

Excavating the Potential Capacity of Self-Supervised Monocular Depth Estimation

Rui Peng Ronggang Wang  Yawen Lai Luyang Tang Yangang Cai
School of Electronic and Computer Engineering, Peking University

{ruipeng, tly926}@stu.pku.edu.cn {alanlyawen, caiyangang}@pku.edu.cn rgwang@pkusz.edu.cn

Abstract

*Self-supervised methods play an increasingly important role in monocular depth estimation due to their great potential and low annotation cost. To close the gap with supervised methods, recent works take advantage of extra constraints, e.g., semantic segmentation. However, these methods will inevitably increase the burden on the model. In this paper, we show theoretical and empirical evidence that the potential capacity of self-supervised monocular depth estimation can be excavated without increasing this cost. In particular, we propose (1) a novel data augmentation approach called data grafting, which forces the model to explore more cues to infer depth besides the vertical image position, (2) an exploratory self-distillation loss, which is supervised by the self-distillation label generated by our new post-processing method - selective post-processing, and (3) the full-scale network, designed to endow the encoder with the specialization of depth estimation task and enhance the representational power of the model. Extensive experiments show that our contributions can bring significant performance improvement to the baseline with even less computational overhead, and our model, named **EPCDepth**, surpasses the previous state-of-the-art methods even those supervised by additional constraints. Code is available at <https://github.com/prstrive/EPCDepth>.*

1. Introduction

Depth estimation has always been a fundamental problem of computer vision, which dominates the performance of various applications, e.g., virtual reality, autonomous driving, robotics, etc. As the cheapest solution, monocular depth estimation (MDE) has made considerable progress due to the evolvement of Convolution Neural Networks [27, 44, 45, 16]. However, most existing state-of-the-art approaches rely on supervised training [7, 6, 8, 29, 1], whose training datasets collection is a cumbersome and formidable challenge. As an alternative, self-supervised methods elimi-

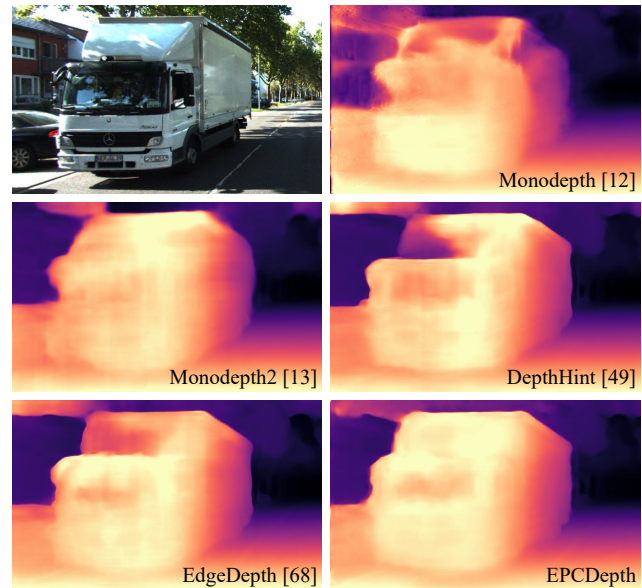


Figure 1. **Depth estimation from a single image.** Our model (**EPCDepth**), trained only on stereo data, performs the best and produces the sharpest and most complete result with the least computational cost.

nate the need for ground-truth depth through recasting depth estimation as the reconstruction problem among stereo images [10, 12, 49, 68], monocular video [67, 2, 42, 32] or a combination of both [62, 13].

In terms of performance alone, recent works have shown that the gap between self-supervision and full-supervision has made a de facto reduction. But on the other hand, this reduction largely benefits from the sophisticated model architecture and extra constraints from external modalities, e.g., semantic segmentation [2, 26, 68, 15], optical flow [57, 40], depth normal [54], etc. Apparently, these factors substantially increase the burden of the model training and run counter to the concept of self-supervision to some extent. In this paper, we show the potential of self-supervised monocular depth estimation even without these additional

constraints from three aspects: **data augmentation, self-distillation, and model architecture.**

Generally, the closer the projection on the image is to the lower boundary, the smaller the depth of the object. This feature of vertical image position has been proven to be the main cue adopted by the MDE model to infer depth, while the apparent size and other cues that humans will rely on are ignored [47]. We conjecture that the reason is that in the traditional training mechanism that takes the entire image as input, the feature of vertical image position exists in almost every training sample, while the number of samples for other cues is relatively small, which leads to a long-tailed distribution on cues. Obviously, this kind of paranoia tends to damage the generalization ability of the model. To solve this, we propose a novel data augmentation method called *Data Grafting*, which breaks this dilemma by vertically grafting a certain proportion from another image to appropriately weaken the relationship between depth and vertical image position. Moreover, there is another fact that the precision of different scales output by the multi-scale network is inconsistent at different pixels, and this motivates us to generate better disparity maps as pseudo-labels to realize the self-distillation of the model. Concretely, we propose *Selective Post-Processing* (SPP) to select the best prediction for each pixel among all scales according to the reconstruction error, which is inspired by the availability of all views during training, and the similar idea has been proven effective in the field of multi-view stereo [55]. Finally, we extend the traditional multi-scale network to the full-scale network by inserting prediction modules not only on the decoder but also on the encoder to advance the specialization of depth prediction from decoder to encoder and absorb the representational power of the model. The superior result of our model is shown in Figure 1.

To summarize, our main contributions are listed below in fourfold:

- We introduce a conceptually simple but empirically efficient data augmentation approach, which enables the model to learn more effective cues besides the vertical image position.
- We apply self-distillation to MDE for the first time without any auxiliary network and generate better pseudo-labels based on our training-oriented selective post-processing method.
- We propose a more efficient full-scale network to strengthen the constraints on the model and enhance the encoder’s specificity of depth estimation.
- Without bells and whistles, we achieve state-of-the-art performance within self-supervised methods even compared to those high-performance models that are trained by extra constraints.

