

Supplementary Materials for Seeing and Seeing Through the Glass: Real and Synthetic Data for Multi-Layer Depth Estimation

Anonymous ICCV submission

Paper ID 2616

001 1. Additional Details

002 1.1. Benchmark

003 All images in our benchmark are released under the CC0
 004 license. We will make the dataset publicly available; how-
 005 ever, a subset of the ground truth annotations for the real-
 006 world images will be withheld for use on a public evalua-
 007 tion server. To validate our approach, we manually anno-
 008 tated 30 synthetic images with known depth ground truth.
 009 Our annotations matched the ground truth in 98% of cases,
 010 demonstrating the reliability. In total, we annotated 1,500
 011 images, with 300 allocated for validation and 1,200 for test-
 012 ing. Our annotations include 5,406 monotonic depth lines
 013 and 38,392 relative depth points across 7 distinct layers,
 014 plus 3,011 fake depth lines and 6,025 fake relative depth
 015 points.

016 1.2. Baseline Training

017 We train all NeWCRFs baseline models from scratch on our
 018 synthetic dataset for 100 epochs. For the Multi-head strat-
 019 egy, we set the number of output depth channels from 1 to
 020 4 to predict multi-layer depth. For the Index Concat strat-
 021 egy, we concatenate the image with an index channel of size
 022 $H \times W$, where all values are set to the corresponding layer
 023 index, and set the number of input channels to 4 (RGB plus

Index). For the Recurrent strategy, we concatenate the im-
 age with the previous depth output and also set the number
 of input channels to 4 (RGB plus Depth). During each train-
 ing step, a random layer is selected as the prediction target.
 To provide richer supervision, we utilize snapped layered
 depth: if a pixel lacks ground-truth depth at layer i , it inher-
 its the depth value from layer $i - 1$. For optimization, we
 use the Scale-Invariant Logarithmic loss [1].

024 1.3. Evaluation and Fine-tuning

025 All methods are evaluated using a single NVIDIA RTX
 026 3090 GPU. When assessed on the synthetic validation
 027 dataset, both the ground-truth values and predictions are
 028 clipped to the range (0.001, 30). Fine-tuning for Metric3D
 029 V2 [2] is performed using the publicly available code, with
 030 training for 100,000 steps.

031 2. Additional Results

032 2.1. Fine-tuning DepthAnything V2

033 The proposed multi-layer depth baseline method can be ap-
 034 plied to any existing depth backbone. To demonstrate how
 035 recent models can benefit from our dataset, we further fine-
 036 tuned DepthAnything V2 [3] for multi-layer depth estima-
 037 tion using the multi-head strategy. To achieve this, we repli-
 038

Method	All			Mixed			Layer 1			Layer 3			Layer 5			Layer 7		
	P	T	Q	P	T	Q	P	T	Q	P	T	Q	P	T	Q	P	T	Q
Multi-head (NeWCRFs)	63.42	42.55	25.97	74.72	46.13	26.38	65.66	39.93	24.21	<u>56.58</u>	<u>43.36</u>	<u>30.77</u>	52.17	38.66	29.94	38.40	35.65	33.10
Index Concat (NeWCRFs)	64.46	44.00	26.00	76.70	48.37	26.46	66.95	41.84	24.45	55.85	41.53	30.06	55.36	42.69	31.57	39.66	45.18	45.08
Recurrent (NeWCRFs)	62.36	41.88	24.64	73.51	45.27	25.35	68.08	44.46	26.12	49.47	31.53	21.25	45.51	30.06	21.29	30.09	34.10	27.59
Multi-head (DA v2)	80.72	70.67	61.68	89.37	72.98	62.23	85.80	74.34	65.62	<u>73.75</u>	<u>62.68</u>	<u>59.19</u>	<u>66.19</u>	<u>59.29</u>	<u>54.75</u>	57.39	50.78	48.09

Table 1. Baseline methods evaluated on our real-world benchmark via tuple-wise accuracy. Best scores are in **bold**. Second best underlined.

