

Channel-Wise Attention-Based Network for Self-Supervised Monocular Depth Estimation

Jiaxing Yan Hong Zhao Penghui Bu YuSheng Jin
 State Key Laboratory for Manufacturing System Engineering,
 School of Mechanical Engineering, Xi'an Jiaotong University, China

{yanjxedu@stu., zhaohong@mail., bph19891027@stu., yushengJ@stu.}@xjtu.edu.cn

Abstract

Self-supervised learning has shown very promising results for monocular depth estimation. Scene structure and local details both are significant clues for high-quality depth estimation. Recent works suffer from the lack of explicit modeling of scene structure and proper handling of details information, which leads to a performance bottleneck and blurry artefacts in predicted results. In this paper, we propose the Channel-wise Attention-based Depth Estimation Network (CADepth-Net) with two effective contributions: 1) The structure perception module employs the self-attention mechanism to capture long-range dependencies and aggregates discriminative features in channel dimensions, explicitly enhances the perception of scene structure, obtains the better scene understanding and rich feature representation. 2) The detail emphasis module re-calibrates channel-wise feature maps and selectively emphasizes the informative features, aiming to highlight crucial local details information and fuse different level features more efficiently, resulting in more precise and sharper depth prediction. Furthermore, the extensive experiments validate the effectiveness of our method and show that our model achieves the state-of-the-art results on the KITTI benchmark and Make3D datasets.

1. Introduction

Accurate depth estimation from a single image is a fundamental task in computer vision. High quality depth information can provide useful cues for various fields, including robotics navigation [4], autonomous driving [33] and augmented reality [34]. Recently, the fully-supervised methods [5, 6, 7, 15] for monocular depth estimation have produced outstanding results, while they need large numbers of accurate ground truth which could only be sparsely collected from expensive LiDAR sensors [9]. As an attractive alternative, self-supervised methods can alleviate this lim-

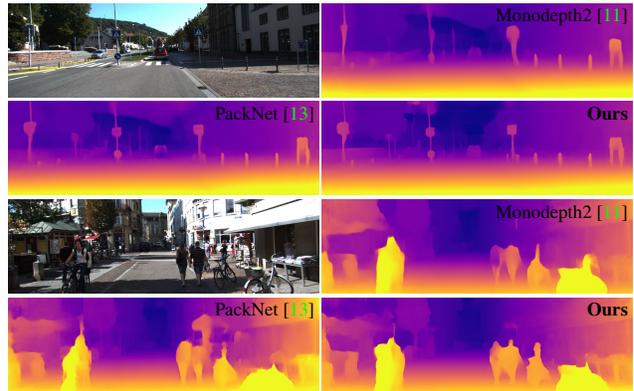


Figure 1. **Depth prediction from a single image.** Our proposed CADepth-Net produces more precise and sharper depth estimation, especially for thin structures *e.g.* road signs and pedestrians.

itation, as they use geometrical constraints on monocular video [56] or synchronized stereo image pairs [10] as the sole source of supervision.

In depth estimation, the most important information is the scene structure aiming at accurately obtaining the overall structure and relative depth information in 3D space. Most previous works [10, 56] just simply use convolutional neural networks to extract semantic features of input images and implicitly learn the structural information of the scene. However, the lack of explicit exploration of the robust representation of 3D scene geometry leads to an incomplete perception of overall layout for the complex scenes.

Local detail is another critical feature that focuses on object boundaries and attempts to generate sharp depth maps. Most depth estimation networks are based on the U-Net [37] framework, and the decoder simply leverages concatenation and a basic convolution to fuse high-level and low-level features. We found that these operations can not preserve sufficient details or precisely recover spatial information, leading to inefficient integration of different levels features and blurry artefacts at the depth discontinuous regions.

To address the above problems and efficiently handle

both the overall structure and local details, we propose a novel Channel-wise Attention-based Depth Estimation Network with *structure perception module* and *detail emphasis module*. To better understand the 3D structure of the whole scene, we perform an in-depth analysis of semantic information in the monocular depth estimation task and conclude that each high-level feature map can be regarded as a region-specific response. Based on this observation, we provide the structure perception module to capture more contextual information of scene geometry and enhance the feature representations. Specifically, we first employ the self-attention mechanism to capture the long-range dependencies between any two channel maps, then each feature map is updated via aggregating features from all channel maps by weighted summation and fuse different local depth responses from non-contiguous regions. To generate sharper object boundaries, instead of using the above naive operations, we propose the detail emphasis module employing channel attention mechanism to re-calibrate channel-wise features and emphasize the specific semantic information. Specifically, we sequentially adopt the detail emphasis module at different scales in the decoding stage to highlight features containing crucial local details (*e.g.* object boundaries information). To summarize our contributions in this paper:

- We introduce a novel Channel-wise Attention-based Depth Network (CADepth-Net) for self-supervised monocular depth estimation employing two channel-wise attention modules to perform the information aggregation and feature re-calibration respectively.
- We propose structure perception module utilizing self-attention mechanism to obtain rich context of scene structure and better feature representation.
- We carefully design the detail emphasis module with channel attention mechanism to efficiently fuse different scale features and emphasize important details for sharper depth estimation.
- We conduct extensive experiments on KITTI and Make3D datasets, demonstrating our model significantly outperforms existing methods and achieve the state-of-the-art results on KITTI benchmark.