2. Related Works

Self-Supervised Monocular Depth Estimation. The depth is predicted as an intermediate in self-supervised MDE to synthesize the reconstructed view from the source view, and the photometric loss between the target view and the reconstructed view is calculated as the target of minimization. There are mainly two kinds of self-supervised methods: trained by synchronized stereo images [10, 12, 38, 49, 68] or monocular video [67, 57, 2, 42]. For the first category, the model with known relative placement only needs to predict the disparity, that is, the inverse of the depth. For the second category, additional predictions of the relative pose of the camera are required. Recently, abundant works have improved the performance of self-supervised MDE through new loss function [10, 49, 13, 42, 68], new architecture [38, 66, 58, 14, 32] and new supervision from extra constraints [54, 57, 40, 2, 26, 68, 15].

In this paper, we further excavate the potential capacity of self-supervised MDE with the realization of training on stereo images.

Self-Distillation. Knowledge distillation is a pioneering work to transfer knowledge from powerful teacher networks to student networks using the softmax output [18], intermediate feature [41, 17], attention [61, 21], relationship [56, 34, 36], *etc.* Self-distillation is a special case where the model itself is used as a teacher. Intuitively, the model can be distilled by the same model trained previously [9], but these approaches are inefficient because they need to train multiple generations synchronously. Therefore, some recent works advocate distilling the model within one generation, which take supervision from prior iterations [52, 24], consistency of distorted data [51], invariance among intra-class [60] and the output of deeper portion [64].

These methods only focus on the self-distillation of the classification task. In this work, we applied self-distillation to the regression task of depth estimation. Different from the method of using the whole network to promote sub-networks in [38], we select the optimal disparity map from all output scales as the self-distillation label to distill the whole network in one generation.

Data Augmentation. For overfitting, data augmentation is an efficient approach to mitigate this drawback by implicitly increasing the total amount of training data and teaching models about the invariance of the data domain. Common data augmentation methods can be summarized into two categories: learnable [46, 4] and parameter learning free [27, 5, 63, 59, 65]. Learnable methods are more universal and work out of the box, while the subsequent methods are easier to be implemented and most of them are tailored to specific datasets.

Motivated by the fact that the monocular depth estimation model mainly relies on the vertical image position and tends to overlook other useful cues, we propose a new pa-

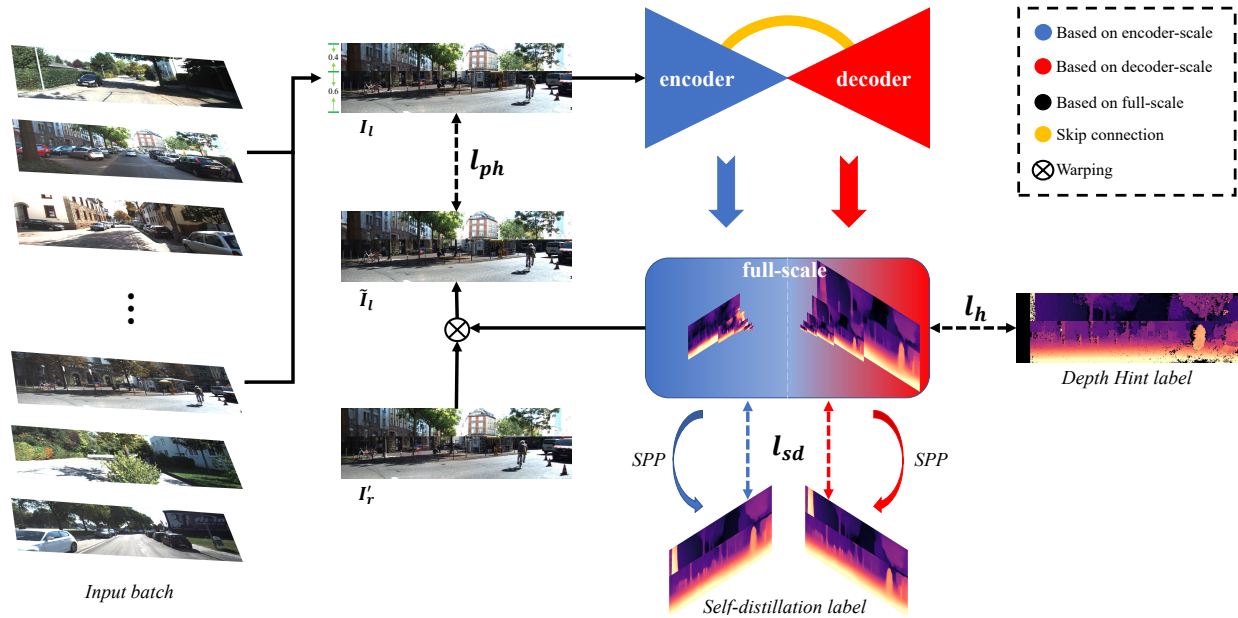


Figure 2. **Framework illustration.** The proposed approach is mainly composed of three procedures. The input batch data is first refactored by data grafting, and here we take the grafting ratio of 0.6 as an example. Immediately after that, the full-scale network will estimate the disparity map at all scales, which means that not only the decoder but also the encoder will infer the disparity. Finally, the full-scale disparity will be used to generate the self-distillation label through selective post-processing for the encoder and decoder scale separately and calculate the loss l_{sd} . Meanwhile, the model will be trained with the assistance of photometric loss l_{ph} and depth hint loss l_h , and it is worth noting that these losses are executed on all scales.

parameter learning free data augmentation method, called data grafting, to force the model to explore more cues.

3. Method

We adopt rectified stereo pairs as the input of our self-supervised model in training, while only a single image is required to infer depth at test time. This kind of self-supervised method is mainly divided into three steps. The model $\mathcal{F} : I \rightarrow d \in \mathbb{R}^{H \times W}$, that will first estimate the disparity map d , which represents the offset of the corresponding pixel between the stereo pair, from the target view $I \in \mathbb{R}^{C \times H \times W}$. Next, the model will be trained iteratively by minimizing the discrepancy between the target view and the view \tilde{I} reconstructed from the source view I' with differentiable warping $f_w(I', d)$. The photometric loss measured with the combination of SSIM [48] and L1 is often adopted to express the discrepancy between the target view and the reconstructed view just as:

$$l_{ph}(d) = l_{ph}(I, \tilde{I}) = \alpha \frac{1 - SSIM(I, \tilde{I})}{2} + \beta |I - \tilde{I}| \quad (1)$$

where $SSIM()$ is computed over a 3×3 pixel window, with $\alpha = 0.85$ and $\beta = 0.15$. Finally, the depth map $z \in \mathbb{R}^{H \times W}$ will be recovered from d , which is outputted by the trained model, with known baseline b and focal length f under formula $z = bf/d$.

