

Supplemental Material

Reproducibility Statement

All results presented in this work are reproducible using the code associated with our research, available at <https://jugit.fz-juelich.de/ias-8/mdeaux>.

Ethics Statement

Our research relies on publicly available datasets, ensuring transparency and reproducibility in our experiments. Additionally, for datasets obtained through agreements, such as Matterport, we adhere to the respective terms and conditions outlined in the agreements [44].

Environmental Impact

Recognizing the environmental impact of computational resources, we are mindful of the compute resources used in our experiments. The experiments, conducted on the aforementioned hardware setup, resulted in an estimated environmental impact of approximately 29.25 million mWh, equivalent to 0.0117 metric tons of carbon dioxide. This is comparable to the emissions from driving 5.1 miles in the average gasoline-powered passenger vehicle in the US [2]. For a comprehensive overview of the environmental impact of compute [49], we refer to metrics provided by platforms like the ML CO2 Impact calculator [28].

A. Training Details

As mentioned in Section 4 of the main paper, we adopt the DINOv2 training procedure [46]. In particular, we use AdamW with initial learning rate of 0.0001 and weight decay of 0.01, and a cosine scheduler with linear warmup for 1/3 of the iterations. In total we train for 38400 steps with a batch size of 4 (2 images per 2 GPUs). For Taskonomy we use a batch size of 16 (across 8 GPUs), whereas for Matterport we train twice longer. When using our method, we duplicate the optimizer and learning rate schedulers, and scale the DPT decoder learning rates by our parameter α .

For the experiment using the Depth Anything backbone, we adapt the Depth Anything training procedure to our scheme. In particular we use an initial learning rate of 0.000161 and a cosine scheduler without linear warmup. We train for 38400 steps with a batch size of 16 (2 images per 8 GPUs).

B. Hardware Details

In our experiments, we leveraged high-performance computing (HPC) nodes equipped with 4 NVIDIA A100 GPUs with 40GB VRAM and 48 CPU cores. For most of the experiments the training process used only 2 of the GPUs, while for some we used 8 GPUs.

C. Additional ablations

C.1. Single-Label Dense Classification

Motivated by the promising outcomes observed with our MLDC task as an auxiliary task for Monocular Depth Estimation, we simplify the problem to Single-Label Dense Classification to investigate the viability of using straight-forward classification datasets as auxiliaries for MDE. This involves extracting the dominant class for each image segmentation mask and output during training and computing the CrossEntropyLoss based on the dominant classes. As depicted in Table 5, the results for single-label and MLDC are comparable. This suggests that classifying the dominant class could potentially suffice as an auxiliary task for MDE. We encourage further exploration with additional single-label classification datasets, such as ImageNet [16], to validate whether they can contribute further improvements in MDE quality.

C.2. Comparison of Single and Multi Source Auxiliary Tasks

We extend our investigation to ascertain if our approach can also be applied within a single dataset, even though this is not the primary focus of our research. For SUN RGBD we use the original provided semantic segmentation labels, whereas for the other datasets we use pseudo labels generated using the approach used in PolyMax [70] for Taskonomy. Note that the quality of the pseudo-labels is not guaranteed to be high, especially when images contain objects unknown to the chosen pseudo-labeler model. This can lead to noisy labels and potentially degrade the performance. The results reported in Table 6 reveal that our method can also improve the MDE quality when only using a single dataset and pre-processing for both tasks, while the quality gains for each MDE dataset are worse compared to using multiple auxiliary data sources. On one hand, this demonstrates the versatility of our method, showing that it can provide improvements in both single and multiple sources scenarios. On the other hand, these results suggest that using multiple and diverse auxiliary sources should be preferable to ensure higher quality gains.

D. Additional qualitative results

We present additional qualitative results for each dataset to demonstrate the effectiveness of our approach. Each row should be considered individually.