Method	Layer 1				Layer 3				Layer 5				Layer 7			
	AbsRel↓	RMS↓	$\delta 1\uparrow$	$\delta 2\uparrow$	AbsRel↓	RMS↓	$\delta 1\uparrow$	$\delta 2\uparrow$	AbsRel↓	RMS↓	$\delta 1\uparrow$	$\delta 2\uparrow$	AbsRel↓	RMS↓	$\delta 1\uparrow$	$\delta 2\uparrow$
Multi-head (NeWCRFs)	17.97	<u>23.41</u>	<u>83.08</u>	93.14	<u>15.90</u>	<u>47.39</u>	<u>80.15</u>	<u>93.89</u>	<u>14.58</u>	<u>47.51</u>	<u>81.48</u>	<u>94.12</u>	<u>16.27</u>	<u>58.78</u>	<u>80.44</u>	<u>93.92</u>
Index Concat (NeWCRFs)	17.26	23.70	83.02	93.27	16.25	48.19	79.63	93.72	15.03	48.11	80.61	93.97	16.49	59.03	80.26	93.76
Recurrent (NeWCRFs)	17.23	25.53	81.89	<u>93.37</u>	17.61	51.79	76.39	92.53	17.25	52.53	76.54	92.20	18.35	61.98	77.04	92.40
Multi-head (DA v2)	8.24	15.21	93.43	97.73	10.68	37.97	88.76	96.83	11.15	40.67	87.51	96.40	13.84	54.63	84.79	95.13

Table 2. Multi-layer baseline methods evaluated on our synthetic validation set. Values are scaled by 100 for clearer comparison. Best scores are in **bold**. Second best underlined.

046 cate the depth head four times, enabling the model to predict
047 multi-layer depth.

048 Results evaluated on the real world benchmark are
049 shown in Tab. 1. Results evaluated on the synthetic vali-
050 dation set are shown in Tab. 2. Qualitative comparison on
051 real benchmark can be found in Fig. 1 and Fig. 2. Qualita-
052 tive comparison on synthetic validation set can be found in
053 Fig. 3. These results show that DepthAnything V2 not only
054 aligns more closely with the training domain but also gener-
055 alizes exceptionally well after fine-tuning solely on our syn-
056 thetic data. This underscores the power of our synthetic data
057 generator for multi-layer depth estimation and transparent-
058 object understanding. Nevertheless, some regions around
059 object boundaries and in deeper layers still exhibit artifacts,
060 leaving rooms for improved model designs in future work.

061 2.2. Additional Qualitative Results

062 To help readers better visually assess the effectiveness of
063 our datasets, we provide additional qualitative results on
064 real world benchmark and synthetic validation set, which
065 can be found in Fig. 4 and Fig. 5, respectively.

066 References

- [1] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 1
067
[2] Mu Hu, Wei Yin, Chi Zhang, Zhipeng Cai, Xiaoxiao Long, Hao Chen, Kaixuan Wang, Gang Yu, Chunhua Shen, and Shaojie Shen. Metric3d v2: A versatile monocular geometric foundation model for zero-shot metric depth and surface normal estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1
068
[3] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2025. 1
069
070
071
072
073
074
075
076
077
078
079
080
081

