

# SphereFusion: Efficient Panorama Depth Estimation via Gated Fusion

## Supplementary Material

### 001 1. Panorama Projections

002 To capture the texture features and avoid distortion and  
003 discontinuity, SphereFusion relies on two panorama pro-  
004 jections: the equirectangular projection and the spherical  
005 mesh. In this section, we describe the details of their con-  
006 version, and how to convert the panorama image in spher-  
007 ical mesh to equirectangular for depth map evaluation or  
008 re-project it to 3D space for visualization.

#### 009 1.1. The E2S

010 Based on the definition of the equirectangular projection  
011 and the spherical mesh, we define the E2S ( equirectangular  
012 projection to spherical projection ) module, which converts  
013 a panorama image from the equirectangular projection to  
014 the spherical projection. Given a triangle center  $(x, y, z)$   
015 from the spherical mesh, we first calculate its position on  
016 the image plane  $(u, v)$  by Eq. 1, and then use a bi-linear  
017 sample to capture value on the image plane, finally assign  
018 corresponding value to the triangle.

$$019 \begin{cases} u = (1 + \operatorname{atan}(y, x)/\pi) \times W/2 \\ v = (0.5 + \operatorname{atan}(z, \sqrt{x^2 + y^2})/\pi) \times H \end{cases} \quad (1)$$

#### 020 1.2. The S2E

021 Meanwhile, we can also define the S2E ( spherical pro-  
022 jection to equirectangular projection ) module, which con-  
023 verts a panorama image from the spherical projection to the  
024 equirectangular projection. Given a pixel  $(u, v)$  on the im-  
025 age plane, we first calculate its position on the sphere sur-  
026 face, then calculate its 3D position to find out the closest  
027 triangle on the spherical mesh, and finally assign the value  
028 of the triangle to the pixel.

$$029 \begin{cases} x = \cos(\operatorname{longitude})\cos(\operatorname{latitude}) \\ y = \sin(\operatorname{longitude})\cos(\operatorname{latitude}) \\ z = \sin(\operatorname{latitude}) \end{cases} \quad (2)$$

#### 030 1.3. Projection To 3D Point Cloud

031 For better comparison, we visualize point clouds using dif-  
032 ferent methods. Although we name the result obtained from  
033 the panorama depth estimation as the depth map, it repre-  
034 sents the distance map, which means the Euclidean distance  
035 between a 3D point and the camera center. Therefore, we  
036 can reproject a panorama depth map to 3D space by Eq. 3  
037 after converting it into the sphere coordinates, where  $d$  is  
038 the distance.

$$\begin{cases} X = \cos(\operatorname{latitude})\cos(\operatorname{longitude}) \times d \\ Y = \cos(\operatorname{latitude})\sin(\operatorname{longitude}) \times d \\ Z = \sin(\operatorname{latitude}) \times d \end{cases} \quad (3) \quad 039$$

### 040 2. Evaluation

041 To compare with other methods, we convert our depth map  
042 in the spherical domain to equirectangular projection. Fol-  
043 lowing BiFuse [7] and SliceNet [5], we use five evalua-  
044 tion metrics, including MAE (mean average error), MRE  
045 (mean relative error), RMSE (root mean square error),  
046 RMSE(log), and  $\delta^n$ . Eq. 4 shows how to calculate them,  
047 where  $gt$  is the ground truth,  $pr$  is the predicted depth,  $V$   
048 is valid pixels, and  $N$  is the number of valid pixels. For  
049 MAE, MRE, RMSE, and RMSE(log), the smaller is better.  
050 For the  $\delta$ , the bigger is better. During the evaluation, we set  
051 the depth range to 0.1 ~ 10 meters.

$$\begin{aligned} MAE &= \sum_{i \in V} |gt_i - pr_i| \\ MRE &= \sum_{i \in V} \frac{|gt_i - pr_i|}{gt_i} \\ RMSE &= \sqrt{\frac{\sum_{i \in V} (gt_i - pr_i)^2}{N}} \\ RMSE_{log} &= \sqrt{\frac{\sum_{i \in V} (\log_{10}(gt_i) - \log_{10}(pr_i))^2}{N}} \\ \delta^n &= \frac{\sum_{i \in V} \max(\frac{gt_i}{pr_i}, \frac{pr_i}{gt_i}) < 1.25^n}{N} \end{aligned} \quad (4) \quad 052$$

### 053 3. Visualization

054 In this section, we add more visualization results on three  
055 datasets and compare SphereFusion (ours) with state-of-  
056 the-art methods.

