

# URCDC-Depth: Uncertainty Rectified Cross-Distillation with CutFlip for Monocular Depth Estimation

Shuwei Shao<sup>1</sup>, Zhongcai Pei<sup>1</sup>, Weihai Chen<sup>1</sup>, Ran Li<sup>1</sup>, Zhong Liu<sup>1</sup>, Zhengguo Li<sup>2</sup>,

<sup>1</sup> School of Automation Science and Electrical Engineering, Beihang University

<sup>2</sup>Institute for Infocomm Research, A\*STAR  
swshao@buaa.edu.cn

## Abstract

This work aims to estimate a high-quality depth map from a single RGB image. Due to the lack of depth clues, making full use of the long-range correlation and the local information is critical for accurate depth estimation. Towards this end, we introduce an uncertainty rectified cross-distillation between Transformer and convolutional neural network (CNN) to learn a unified depth estimator. Specifically, we use the depth estimates from the Transformer branch and the CNN branch as pseudo labels to teach each other. Meanwhile, we model the pixel-wise depth uncertainty to rectify the loss weights of noisy pseudo labels. To avoid the large capacity gap induced by the strong Transformer branch deteriorating the cross-distillation, we transfer the feature maps from Transformer to CNN and design coupling units to assist the weak CNN branch to leverage the transferred features. Furthermore, we propose a surprisingly simple yet highly effective data augmentation technique CutFlip, which enforces the model to exploit more valuable clues apart from the vertical image position for depth inference. Extensive experiments demonstrate that our model, termed **URCDC-Depth**, exceeds previous state-of-the-art methods on the KITTI, NYU-Depth-v2 and SUN RGB-D datasets, even with no additional computational burden at inference time. The source code is publicly available at <https://github.com/ShuweiShao/URCDC-Depth>.

## Introduction

Monocular depth estimation is a fundamental research topic in the computer vision community, with applications ranging from scene understanding, 3D reconstruction through to augmented reality. Benefiting from the advances in convolutional neural networks (CNNs) (He et al. 2016; Tan and Le 2019), recent studies (Lee et al. 2019; Bhat, Alhashim, and Wonka 2021) achieve promising depth results. Due to the lack of depth cues, fully exploiting the long-range correlation (i.e., inter-object distance relationship) and the local information (i.e., intra-object consistency), is crucial for accurate depth estimation (Saxena et al. 2005). However, the convolution operator with a limited receptive field is hard to capture the long-range correlation, which becomes a potential bottleneck of current CNN-based depth estimation methods (Bhat, Alhashim, and Wonka 2021).

There are extensive works dedicated to alleviating the above limitation of CNN, which can be roughly divided into two categories: manipulating the convolution operation

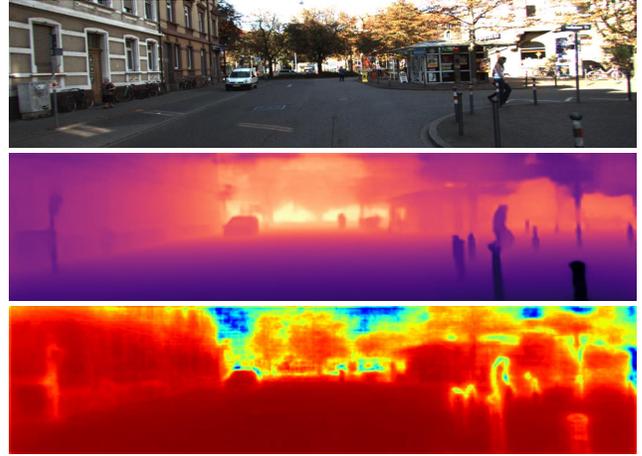


Figure 1: **Illustration of the depth and uncertainty maps from UR CDC-Depth.** *Top: input image; Middle: estimated depth map; Bottom: pixel-wise depth uncertainty (red: low uncertainty; yellow/blue: high/highest uncertainty).*

and introducing the attention mechanism (Vaswani et al. 2017). The former leverages atrous spatial pyramid pooling (Chen et al. 2017), coarse-to-fine fusion (Lin et al. 2017) and densely connecting (Zhang et al. 2020) to enhance the efficacy of convolution operator. The latter integrates the attention module to establish the long-distance dependency in the feature map (Zhou et al. 2019; Johnston and Carneiro 2020). In addition, several general methods adopt both of these strategies (Huynh et al. 2020; Bhat, Alhashim, and Wonka 2021). Despite the considerable improvements in performance, the dilemma remains.

Recently, visual Transformer has been demonstrated as a promising alternative to the CNN (Dosovitskiy et al. 2021; Liu et al. 2021; Peng et al. 2021). Building upon the attention mechanism, the Transformer with a global receptive field is more proficient in capturing the long-range correlation. Nevertheless, local feature details are prone to be ignored by it due to the lack of spatial inductive bias, resulting in unsatisfactory performance. A few depth estimation methods try to overcome the drawback of Transformer by utilizing additional CNN branch (Li et al. 2022; Shao et al. 2021).