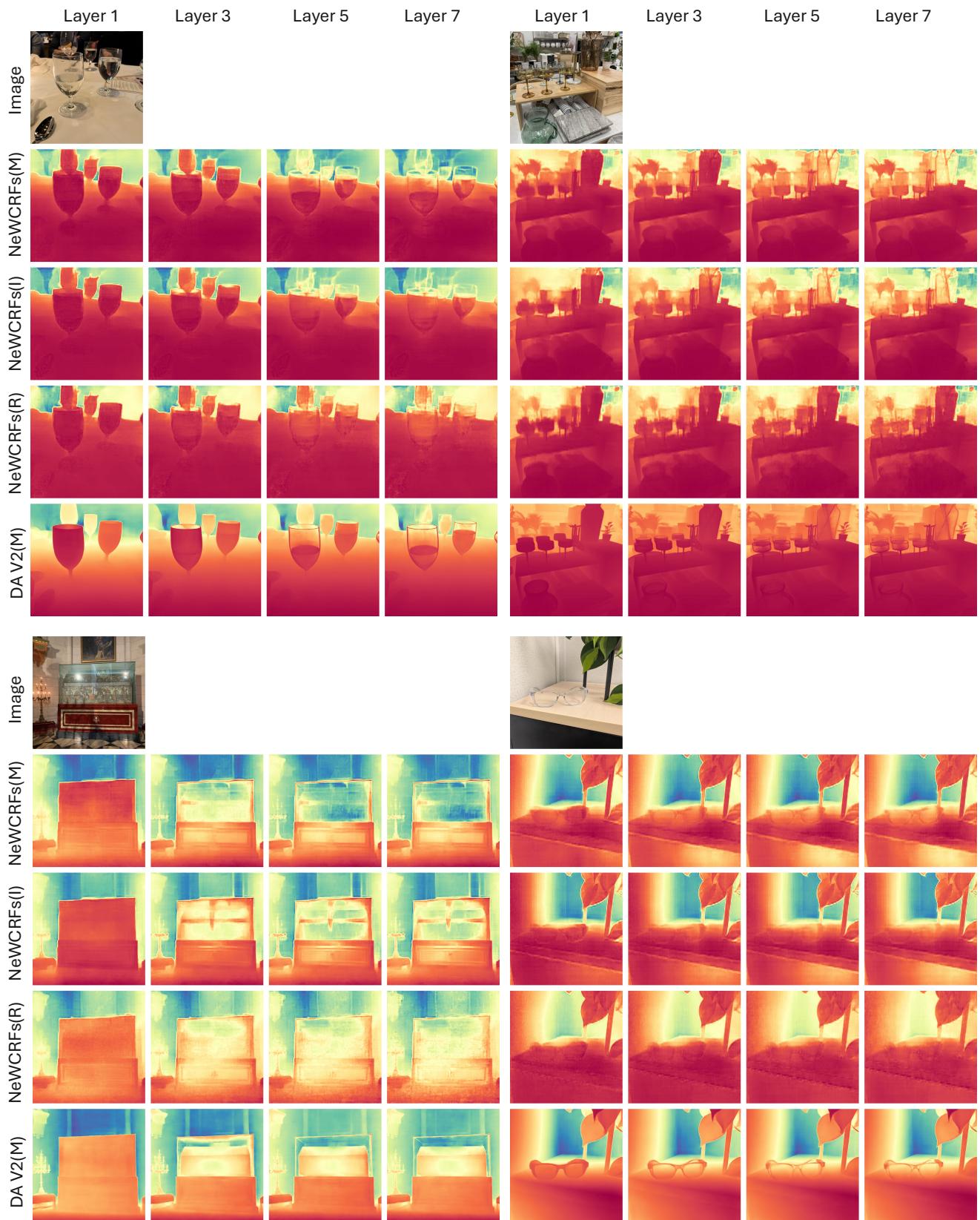


Figure 1. Additional qualitative comparison of multi-layer depth baselines on real benchmark.