057 On 360D [9], we compare our method with SphereDepth  
058 [8] and UniFuse [3]. Figure 1 shows depth maps and point  
059 clouds generated by different methods. Our method uses  
060 texture features extracted by the 2D encoder to enhance the  
061 mesh encoder in the spherical domain and combines the  
062 strengths of two encoders. Results show that our method  
063 captures more features from the scene and can reconstruct  
064 more details in the scene, such as doors, tables, and walls.

065 On Matterport3D [2] and Stanford2D3D [1], we com-  
066 pare our method with OmniFusion [8] and PanoFormer [3].  
067 Figure 2 shows depth maps generated by different methods  
068 on Matterport3D and corresponding point clouds. Figure

069 3 shows results generated by different methods on Stan-  
070 ford2D3D. Compared with tangent patches, our method di-  
071 rectly estimates the panorama depth in the spherical domain  
072 and does not need any special mechanism to fuse these  
073 patches. The results of OmniFusion have obvious patch  
074 gaps. Meanwhile, the results of PanoFormer are smoother  
075 but also lose some details. Our method reconstructs more  
076 details but suffers from imperfect ground truth [3]. Depth  
077 maps and point clouds show that our method achieves com-  
078 petitive results with state-of-the-art methods with a lighter  
079 network and higher efficiency.

## 080 References

- 081 [1] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese.  
082 Joint 2d-3d-semantic data for indoor scene understanding.  
083 *arXiv preprint arXiv:1702.01105*, 2017. 1
- 084 [2] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej  
085 Halber, Matthias Niessner, Manolis Savva, Shuran Song,  
086 Andy Zeng, and Yinda Zhang. Matterport3d: Learning  
087 from rgb-d data in indoor environments. *arXiv preprint*  
088 *arXiv:1709.06158*, 2017. 1
- 089 [3] Hualie Jiang, Zhe Sheng, Siyu Zhu, Zilong Dong, and Rui  
090 Huang. Unifuse: Unidirectional fusion for 360 panorama  
091 depth estimation. *IEEE Robotics and Automation Letters*, 6  
092 (2):1519–1526, 2021. 1, 2, 3
- 093 [4] Yuyan Li, Yuliang Guo, Zhixin Yan, Xinyu Huang, Ye Duan,  
094 and Liu Ren. Omnifusion: 360 monocular depth estimation  
095 via geometry-aware fusion. In *Proceedings of the IEEE/CVF*  
096 *Conference on Computer Vision and Pattern Recognition*,  
097 pages 2801–2810, 2022. 4, 5
- 098 [5] Giovanni Pintore, Marco Agus, Eva Almansa, Jens Schneider,  
099 and Enrico Gobbetti. Slicenet: deep dense depth estimation  
100 from a single indoor panorama using a slice-based representa-  
101 tion. In *Proceedings of the IEEE/CVF Conference on Com-*  
102 *puter Vision and Pattern Recognition*, pages 11536–11545,  
103 2021. 1
- 104 [6] Zhijie Shen, Chunyu Lin, Kang Liao, Lang Nie, Zishuo  
105 Zheng, and Yao Zhao. Panoformer: Panorama transformer  
106 for indoor 360  $\{\deg\}$  depth estimation. *arXiv preprint*  
107 *arXiv:2203.09283*, 2022. 4, 5
- 108 [7] Fu-En Wang, Yu-Hsuan Yeh, Min Sun, Wei-Chen Chiu, and  
109 Yi-Hsuan Tsai. Bifuse: Monocular 360 depth estimation via  
110 bi-projection fusion. In *Proceedings of the IEEE/CVF Con-*  
111 *ference on Computer Vision and Pattern Recognition*, pages  
112 462–471, 2020. 1
- 113 [8] Qingsong Yan, Qiang Wang, Kaiyong Zhao, Bo Li, Xiaowen  
114 Chu, and Fei Deng. Spheredepth: Panorama depth estimation  
115 from spherical domain. In *2022 Tenth international confer-*  
116 *ence on 3D vision (3DV)*, 2022. 1, 3
- 117 [9] Nikolaos Zioulis, Antonis Karakottas, Dimitrios Zarpalas, and  
118 Petros Daras. Omnidepth: Dense depth estimation for indoors  
119 spherical panoramas. In *Proceedings of the European Con-*  
120 *ference on Computer Vision (ECCV)*, pages 448–465, 2018.  
121 1