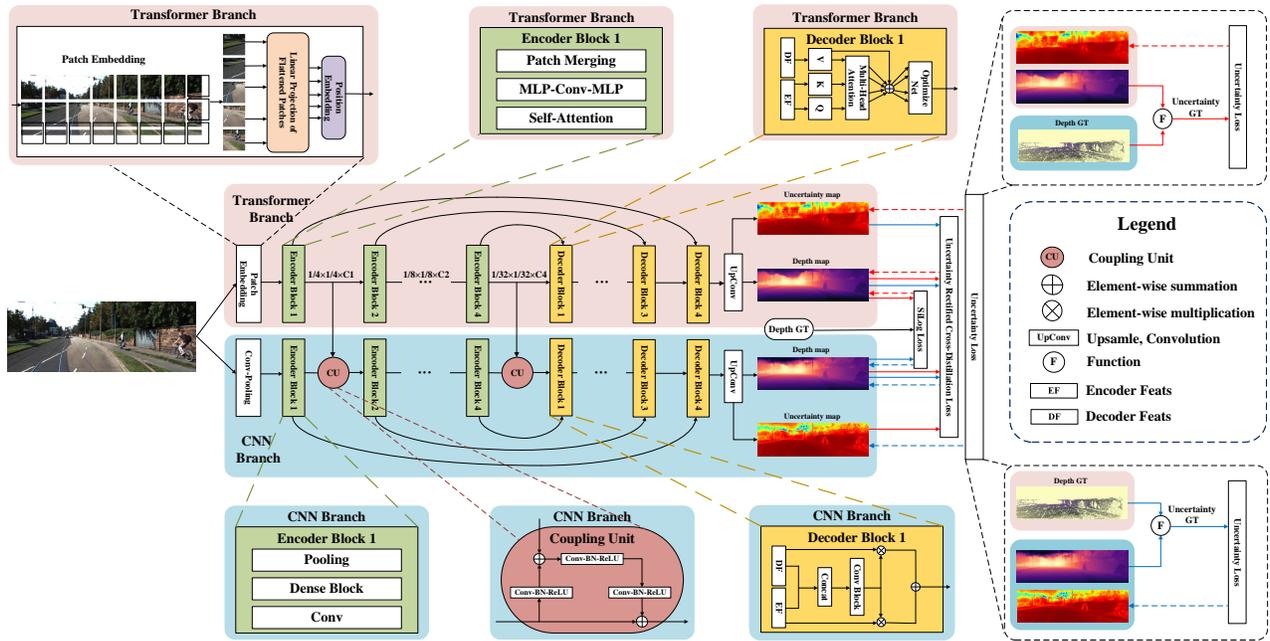


Figure 2: **Overview of the developed UR CDC-Depth.** Our UR CDC-Depth in the training phase consists of two branches, a Transformer branch and a CNN branch. During the evaluation phase, we only leverage the Transformer branch to generate the depth map.

However, these frameworks also rely on the CNN branch in the evaluation phase, increasing the computational cost at inference time.

In this work, we introduce a novel monocular depth estimation framework, termed **UR CDC-Depth** (Fig. 2), which integrates the strengths from both the Transformer and CNN using **cross-distillation** to enhance the performance. The core idea of UR CDC-Depth lies in that the Transformer branch establishes the long-range correlation while the CNN branch focuses on the local information, so cross-distillation between these two branches can help learn a unified depth estimator with both properties. To be specific, we use the derived depth estimates as pseudo labels to teach their counterparts. To alleviate the negative impact of noisy labels, we model the **pixel-wise depth uncertainty** to rectify the loss weights in these regions. The uncertainty map is predicted jointly with the corresponding depth map. In addition, we transfer the feature maps from the Transformer to the CNN so as to bridge the large performance gap induced by the strong Transformer branch and design **coupling units** to assist the weak CNN branch to utilize the transferred feature maps, which contribute to boost the performance of cross-distillation.

Furthermore, we train the UR CDC-Depth with a very simple yet effective data augmentation technique, **CutFlip**, based on the observation that monocular depth estimation model relies heavily on the vertical image position to infer depth, while other clues such as apparent sizes are ignored, deteriorating the model generalization ability (Dijk and Croon 2019). Generally, the feature of vertical image position is in that the closer the projection on the image is

to the lower boundary, the smaller the depth of the scene point. The traditional training mechanism allows the clue of vertical image position to exist in almost all training samples. In contrast, other cues are much less numerous. To resolve this, we vertically cut the training sample into upper and lower parts, and flip these two parts along the vertical direction with a certain probability, weakening the relationship between depth and vertical image position. In such case, the accuracy of predicted depth is significantly improved.

To summarize, the main contributions of this work are listed as:

- We introduce a novel monocular depth estimation model equipped with uncertainty rectified cross-distillation to exploit both the long-range correlation and the local information. Besides, the model has no additional computational burden in the evaluation phase thanks to the cross-distillation paradigm.
- We design a simple yet effective data augmentation strategy, which enforces the model to focus on more valuable cues for depth estimation, not just the clue of vertical image position.
- Detailed experiments and analysis indicate the efficacy of our developed components in improving the depth accuracy. The proposed approach achieves state-of-the-art performance on the KITTI (Geiger et al. 2013) and NYU-Depth-v2 (Silberman et al. 2012) datasets.

## Related work

**Monocular depth estimation** attempts to regress depth map from a single RGB image. As a seminal work, Saxena et al.

(2005) used a Markov random field to predict depth. Later, benefiting from the encoded features of CNNs that generalize well across diverse tasks, many follow-up works have achieved drastic performance improvement (Eigen, Puhrsch, and Fergus 2014; Qi et al. 2018; Fu et al. 2018). Recently, Lee et al. (2019) introduced local planar guidance layers to infer plane coefficients in the decoding stage, which were leveraged to recover the full resolution depth map. Bhat, Alhashim, and Wonka (2021) revisited the ordinal regression network (Fu et al. 2018) and proposed to calculate adaptive bins based on the image content.

**Transformer** has attracted a widespread attention owing to its effectiveness in natural language processing (Vaswani et al. 2017). In terms of computer vision, Dosovitskiy et al. (2021) introduced Vision Transformer (ViT) and indicated its feasibility on the image classification task. The success of the ViT accelerates the application of the Transformer to other tasks. Zheng et al. (2021) developed one of the first attempts at dense prediction tasks by using the ViT as the backbone.

There have been some attempts at applying Transformer to monocular depth estimation. Bhat, Alhashim, and Wonka (2021) utilized a minimized version of ViT to calculate bin width adaptively. Yang et al. (2021); Ranftl, Bochkovskiy, and Koltun (2021); Kim et al. (2022); Yuan et al. (2022) used the Transformer as an encoder to attain a global receptive field. As demonstrated in (Li et al. 2022), the model with a Transformer encoder tends to lose local depth details, *e.g.*, sharp edges. Li et al. (2022) proposed to use the Transformer encoder and an additional CNN encoder so that the model can enjoy the desired properties from both networks. However, the ensembled model increases the computational complexity. By contrast, we leverage the cross-distillation between the Transformer and CNN to construct a unified depth estimator, which allows our model to use only the Transformer branch at inference time with no additional computational burden.

**Knowledge distillation** is a learning paradigm targeting to transfer the learned knowledge from a teacher model to a lower-capacity student model, which is initially proposed on image recognition (Hinton et al. 2015). Since then, numerous knowledge distillation variants have been proposed, either working to improve its effectiveness (Yim et al. 2017; Tian, Krishnan, and Isola 2019; Sun et al. 2019) or applying it to other tasks (Garcia, Morerio, and Murino 2018; Hafner et al. 2018). Cross-distillation is a special case where models reach a consensus by simultaneously teaching each other, which is similar to the mutual learning (Zhang et al. 2018). A few studies also apply knowledge distillation to enhance monocular depth estimation (Pilzer et al. 2019; Aleotti et al. 2020). Unlike these methods, we introduce an uncertainty rectified cross-distillation for accurate depth estimation.