2. Related Work

2.1. Supervised Depth Estimation

Estimating depth from a single image is an inherently ill-posed problem as pixels in the image can have multiple plausible depths. Recently, fully supervised methods had shown the capacity of fitting predictive models to estimate depth from color input images correctly. Eigen *et al.* [6] produced dense pixel depth estimates by utilizing the multi-scale neural networks, one that estimated a coarse global

depth prediction and another locally refined prediction produced by the first network. Rapidly, various fully supervised methods based on deep learning had been continuously explored [7, 5, 24, 27]. However, all the above methods required high-quality ground truth depth, which can be costly to obtain.

2.2. Self-supervised Monocular Depth Estimation

To overcome the limitation of supervised approaches, self-supervised methods unified depth estimation and ego-motion estimation into one framework with view synthesis as supervision signal. SfMLearner introduced by Zhou *et al.* [56] simultaneously learned depth and ego-motion from monocular video by training a depth estimation network along with a separate pose network. Furthermore, Yin *et al.* [50] decomposed scene motion into rigid and non-rigid parts to account for object motion. Wang *et al.* [44] incorporated Direct Visual Odometry to estimate the relative camera pose. [31] proposed 3D constraints loss to enforce consistency of the estimated depth and ego-motion across consecutive frames. Guizilini *et al.* [13] learned to compress and decompress detail-preserving representations by symmetrical packing and unpacking blocks. Other published methods were based upon edge and normal [48, 49], Competitive Collaboration [36], semantic segmentation [14, 22] and feature representations learning [41]. A state-of-the-art framework was Monodepth2 proposed by Godard *et al.* [11], which introduced a minimum re-projection loss to deal with occlusions and auto-masking scheme removing invalid pixels robustly. Our model is based on Monodepth2 extended with our contributions.

2.3. Self-attention Mechanism

Self-attention mechanism had been widely used to capture long-range dependencies in various tasks. The Transformer [43] was the first work that proposed the self-attention mechanism to handle long-range dependencies between words in machine translation. Wang *et al.* [45] modeled the spatial-temporal dependencies in video sequences and images via aggregating query-specific global context to each query position. Zhang *et al.* [52] incorporated the self-attention mechanism into the GAN framework and learned a better image generator. Fu *et al.* [8] enhanced the ability of feature representations for scene segmentation by designing two types of attention modules. For the monocular depth estimation task, Johnston *et al.* [18] captured the context of similar disparity values at non-contiguous regions by exploring the feature similarity at spatial dimensions. Unlike previous works, we demonstrate that capturing global dependencies along the channel dimensions and aggregating discriminative features will achieve better performance for depth estimation, as each channel map gains more relative depth information from the distant regions.