In this section, we will introduce the main contributions of this paper in detail. The framework pipeline is just shown in Figure 2.

3.1. Data Grafting

Lack of data in both quantity and diversity is the first tricky obstacle faced by monocular depth estimation, which will damage the generalization ability of the model. One of the significant overfitting risks in MDE is the excessive dependence on the vertical image position as described in Sec. 1. Although data augmentation is the most cost-effective and ubiquitous solution, there is almost no relevant research on existing self-supervised MDE methods, and only some simple data perturbations such as horizontal flipping are used. The reason mainly lies in that self-supervised MDE methods generate supervisory signals based on the degree of matching between views, which requires strict pixel correspondence (epipolar constraint) to ensure that the matching error only comes from the estimated disparity. Obviously, the traditional data augmentation method will break this correspondence, thereby damaging the performance of the model as shown in our experiments in Sec. 5.2.

However, we note that this restriction is relaxed in the category with stereo pairs as input. Because the two views were taken with parallel cameras and rectified, the match between them will only occur in the horizontal direction,

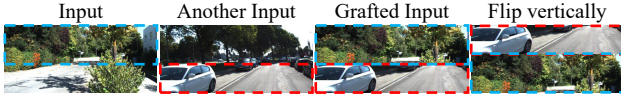


Figure 3. Illustration of data grafting.

e.g., panning left or right. Therefore, we can do perturbation in the vertical direction to augment our data.

To this end, we found that grafting two images with different semantics together can effectively alleviate the overfitting risk in MDE and encourage the model to better utilize the full context of the input without destroying the epipolar constraint. We conduct the data grafting within a mini-batch, and it is determined by two hyper-parameters: the grafting ratio r and the corresponding uniform probability p . We reconstruct the input by vertically grafting an area with a proportion of r from another input with the probability of p , and randomly flip these two parts vertically, as shown in Figure 3. Meanwhile, grafting is not only for the target view, but also for its corresponding Depth Hint, which will be introduced in Sec. 3.4, and the source view. But each grafting operation can only be performed between the same category, *e.g.*, both are target views. And the grafting config of all inputs in a batch is the same. The grafting detail of a single input is shown in Algorithm 1.

Algorithm 1: Data Grafting

Input: Input I^1 ; Another input of the same category randomly sampled from the same batch I^2 ; Shape of input (c, h, w) ; Random vertical flip factor $flip$.

Output: Grafted input I^1 .

- 1 Random sampling r from $\{0, 0.2, 0.4, 0.6, 0.8\}$ with the uniform probability of 0.2;
- 2 **if** $r = 0$ **then**
- 3 | **return** I^1 .
- 4 **else**
- 5 | $graft_h = Ceil(h \times r)$;
- 6 | $I^1[:, graft_h :, :] \leftarrow I^2[:, graft_h :, :]$;
- 7 | **if** $flip < 0.5$ **then**
- 8 | | $T = I^1$;
- 9 | | $I^1[:, h - graft_h :, :] \leftarrow T[:, : graft_h, :]$;
- 10 | | $I^1[:, : h - graft_h, :] \leftarrow T[:, graft_h :, :]$;
- 11 | **end**
- 12 **end**
- 13 **return** I^1 .

3.2. Full-scale Network

The coarse-to-fine strategy has been proven effective in MDE which continuously refines the estimation with iterative warping [10, 12, 13, 49]. The common practice is to output multi-scale disparity prediction in the decoder,

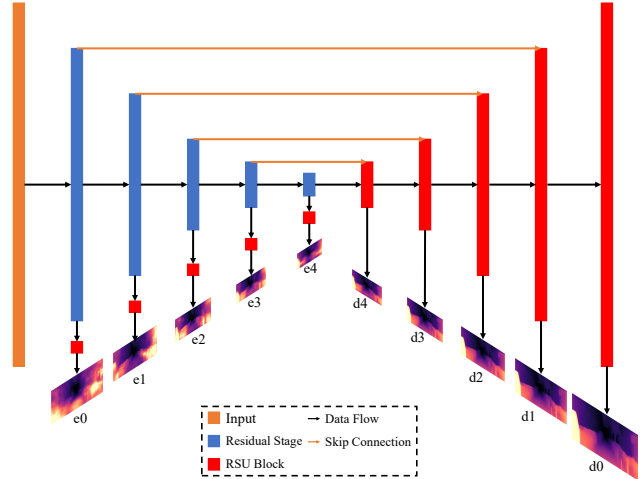


Figure 4. Full-scale network. The “ $e_0 \sim e_4$ ” represents the scale in the encoder and “ $d_0 \sim d_4$ ” represents the scale in the decoder. The spatial size of each scale increases with the decrease of serial number.

whose spatial size is incremental. In this scenario, the knowledge learned by the encoder is more abstract and general, while that in the decoder is more specific to the depth estimation task.

Intuitively, advancing the specialization of depth estimation to the encoder can give stronger constraints to the model and further improve its performance. Therefore, we extend the traditional multi-scale to full-scale, which means that we also add the multi-scale disparity prediction block to the encoder. Meanwhile, we insert a residual block, or more precisely an RSU block [39], between the prediction block and the residual stage in the encoder as the bridge to mitigate the impacts between different scales.

Furthermore, just as depicted in Figure 4, we adopt the RSU block, which is more powerful and more lightweight, to construct the decoder to draw the representational capacity of our full-scale network. After training, we can discard the encoder-scale or even part of the decoder-scale, and only retain the largest scale of the decoder, which means that the full-scale network will not bring more parameters or computation than the traditional network.

3.3. Self-distillation

Self-distillation is an effective way to generate more supervised signals for the model, and it is particularly important for self-supervised learning. Here, we propose selective post-processing to generate the self-distillation label, and with which we create a new loss, termed Self-Distillation Loss l_{sd} , for the model.

