

# Transformer-based Monocular Depth Estimation with Attention Supervision - Supplementary Material

Wenjie Chang  
changwj@mail.ustc.edu.cn

Yueyi Zhang\*  
zhyueyi@ustc.edu.cn

Zhiwei Xiong  
zwxiong@ustc.edu.cn

Department of Electronic Engineering  
and Information Science,  
University of Science and Technology of  
China, Hefei, China

## 1 DUB Architecture

In this section, we show the difference of Attention-based Up-sample Block (AUB) and Direct Up-sample Block (DUB). Fig. 1(a) and Fig. 1(b) show the architectures of AUB and DUB, respectively. Compared with AUB, DUB directly concatenates feature maps from Image Branching and the outputs from Transformer Layer at the same scale.

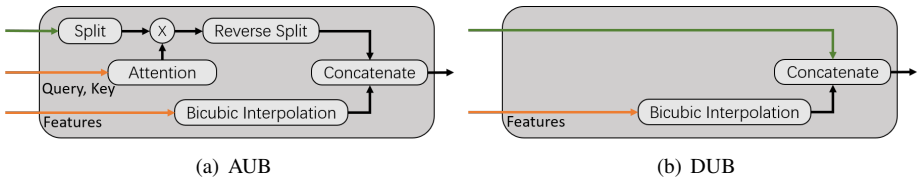


Figure 1: Architecture comparison between AUB and DUB. The orange route denotes elements from Transformer Layer. The green route represents the feature maps from Image Branching.

## 2 Influence of the factor $\lambda$

We calculate attention maps used in AS as follows:

$$\hat{A}_{i,j} = \text{SOFTMAX}(-\lambda |\hat{D} - \hat{d}_{i,j}|), \quad (1)$$

where  $\hat{D}$  represents the depth map which is  $16 \times$  down-sampled from the ground-truth depth map and normalized to  $[0, 1]$ ,  $\hat{d}_{i,j}$  represents the depth value at position  $(i, j)$ ,  $\lambda$  is a hyper-parameter. Fig. 2 shows attention maps with different  $\lambda$  values. It can be seen that with the

increase of the  $\lambda$  value, small regions are be paid attention to. In our experiment, we set  $\lambda$  to 8.

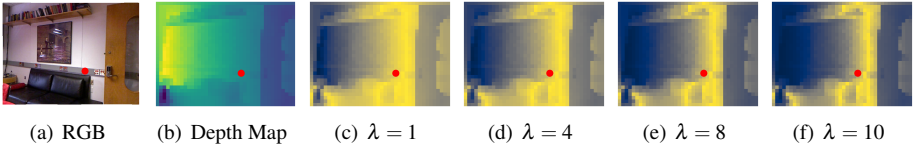


Figure 2: Attention maps (c) (d) (e) (f) shows the correlations between the red point and other pixels.

### 3 Additional Qualitative Results

We present the qualitative comparison results with the state-of-the-art methods in Fig. 3 and Fig. 4. Fig. 3 shows the depth results on the KITTI dataset. Our proposed method is compared with VNL [10] and LapDepth [10]. It can be observed that in the areas with complex lighting situations, such as car windows (Row 1, 3, 5 in Fig. 3) and shadow regions (Row 4, 8 in Fig. 3), our method always provides better depth results. Fig. 4 shows the depth results on the NYU Depth V2 dataset. Comparing with SARPNet [10], LapDepth [10] and our network without Attention Supervision (AS), our method demonstrates more accurate depth results around the planes (Row 4, 9 in Fig. 4) and overcomes the artifacts coming from shadows (Row 5, 6 in Fig. 4). By comparing Fig. 4(e) and Fig. 4(f), it is obvious that after adding the Attention-based Up-Sample Block (AUB), our method predicts depth maps with sharper boundaries.

