7.700	8.300	8.900
CVXNet at 2018 CVPR	1.050	1.650	2.250	2.850	3.450	4.050	4.650	5.250	5.850	6.450	7.050	7.650	8.250	8.850
SurfaceNet at 2018 CVPR	1.000	1.600	2.200	2.800	3.400	4.000	4.600	5.200	5.800	6.400	7.000	7.600	8.200	8.800
Ours	0.950	1.550	2.150	2.750	3.350	3.950	4.550	5.150	5.750	6.350	6.950	7.550	8.150	8.750

**Qualitative compare to voxel-based method [2]**

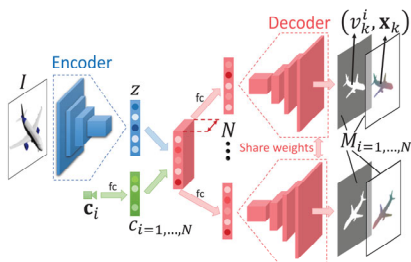
**Reconstruction results on real word data**

## Network Architecture

The MVPNet is an encoder-decoder generative network incorporating camera parameters to generate view-dependent point clouds.

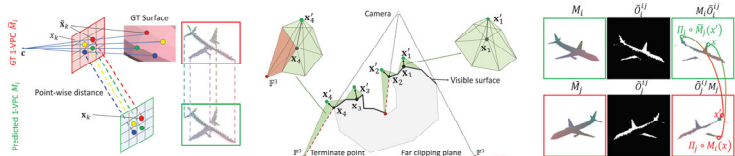
- The **encoder** maps an image  $I$  to an embedding space to obtain a feature  $Z$ . Each camera matrix  $c_i$  is first mapped to a feature  $c_i$ , serving as a view indicator, and then is concatenated with  $Z$  to get  $(Z, c_i)$ .
- The **decoder** converting  $(Z, c_i)$  to a 1-VPC  $M_i$  indicated by  $c_i$  learns the projective transformation and space completion. The decoder shares weights among  $N$  branches.

The output MVPC is of shape  $N \times H \times W \times 4$ . The last channel corresponds to a 3D coordinate  $x_k = (x_k, y_k, z_k)$  and visibility  $v_k$  of a point  $x_k$ .



## Geometric Loss

We propose a geometric loss (GeoLoss) that is able to capture variances over 3D surfaces rather than over sparse point sets or 2D projective planes. The GeoLoss is made up of three components:



**Point-wise distance term.** The points in ground truth and predicted 1-VPC have a one-to-one mapping, since 2D pixels with equal 2D coordinates are defined to store the same surface point induced by the same viewpoint. The sum of point-wise distances for ground truth and predicted 1-VPC is the L2 loss.

$$\mathcal{L}_{ptd} = \sum_{x \in M_i} \|M_i(x) - \tilde{M}_i(x)\|_2$$

**Quasi-volume term.** Inspired by the volume-preserving constraints used in variational surface deformation, we propose a quasi-volume discrepancy metric to describe the surface discrepancy, characterizing details and handling occluding contours.

$$\mathcal{L}_{vol} = \int_{\mathcal{S}} (\mathbf{x} - \tilde{\mathbf{x}}) \cdot \mathbf{n} dx$$

It is discretized in the MVPC's mesh as:

$$\mathcal{L}_{vol} = \sum_{x \in M_i} \tilde{V}_i(x)(M_i(x) - \tilde{M}_i(x)) \cdot \tilde{N}_i(x)$$

**Multi-view consistency term.** Some 1-VPCs may have overlap. Corresponding points should be close. We minimize sum of two distances between stored 3D points in two corresponding pixels and their reprojected pixels in another 1-VPC.

$$\mathcal{L}_{mvo} = \sum_{x \in \mathcal{O}_i^j} \|M_i(x) - \tilde{M}_j(\Pi_j \circ M_i(x))\|_2 + \sum_{x \in \mathcal{O}_j^i} \|\tilde{M}_j(x) - M_i(\Pi_i \circ \tilde{M}_j(x))\|_2$$

## Application

We show the generative representation of the learned features with interpolation, arithmetic, classification and clustering.

**Arithmetic.** The shapes of last column are obtained by decoding the feature  $(A+B-C)$ .

**Interpolation.** Reconstructions from linear interpolation of two learned features.

**Classification.**

**Clustering.**

## Reference

- [1] Fan, H.; Su, H.; and Guibas, L. 2017. A point set generation network for 3d object reconstruction from a single image. In ICCV.
- [2] Choy, C. B.; Xu, D.; Gwak, J.; Chen, K.; and Savarese, S. 2016. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In ECCV.

## Website

