

# Pano3D: A Holistic Benchmark and a Solid Baseline for 360° Depth Estimation – Supplementary Material

Georgios Albanis<sup>\*1</sup> Nikolaos Zioulis<sup>\*1,2</sup> Petros Drakoulis<sup>1</sup> Vasileios Gkitsas<sup>1</sup> Vladimiros Strezenenko<sup>1</sup>

Federico Alvarez<sup>2</sup> Dimitrios Zarpalas<sup>1</sup> Petros Daras<sup>1</sup>

<sup>1</sup> Centre for Research and Technology Hellas, Thessaloniki, Greece

<sup>2</sup> Universidad Politécnica de Madrid, Madrid, Spain

{galbanis, nzioulis, petros.drakoulis, gkitsasv, vladster}@iti.gr

fag@gatv.ssr.upm.es {zarpalas, daras}@iti.gr

vcl3d.github.io/Pano3D

## 1. Introduction

This supplementary material complements our original manuscript with additional results, supporting further ablation experiments, providing qualitative results on real data and comparisons between the different architectures.

## 2. Supplementary Results

Table 1 complements Table 1 of the main document, presenting the performance of all remaining metrics, namely the spherical direct depth metrics, the boundary preservation metrics, and the smoothness metrics. In addition, Figure 1 presents the different models’ performance in terms of three indicators, one for each trait. These indicators combine an error and an accuracy metric:

$$i_d = \frac{1}{(1 - \delta_{1.25}) \times RMSE}, \quad (1)$$

$$i_b = \frac{1}{(1 - (rec_{0.25} + rec_{0.5} + rec_{1.0})/3) \times dbeacc}, \quad (2)$$

$$i_s = \frac{1}{(1 - (\alpha_{11.25^\circ} + \alpha_{22.5^\circ} + \alpha_{30^\circ})/3) \times RMSE^\circ}, \quad (3)$$

with  $i_d$ ,  $i_b$ , and  $i_s$  the depth, boundary and smoothness performance indicators. Evidently, UNet performs significantly better than the other models, especially in the boundary consistency metrics, while all models benefit of the addition of extra losses. The addition, of skip connections in a common ResNet architecture offers better performance. While  $\mathcal{L}_{grad}$  offers better depth performance for ResNet<sub>skip</sub>, the variant trained with  $\mathcal{L}_{comb}$  offers higher performance across the two secondary traits.

<sup>\*</sup>Indicates equal contribution.

In addition, we complement the main’s paper spherical metrics Table 2 by collating the traditional ones for a straightforward comparison.

Finally, Table 3 reproduces the grounds upon our methodology was designed, namely the efficacy of pre-trained models [6] and the L1 loss [2]. We use the DenseNet and Pnas models with the encoders initialized using weights pre-trained on ImageNet. Both claims stand, with all pre-trained models achieving better performance than the model trained from scratch. In addition, the L1 loss outperforms both berHu [5] and log loss. Interestingly, the performance drops significantly in DenseNet when trained with other losses, while for Pnas the performance gap is smaller. Therefore, when benchmarking different models, this needs to be taken into account as well. Only through consistent experimentation across different aspects measurable and explainable progress will be possible.

## 3. Qualitative Results

Finally we present additional qualitative results for different models. Apart from the collation of the predicted depth maps between the different models, we provide an advantage visualisation technique similar to that presented in HoHoNet [7]. The visualisation is the MAE difference between two comparable models.

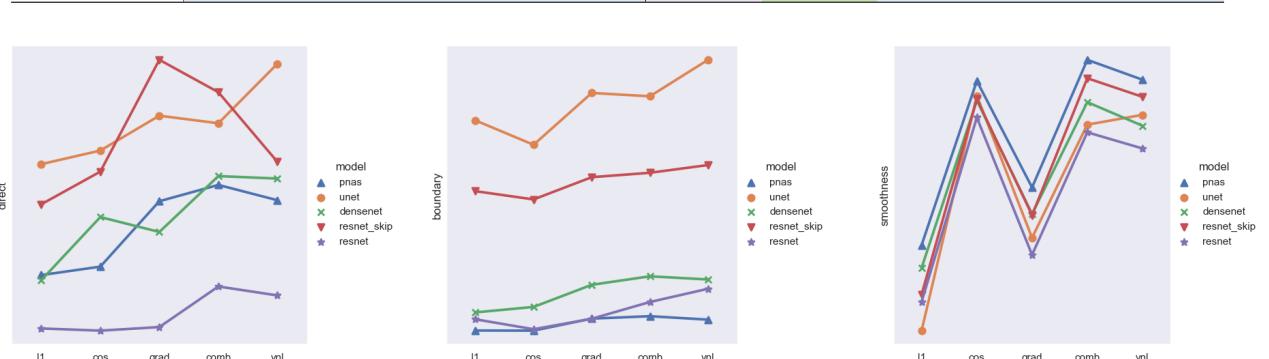
To that end, Figure 2 demonstrates the comparison of ResNet and ResNet<sub>skip</sub> architectures, Figure 3 that of the UNet and Pnas architectures, and, finally Figure 4 presents the differences between the Unet and ResNet<sub>skip</sub> architectures.

Additionally, Figure 5 presents comparative results regarding the boundary preservation performance across models. Once again, UNet is able to capture finer-grained details while the Pnas model produces smoother results.

Table 1: Three axis depth metrics performance across models and supervision schemes. Best three performers are denoted with bold faced **light green** (1<sup>st</sup>), **light blue** (2<sup>nd</sup>) and **light purple** (3<sup>rd</sup>) respectively following the ranking order. Same scheme applies to all tables.