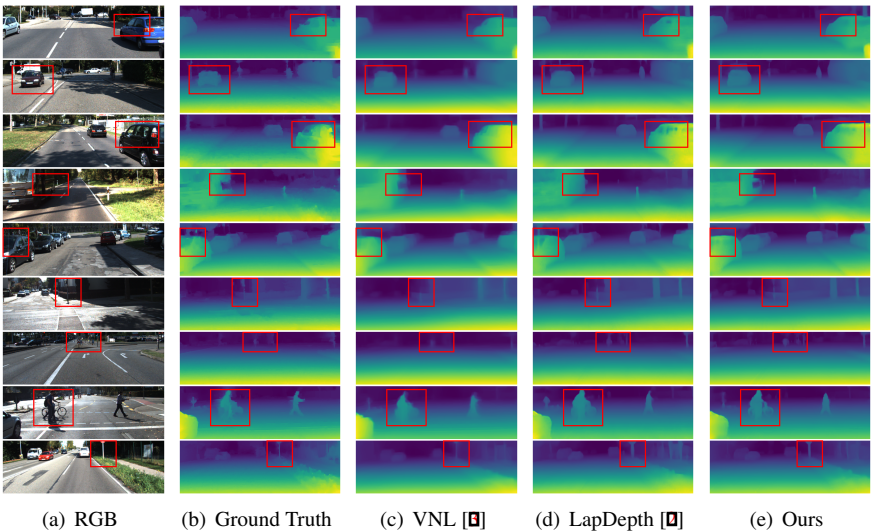


Figure 3: Qualitative comparisons with other methods on the KITTI dataset.

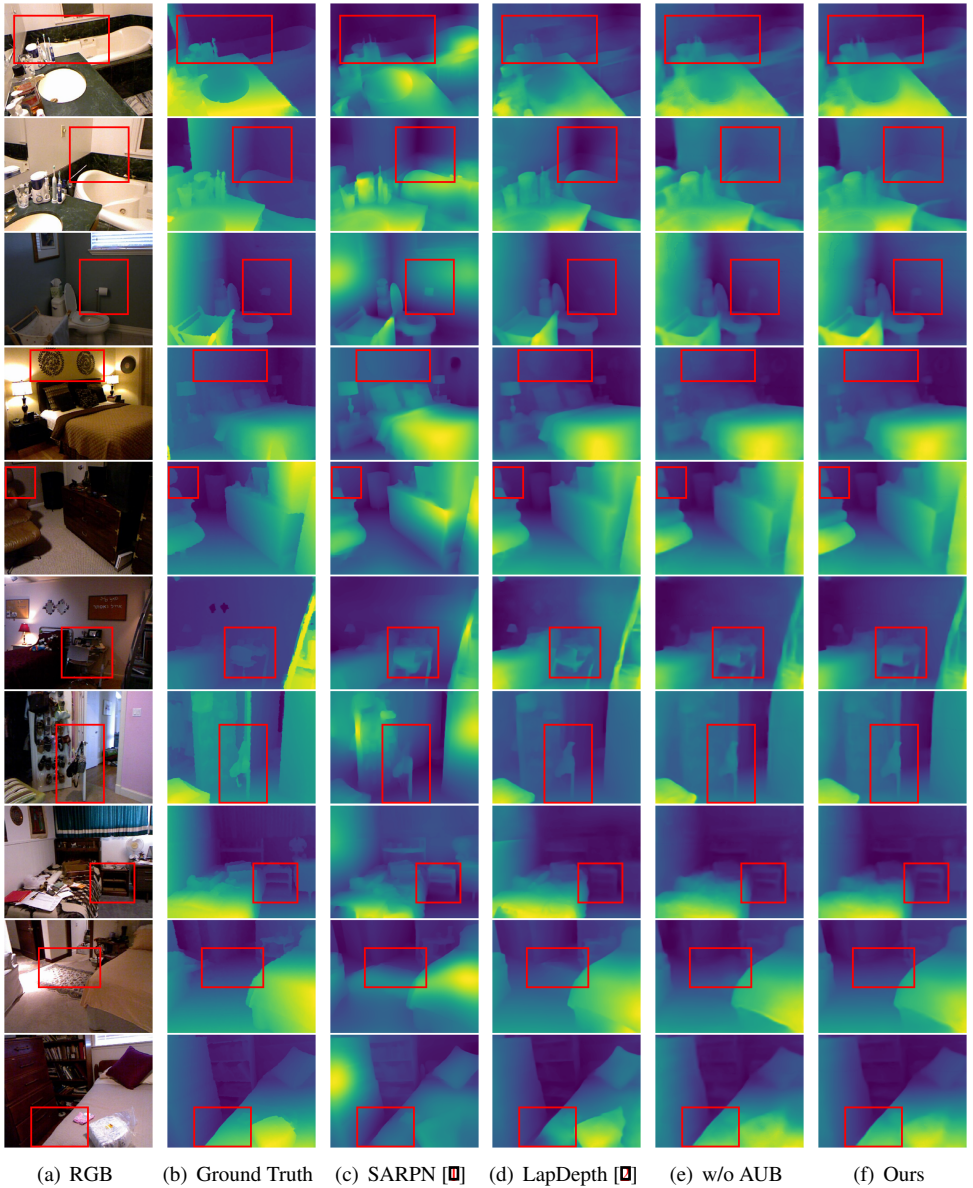
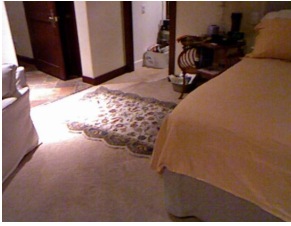