Selective Post-Processing aims to filter out the optimal disparity at each pixel from multiple disparity scales. Actually, the largest disparity map in the decoder that we of-

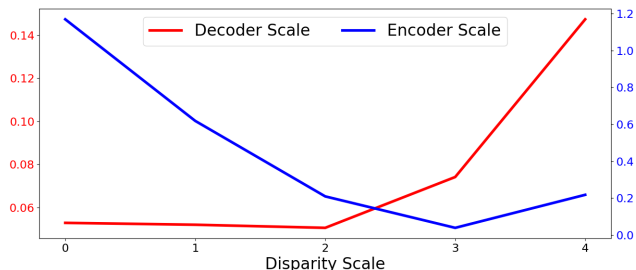


Figure 5. **Precision improvement statistics of SPP result on each scale for all test samples in Eigen split [6].**

ten output is not always the best at all pixels, as shown in Table 5. Maybe the “d0” scale is better at pixel a but the “d3” scale is better at pixel b . Hence, to distinguish the optimal scale at each pixel, we adopt the reconstruction error or the photometric loss as our criterion, which is inspired by [49]. Given the full-scale disparity maps $D = [d_{d0}, \dots, d_{d4}, d_{e0}, \dots, d_{e4}]$, we will calculate a reconstruction error map for each scale according to Equation (1). Then, the self-distillation label of encoder y_e and decoder y_d can be constructed based on the assumption that the smaller the error, the better the predicted disparity. The detailed procedure of SPP, which is the same between the encoder-scale and decoder-scale, is shown in Algorithm 2. The statistic result in Figure 5 shows that the SPP can get the most precise results.

Algorithm 2: Selective Post-Processing

Input: The target view I ; The source view I' ;
Multi-scale disparity maps D' .
Output: Self-distillation label y .

- 1 **Initialization:** $e_{min} = None$;
- 2 **for** d in D' **do**
- 3 Upsample d to the same size as I ;
- 4 Reconstruct target view $\tilde{I} = f_w(I', d)$;
- 5 Calculate the reconstruction error $e = l_{ph}(I, \tilde{I})$;
- 6 **if** $d = D'[0]$ **then**
- 7 $y = d$;
- 8 $e_{min} = e$;
- 9 **else**
- 10 Find all the pixels where $e < e_{min}$;
- 11 Update y with d at these pixels;
- 12 Update e_{min} with e at these pixels;
- 13 **end**
- 14 **end**
- 15 **return** y .

Self-Distillation Loss is the difference between the disparity map and the self-distillation label for each scale, and it can be modeled as:

$$l_{sd}(d) = \log(|y_{c(d)} - d| + 1) \quad (2)$$

where $c(\cdot)$ is used to determine whether d belongs to the decoder-scale or the encoder-scale.

3.4. Training Loss

Following [49], we incorporate the hint loss that has been proven effective for thin structures into our model. The Depth Hint h is generated by the Semi-Global Matching (SGM) algorithm [19, 20] and be consulted only when the reconstruction error can be improved upon. It can be formulated for pixel i in each scale as:

$$l_h(d_i) = \begin{cases} \log(|h_i - d_i| + 1), & \text{if } l_{ph}(I, \tilde{I}_h)_i < l_{ph}(I, \tilde{I})_i \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

where \tilde{I}_h denotes the reconstructed view with Hint h .

Therefore, the final training loss is composed of the average of the three items of photometric loss, self-distillation loss and hint loss at each scale:

$$l = \frac{1}{|D|} \sum_{d \in D} (l_{ph}(d) + l_{sd}(d) + l_h(d)) \quad (4)$$

4. Implementation Details

We implement our model in PyTorch [35]. The procedure of calculating Depth Hint is the same as that of [49]. We use Adam [25] optimizer with the base learning rate of $1e-4$ and train the joint loss for 20 epochs. Besides our new data augmentation approach, we adopted the preprocessing techniques in [13]. In data grafting, we found that the grafting ratio $r = 0.2 \times n$, where $n \in \mathbb{N}$ and $r < 1$, can get the best effect, as shown in Algorithm 1. Unless otherwise specified, we take ResNet-18 which is pre-trained on ImageNet [23] as the encoder and resize the input to 320×1024 . As for the RSU block [39], we remove the Batch Normalization layer [22] and replace the ReLU [33] with ELU [3] activation. More specifically, we take RSU3 \sim RSU7 to construct the decoder’s layers and the encoder’s bridges from minimum scale to maximum scale respectively.

5. Experiments

We first verify the performance of our model on the KITTI dataset [11], and perform a comprehensive ablation study on each component. Finally, the generalization ability of our model is validated on the NYU-Depth-v2 dataset [43].

KITTI Stereo was recorded from a driving vehicle and contains 42,382 rectified stereo pairs from 61 scenes. To ensure the objectivity of comparison, we utilize the Eigen split [6], which is composed of 22,600 training image pairs in 32 scenes, and 697 test pairs in other 29 scenes. We report all seven of the standard metrics [7] with Garg’s crop [10] and a standard distance cap of 80 meters [12].