Table 5. AbsRel $\times 10^4$ (\downarrow) scores of the best task w.r.t. the DINOv2 baseline. The last column shows the result of Single-Label Dense Classification, while others are taken from Table 3. The indoor-outdoor dominated MIX6 dataset with $\alpha = 0.9$ is used for all the experiments. The best and second-best results are highlighted in **bold** and *italic*, respectively.

MDE Datasets	In MIX6	Aux Tasks						Gain %
		DINOv2	Classification	Segmentation	Reconstruction	S-Classification		
NYUv2	\times	809 \pm 10	696 \pm 3	755 \pm 6	838 \pm 7	727 \pm 8	13.9	
SUN RGBD	\checkmark	1128 \pm 11	<i>1024 \pm 10</i>	1069 \pm 9	1111 \pm 5	1007 \pm 11	10.3	
Matterport3D	\times	1874 \pm 19	<i>1728 \pm 9</i>	1793 \pm 13	1805 \pm 6	1720 \pm 5	8.2	
Taskonomy	\times	1506 \pm 13	1481 \pm 4	1567 \pm 12	1543 \pm 6	<i>1491 \pm 20</i>	1.7	
DIODE In	\times	3588 \pm 19	3239 \pm 26	3451 \pm 47	3368 \pm 31	<i>3289 \pm 59</i>	9.7	
DIODE Out	\times	5820 \pm 223	4965 \pm 162	5085 \pm 172	4530 \pm 91	<i>4926 \pm 121</i>	22.2	

Table 6. AbsRel $\times 10^4$ (\downarrow) scores on various depth datasets using different auxiliary tasks and percentage gain of the best task w.r.t. the DINOv2 baseline. For every dataset, the same dataset is used as auxiliary with $\alpha = 0.9$, either with original or pseudo labels. The best and second-best results are highlighted in **bold** and *italic*, respectively.

MDE Datasets	Pseudo Labels	Aux Tasks – Same Dataset Only					Gain %
		DINOv2	Classification	Segmentation	Reconstruction		
NYUv2	\checkmark	809 \pm 10	762 \pm 25	838 \pm 32	846 \pm 29	5.8	
SUN RGBD	\times	<i>1128 \pm 11</i>	1105 \pm 32	1161 \pm 7	1167 \pm 12	2	
Matterport3D	\checkmark	1874 \pm 19	1841 \pm 21	<i>1856 \pm 17</i>	1864 \pm 15	1.8	
Taskonomy	\checkmark	1506 \pm 13	<i>1576 \pm 16</i>	1606 \pm 5	1607 \pm 8	-4.6	
DIODE In	\checkmark	3588 \pm 19	3447 \pm 67	<i>3554 \pm 90</i>	3620 \pm 26	3.1	
DIODE Out	\checkmark	5820 \pm 223	5918 \pm 88	6040 \pm 60	5682 \pm 96	2.4	