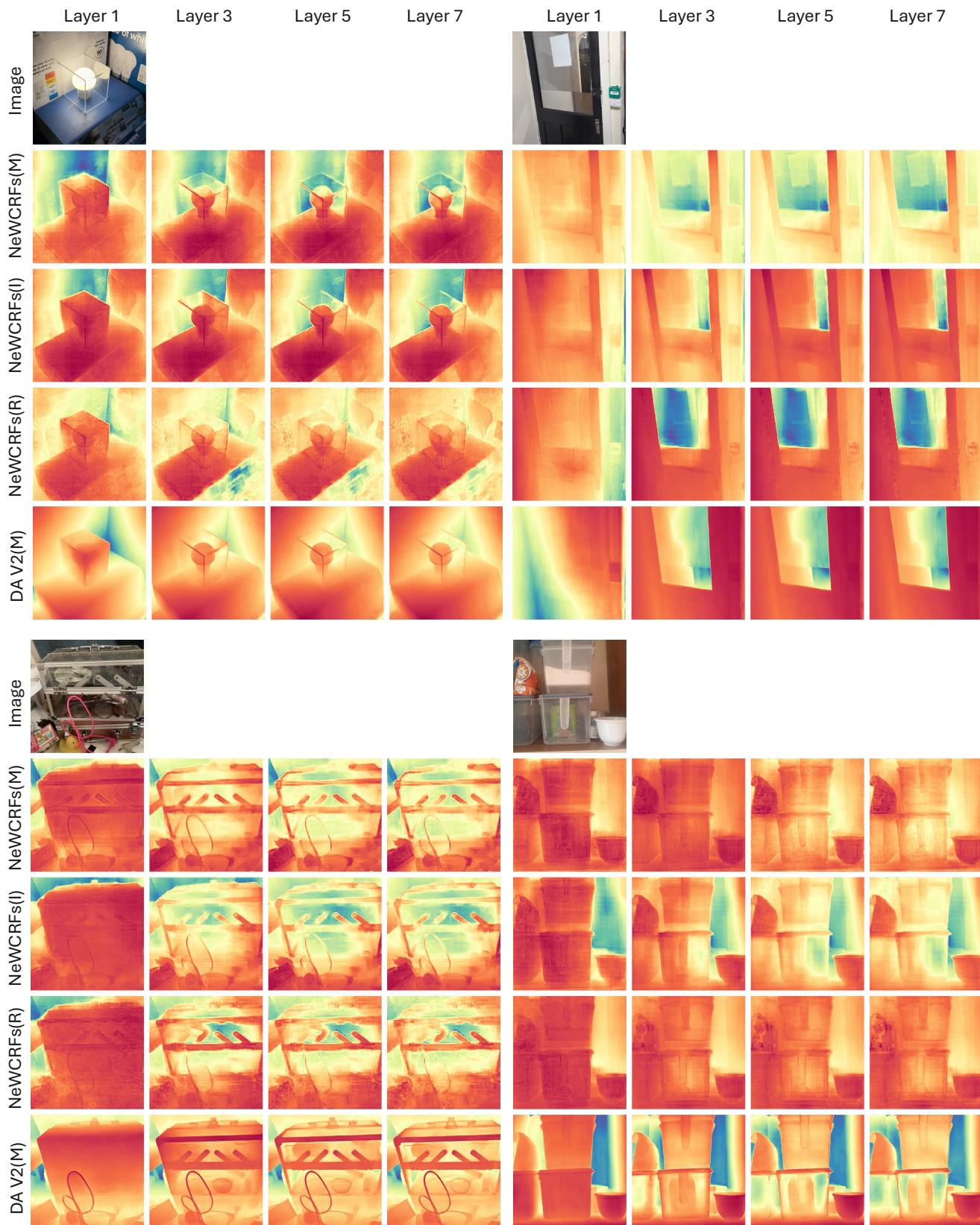


Figure 2. Additional qualitative comparison of multi-layer depth baselines on real benchmark.