**Data augmentation** is a powerful technique in mitigating overfitting by increasing the effective amount of training samples. Therefore, common data augmentation techniques such as color jitter, crop, rotation are used in various tasks to improve model performance. Besides, there are methods tailored for monocular depth estimation, CutDepth (Ishii and Yamashita 2021) and DataGrafting (Peng et al. 2021). The

motivation of CutDepth differs significantly from our CutFlip. Concretely, the CutDepth aims to shorten the distance between the RGB image and the depth map in the latent space by replacing part of the RGB image with the corresponding depth ground-truth. While DataGrafting also aims to mitigate the overfitting risk for vertical image position, it relies on grafting together two training samples with different semantics. In contrast, our CutFlip is simpler and easier to implement as it only requires one training sample.

## Methodology

In this section, we elaborate on the main contributions of this work, namely uncertainty rectified cross-distillation and CutFlip. An overview of URDCD-Depth is presented in Fig. 2.

### Uncertainty Rectified Cross-Distillation

**Network architecture.** The proposed Transformer branch shares a same network architecture with NeWCRFs (Yuan et al. 2022) apart from the final prediction layer, which generates not only the depth map, but also the pixel-wise depth uncertainty. The encoder uses the Swin Transformer (Liu et al. 2021) to extract hierarchical feature representations. The decoder is composed of four neural window fully-connected conditional random fields (CRFs) modules.

The proposed CNN branch is also based on an encoder-decoder structure, where the encoder is DenseNet (Huang et al. 2017) and the decoder is similar to (Kim et al. 2022). The CNN branch is only used for complementary training, and will be discarded once training process is complete.

**Cross-distillation between Transformer and CNN.** We make use of the cross-distillation to construct a unified depth estimator that fully exploits the long-range correlation and the local information. For an input RGB image  $\mathbf{r}_n(\mathbf{p})$ , where  $\mathbf{p}$  denotes the pixel coordinate, our model generates two depth predictions,

$$\mathbf{d}_n^t(\mathbf{p}) = f_\theta^t(\mathbf{r}_n(\mathbf{p})); \mathbf{d}_n^c(\mathbf{p}) = f_\theta^c(\mathbf{r}_n(\mathbf{p})), \quad (1)$$

where  $\mathbf{d}_n^t(\mathbf{p})$  and  $\mathbf{d}_n^c(\mathbf{p})$  denotes the predictions from Transformer branch  $f_\theta^t(\cdot)$  and CNN branch  $f_\theta^c(\cdot)$ , respectively. As mentioned before, the Transformer and CNN are asymmetric learning networks, where the Transformer relies on the long-range self-attention while the CNN is built upon the local convolution operator, so the predictions  $\mathbf{d}_n^t(\mathbf{p})$  and  $\mathbf{d}_n^c(\mathbf{p})$  have inherently diverse properties and are used as pseudo labels to guide their counterparts towards the correct depth.

**Uncertainty-based rectification.** However, the pseudo labels contain heavy noises, particularly at the beginning of training, which inevitably damage the entire training process and enforce wrong predictions. To alleviate the negative impact of depth noises, we introduce an uncertainty rectified cross-distillation loss, defined as

$$\begin{aligned} \mathcal{L}_{urcd} = & \sum_{\mathbf{p}} (1 - \overline{\mathbf{u}}_n^c(\mathbf{p})) \odot \left| \mathbf{d}_n^t(\mathbf{p}) - \overline{\mathbf{d}}_n^c(\mathbf{p}) \right| \\ & + \sum_{\mathbf{p}} \left( 1 - \overline{\mathbf{u}}_n^t(\mathbf{p}) \right) \odot \left| \mathbf{d}_n^c(\mathbf{p}) - \overline{\mathbf{d}}_n^t(\mathbf{p}) \right|, \end{aligned} \quad (2)$$



Figure 3: **Illustration of the technique CutFlip.** The CutFlip contains two key steps: vertical cut and vertical flip.

where  $\bar{\cdot}$  is the gradient stopping operation,  $\odot$  is the element-wise multiplication, and  $\mathbf{u}_n^t(\mathbf{p})$  and  $\mathbf{u}_n^c(\mathbf{p})$  are the uncertainty predictions from Transformer branch  $f_\theta^t(\cdot)$  and CNN branch  $f_\theta^c(\cdot)$ , respectively, which are used to downweight the relevant pixels to alleviate the negative impact on regions with high uncertainty. The values of uncertainty map are ranging from 0 to 1.

Since there is no ground-truth for the uncertainty prediction, we model it with a function inspired by the probability density function of Laplace distribution, mathematically,

$$\mathbf{u}_n^*(\mathbf{p}) = 1 - \exp\left(-\frac{|\mathbf{d}_n(\mathbf{p}) - \mathbf{d}_n^*(\mathbf{p})|}{b(\mathbf{d}_n(\mathbf{p}) + \mathbf{d}_n^*(\mathbf{p}))}\right), \mathbf{p} \in \mathbf{T}, \quad (3)$$

where  $\mathbf{d}_n(\mathbf{p})$  stands for the predicted depth map,  $\mathbf{d}_n^*(\mathbf{p})$  denotes the ground-truth depth map,  $b$  is a coefficient that controls the tolerance for error and is set as 0.2 in this work, and  $\mathbf{T}$  denotes a set of pixels with valid ground-truth depth values. Here, instead of using the absolute difference between  $\mathbf{d}_n(\mathbf{p})$  and  $\mathbf{d}_n^*(\mathbf{p})$  directly, we normalize it by their sum due to the fact that a unit depth difference (e.g., 1m) represents the different uncertainty between distant and nearby points in a scene, and it should be higher on nearby points and less on distant points. We apply  $\mathcal{L}_u$  to enforce the uncertainty prediction to approximate  $\mathbf{u}_n^*(\mathbf{p})$ ,

$$\mathcal{L}_u = \sum_{\mathbf{p}} |\mathbf{u}_n^t(\mathbf{p}) - \mathbf{u}_n^{t*}(\mathbf{p})| + \sum_{\mathbf{p}} |\mathbf{u}_n^c(\mathbf{p}) - \mathbf{u}_n^{c*}(\mathbf{p})|, \mathbf{p} \in \mathbf{T}. \quad (4)$$

**Coupling unit.** As presented in recent studies (Li et al. 2022; Yuan et al. 2022; Lee et al. 2019), the Swin Transformer-based models perform much better than CNN-based ones for depth estimation. Moreover, the large capacity gap between teacher and student tends to cause poor performance of knowledge distillation (Hu et al. 2021). To bridge it, we transfer the feature maps encoded by the Transformer branch to the CNN branch, and design **coupling units** to fuse these two types of features. First, we align the channel dimension of the feature maps via a  $1 \times 1$  convolution and add them together. Second, we use a  $3 \times 3$  convolution to adaptively fuse the added features. Finally, we adjust the feature channel dimension with a  $1 \times 1$  convolution to form a residual connection. Meanwhile, BatchNorm (Ioffe and Szegedy 2015) and ReLU activation are utilized to regularize features. Note that the feature transfer operation only exists in the encoder that mainly determines the model performance.