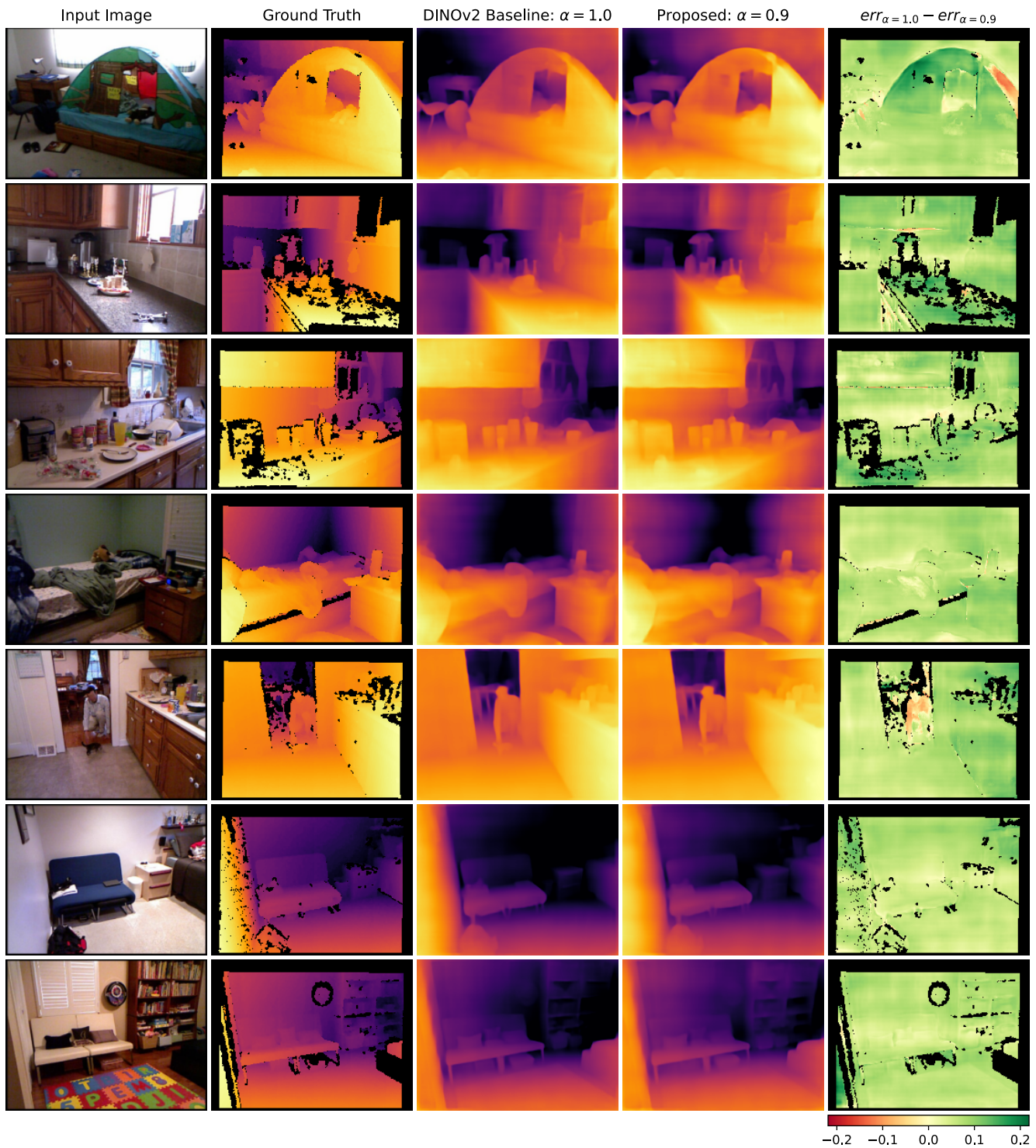


Figure 6. Results on NYU with MIX6 auxiliary MLDC task. From left to right: image and respective ground truth, baseline and our method predictions, and error difference between the last two w.r.t. to the ground truth.