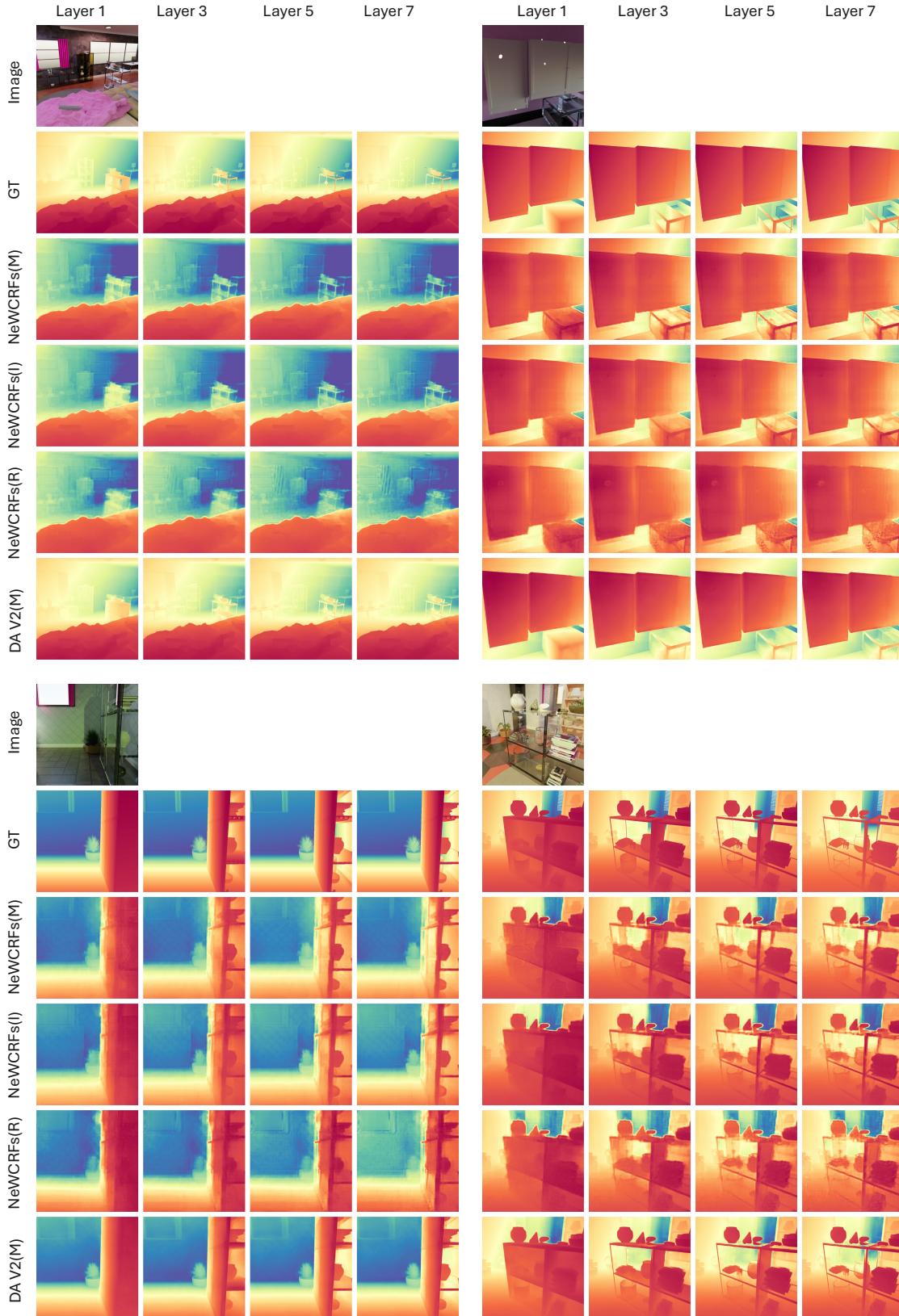


Figure 3. Additional qualitative comparison of multi-layer depth baselines on synthetic validation set.