**Overall loss.** The total optimization objective to train the UR CDC-Depth is summarized as follows

$$\mathcal{L}_{total} = \mathcal{L}_{ssi} + \lambda_1 \mathcal{L}_{urcd} + \lambda_2 \mathcal{L}_u, \quad (5)$$

---

### Algorithm 1: CutFlip

---

**Input:** RGB image  $\mathbf{r}_n$ ; Ground-truth depth map  $\mathbf{d}_n^*$ ; Height of input  $h$ .

**Output:** Transformed  $\mathbf{r}_n$ ; Transformed  $\mathbf{d}_n^*$ .

- 1: Random sampling  $p$  from the uniform distribution  $\mathbf{U}(0, 1)$ .
  - 2: **if**  $p < 0.5$  **then**
  - 3:   **return**  $\mathbf{r}_n; \mathbf{d}_n^*$ .
  - 4: **else**
  - 5:   Random sampling  $\varsigma$  from  $[\text{floor}(0.2h), \text{floor}(0.8h)]$ ;
  - 6:    $\mathbf{r} = \mathbf{r}_n.\text{copy}()$ ;  $\mathbf{d}^* = \mathbf{d}_n^*.\text{copy}()$ ;
  - 7:    $\mathbf{r}_n[h - \varsigma, :, :, :] = \mathbf{r}[\varsigma, :, :, :]$ ;  $\mathbf{d}_n^*[h - \varsigma, :, :, :] = \mathbf{d}^*[\varsigma, :, :, :]$ ;
  - 8:    $\mathbf{r}_n[h - \varsigma, :, :, :] = \mathbf{r}[\varsigma, :, :, :]$ ;  $\mathbf{d}_n^*[h - \varsigma, :, :, :] = \mathbf{d}^*[\varsigma, :, :, :]$ ;
  - 9:   **return**  $\mathbf{r}_n; \mathbf{d}_n^*$ .
  - 10: **end if**
- 

with

$$\begin{aligned} \mathcal{L}_{ssi} = & \kappa \sqrt{\frac{1}{|\mathbf{T}|} \sum_{\mathbf{p}} (\mathbf{g}_n^t(\mathbf{p}))^2 - \frac{\eta}{|\mathbf{T}|^2} \left( \sum_{\mathbf{p}} \mathbf{g}_n^t(\mathbf{p}) \right)^2} \\ & + \kappa \sqrt{\frac{1}{|\mathbf{T}|} \sum_{\mathbf{p}} (\mathbf{g}_n^c(\mathbf{p}))^2 - \frac{\eta}{|\mathbf{T}|^2} \left( \sum_{\mathbf{p}} \mathbf{g}_n^c(\mathbf{p}) \right)^2}, \mathbf{p} \in \mathbf{T}, \end{aligned} \quad (6)$$

where  $\mathcal{L}_{ssi}$  denotes the scaled scale-invariant loss introduced by (Lee et al. 2019),  $\mathbf{g}_n(\mathbf{p}) = \log \mathbf{d}_n(\mathbf{p}) - \log \mathbf{d}_n^*(\mathbf{p})$ ,  $\kappa$  and  $\eta$  are set as 10 and 0.85 based on (Lee et al. 2019),  $\lambda_1$  and  $\lambda_2$  are empirically set as 0.1 and 0.5.

### CutFlip

The quantity and diversity of training data are critical for deep learning-based models, which however are hard to satisfy in depth estimation because data acquisition is extremely expensive and laborious. Lack of data deteriorates the model generalization ability, and one of the serious overfitting threats is the heavy reliance on the vertical image position (Dijk and Croon 2019). To enforce the model to focus on more valuable clues, we propose a surprisingly simple yet highly effective data augmentation technique, **CutFlip**. As illustrated in Fig. 3, we vertically cut the input sample into upper and lower parts, highlighted by the orange box and the green box, respectively, and flip these two parts along the vertical direction to weaken the relationship of depth and vertical image position. The details are shown in Algorithm 1. The CutFlip is performed with a probability of 0.5, and the vertical position to cut is randomly sampled, which allows the model to greatly adapt to various types of data.

## Experiments

### Datasets

**KITTI dataset** is captured from outdoor scenes with equipment placed on a moving vehicle (Geiger et al. 2013). The image resolution is around  $1241 \times 376$  pixels. We use two commonly used splits for monocular depth estimation. One

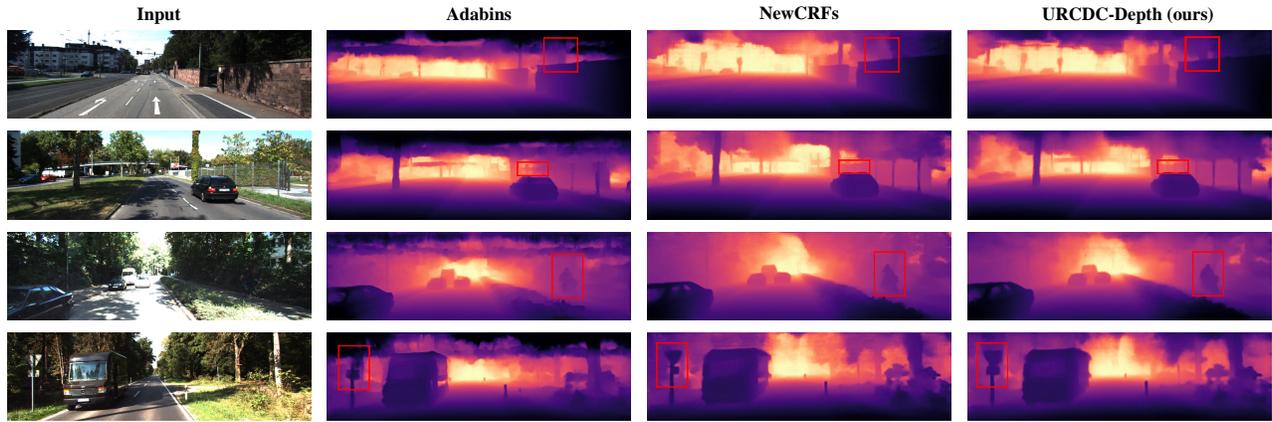


Figure 4: **Qualitative depth results on the Eigen split of KITTI dataset.** The red boxes indicate the regions to emphasize.