Method	PP	Data	$H \times W$	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Eigen <i>et al.</i> [7]		D	184 × 612	0.203	1.548	6.307	0.282	0.702	0.890	0.890
Kuznetsov <i>et al.</i> [28]		DS	187 × 621	0.113	0.741	4.621	0.189	0.862	0.960	0.986
Yang <i>et al.</i> [53]	✓	D [†] S	256 × 512	0.097	0.734	4.442	0.187	0.888	0.958	0.980
Luo <i>et al.</i> [31]		D*DS	192 × 640 crop	0.094	0.626	4.252	0.177	0.891	0.965	0.984
Fu <i>et al.</i> [8]		D	385 × 513 crop	0.099	0.593	3.714	0.161	0.897	0.966	0.986
Lee <i>et al.</i> [30]		D	352 × 1216	0.091	0.555	4.033	0.174	0.904	0.967	0.984
Zhan <i>et al.</i> [62]		MS	160 × 608	0.135	1.132	5.585	0.229	0.820	0.933	0.971
Godard <i>et al.</i> [13]	✓	MS	320 × 1024	0.104	0.775	4.562	0.191	0.878	0.959	0.981
Watson <i>et al.</i> [49]	✓	MS	320 × 1024	0.098	0.702	4.398	0.183	0.887	0.963	0.983
Shu <i>et al.</i> [42]		MS	320 × 1024	0.099	0.697	4.427	0.184	0.889	0.963	0.982
Lyu <i>et al.</i> [32]		MS	320 × 1024	0.101	0.716	4.395	0.179	0.899	0.966	0.983
Garg <i>et al.</i> [10]		S	188 × 620	0.169	1.080	5.104	0.273	0.740	0.904	0.962
Godard <i>et al.</i> [12]	✓	S	256 × 512	0.138	1.186	5.650	0.234	0.813	0.930	0.969
Wong <i>et al.</i> [50]		S	256 × 512	0.133	1.126	5.515	0.231	0.826	0.934	0.969
Pilzer <i>et al.</i> [38] Teacher		S	256 × 512	0.098	0.831	4.656	0.202	0.882	0.948	0.973
Chen <i>et al.</i> [2]	✓	SC	256 × 512	0.118	0.905	5.096	0.211	0.839	0.945	0.977
Godard <i>et al.</i> [13]	✓	S	192 × 640	0.108	0.842	4.891	0.207	0.866	0.949	0.976
Watson <i>et al.</i> [49]	✓	S	192 × 640	0.106	0.780	4.695	0.193	0.875	0.958	0.980
Ours	✓	S	192 × 640	0.099	0.754	4.490	0.183	0.888	0.963	0.982
Pillai <i>et al.</i> [37]	✓	S	384 × 1024	0.112	0.875	4.958	0.207	0.852	0.947	0.977
Godard <i>et al.</i> [13]	✓	S	320 × 1024	0.105	0.822	4.692	0.199	0.876	0.954	0.977
Watson <i>et al.</i> [49]	✓	S	320 × 1024	0.099	0.723	4.445	0.187	0.886	0.962	0.981
Zhu <i>et al.</i> [68] Finetuned	✓	SC [†]	320 × 1024	0.097	0.675	4.350	0.180	0.890	0.964	0.983
Ours	✓	S	320 × 1024	0.093	0.671	4.297	0.178	0.899	0.965	0.983
Watson <i>et al.</i> [49] ResNet50	✓	S	320 × 1024	0.096	0.710	4.393	0.185	0.890	0.962	0.981
Zhu <i>et al.</i> [68] Finetuned ResNet50	✓	SC [†]	320 × 1024	0.091	0.646	4.244	0.177	0.898	0.966	0.983
Ours ResNet50	✓	S	320 × 1024	0.091	0.646	4.207	0.176	0.901	0.966	0.983

Table 1. **Quantitative results on the KITTI dataset [11] using the split of Eigen *et al.* [6].** Best results in each category are in **bold**. For **red** metrics, lower is better. And higher is better for **blue** metrics. Abbreviation in Data column: D refers to methods that are supervised by the ground truth depth, D[†] use auxiliary depth supervision from SLAM, D* use auxiliary depth supervision from synthetic depth labels, C for supervision from segmentation labels, C[†] for supervision from predicted segmentation labels, S refers to the supervision from stereo images and M for models trained by monocular video. PP represents post-processing [12]. The underlined model is our baseline. We annotate all the methods that use extra tricks, *e.g.*, fine-tuning and teacher model.

NYU-Depth-v2 was captured with a Microsoft Kinect sensor and consists of a total 582 indoor scenes. We validate our model on the official test set using the same standard metrics as in KITTI.

5.1. Depth Estimation Performance

We conduct a comprehensive comparison with multifarious methods on the KITTI benchmark to verify our depth estimation performance. First of all, we need to emphasize that our model is only trained on KITTI stereo data and is trick-free. We compare our approach with the recent self-, semi- and fully-supervised monocular depth estimation methods in Table 1. And the results show that our approach outperforms all existing self-supervised methods on all metrics and even some of the fully-supervised methods. Our approach of training only on stereo pairs improves 0.013 on the $\delta < 1.25$ compared to our baseline model [49], and this improvement is 225% ($= \frac{0.899-0.886}{0.890-0.886} - 1$) higher than that of [68], which has finetuned the model and used additional constraints. Furthermore, our method is not only outstanding in the category trained with stereo images, but also has a major advantage in the category of methods trained with stereo video (MS). Even if compared with the best score of each metric in the MS category, our approach won out in

most metrics. Moreover, we have done more experiments on low-resolution and complex backbones to demonstrate the generality and robustness of our model, and the consistent performance improvement obtained just proves it. It’s worth noting that we have further reduced the gap between full-supervision and self-supervision by nearly 79% ($= 1 - \frac{0.904-0.901}{0.904-0.890}$) compared to our baseline [49]. Besides, the qualitative results in Figure 6 show that our model predicts more accurately in challenging areas.

While our model significantly improves the performance of the baseline, it also retains the advantages of simple implementation. Each plug-and-play improvement can be easily integrated into other models, which is critical for future in-depth studies of monocular depth estimation.

5.2. Ablation Studies

We perform ablation analysis on the KITTI. The results in Table 6 show that our full model combining all components has leading performance, and the baseline model, without any of our contributions, performs the worst.

Benefits of data grafting. With data grafting, we can implicitly increase the amount of data by $1/p$ times on the basis of our baseline. The results in Table 2 show that only 20% of the data is used to obtain competitive performance