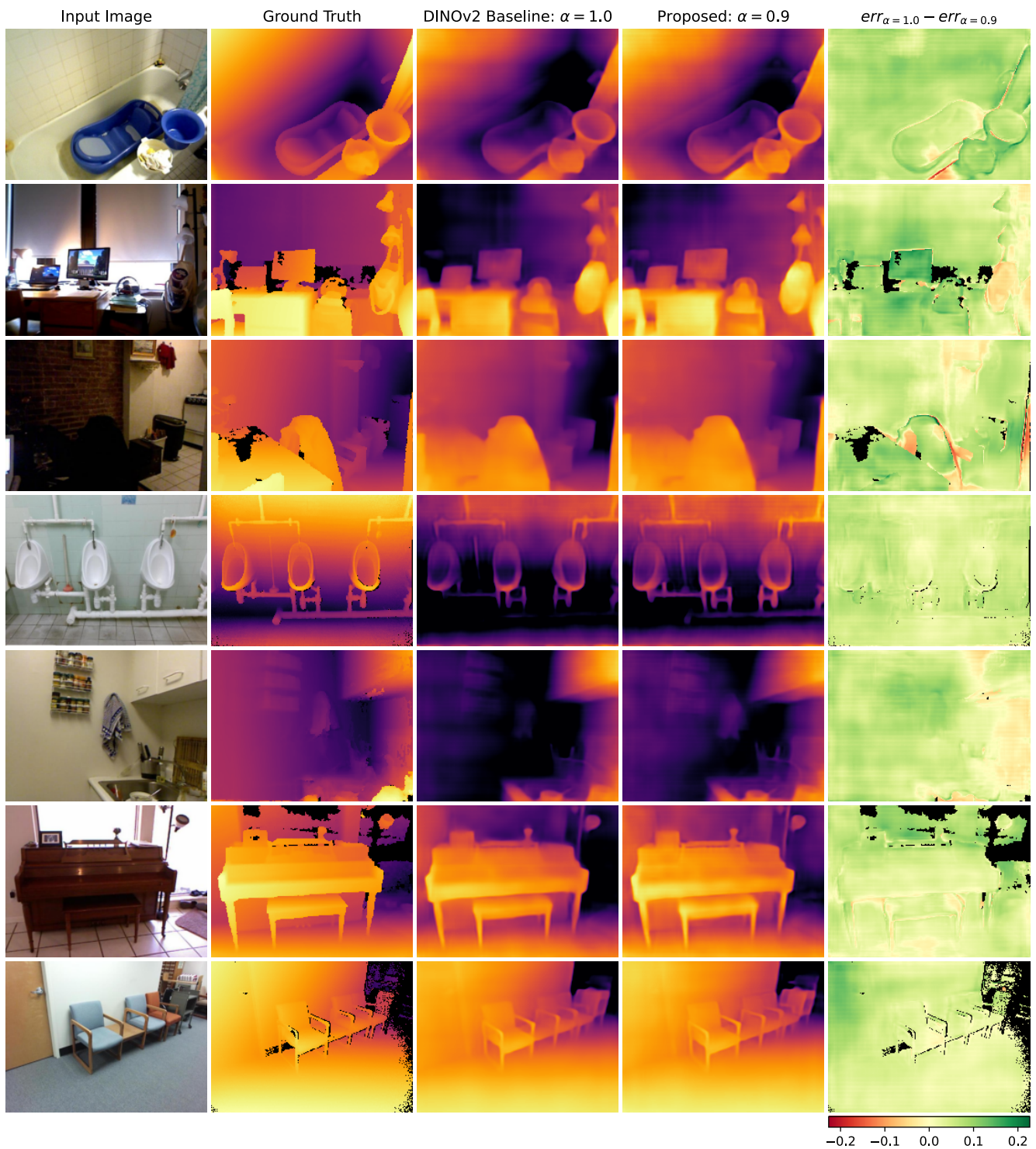


Figure 7. Results on SUNRGBD with MIX6 auxiliary MLDC task. From left to right: image and respective ground truth, baseline and our method predictions, and error difference between the last two w.r.t. to the ground truth.