Method	Cap	Abs Rel ↓	Sq Rel ↓	RMSE ↓	RMSE log ↓	$\delta < 1.25 \uparrow$	$\delta < 1.25^2 \uparrow$	$\delta < 1.25^3 \uparrow$
Eigen <i>et al.</i> (Eigen and Fergus 2015)	0-80m	0.203	1.548	6.307	0.282	0.702	0.898	0.967
Fu <i>et al.</i> (Fu et al. 2018)	0-80m	0.072	0.307	2.727	0.120	0.932	0.984	0.994
VNL (Yin et al. 2019)	0-80m	0.072	-	3.258	0.117	0.938	0.990	0.998
BTS (Lee et al. 2019)	0-80m	0.061	0.261	2.834	0.099	0.954	0.992	0.998
PWA (Lee et al. 2021)	0-80m	0.060	0.221	2.604	0.093	0.958	0.994	<b>0.999</b>
TransDepth (Yang et al. 2021)	0-80m	0.064	0.252	2.755	0.098	0.956	0.994	<b>0.999</b>
Adabins (Bhat, Alhashim, and Wonka 2021)	0-80m	0.058	0.190	2.360	0.088	0.964	0.995	<b>0.999</b>
P3Depth (Patil et al. 2022)	0-80m	0.071	0.270	2.842	0.103	0.953	0.993	0.998
DepthFormer (Li et al. 2022)	0-80m	0.052	0.158	2.143	0.079	0.975	<b>0.997</b>	<b>0.999</b>
NeWCRFs (Yuan et al. 2022)	0-80m	0.052	0.155	2.129	0.079	0.974	<b>0.997</b>	<b>0.999</b>
<b>URCDC-Depth (ours)</b>	0-80m	<b>0.050</b>	<b>0.142</b>	<b>2.032</b>	<b>0.076</b>	<b>0.977</b>	<b>0.997</b>	<b>0.999</b>
Fu <i>et al.</i> (Fu et al. 2018)	0-50m	0.071	0.268	2.271	0.116	0.936	0.985	0.995
BTS (Lee et al. 2019)	0-50m	0.058	0.183	1.995	0.090	0.962	0.994	0.999
PWA (Lee et al. 2021)	0-50m	0.057	0.161	1.872	0.087	0.965	0.995	0.999
TransDepth (Yang et al. 2021)	0-50m	0.061	0.185	1.992	0.091	0.963	0.995	0.999
P3Depth (Patil et al. 2022)	0-50m	0.055	0.130	1.651	0.081	0.974	0.997	0.999
<b>URCDC-Depth (ours)</b>	0-50m	<b>0.049</b>	<b>0.108</b>	<b>1.528</b>	<b>0.072</b>	<b>0.981</b>	<b>0.998</b>	<b>1.000</b>

Table 1: **Quantitative depth comparison on the Eigen split of KITTI dataset.** Note that the backbones of DepthFormer, NeWCRFs and our UR CDC-Depth at inference time are Swin-Large and ResNet-50, Swin-Large and Swin-Large, respectively. “-” indicates not applicable. The best results are highlighted in **bold**.

Method	SILog ↓	sqErrRel ↓	absErrRel ↓	iRMSE ↓
Fu <i>et al.</i> (Fu et al. 2018)	11.77	8.78	2.23	12.98
BTS (Lee et al. 2019)	11.67	9.04	2.21	12.23
BA-Full (Aich et al. 2021)	11.61	9.38	2.29	12.23
PackNet-SAN (Guizilini et al. 2021)	11.54	9.12	2.35	12.38
PWA (Lee et al. 2021)	11.45	9.05	2.30	12.32
NeWCRFs (Yuan et al. 2022)	10.39	8.37	1.83	11.03
<b>URCDC-Depth (ours)</b>	<b>10.03</b>	<b>8.24</b>	<b>1.74</b>	<b>10.71</b>

Table 2: **Quantitative depth comparison on the official split of KITTI dataset.** The results are available from the online server.

is the Eigen split (Eigen, Puhrsch, and Fergus 2014) including 23488 training image pairs and 697 testing images. The other one is the official split (Geiger et al. 2013) including

42949 training image pairs, 1000 validation images and 500 testing images. The evaluation results on the official split are generated by the online server.

**NYU-Depth-v2 dataset** provides RGB images and depth maps collected from indoor scenes at a resolution of  $640 \times 480$  pixels (Silberman et al. 2012). Following prior works, we adopt the official split and the dataset processed by Lee et al. (2019), which contains 24231 training images and 654 testing images.

### Implementation Details

The UR CDC-Depth is implemented in PyTorch (Paszke et al. 2017) and trained on NVIDIA RTX A5000 GPUs. We optimize it using the Adam optimizer (Kingma and Ba 2015) where  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . The training process runs a total number of 20 epochs with a batch size of 8 and a learning