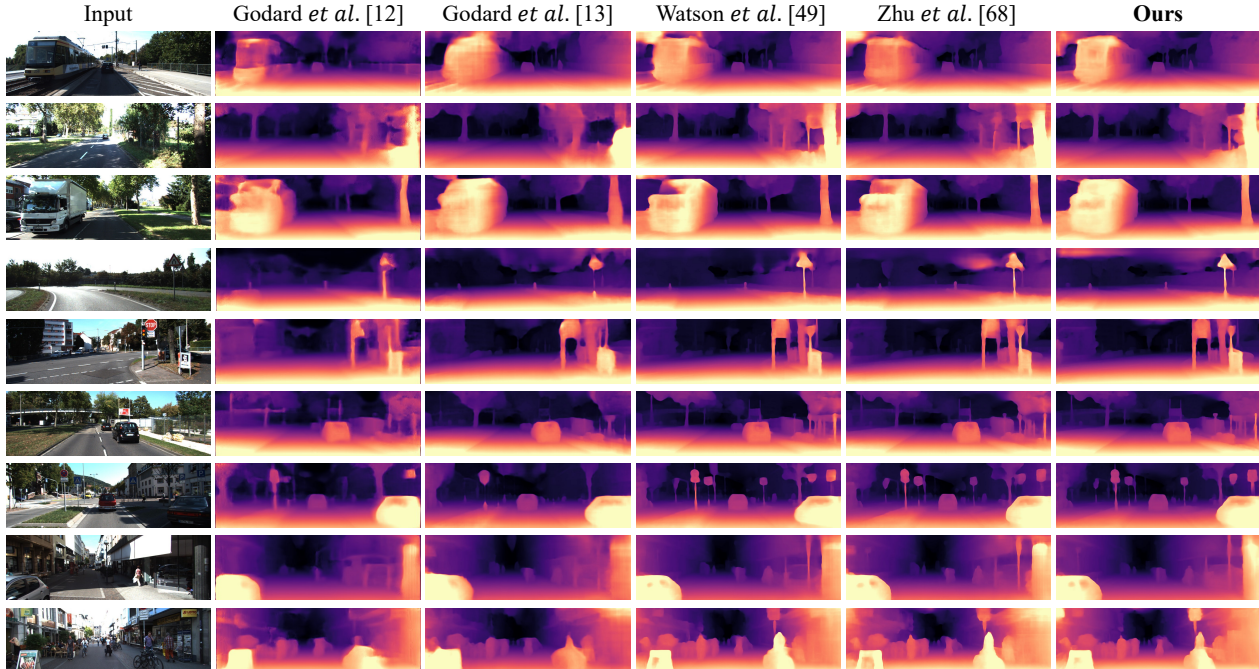


Figure 6. **Qualitative results.** Our model (**EPCDepth**) in the last column produces the most accurate and sharpest results, especially in challenging areas, *e.g.*, tree trunks, cars, *etc.*

Data Amount	Abs Rel	Sq Rel	RMSE	$\delta < 1.25$
w/o DG 100%	0.096	0.696	4.368	0.892
Full 20%	0.098	0.696	4.344	0.890
Full 50%	0.096	0.683	4.305	0.896
Full 100%	0.093	0.671	4.297	0.899

Table 2. **Ablation study on training data amount.** DG refers to data grafting. And the % means the percentage of data amount.

Augmentation	Abs Rel	Sq Rel	RMSE	$\delta < 1.25$
RandErasing [65]	0.115	0.992	4.987	0.858
Cutout [5]	0.106	0.830	4.753	0.874
CutMix [59]	0.105	0.831	4.752	0.876
DataGrafting	0.102	0.782	4.581	0.883

Table 3. **Comparison against other similar augmentation methods.** And the input size is 192×640 .

to the model without data grafting under 100% of the data, which just verifies the strong generalization ability of our model. Moreover, we make a comparison with other similar augmentation methods to demonstrate our effectiveness in Table 3. The result just shows that breaking the relationship between the depth and the vertical image position with a certain probability, which is the uniqueness of data grafting, can allow the model to potentially grasp more effective cues. The unsatisfaction of other methods may lie in the lack of regularization ability for the vertical image position

Source	Abs Rel	Sq Rel	RMSE	$\delta < 1.25$
PP	0.094	0.680	4.320	0.898
SPP	0.094	0.675	4.312	0.899
SPP separate	0.093	0.671	4.297	0.899

Table 4. **Ablation study on distillation source.** PP refers to the post-processing result of the largest scale in the decoder.

Scale	Abs Rel	Sq Rel	RMSE	$\delta < 1.25$	$\delta < 1.25^2$
d0	0.0925	0.671	4.297	0.899	0.965
d1	0.0922	0.668	4.292	0.899	0.965
d2	0.092	0.655	4.268	0.898	0.965

Table 5. **Quantitative results of different scales.**

and the damage of the epipolar constraint between views at the edge of the hole. Meanwhile, we conduct a sensitivity experiment on the grafting ratio. The results in Figure 7 show that the odd setting (*e.g.* $n/3$) is generally better than even setting (*e.g.* $n/2$), and performs best when $r = n/5$, which indicates that the grafting result holding a piece of dominant semantic information is more effective.

Benefits of self-distillation. From Table 5, we expect to select the optimal scale for each pixel through selective post-processing. The comparison results in Table 4 between different label generation methods show that the SPP can get more stable improvement and distilling encoder and decoder separately is more effective. Meanwhile, we no-

Method	DG	SD	FS	HR	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Baseline					0.107	0.848	4.745	0.194	0.875	0.957	0.980
Baseline + DG	✓				0.102	0.782	4.581	0.188	0.883	0.960	0.981
Baseline + SD		✓			0.105	0.822	4.708	0.193	0.876	0.958	0.981
Baseline + FS			✓		0.103	0.785	4.628	0.189	0.881	0.960	0.981
Baseline HR				✓	0.101	0.758	4.497	0.187	0.886	0.962	0.982
Baseline HR + DG	✓			✓	0.098	0.694	4.371	0.182	0.890	0.963	0.983
Baseline HR + SD		✓		✓	0.099	0.744	4.465	0.186	0.888	0.962	0.982
Baseline HR + FS			✓	✓	0.097	0.701	4.364	0.182	0.892	0.963	0.982
Full HR w/o FS	✓	✓		✓	0.098	0.702	4.377	0.184	0.888	0.963	0.983
Full HR w/o SD	✓		✓	✓	0.094	0.678	4.312	0.180	0.898	0.965	0.982
Full HR w/o DG		✓	✓	✓	0.096	0.696	4.368	0.182	0.892	0.963	0.982
Full HR	✓	✓	✓	✓	0.093	0.671	4.297	0.178	0.899	0.965	0.983

Table 6. Ablation results amongst variants of our model (EPCDepth) on the KITTI dataset. DG refers to data grafting, SD refers to self-distillation, FS refers to full-scale and HR refers to high resolution.