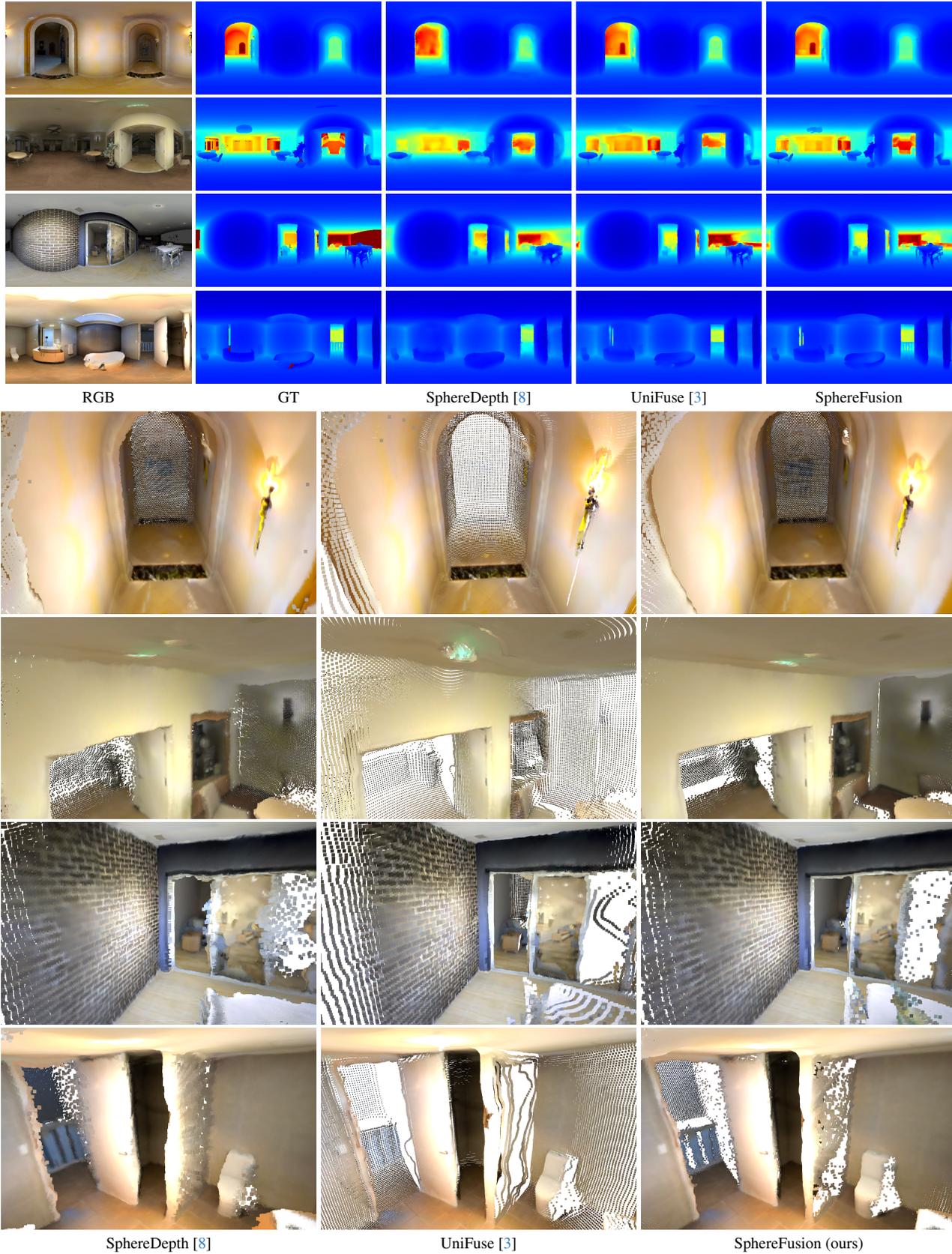


Figure 1. Depth maps and point clouds of 360D. Invalid parts of the depth map are set to red. Our method SphereFusion reconstructs more details of the scene.