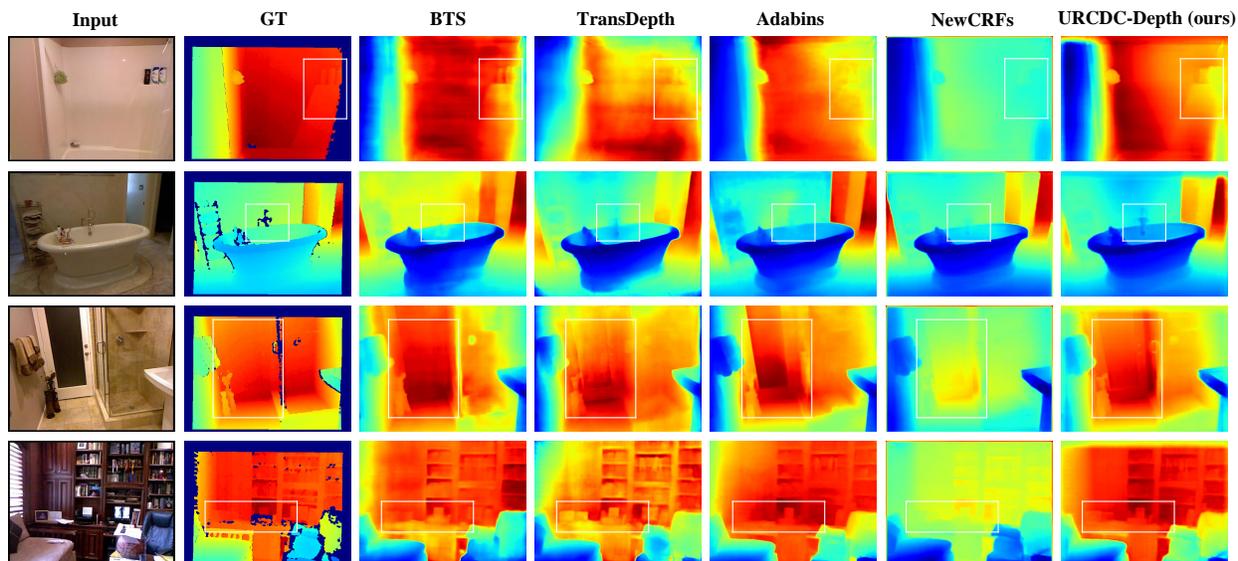


Figure 5: **Qualitative depth results on the NYU-Depth-v2 dataset.** The white boxes indicate the regions to emphasize.

Method	Cap	Abs Rel ↓	RMSE ↓	$\log_{10}$ ↓	$\delta < 1.25$ ↑	$\delta < 1.25^2$ ↑	$\delta < 1.25^3$ ↑
Eigen <i>et al.</i> (Eigen, Puhrsch, and Fergus 2014)	0-10m	0.158	0.641	-	0.769	0.950	0.988
Fu <i>et al.</i> (Fu <i>et al.</i> 2018)	0-10m	0.115	0.509	0.051	0.828	0.965	0.992
VNL (Yin <i>et al.</i> 2019)	0-10m	0.108	0.416	0.048	0.875	0.976	0.994
BTS (Lee <i>et al.</i> 2019)	0-10m	0.113	0.407	0.049	0.871	0.977	0.995
DAV (Huynh <i>et al.</i> 2020)	0-10m	0.108	0.412	-	0.882	0.980	0.996
PWA (Lee <i>et al.</i> 2021)	0-10m	0.105	0.374	0.045	0.892	0.985	0.997
Long <i>et al.</i> (Long <i>et al.</i> 2021)	0-10m	0.101	0.377	0.044	0.890	0.982	0.996
TransDepth (Yang <i>et al.</i> 2021)	0-10m	0.106	0.365	0.045	0.900	0.983	0.996
Adabins (Bhat, Alhashim, and Wonka 2021)	0-10m	0.103	0.364	0.044	0.903	0.984	0.997
P3Depth (Patil <i>et al.</i> 2022)	0-10m	0.104	0.356	0.043	0.898	0.981	0.996
DepthFormer (Li <i>et al.</i> 2022)	0-10m	0.096	0.339	0.041	0.921	0.989	0.998
NeWCRFs (Patil <i>et al.</i> 2022)	0-10m	0.095	0.334	0.041	0.922	<b>0.992</b>	<b>0.998</b>
<b>URCDC-Depth (ours)</b>	0-10m	<b>0.088</b>	<b>0.316</b>	<b>0.038</b>	<b>0.933</b>	<b>0.992</b>	<b>0.998</b>

Table 3: **Quantitative depth comparison on the NYU-Depth-v2 dataset.**

rate scheduled via polynomial decay from  $1e-4$  to  $1e-5$ . We use the standard data augmentation techniques and evaluation metrics following previous works (Yuan *et al.* 2022; Lee *et al.* 2019).

### Comparison to State-of-the-Arts

**KITTI.** We first conduct comparison with the leading methods on the Eigen split. Table 1 shows the results, indicating that our UR CDC-Depth exceeds previous methods. It is worth noting that although UR CDC-Depth and NewCRFs share the almost identical network structure in the evaluation phase, it improves the NeWCRFs by 8.4% and 4.6% on the Sq Rel and RMSE. Fig. 4 presents qualitative depth comparisons. As we can see, the NeWCRFs is struggle with thinner structures, *e.g.*, posts and difficult object boundaries such as human boundary, while our UR CDC-Depth is capable of es-

timating these small details, which supports our standpoint that the cross-distillation helps to learn a unified depth estimator with both desired properties from the Transformer and CNN.

We then compare our UR CDC-Depth against the competing methods on the official split. The results are generated by the online server and reported in Table 2. Here we can see that our UR CDC-Depth outperforms previous methods again. A notable phenomenon is that the main ranking metric SILog, from the model proposed by Fu *et al.* to the PWA, has only increased by 2.5% in three years. With the advent of visual Transformer, the NeWCRFs makes a significant breakthrough on this metric via the designed neural CRFs. Our UR CDC-Depth further improves the NeWCRFs by 3.5% on the SILog, even with no additional computational burden at inference time.

ID	CD	UP	CU	CF	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta < 1.25 \uparrow$
1					0.052	0.155	2.129	0.974
2	✓				0.055	0.155	2.086	0.975
3	✓	✓			0.051	0.147	2.076	0.976
4	✓		✓		0.051	0.148	2.078	0.975
5	✓	✓	✓		0.051	0.147	2.062	0.976
6				✓	0.051	0.144	2.056	<b>0.977</b>
7	✓	✓	✓	✓	<b>0.050</b>	<b>0.142</b>	<b>2.032</b>	<b>0.977</b>

Table 4: **Ablation study of the proposed UR CDC-Depth on the KITTI dataset.** CD: cross-distillation; UP: uncertainty map; CU: coupling unit; CF: CutFlip.