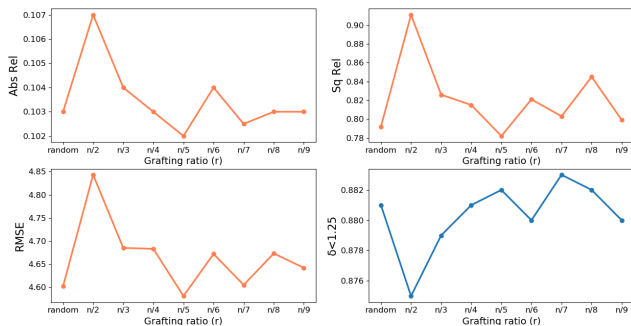


Figure 7. Sensitivity analysis of grafting ratio r . The smaller the value, the better in the red line chart, and the worse in the blue.

Full-Scale	Abs Rel	Sq Rel	RMSE	$\delta < 1.25$
+ Encoder Scale	0.105	0.811	4.668	0.877
+ Bridges	0.104	0.798	4.655	0.878
+ RSU	0.103	0.785	4.628	0.881

Table 7. Ablation study on full-scale network. Conducted by continuously accumulating each module with input size 192×640 .

ticed that the magnitude of its performance improvement is minimal and is affected by the capacity of the model. But we hope that our exploration can open the door to self-distillation in this regression task.

Benefits of full-scale network. Our full-scale network draws on some advantages of the multi-generation strategy, that is to impose more constraints on the model, and the results in Table 6 just prove its power. Furthermore, we explore the effectiveness of the encoder scale, RSU blocks and the encoder’s bridges respectively, by ablating their effects in Table 7. Note that each experiment is carried out on the basis of the previous experiment. The continuous performance improvement of each module proves their effectiveness. Meanwhile, our full-scale network achieves superior performance with **9.88 GFLOPS** at test time, compared to **10.1 GFLOPS** of the traditional network [13, 49, 68].

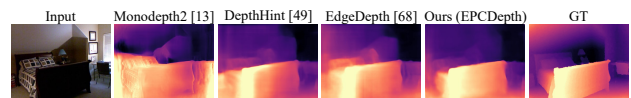


Figure 8. Qualitative results on the NYU-Depth-v2 dataset.

Method	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Monodepth2 [13]	0.355	0.673	1.252	0.373	0.485	0.771	0.907
DepthHint [49]	0.298	0.457	1.043	0.331	0.539	0.821	0.937
EdgeDepth [68]	0.292	0.437	1.018	0.319	0.563	0.834	0.941
Ours (EPCDepth)	0.247	0.277	0.818	0.285	0.605	0.869	0.961

Table 8. Quantitative results on the NYU-Depth-v2 dataset.

5.3. Generalizing to NYU-Depth-v2

Since there are no stereo pairs in NYU-Depth-v2 dataset, we train on the KITTI dataset and then test on it just as Monodepth [12] did on Make3D. The preprocessing strategy we adopt is the same as that of [58], and median scaling is applied for all models. The results shown in Table 8 and Figure 8 just verify our strong generalization ability.

6. Conclusion

We extracted the potential capacity of self-supervised monocular depth estimation through our novel data augmentation method, exploratory self-distillation and efficient full-scale network. The experiments demonstrate that our model (EPCDepth) can yield the best performance with the least computational cost. In future work, we will try to further improve the performance of self-distillation by exploring more accurate label generation methods. Besides, applying our contributions to other categories, *e.g.*, M, MS and even supervised method, is also a potential direction.

Acknowledgements. Thanks to National Natural Science Foundation of China 61672063 and 62072013, Shenzhen Research Projects of JCYJ20180503182128089, 201806080921419290 and RCJC20200714114435057. In addition, we thank the anonymous reviewers for their valuable comments.

References

- [1] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *CVPR*, pages 4009–4018, 2021.
- [2] Po Yi Chen, Alexander H. Liu, Yen Cheng Liu, and Yu Chiang Frank Wang. Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation. In *CVPR*, pages 2624–2632, 2019.
- [3] Djork Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (ELUs). In *ICLR*, 2016.
- [4] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *CVPR*, pages 113–123, 2019.
- [5] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- [6] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *ICCV*, pages 2650–2658, 2015.
- [7] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NeurIPS*, pages 2366–2374, 2014.
- [8] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *CVPR*, pages 2002–2011, 2018.
- [9] Tommaso Furlanello, Zachary C Lipton, Michael Tschanen, Laurent Itti, and Anima Anandkumar. Born-again neural networks. In *ICML*, pages 1607–1616, 2018.
- [10] Ravi Garg, B. G. Vijay Kumar, Gustavo Carneiro, and Ian Reid. Unsupervised CNN for single view depth estimation: Geometry to the rescue. In *ECCV*, pages 740–756, 2016.
- [11] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *CVPR*, pages 3354–3361, 2012.
- [12] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, pages 270–279, 2017.
- [13] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel Brostow. Digging into self-supervised monocular depth estimation. In *ICCV*, pages 3828–3838, 2019.
- [14] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raveentos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *CVPR*, pages 2485–2494, 2020.
- [15] Vitor Guizilini, Rui Hou, Jie Li, Ambrus Rares, and Adrien Gaidon. Semantically-guided representation learning for self-supervised monocular depth. In *ICLR*, pages 1–14, 2020.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [17] Byeongho Heo, Jeesoo Kim, Sangdoon Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation. In *ICCV*, pages 1921–1930, 2019.
- [18] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv*, pages 1–9, 2015.
- [19] Heiko Hirschmüller. Accurate and efficient stereo processing by semi-global matching and mutual information. In *CVPR*, pages 807–814, 2005.
- [20] Heiko Hirschmüller. Stereo processing by semiglobal matching and mutual information. *IEEE TPAMI*, pages 328–341, 2008.
- [21] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning lightweight lane detection CNNs by self-attention distillation. In *ICCV*, pages 1013–1021, 2019.
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456, 2015.
- [23] Jia Deng, Wei Dong, R. Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.
- [24] Kyungyul Kim, Byeong Moon Ji, Doyoung Yoon, and Sangheum Hwang. Self-knowledge distillation: A simple way for better generalization. In *NeurIPS*, pages 1–15, 2020.
- [25] Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [26] Marvin Klingner, Jan-Aike Termöhlen, Jonas Mikolajczyk, and Tim Fingscheidt. Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance. In *ECCV*, pages 582–600, 2020.
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, pages 1097–1105, 2012.
- [28] Yevhen Kuznietsov, Jörg Stückler, and Bastian Leibe. Semi-supervised deep learning for monocular depth map prediction. In *CVPR*, pages 6647–6655, 2017.
- [29] Katrin Lasinger, Rene Ranftl, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE TPAMI*, pages 1–14, 2020.
- [30] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv preprint arXiv:1907.10326*, 2019.
- [31] Yue Luo, Jimmy Ren, Mude Lin, Jiahao Pang, Wenxiu Sun, Hongsheng Li, and Liang Lin. Single view stereo matching. In *CVPR*, pages 155–163, 2018.
- [32] Xiaoyang Lyu, Liang Liu, Mengmeng Wang, Xin Kong, Lina Liu, Yong Liu, Xinxin Chen, and Yi Yuan. HR-Depth: High resolution self-supervised monocular depth estimation. In *AAAI*, pages 2294–2301, 2021.
- [33] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve Restricted Boltzmann machines. In *ICML*, pages 807–814, 2010.
- [34] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *CVPR*, pages 3967–3976, 2019.
- [35] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.