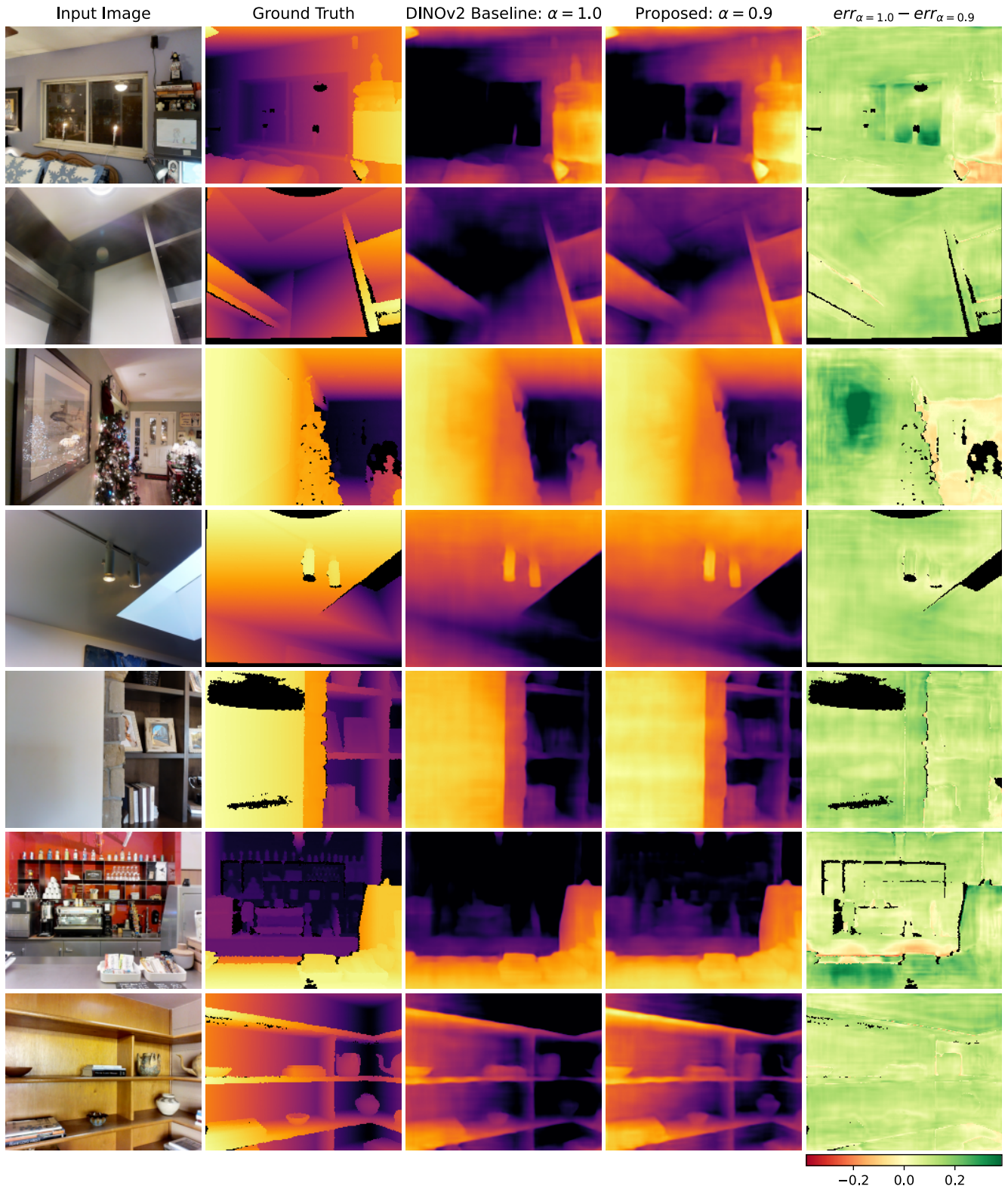


Figure 8. Results on Matterport with MIX6 auxiliary MLDC task. From left to right: input image and respective ground truth, baseline and our method predictions, and error difference between the last two w.r.t. to the ground truth.