ID	CD	UP	CU	CF	Abs Rel ↓	RMSE ↓	$\log_{10} \downarrow$	$\delta < 1.25 \uparrow$
1					0.095	0.334	0.041	0.922
2	✓				0.095	0.329	0.040	0.923
3	✓	✓			0.091	0.323	0.039	0.928
4	✓		✓		0.091	0.326	0.039	0.926
5	✓	✓	✓		0.089	0.319	<b>0.038</b>	0.931
6				✓	0.091	0.322	0.039	<b>0.934</b>
7	✓	✓	✓	✓	<b>0.088</b>	<b>0.316</b>	<b>0.038</b>	0.933

Table 5: **Ablation study of the proposed UR CDC-Depth on the NYU-Depth-v2 dataset.**

Augmentation method	Abs Rel ↓	Sq Rel ↓	RMSE ↓	$\delta < 1.25 \uparrow$
CutMix (Yun et al. 2019)	0.054	0.154	2.093	0.974
CutOut (DeVries and Taylor 2017)	0.052	0.148	2.078	0.975
CutDepth (Ishii and Yamashita 2021)	0.052	0.150	2.076	0.975
DataGrafting (Peng et al. 2021)	0.051	0.146	2.051	0.976
<b>CutFlip</b>	<b>0.050</b>	<b>0.142</b>	<b>2.032</b>	<b>0.977</b>

Table 6: **Comparison of data augmentation techniques on the KITTI dataset.**

**NYU-Depth-v2.** To demonstrate the competitiveness of our UR CDC-Depth in the indoor scenario, we also evaluate it on the NYU-Depth-v2 dataset. The results are reported in Table 3, which indicates that our method greatly boosts the performance on most metrics, such as Abs Rel and RMSE. This emphasizes our contributions in improving the results. We display qualitative depth comparisons in Fig. 5. As can be seen, our UR CDC-Depth preserves small details *e.g.*, handle and predicts sharp depth edges even in scenes with extremely scarce texture (top row).

### Ablation Study

To better inspect how the proposed components in UR CDC-Depth affect the performance, we present detailed ablation studies on the KITTI and NYU-Depth-v2 datasets, which are shown in Table 4 and Table 5, respectively.

**Cross-distillation.** We start from the baseline NeWCRFs (ID 1). By directly introducing the cross-distillation, we observe a slight performance improvement on most metrics (ID 2). The cross-distillation suffers from the heavy depth noises from pseudo labels. Besides, the large capacity gap between Transformer branch and CNN branch limits the performance

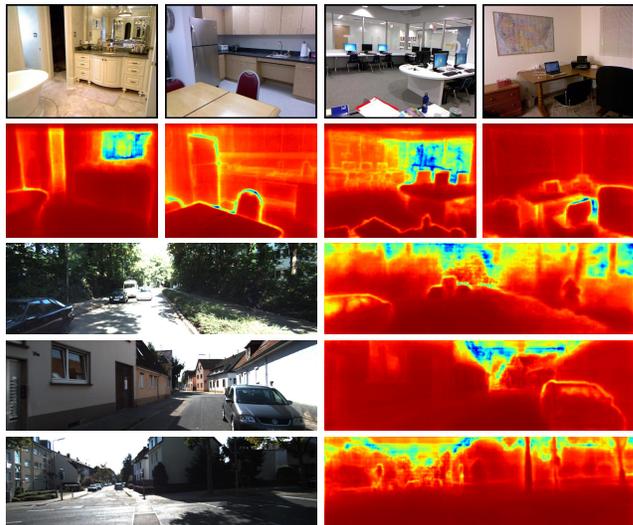


Figure 6: **Pixel-wise depth uncertainty from our UR CDC-Depth.** Red indicates areas of low uncertainty, yellow/blue indicates areas of high/highest uncertainty.

gain from the cross-distillation.

**Uncertainty map.** ID 3 presents the addition of the uncertainty map, which greatly boosts the performance. Especially notable are the results on Sq Rel sensitive to the large depth errors in Table 4. The uncertainty map helps mitigate the negative impact of depth noises in the training process, hence resulting in a considerably lower Sq Rel result. Fig. 6 presents a visualization of uncertainty maps. Unlike the prior method (Johnston and Carneiro 2020) that is also capable of estimating uncertainty maps showing a clear trend where uncertainty increases with distance, our uncertainty maps focus on the difficult regions, such as object boundaries and vanishing points. We attribute this to the normalization operation when modeling the ground-truth of uncertainty map.

**Coupling unit.** To bridge the large performance gap between the Transformer branch and CNN branch, we transfer the feature maps and use the coupling units to fuse the transferred features from the Transformer branch and the features in the CNN branch. The results after adding the coupling units are in IDs 4 and 5, which contributes to the performance.

**CutFlip.** With the CutFlip data augmentation, we can see that the performance is improved significantly (IDs 6 and 7). Besides, we compare the CutFlip against other similar augmentation methods to further demonstrate its efficacy in Table 6. To make a fair comparison, we remain all other configurations the same except for these specially crafted data augmentation techniques. The inferior performance of CutMix, CutOut and CutDepth may lie in the lack of constraint on the vertical image position. DataGrafting takes the overfitting risk of vertical image position into account. However, grafting together two training samples with different semantics increases the learning burden of network.

## Conclusion

In this work, we introduce a novel monocular depth estimation framework URDC-Depth, which leverages uncertainty rectified cross-distillation to fully exploit the long-range correlation and the local information. The paradigm allows our framework to generate precisely estimated depth maps with no additional computational burden at inference time. In addition, we propose a simple yet effective data augmentation technique CutFlip to enforce the model to emphasize more valuable depth reasoning clues apart from the vertical image position. We conduct comprehensive experiments on the KITTI, NYU-Depth-v2 and SUN RGB-D datasets, and the experimental results verify the efficacy of the proposed URDC-Depth.