- [36] Baoyun Peng, Xiao Jin, Dongsheng Li, Shunfeng Zhou, Yichao Wu, Jiaheng Liu, Zhaoning Zhang, and Yu Liu. Correlation congruence for knowledge distillation. In *ICCV*, pages 5007–5016, 2019.
- [37] Sudeep Pillai, Rareş Ambruş, and Adrien Gaidon. SuperDepth: Self-supervised, super-resolved monocular depth estimation. In *ICRA*, pages 9250–9256, 2019.
- [38] Andrea Pilzer, Stephane Lathuiliere, Nicu Sebe, and Elisa Ricci. Refine and distill: Exploiting cycle-inconsistency and knowledge distillation for unsupervised monocular depth estimation. In *CVPR*, pages 9768–9777, 2019.
- [39] Xuebin Qin, Zichen Zhang, Chenyang Huang, Masood Dehghan, Osmar R. Zaiane, and Martin Jagersand. U2-Net: Going deeper with nested U-structure for salient object detection. *PR*, page 107404, 2020.
- [40] Anurag Ranjan, Varun Jampani, Lukas Balles, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J. Black. Competitive collaboration: Joint unsupervised learning of depth, camera motion, optical flow and motion segmentation. In *CVPR*, pages 12240–12249, 2019.
- [41] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. FitNets: Hints for thin deep nets. In *ICLR*, pages 1–13, 2015.
- [42] Chang Shu, Kun Yu, Zhixiang Duan, and Kuiyuan Yang. Feature-metric loss for self-supervised learning of depth and egomotion. In *ECCV*, pages 572–588, 2020.
- [43] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, pages 746–760, 2012.
- [44] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [45] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, pages 1–9, 2015.
- [46] Toan Tran, Trung Pham, Gustavo Carneiro, Lyle Palmer, and Ian Reid. A Bayesian data augmentation approach for learning deep models. In *NeurIPS*, pages 1–10, 2017.
- [47] Tom van Dijk and Guido C. H. E. de Croon. How do neural networks see depth in single images? In *ICCV*, pages 2183–2191, 2019.
- [48] Zhou Wang, Alan Conrad Bovik, Hamid Rahim Sheikh, and Eero P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE TIP*, pages 600–612, 2004.
- [49] Jamie Watson, Michael Firman, Gabriel J. Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *ICCV*, pages 2162–2171, 2019.
- [50] Alex Wong and Stefano Soatto. Bilateral cyclic constraint and adaptive regularization for unsupervised monocular depth prediction. In *CVPR*, pages 5644–5653, 2019.
- [51] Ting-bing Xu and Cheng-lin Liu. Data-distortion guided self-distillation for deep neural networks. In *AAAI*, pages 5565–5572, 2019.
- [52] Chenglin Yang, Lingxi Xie, Chi Su, and Alan L. Yuille. Snapshot distillation: Teacher-student optimization in one generation. In *CVPR*, pages 2859–2868, 2019.
- [53] Nan Yang, Rui Wang, Jörg Stückler, and Daniel Cremers. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *ECCV*, pages 817–833, 2018.
- [54] Zhenheng Yang, Peng Wang, Wei Xu, Liang Zhao, and Ramakant Nevatia. Unsupervised learning of geometry from videos with edge-aware depth-normal consistency. In *AAAI*, pages 7493–7500, 2018.
- [55] Hongwei Yi, Zizhuang Wei, Mingyu Ding, Runze Zhang, Yisong Chen, Guoping Wang, and Yu-Wing Tai. Pyramid multi-view stereo net with self-adaptive view aggregation. In *ECCV*, pages 766–782, 2020.
- [56] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *CVPR*, pages 4133–4141, 2017.
- [57] Zhichao Yin and Jianping Shi. GeoNet: Unsupervised learning of dense depth, optical flow and camera pose. In *CVPR*, pages 1983–1992, 2018.
- [58] Zehao Yu, Lei Jin, and Shenghua Gao. P2Net: Patch-match and plane-regularization for unsupervised indoor depth estimation. In *ECCV*, pages 206–222, 2020.
- [59] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. CutMix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, pages 6023–6032, 2019.
- [60] Sukmin Yun, Jongjin Park, Kimin Lee, and Jinwoo Shin. Regularizing class-wise predictions via self-knowledge distillation. In *CVPR*, pages 13876–13885, 2020.
- [61] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR*, pages 1–13, 2017.
- [62] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian M. Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *CVPR*, pages 340–349, 2018.
- [63] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. Mixup: Beyond empirical risk minimization. In *ICLR*, pages 1–13, 2017.
- [64] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *ICCV*, pages 3713–3722, 2019.
- [65] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, pages 13001–13008, 2020.
- [66] Junsheng Zhou, Yuwang Wang, Kaihuai Qin, and Wenjun Zeng. Unsupervised high-resolution depth learning from videos with dual networks. In *ICCV*, pages 6872–6881, 2019.
- [67] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, pages 1851–1858, 2017.
- [68] Shengjie Zhu, Garrick Brazil, and Xiaoming Liu. The edge of depth: Explicit constraints between segmentation and depth. In *CVPR*, pages 13116–13125, 2020.