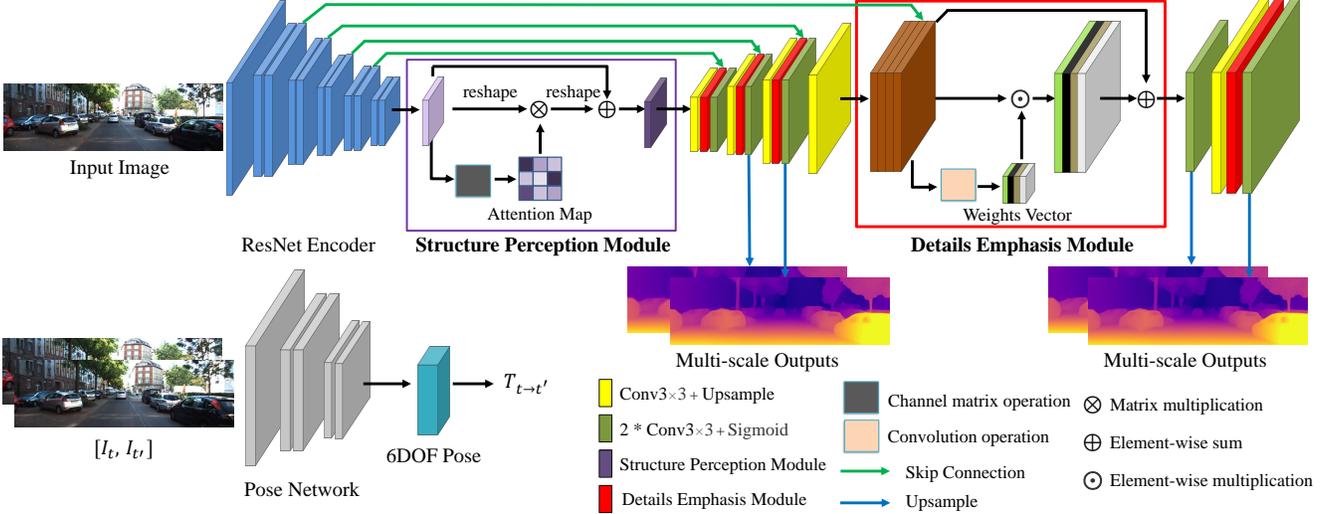


Figure 2. **Overview of Framework.** Our proposed CADepth-Net is a fully convolutional U-Net architecture. We first use a ResNet encoder to extract semantic features and input them to the **Structure Perception Module**. We perform the matrix multiplication between input features and attention maps to generate aggregated features. The low-resolution feature maps are passed through successive blocks of UpConv(upsample + convolution), as well as the **Detail Emphasis Module** which computes a weight vector to re-calibrate channel-wise features. Finally, we upsample the predicted disparities at multiple scales to original input resolutions. Besides, the pose network takes temporally adjacent images $I_t, I_{t'}$ as input and outputs relative pose $T_{t \rightarrow t'}$.

3. Method

In this section, we firstly review the training methods for self-supervised monocular depth estimation, then introduce the architecture of our Channel-wise Attention-based Network, finally describe our main contributions, the structure perception module and detail emphasis module.

3.1. Self-Supervised Training

The goal of self-supervised monocular depth estimation is to predict the depth map from a single RGB image without ground truth. Specifically, given a single input image I_t , the depth network predicts its corresponding depth map D_t , then the pose network takes temporally adjacent images as input and predicts relative pose $T_{t \rightarrow t'}$ between the target image I_t and source images $I_{t'}, t' \in \{t-1, t+1\}$, finally we use the predicted D_t and $T_{t \rightarrow t'}$ to perform view synthesis as the supervisory signal.

At training time, both the depth network and pose network are optimized jointly by minimizing the per-pixel minimum photometric re-projection error L_p [11]

$$L_p = \min_{t'} pe(I_t, I_{t' \rightarrow t}), \quad (1)$$

where $pe()$ denotes the photometric error which consists of L1 and the Structural Similarity (SSIM) [46], $I_{t' \rightarrow t}$ is the warped result from $I_{t'}$ to I_t as in

$$pe = \frac{\alpha}{2} (1 - \text{SSIM}(I_t, I_{t' \rightarrow t})) + (1 - \alpha) \|I_t - I_{t' \rightarrow t}\|_1, \quad (2)$$

$$I_{t' \rightarrow t} = I_{t'} \langle \text{proj}(D_t, T_{t \rightarrow t'}, K) \rangle, \quad (3)$$

where $\text{proj}()$ represents the resulting 2D coordinates of the projected depths D_t in $I_{t'}$ and $\langle \rangle$ is the sampling operator. We use the differentiable bilinear sampling mechanism proposed in the STN [17] to sample the source images.

In a real-world scenario, situations like stationary camera and moving objects will break down the assumptions of a moving camera and a static scene. To handle this issue, we apply auto-masking method [11] to filter out stationary pixels that remain with the same appearance between two frames in a sequence. Since the binary mask μ is computed in this form on the forward pass

$$\mu = \left[\min_{t'} pe(I_t, I_{t' \rightarrow t}) < \min_{t'} pe(I_t, I_{t'}) \right], \quad (4)$$

where $[\]$ is the Iverson bracket.

In addition, in order to regularize the disparities in texture-less regions, an edges-aware smoothness regularization term L_s is used

$$L_s = |\partial_x d_t^*| e^{-|\partial_x I_t|} + |\partial_y d_t^*| e^{-|\partial_y I_t|}, \quad (5)$$

where $d_t^* = d_t / \bar{d}_t$ is the mean-normalized inverse depth from [44] to discourage shrinking of the estimated depth.

The final loss L is computed as the combination of photometric loss L_p and smoothness loss L_s at multiple scales

$$L = \frac{1}{S} \sum_i (\mu L_p^i + \lambda L_s^i), \quad (6)$$

where S is the number of scales, and λ is the weighting for the smoothness regularization term.