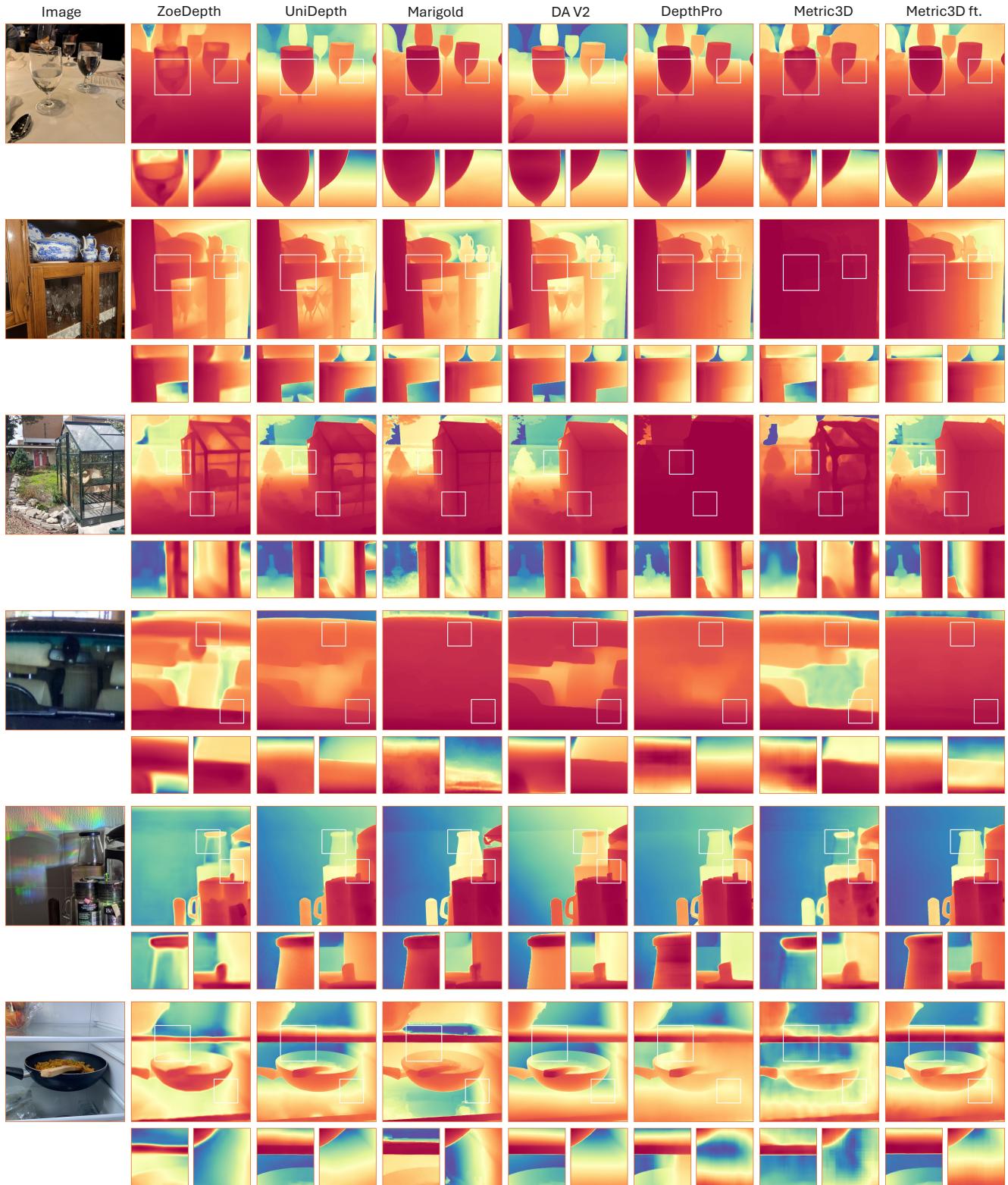


Figure 4. Additional qualitative comparison of state-of-the-art single layer depth methods on our benchmark.