Model	Direct Depth										Depth Discontinuity										Depth Smoothness				
	Error ↓				Accuracy ↑				Error ↓				Accuracy ↑				Error ↓				$RMSE^v$				$\alpha_{11.25^*}$
	wRMSE	wRMSLE	wAbsRel	wSqRel	$\delta_{1.05}^{ico^6}$	$\delta_{1.1}^{ico^6}$	$\delta_{1.25}^{ico^6}$	$\delta_{1.25^2}^{ico^6}$	$\delta_{1.25^3}^{ico^6}$	$dbe_{acc}$	$dbe_{comp}$	$prec_{0.25}$	$prec_{0.5}$	$prec_1$	$rec_{0.25}$	$rec_{0.5}$	$rec_1$	$RMSE^v$	$\alpha_{11.25^*}$	$\alpha_{22.5^*}$	$\alpha_{30^*}$				
Pnas	$L_1$	0.5602	0.0854	0.1326	0.1196	32.69%	56.94%	85.12%	95.38%	97.95%	2,654.2	5,730.3	38.73%	30.26%	23.58%	18.74%	10.48%	8.48%	20.12	53.88%	69.81%	75.65%			
	$L_{cosine}$	0.5622	0.0858	0.1336	0.1317	34.49%	58.06%	85.52%	95.44%	97.88%	2,719.4	5,496.4	36.16%	27.76%	22.22%	21.48%	13.55%	10.02%	15.70	67.14%	77.37%	81.05%			
	$L_{grad}$	0.5374	0.0822	0.1276	<b>0.1146</b>	35.68%	59.54%	86.41%	95.72%	98.04%	<b>2,506.8</b>	5,454.8	<b>40.91%</b>	<b>32.05%</b>	25.27%	22.49%	12.17%	9.07%	18.12	59.56%	73.24%	78.30%			
	$L_{comb}$	<b>0.5367</b>	<b>0.0811</b>	<b>0.1259</b>	0.1153	<b>36.44%</b>	<b>60.52%</b>	<b>86.80%</b>	<b>95.83%</b>	<b>98.11%</b>	2,511.9	<b>5,350.1</b>	39.83%	31.59%	<b>27.01%</b>	<b>23.53%</b>	<b>14.42%</b>	<b>10.98%</b>	<b>15.26</b>	<b>67.73%</b>	<b>77.99%</b>	<b>81.67%</b>			
	$L_{ent}$	0.5403	0.0815	0.1280	0.1183	35.43%	59.72%	86.58%	95.79%	98.11%	2,514.1	5,389.3	40.14%	31.77%	24.47%	22.14%	12.69%	8.74%	15.57	66.61%	77.34%	81.23%			
Unet	$L_1$	0.4834	0.2361	0.1211	0.0913	35.18%	58.24%	86.80%	96.45%	98.43%	1,401.1	4,315.2	57.59%	58.00%	53.85%	38.74%	31.57%	24.31%	24.66	36.80%	60.60%	69.73%			
	$L_{cosine}$	0.4736	<b>0.0906</b>	0.1217	0.0891	32.65%	58.04%	87.40%	96.68%	98.61%	1,451.3	5,045.5	55.35%	52.16%	46.01%	39.36%	30.01%	21.69%	<b>15.80</b>	<b>63.10%</b>	<b>77.60%</b>	<b>82.38%</b>			
	$L_{grad}$	0.4659	0.5186	0.1209	0.0833	35.25%	58.79%	87.33%	96.56%	98.45%	1,330.5	4,058.2	<b>63.13%</b>	<b>56.54%</b>	40.39%	32.47%	23.37%	19.52	52.23%	70.40%	76.91%				
	$L_{comb}$	0.4630	0.1690	0.1222	0.0847	34.79%	58.21%	87.08%	96.63%	98.66%	1,307.7	4,208.0	<b>63.31%</b>	61.74%	54.96%	39.38%	30.27%	22.00%	16.19	61.01%	76.18%	81.45%			
	$L_{ent}$	<b>0.4520</b>	0.1300	<b>0.1147</b>	<b>0.0811</b>	<b>36.68%</b>	<b>60.59%</b>	<b>88.31%</b>	<b>96.96%</b>	<b>98.73%</b>	<b>1,269.9</b>	<b>3,887.6</b>	58.97%	57.54%	51.85%	<b>43.96%</b>	<b>36.69%</b>	<b>28.59%</b>	16.02	61.80%	76.58%	81.70%			
DenseNet	$L_1$	0.5441	0.6872	0.1349	0.1144	34.34%	57.10%	84.73%	95.28%	97.69%	2,369.0	5,513.5	40.40%	36.07%	28.78%	20.45%	11.54%	8.05%	21.08	49.98%	68.29%	74.78%			
	$L_{cosine}$	0.5361	<b>0.0822</b>	0.1239	0.1034	34.98%	59.34%	86.36%	95.94%	<b>98.13%</b>	2,348.6	5,370.2	41.01%	35.45%	29.10%	22.80%	14.19%	9.39%	<b>15.97</b>	<b>64.92%</b>	<b>76.91%</b>	81.15%			
	$L_{grad}$	<b>0.5202</b>	0.4655	0.1304	0.1045	32.68%	57.59%	85.69%	95.85%	98.06%	2,078.9	5,215.0	47.01%	40.61%	33.32%	23.68%	13.71%	9.35%	18.90	56.86%	71.79%	77.23%			
	$L_{comb}$	0.5209	0.1982	<b>0.1209</b>	<b>0.1013</b>	35.97%	<b>60.41%</b>	<b>87.02%</b>	<b>95.96%</b>	98.09%	2,062.8	5,097.7	<b>47.16%</b>	<b>40.77%</b>	<b>35.20%</b>	<b>26.09%</b>	<b>16.87%</b>	<b>12.21%</b>	15.98	64.58%	76.86%	<b>81.20%</b>			
	$L_{ent}$	0.5232	0.7560	0.1258	0.1030	<b>36.28%</b>	60.04%	86.61%	95.66%	97.74%	<b>2,052.5</b>	<b>5,093.1</b>	44.81%	40.14%	32.30%	25.21%	15.71%	10.33%	16.51	63.43%	76.02%	80.53%			
ResNet <sub>skip</sub>	$L_1$	0.5500	0.1922	0.1394	0.1186	30.59%	54.17%	84.07%	95.47%	98.03%	2,438.6	5,768.8	39.10%	31.69%	23.28%	20.92%	10.24%	6.32%	22.83	44.68%	64.51%	72.02%			
	$L_{cosine}$	0.5435	<b>0.0864</b>	0.1364	0.1194	<b>34.77%</b>	56.32%	84.29%	95.64%	98.11%	2,691.8	5,792.8	38.35%	32.13%	26.82%	21.88%	12.61%	8.71%	<b>16.37</b>	<b>64.24%</b>	<b>76.30%</b>	<b>80.63%</b>			
	$L_{grad}$	0.5475	0.2976	0.1387	0.1151	32.43%	54.46%	83.76%	95.37%	97.97%	2,411.2	5,759.5	41.87%	33.23%	21.60%	21.31%	9.27%	4.95%	20.50	52.77%	68.97%	75.00%			
	$L_{comb}$	<b>0.5294</b>	0.1365	0.1374	0.1127	32.03%	55.31%	84.74%	95.81%	<b>98.21%</b>	2,239.3	5,379.6	44.10%	32.44%	22.91%	12.23%	7.20%	16.63	63.09%	75.70%	80.20%				
	$L_{ent}$	0.5324	0.3320	<b>0.1301</b>	<b>0.1070</b>	33.60%	<b>57.50%</b>	<b>85.20%</b>	<b>95.83%</b>	98.07%	<b>2,133.5</b>	<b>5,186.6</b>	<b>45.00%</b>	<b>38.70%</b>	<b>30.85%</b>	<b>24.88%</b>	<b>14.43%</b>	<b>9.28%</b>	17.07	61.99%	75.22%	79.91%			
ResNet <sub>skip</sub>	$L_1$	0.5041	0.2924	0.1259	0.0977	34.10%	57.64%	86.05%	96.13%	98.30%	1,546.2	4,764.0	49.48%	47.23%	43.31%	32.86%	23.57%	16.63%	22.30	44.07%	65.82%	73.55%			
	$L_{cosine}$	0.5024	0.1207	0.1208	0.0958	37.15%	59.61%	87.03%	96.34%	98.35%	1,601.2	4,707.8	52.83%	49.23%	41.05%	32.03%	23.82%	16.75%	15.76	63.32%	77.05%	81.83%			
	$L_{grad}$	<b>0.4754</b>	0.3274	0.1183	0.0905	36.23%	60.44%	87.96%	<b>96.62%</b>	98.45%	1,501.3	4,483.1	56.27%	<b>54.26%</b>	<b>47.88%</b>	33.96%	23.52%	16.07%	18.72	55.00%	71.76%	77.82%			
	$L_{comb}$	0.4788	0.0927	<b>0.1166</b>	<b>0.0893</b>	36.20%	<b>60.64%</b>	<b>87.99%</b>	96.62%	<b>98.49%</b>	1,488.3	4,534.6	<b>57.34%</b>	54.11%	47.57%	33.99%	24.30%	16.37%	<b>15.27</b>	<b>64.18%</b>	<b>77.57%</b>	<b>82.27%</b>			
	$L_{ent}$	0.4923	0.1095	0.1197	0.0941	<b>37.55%</b>	60.43%	87.23%	96.42%	98.46%	<b>1,462.9</b>	<b>4,140.8</b>	54.99%	51.98%	45.40%	<b>35.29%</b>	<b>25.22%</b>	<b>17.68%</b>	15.67	63.28%	77.05%	81.94%			