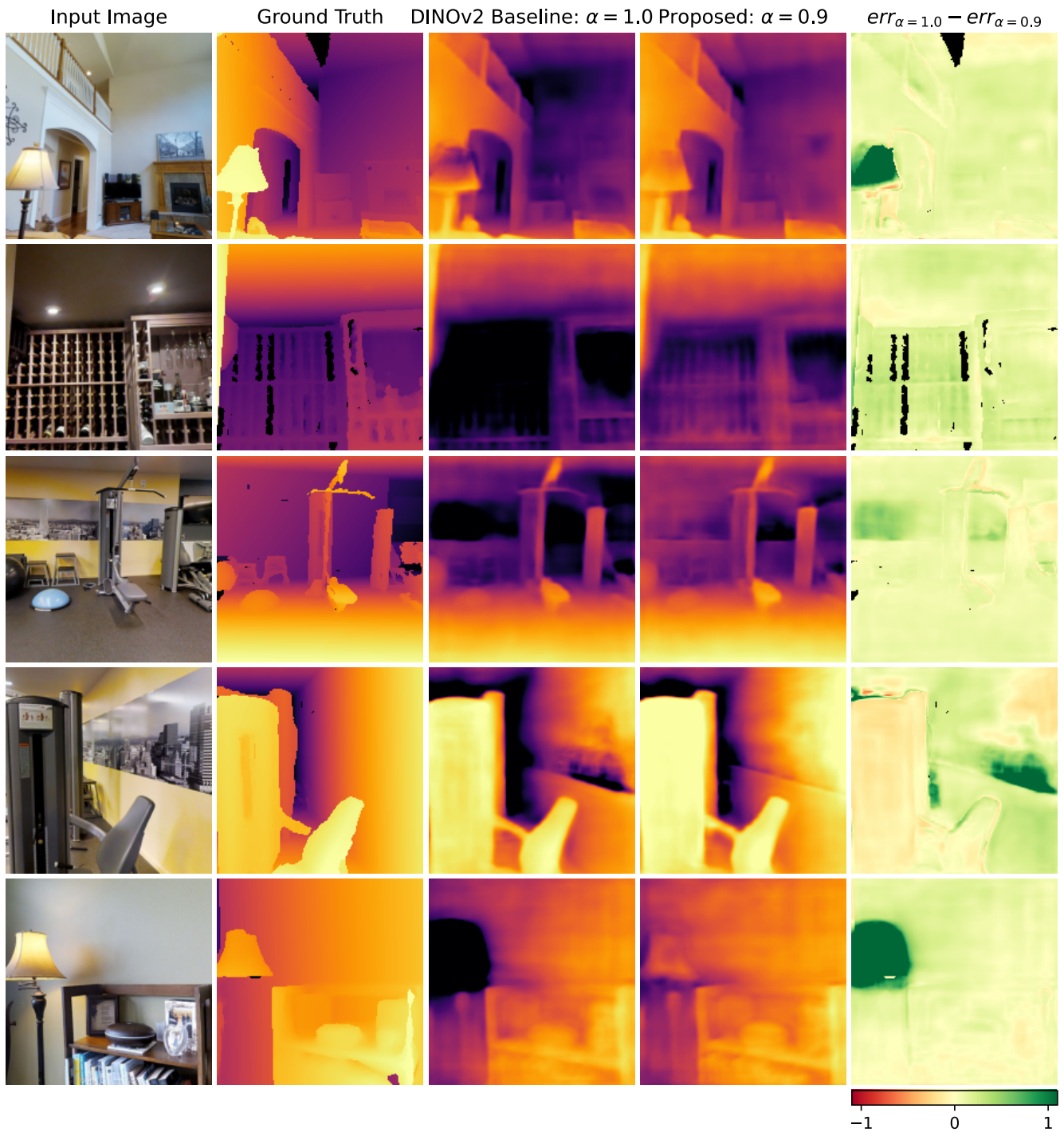


Figure 9. Results on Taskonomy with MIX6 auxiliary MLDC task. From left to right: input image and respective ground truth, baseline and our method predictions, and error difference between the last two w.r.t. to the ground truth.