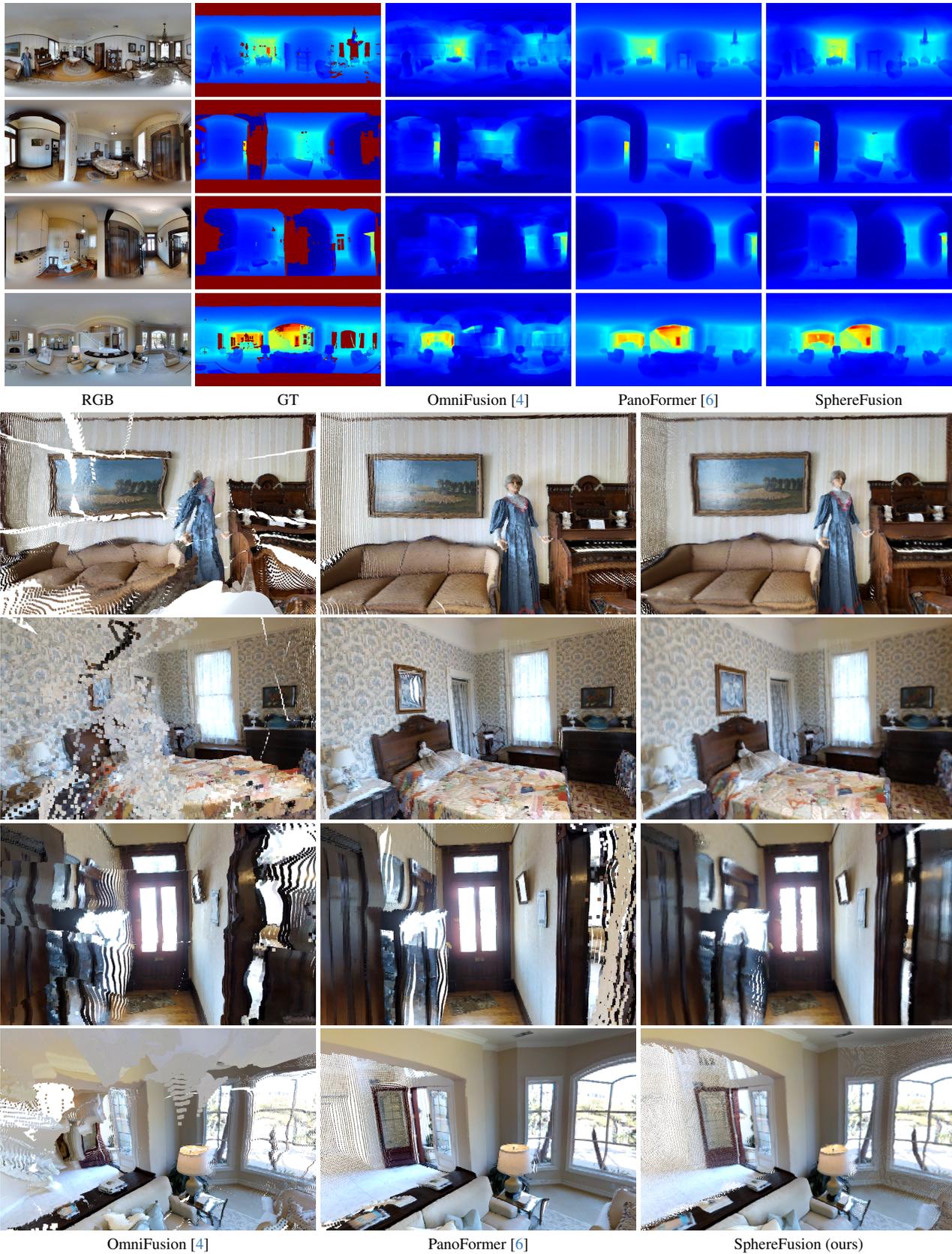


Figure 2. Depth maps and point clouds of Matterport3D. Invalid parts of the depth map are set to red. Our method has less noise and maintains the structure of the scene.