Table 2: Direct depth performance using spherical and conventional metrics. Bottom part results are the same as those presented in Table 2 of the original document. Top part are the corresponding results from Table 1 of the original manuscript.



Similarly, the differences between ResNet and ResNet<sub>skip</sub>, attributed to the addition of the skip connections are apparent across all samples.

Nonetheless, Pnas better captures the global context as seen in Figure 6, where the scene's dominant planar surfaces are better preserved by it than UNet.

Table 3: Direct depth performance metrics across different variations of DenseNet and Pnas.

				Depth Error ↓				Depth Accuracy ↑
model	pretrained	$\mathcal{L}$		RMSE	RMSLE	AbsRel	SqRel	$\delta_{1.25}$
DenseNet	✗	$\mathcal{L}_1$		0.4672	0.5580	0.1223	0.0896	86.72%
	✓	$\mathcal{L}_1$		<b>0.4072</b>	<b>0.3194</b>	<b>0.1140</b>	<b>0.0694</b>	<b>88.91%</b>
	✓	$\mathcal{L}_{log}$		0.5597	0.5720	0.1528	0.4475	80.48%
	✓	$\mathcal{L}_{berHu}$		0.4532	0.3754	0.1228	0.0852	86.68%
Pnas	✗	$\mathcal{L}_1$		0.4817	0.0780	0.1213	0.0933	87.25%
	✓	$\mathcal{L}_1$		<b>0.3998</b>	<b>0.0634</b>	<b>0.0975</b>	<b>0.0661</b>	<b>91.91%</b>
	✓	$\mathcal{L}_{log}$		0.4135	0.0656	0.0999	0.0697	91.09%
	✓	$\mathcal{L}_{berHu}$		0.4059	0.0643	0.0992	0.0666	91.56%

Figures 7, 8, 9 demonstrate qualitative results in GV2 *tiny* split for the UNet, Pnas, and ResNet<sub>skip</sub> architectures respectively. Apart from the predicted point cloud we visualise the *c2c* error on the ground truth point cloud, with a blue-green-red colormap denoting the error’s magnitude.