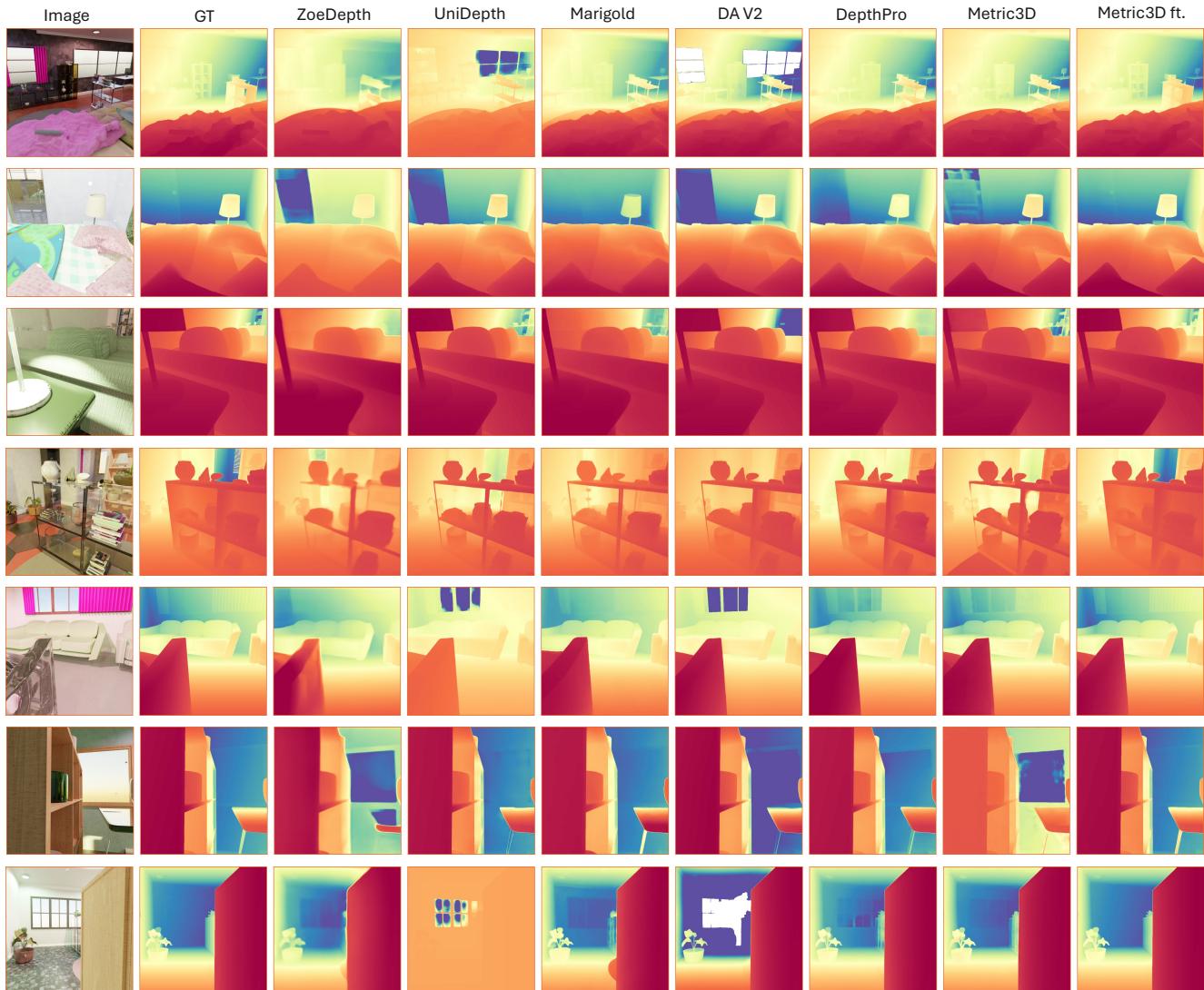


Figure 5. Additional qualitative comparison of state-of-the-art single layer depth methods on our synthetic validation set.