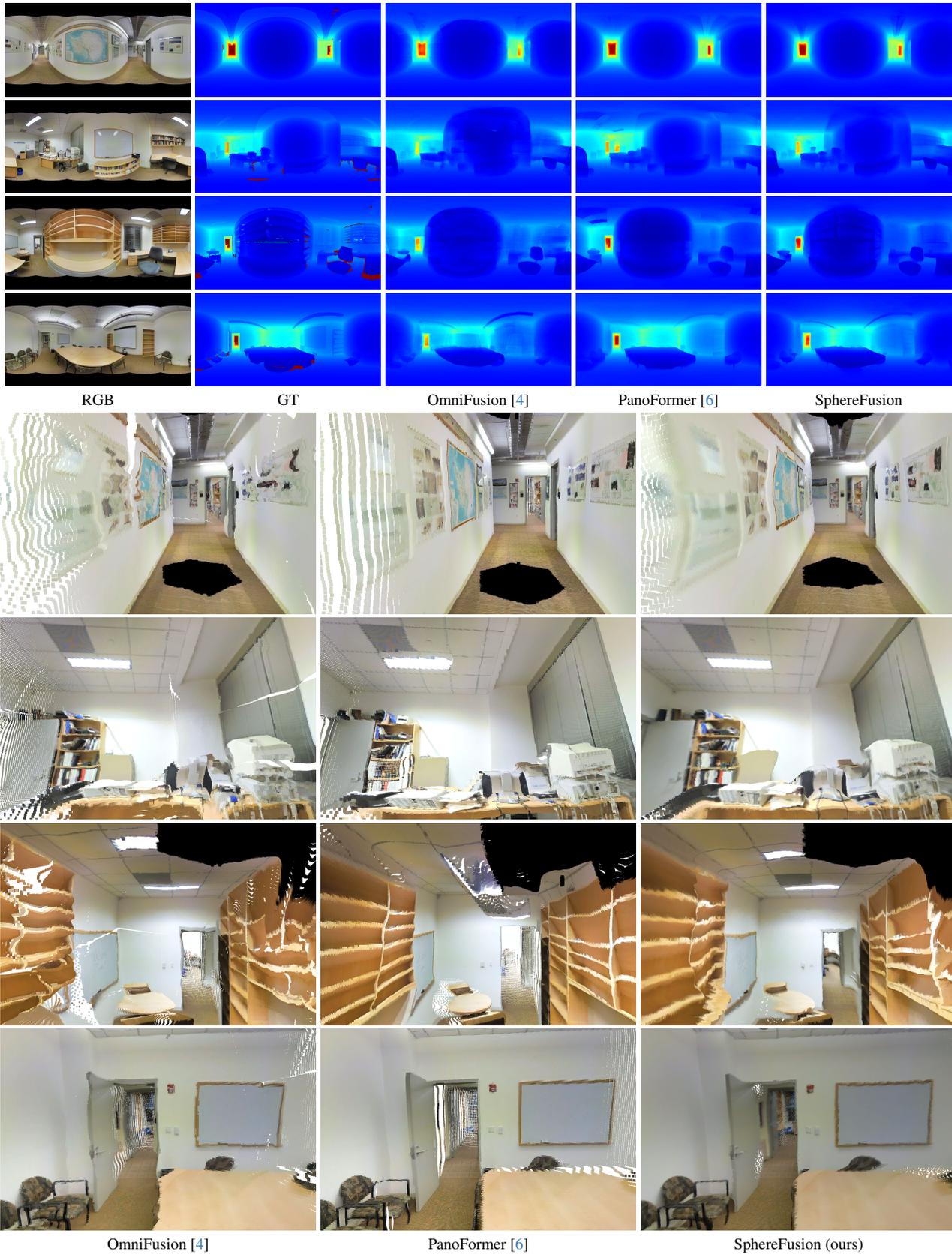


Figure 3. Depth maps and point clouds of Stanford2D3D. Invalid parts of the depth map are set to red. Our method does not suffer from discontinuity.