Finally, Figures 10 and 11 offer qualitative results of our best performing method in real world, in-the-wild, data captures. We also qualitatively compare our predictions with a state-of-the-art 360° depth estimation model (i.e. BiFuse [8]). It is worth highlighting that even the two of the three 360° images are captured by a panorama camera, the last two images are captured by a smartphone camera, and as such there are artifacts. Yet, it seems that this does not greatly affect the performance of models. The UNet model produces higher quality depth estimates than BiFuse, albeit trained only on the train split of M3D, while the publicly available BiFuse model, as reported in UniFuse [3], has been trained on the *entire* M3D dataset.

## References

- [1] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, (6):679–698, 1986. 5
- [2] Marcela Carvalho, Bertrand Le Saux, Pauline Trouve-Peloux, Andres Almansa, and Frederic Champagnat. On regression losses for deep depth estimation. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 2915–2919. IEEE, 2018. 1
- [3] Hualie Jiang, Zhe Sheng, Siyu Zhu, Zilong Dong, and Rui Huang. Unifuse: Unidirectional fusion for 360 panorama depth estimation. *arXiv preprint arXiv:2102.03550*, 2021. 3
- [4] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 6
- [5] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth*

*international conference on 3D vision (3DV)*, pages 239–248. IEEE, 2016. 1

- [6] Rene Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 1
- [7] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Hohonet: 360 indoor holistic understanding with latent horizontal features. *arXiv preprint arXiv:2011.11498*, 2020. 1
- [8] Fu-En Wang, Yu-Hsuan Yeh, Min Sun, Wei-Chen Chiu, and Yi-Hsuan Tsai. Bifuse: Monocular 360 depth estimation via bi-projection fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 462–471, 2020. 3, 11

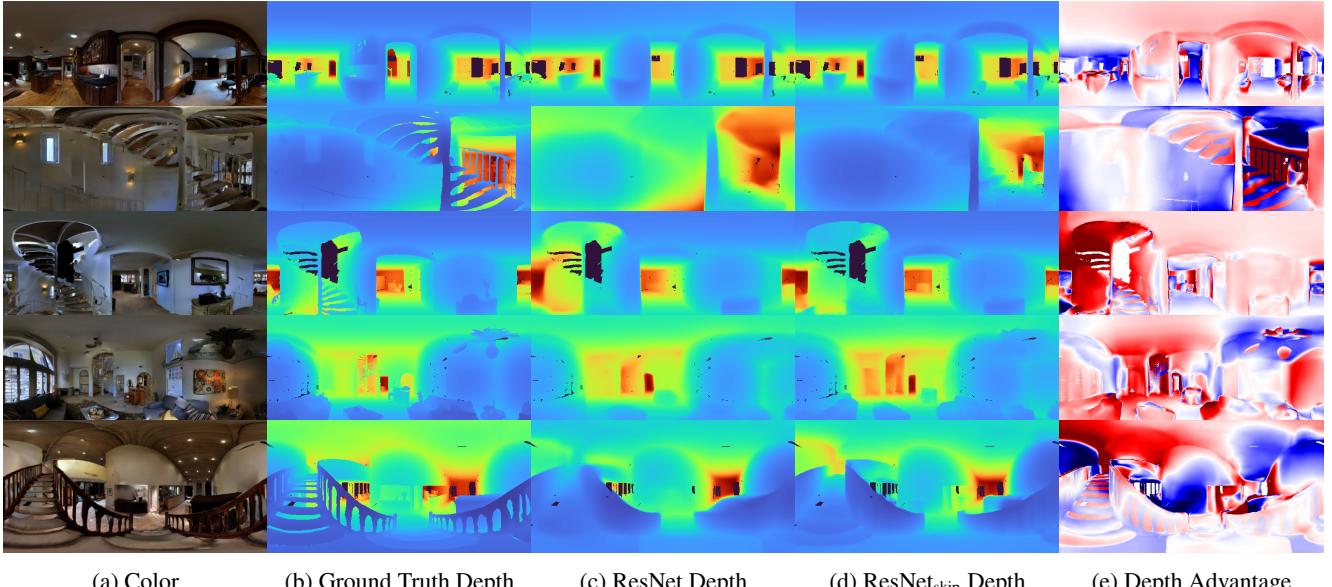


Figure 2: Qualitative comparison between the ResNet and ResNet<sub>skip</sub> architectures. On the right the advantage visualization shows with **blue** color the areas where the former performs better, and with **red** color the areas where the latter performs better. The color magnitude corresponds to the MAE difference between the two models, illustrating the performance deviation between the two models. The addition of skip connections allows ResNet<sub>skip</sub> to capture finer-grained details.