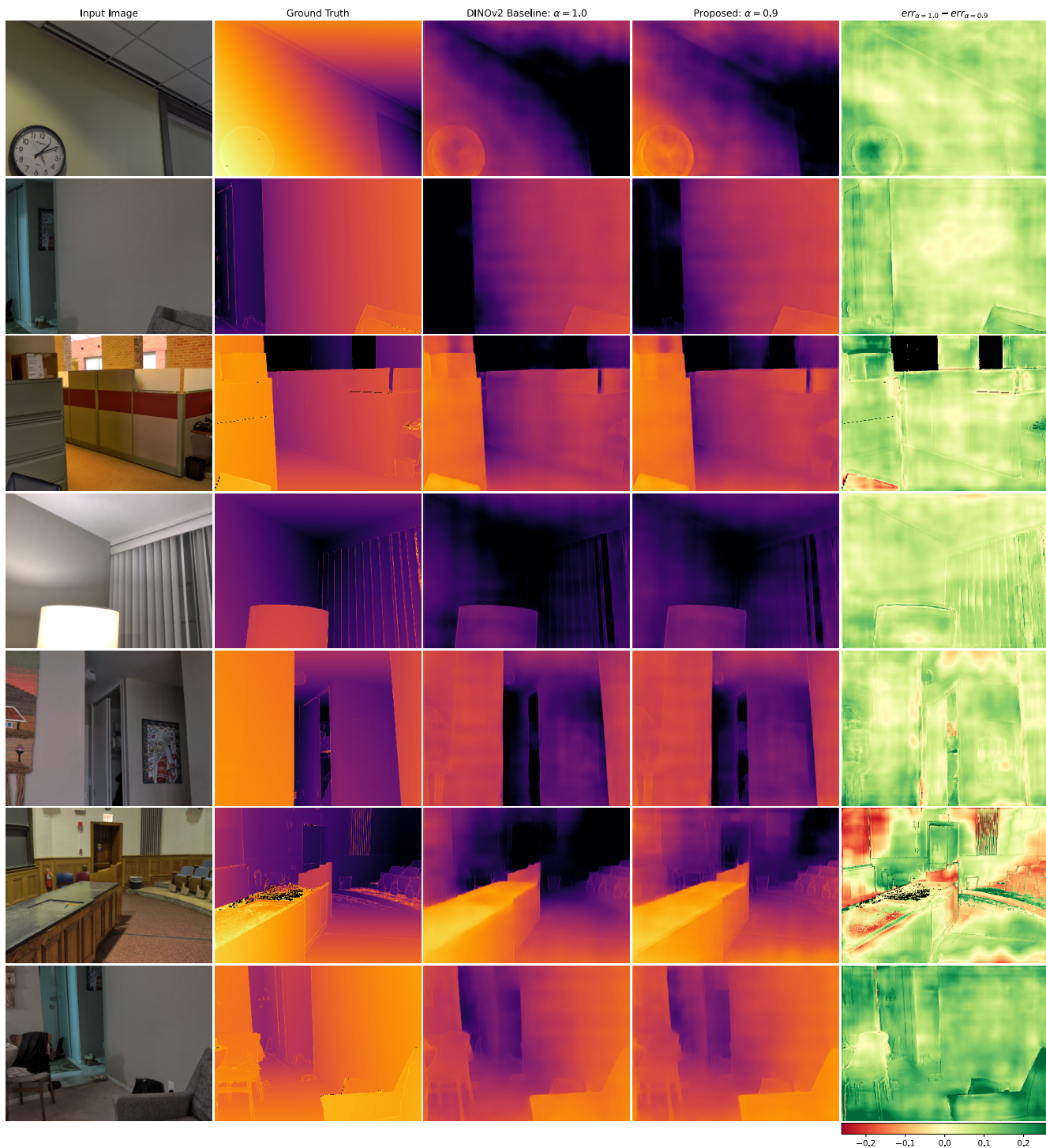


Figure 10. Results on DIODE Indoor with MIX6 auxiliary MLDC task. From left to right: input image and respective ground truth, baseline and our method predictions, and error difference between the last two w.r.t. to the ground truth.