## References

- Aich, S.; Vianney, J. M. U.; Islam, M. A.; and Liu, M. K. B. 2021. Bidirectional attention network for monocular depth estimation. In *2021 IEEE International Conference on Robotics and Automation*, 11746–11752. IEEE.
- Aleotti, F.; Zaccaroni, G.; Bartolomei, L.; Poggi, M.; Tosi, F.; and Mattoccia, S. 2020. Real-time single image depth perception in the wild with handheld devices. *Sensors*, 21(1): 15.
- Bhat, S. F.; Alhashim, I.; and Wonka, P. 2021. Adabins: Depth estimation using adaptive bins. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4009–4018.
- Chen, L.-C.; Papandreou, G.; Kokkinos, I.; Murphy, K.; and Yuille, A. L. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4): 834–848.
- DeVries, T.; and Taylor, G. W. 2017. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*.
- Dijk, T. v.; and Croon, G. d. 2019. How do neural networks see depth in single images? In *Proceedings of the IEEE International Conference on Computer Vision*, 2183–2191.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*.
- Eigen, D.; and Fergus, R. 2015. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, 2650–2658.
- Eigen, D.; Puhersch, C.; and Fergus, R. 2014. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems*, 2366–2374.
- Fu, H.; Gong, M.; Wang, C.; Batmanghelich, K.; and Tao, D. 2018. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2002–2011.
- Garcia, N. C.; Morerio, P.; and Murino, V. 2018. Modality distillation with multiple stream networks for action recognition. In *Proceedings of the European Conference on Computer Vision*, 103–118.
- Geiger, A.; Lenz, P.; Stiller, C.; and Urtasun, R. 2013. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11): 1231–1237.
- Guizilini, V.; Ambrus, R.; Burgard, W.; and Gaidon, A. 2021. Sparse auxiliary networks for unified monocular depth prediction and completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 11078–11088.
- Hafner, F.; Bhuiyan, A.; Kooij, J. F.; and Granger, E. 2018. A cross-modal distillation network for person re-identification in rgb-depth. *arXiv preprint arXiv:1810.11641*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Hinton, G.; Vinyals, O.; Dean, J.; et al. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7).
- Hu, J.; Fan, C.; Jiang, H.; Guo, X.; Gao, Y.; Lu, X.; and Lam, T. L. 2021. Boosting light-weight depth estimation via knowledge distillation. *arXiv preprint arXiv:2105.06143*.
- Huang, G.; Liu, Z.; Van Der Maaten, L.; and Weinberger, K. Q. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4700–4708.
- Huynh, L.; Nguyen-Ha, P.; Matas, J.; Rahtu, E.; and Heikkilä, J. 2020. Guiding monocular depth estimation using depth-attention volume. In *European Conference on Computer Vision*, 581–597. Springer.
- Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, 448–456. PMLR.
- Ishii, Y.; and Yamashita, T. 2021. CutDepth: edge-aware data augmentation in depth estimation. *arXiv preprint arXiv:2107.07684*.
- Johnston, A.; and Carneiro, G. 2020. Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume. In *Proceedings of the IEEE International Conference on Computer Vision*, 4756–4765.
- Kim, D.; Ga, W.; Ahn, P.; Joo, D.; Chun, S.; and Kim, J. 2022. Global-local path networks for monocular depth estimation with vertical CutDepth. *arXiv preprint arXiv:2201.07436*.
- Kingma, D. P.; and Ba, J. 2015. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*.
- Lee, J. H.; Han, M.-K.; Ko, D. W.; and Suh, I. H. 2019. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv preprint arXiv:1907.10326*.
- Lee, S.; Lee, J.; Kim, B.; Yi, E.; and Kim, J. 2021. Patch-wise attention network for monocular depth estimation. In

- Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 1873–1881.
- Li, Z.; Chen, Z.; Liu, X.; and Jiang, J. 2022. DepthFormer: Exploiting long-range correlation and local information for accurate monocular depth estimation. *arXiv preprint arXiv:2203.14211*.
- Lin, T.-Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; and Belongie, S. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2117–2125.
- Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021. Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. In *Proceedings of the IEEE International Conference on Computer Vision*, 10012–10022.
- Long, X.; Lin, C.; Liu, L.; Li, W.; Theobalt, C.; Yang, R.; and Wang, W. 2021. Adaptive surface normal constraint for depth estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, 12849–12858.
- Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic differentiation in pytorch. In *Advances in Neural Information Processing Systems Workshop Autodiff*.
- Patil, V.; Sakaridis, C.; Liniger, A.; and Van Gool, L. 2022. P3Depth: Monocular depth estimation with a piecewise planarity prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1610–1621.
- Peng, Z.; Huang, W.; Gu, S.; Xie, L.; Wang, Y.; Jiao, J.; and Ye, Q. 2021. Conformer: local features coupling global representations for visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, 367–376.
- Pilzer, A.; Lathuiliere, S.; Sebe, N.; and Ricci, E. 2019. Refine and distill: Exploiting cycle-inconsistency and knowledge distillation for unsupervised monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9768–9777.
- Qi, X.; Liao, R.; Liu, Z.; Urtasun, R.; and Jia, J. 2018. Geonet: Geometric neural network for joint depth and surface normal estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 283–291.
- Ranftl, R.; Bochkovskiy, A.; and Koltun, V. 2021. Vision transformers for dense prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, 12179–12188.
- Saxena, A.; Chung, S. H.; Ng, A. Y.; et al. 2005. Learning depth from single monocular images. In *Advances in Neural Information Processing Systems*, volume 18, 1–8.
- Shao, S.; Li, R.; Pei, Z.; Liu, Z.; Chen, W.; Zhu, W.; Wu, X.; and Zhang, B. 2021. NENet: Monocular depth estimation via neural ensembles. *arXiv preprint arXiv:2111.08313*.
- Silberman, N.; Hoiem, D.; Kohli, P.; and Fergus, R. 2012. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision*, 746–760. Springer.
- Sun, D.; Yao, A.; Zhou, A.; and Zhao, H. 2019. Deeply-supervised knowledge synergy. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6997–7006.
- Tan, M.; and Le, Q. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 6105–6114. PMLR.
- Tian, Y.; Krishnan, D.; and Isola, P. 2019. Contrastive representation distillation. In *International Conference on Learning Representations*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, 5998–6008.
- Yang, G.; Tang, H.; Ding, M.; Sebe, N.; and Ricci, E. 2021. Transformer-based attention networks for continuous pixel-wise prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, 16269–16279.
- Yim, J.; Joo, D.; Bae, J.; and Kim, J. 2017. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4133–4141.
- Yin, W.; Liu, Y.; Shen, C.; and Yan, Y. 2019. Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, 5684–5693.
- Yuan, W.; Gu, X.; Dai, Z.; Zhu, S.; and Tan, P. 2022. Neural window fully-connected CRFs for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3916–3925.
- Yun, S.; Han, D.; Oh, S. J.; Chun, S.; Choe, J.; and Yoo, Y. 2019. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE International Conference on Computer Vision*, 6023–6032.
- Zhang, J.; Yue, H.; Wu, X.; Chen, W.; and Wen, C. 2020. Densely connecting depth maps for monocular depth estimation. In *European Conference on Computer Vision*, 149–165. Springer.
- Zhang, Y.; Xiang, T.; Hospedales, T. M.; and Lu, H. 2018. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4320–4328.
- Zheng, S.; Lu, J.; Zhao, H.; Zhu, X.; Luo, Z.; Wang, Y.; Fu, Y.; Feng, J.; Xiang, T.; Torr, P. H.; et al. 2021. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6881–6890.
- Zhou, J.; Wang, Y.; Qin, K.; and Zeng, W. 2019. Unsupervised high-resolution depth learning from videos with dual networks. In *Proceedings of the IEEE International Conference on Computer Vision*, 6872–6881.