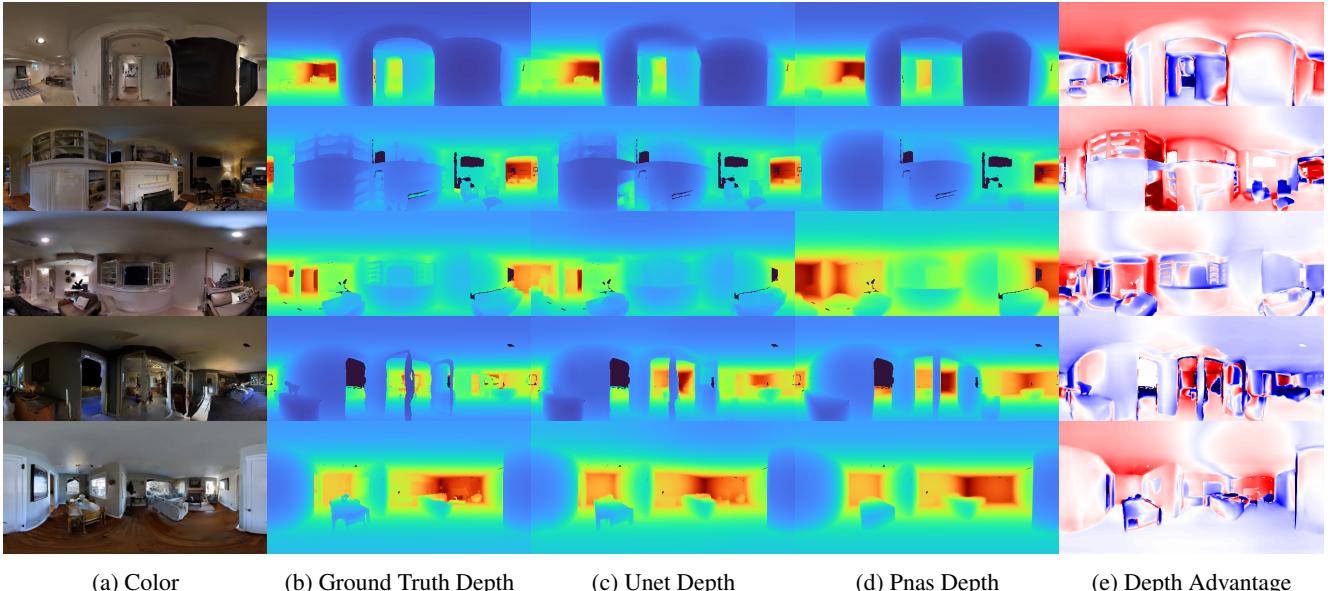


Figure 3: Qualitative comparison between the UNet and Pnas architectures. On the right the advantage visualization shows with **blue** color the areas where the former performs better, and with **red** color the areas where the latter performs better. The color magnitude corresponds to the MAE difference between the two models, illustrating the performance deviation between the two models. Pnas provides smoother results while it is clear that UNet is able to capture finer-grained details.

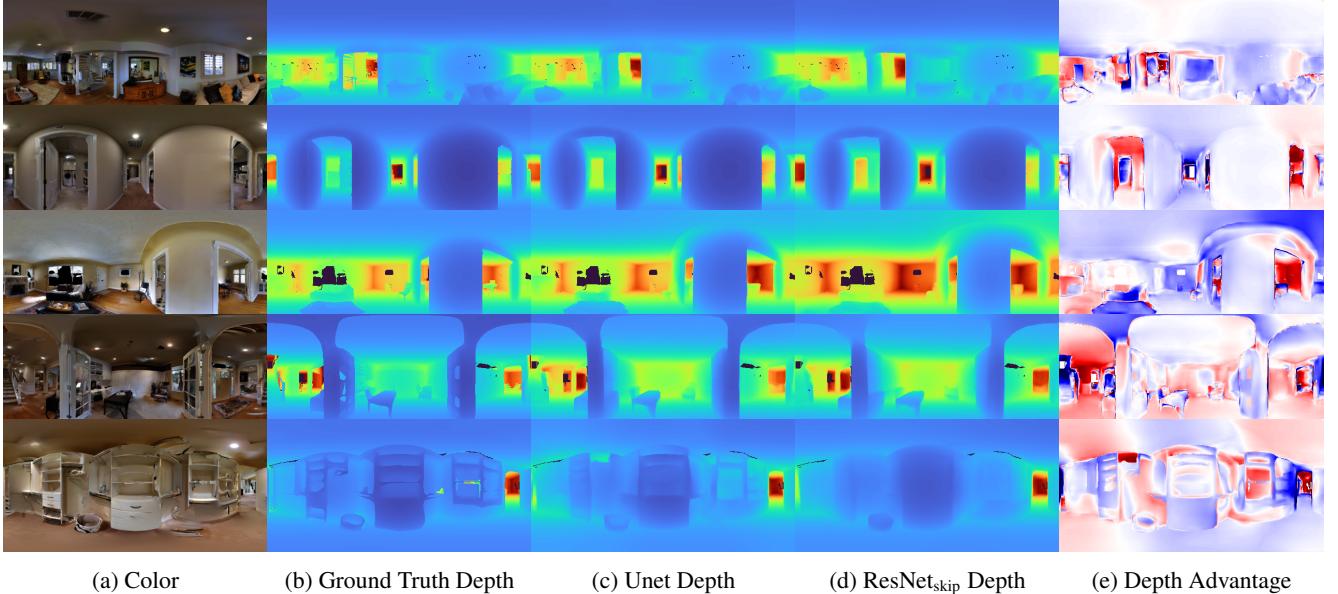


Figure 4: Qualitative comparison between the UNet and ResNet<sub>skip</sub> architectures. On the right the advantage visualization shows with **blue** color the areas where the former performs better, and with **red** color the areas where the latter performs better. The color magnitude corresponds to the MAE difference between the two models, illustrating the performance deviation between the two models.