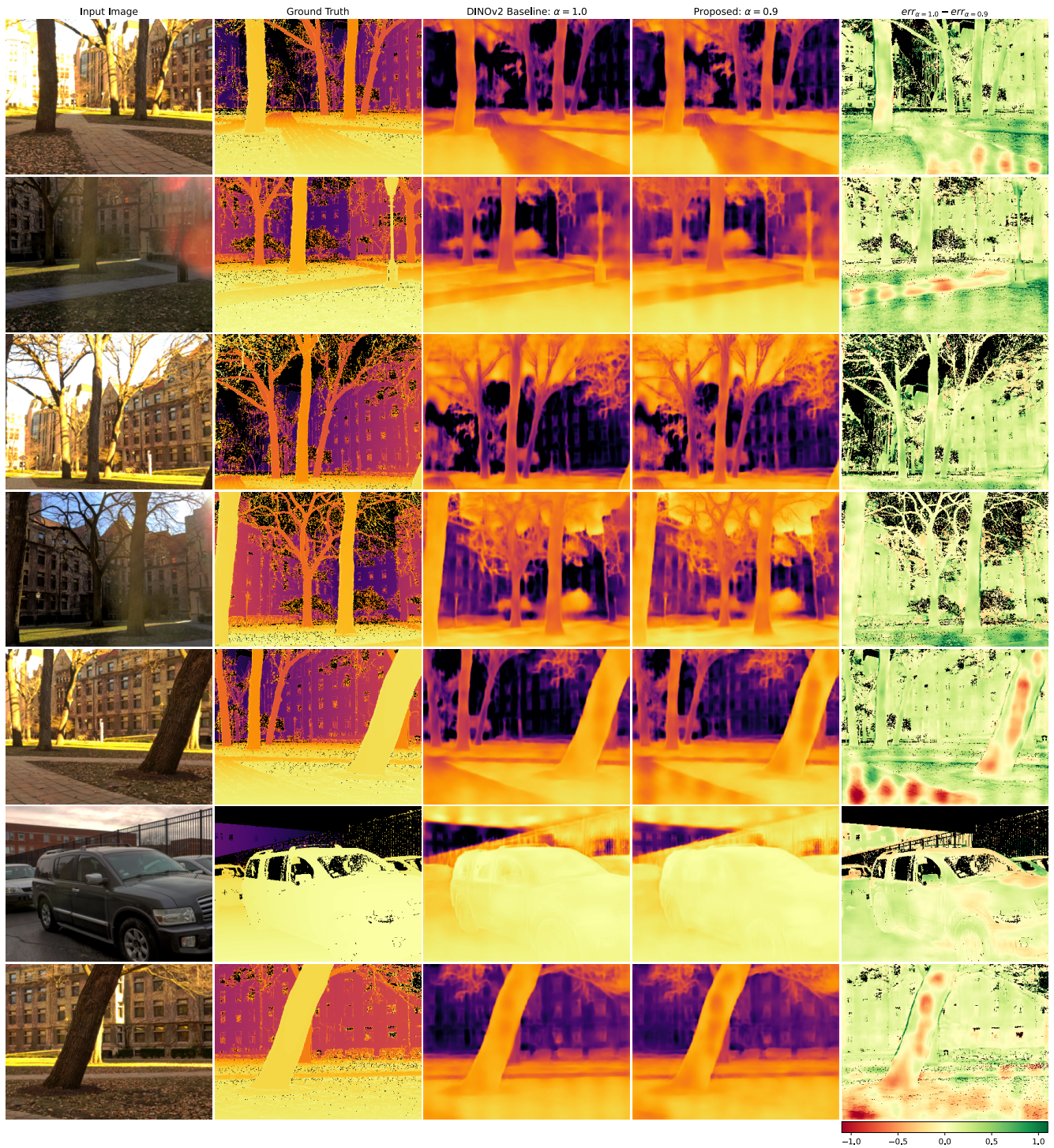


Figure 11. Results on DIODE Outdoor with MIX6 auxiliary MLDC task. Left to right: input image and ground truth; baseline and our method predictions; error difference relative to ground truth highlighting visible improvements in the facade behind the trees.