Figure 4: Qualitative comparisons with other methods on the NYU Depth V2 dataset.



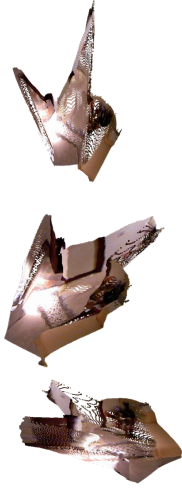
(a) Scene 1



(b) Scene 2



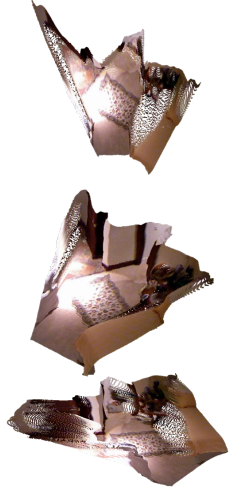
(c) Ground Truth



(d) LapDepth [10]



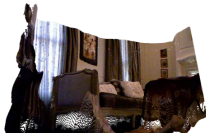
(e) Ours w/o AS



(f) Ours



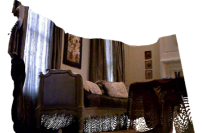
(g) Ground Truth



(h) LapDepth [10]



(i) Ours w/o AS



(j) Ours

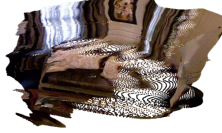


Figure 5: Point clouds comparison with other methods on the NYU Depth V2 dataset.



## 4 Point Cloud Reconstructions

The point clouds shown in Fig. 5 are rendered from depth maps in the same way as [4]. Point clouds in same row are captured from the same viewpoint. We compare our method with LapDepth [4] and our method without AS. It can be observed that point clouds generated from our proposed method is closer to that of ground-truth, especially around the wall and floor.

## References

- [1] Xiaotian Chen, Xuejin Chen, and Zheng-Jun Zha. Structure-aware residual pyramid network for monocular depth estimation. In *International Joint Conferences on Artificial Intelligence*, page 694–700, 2019.
- [2] M. Song, S. Lim, and W. Kim. Monocular depth estimation using laplacian pyramid-based depth residuals. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [3] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5684–5693, 2019.
- [4] Wei Yin, Yifan Liu, and Chunhua Shen. Virtual normal: Enforcing geometric constraints for accurate and robust depth prediction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.