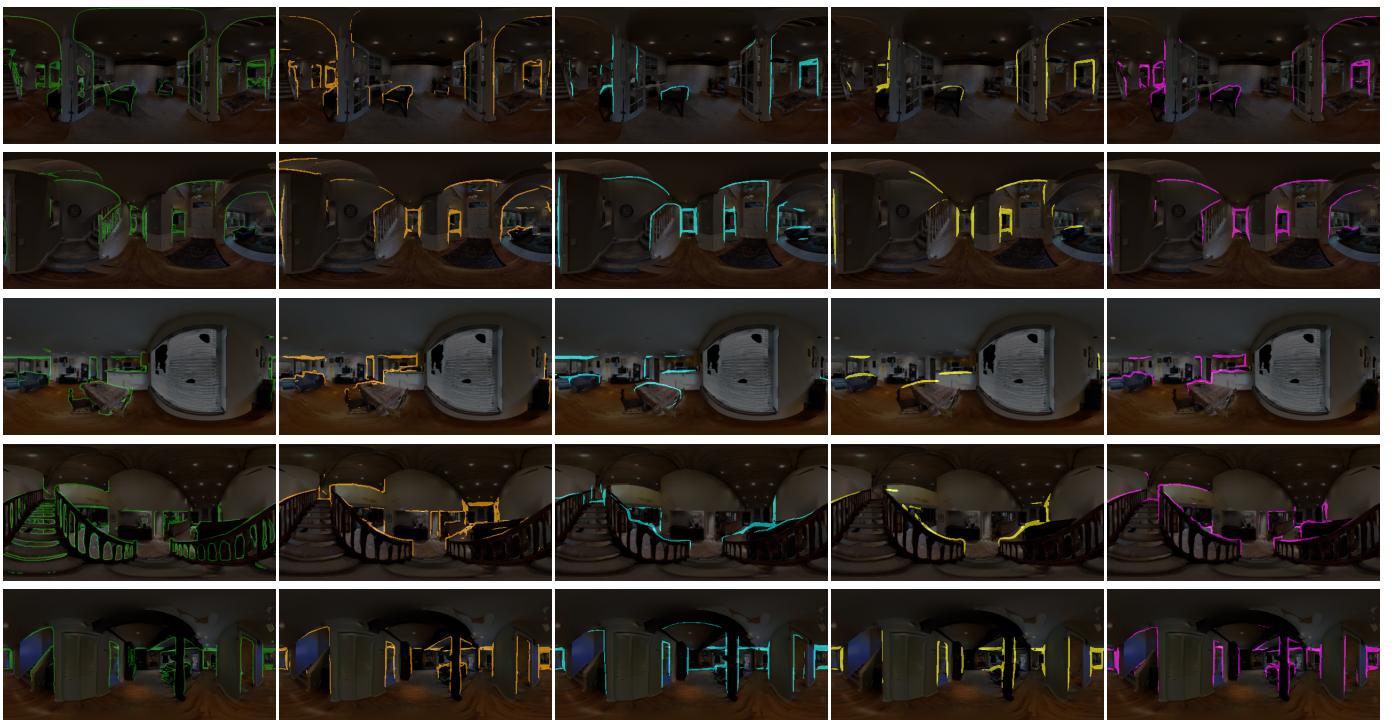


Figure 5: Boundary preservation qualitative comparison between the UNet, Pnas, ResNet and ResNet<sub>skip</sub> models. Boundaries are extracted by applying a Canny edge detector [1] with predefined thresholds on normalized predicted depth maps, and then are blended with the original color panorama. From left to right: **i**) GT depth (**green**), **ii**) UNet (**orange**), **iii**) Pnas (**cyan**), **iv**) ResNet (**yellow**), and **v**) ResNet<sub>skip</sub> (**magenta**).

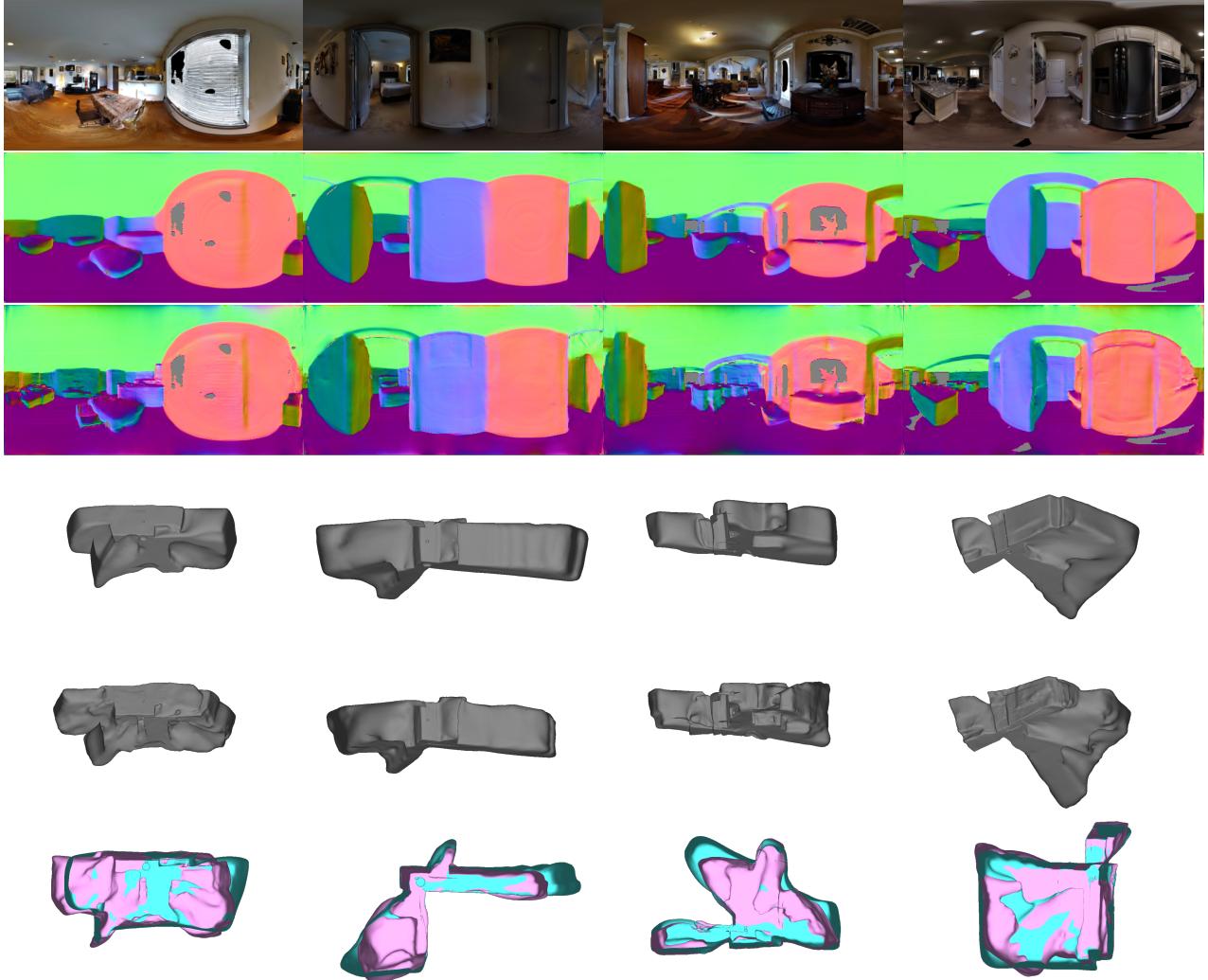


Figure 6: Qualitative comparison of the Pnas and UNet models in surface reconstruction. From top to bottom: **i**) input color panorama, **ii**) Pnas normal map from the estimated depth map, **iii**) UNet normal map, **iv**) Pnas Screened Poisson Surface Reconstruction [4] 3D surface reconstruction, **v**) UNet 3D surface reconstruction, **vi**) overlaid Pnas (cyan) and UNet (pink) 3D surface reconstructions from birds eye view.

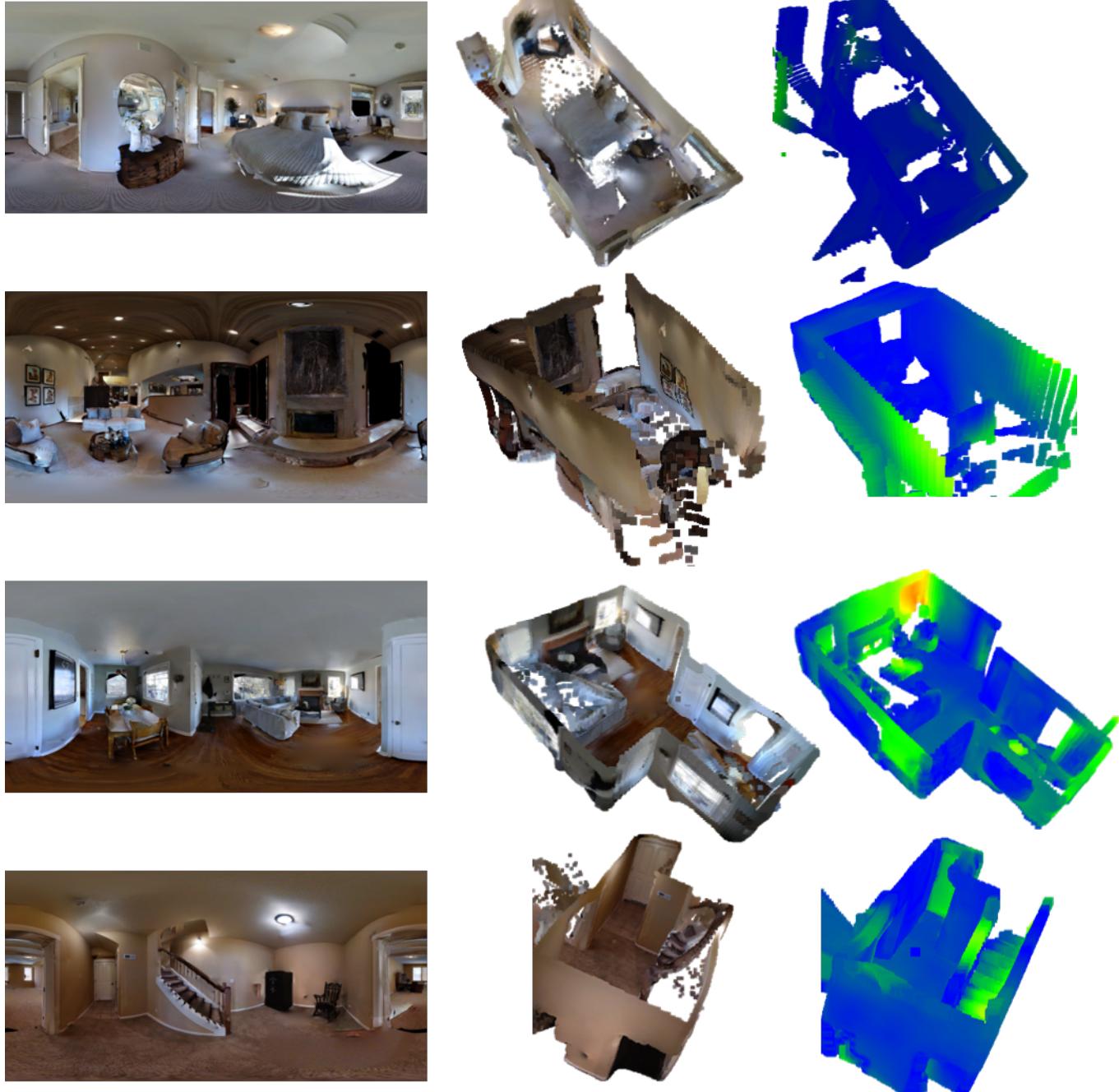


Figure 7: UNet qualitative results. From left to right: **i)** Input color panorama, **ii)** colored predicted point cloud, and **iii)** heatmap visualization of the  $c2c$  error on the ground truth point cloud.



Figure 8: Pnas qualitative results. From left to right: **i**) Input color panorama, **ii**) colored predicted point cloud, and **iii**) heatmap visualization of the  $c2c$  error on the ground truth point cloud.



Figure 9: ResNet<sub>skip</sub> qualitative results. From left to right: **i)** Input color panorama, **ii)** colored predicted point cloud, and **iii)** heatmap visualization of the  $c2c$  error on the ground truth point cloud.

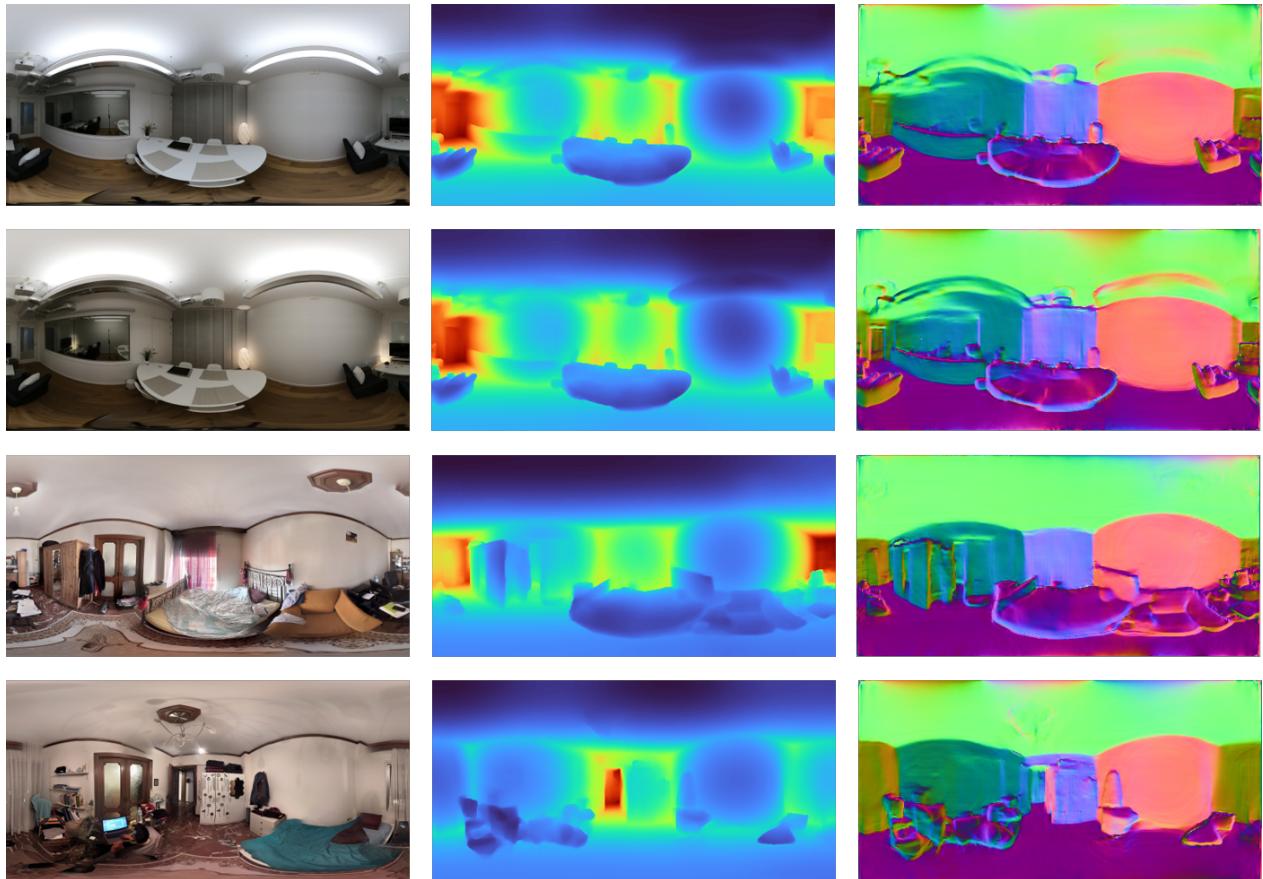


Figure 10: Qualitative results using the UNet model applied to in-the-wild real data captures. The top two rows are captures with a  $360^\circ$  camera, while the bottom two rows are stitched panoramas from a mobile phone. From left to right: **i**) Input color panorama, **ii**) predicted depth, and **iii**) normals derived from the predicted depth.

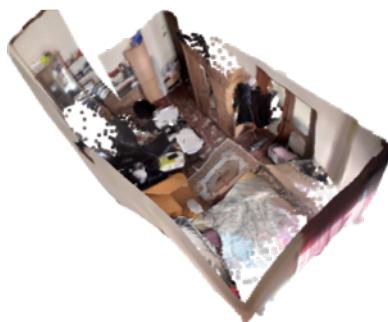
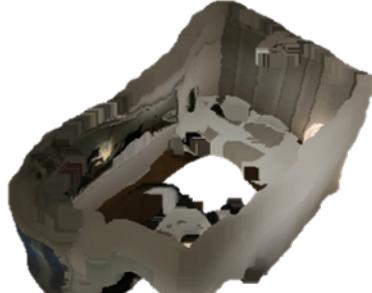


Figure 11: Qualitative results using in-the-wild data. On the left the input color panoramas are depicted. The two top rows are captured with a  $360^\circ$  camera, while the bottom two rows are stitched panoramas from a mobile phone. The colored point clouds of the predicted depths from our UNet model (middle) and BiFuse [8] (right). Ceilings have been removed for visualization purposes.