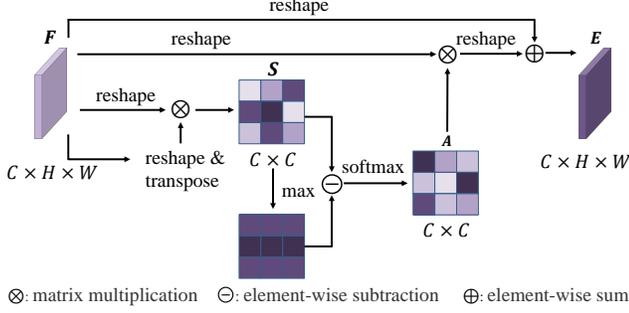


Figure 3. Details of the structure perception module.

3.2. Channel-wise Attention-based Network

As shown in Fig. 2, our CADepth-Net is a fully convolutional U-Net architecture. We adopt a pretrained residual network as the backbone to extract semantic features. Then these features would be fed into the structure perception module and generate new features to explicitly enhance the perception of scene structure. Moreover, we gradually recover the spatial resolution at decoder stage, with skip-connection to facilitate the flow of gradients and information throughout the model, and sequentially employ our detail emphasis module to generate sharp edges and finer details. Finally, we successively upsample the predicted inverse depth maps until original input resolutions using nearest neighbors interpolation at multiple scales, and compute the training loss at this higher input resolution.

3.3. Structure Perception Module

In depth estimation, each high-level feature map can be regarded as a region-specific response as shown in Fig. 6 (b), and different region responses are associated with each other. If each channel map captures more different region responses from all the other channel maps, as shown in Fig. 6 (c), it will obtain more relative depth information from distant regions and significantly enhance the perception of scene structure. Therefore, we propose a self-attention module to model interdependencies between channels and aggregate different region responses.

The first step is to generate an attention matrix which models the relationship between any two channel maps. As illustrated in Fig. 3, given the feature map $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$ produced by ResNet encoder, we firstly reshape \mathbf{F} to $\mathbb{R}^{C \times N}$, where $N = H \times W$ is the number of pixels, then perform a matrix multiplication between \mathbf{F} and the transpose of \mathbf{F} to compute the feature similarity $\mathbf{S} \in \mathbb{R}^{C \times C}$

$$S_{ij} = F_i \cdot F_j^T. \quad (7)$$

The similarity between channel maps indicates the spatial relationship of region responses *i.e.* any two feature maps have higher similarity means that they also have strong responses to the same region. As we need to fuse

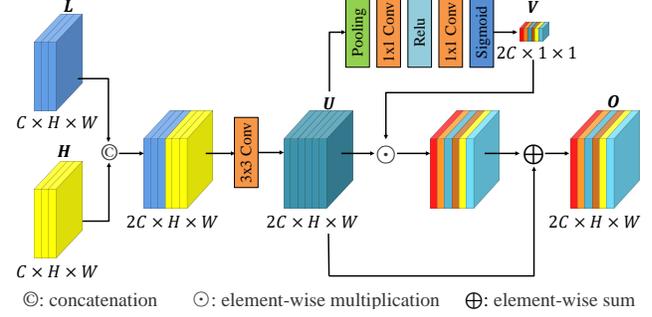


Figure 4. Design of the detail emphasis module.

more responses from *different regions*, we convert the similarity \mathbf{S} to discrimination $\mathbf{D} \in \mathbb{R}^{C \times C}$ by performing the element-wise subtraction

$$D_{ij} = \max_i(S) - S_{i,j}, \quad (8)$$

where D_{ij} measures the j^{th} channel's impact on the i^{th} channel. For each channel map, other channels with discriminative features (*i.e.* different region response) will get higher scores D_{ij} during feature aggregation. Then we apply a softmax layer to obtain the attention map $\mathbf{A} \in \mathbb{R}^{C \times C}$

$$A_{ij} = \frac{\exp(D_{ij})}{\sum_{j=1}^C \exp(D_{ij})}. \quad (9)$$

In addition, we perform a matrix multiplication between the transpose of \mathbf{A} and \mathbf{F} and reshape the result to $\mathbb{R}^{C \times H \times W}$. Finally we perform an element-wise sum operation between \mathbf{F} and the result to obtain the final output $\mathbf{E} \in \mathbb{R}^{C \times H \times W}$ as follows

$$E_i = \sum_{j=1}^C (A_{ij} F_j) + F_i. \quad (10)$$

The Eq. 10 shows that the final feature of each channel is the weighted sum of the features from all channels and the original feature. By capturing the long-range dependencies between feature maps, we obtain the aggregated features encoding rich context information of scene structure.

3.4. Detail Emphasis Module

The decoder recovers the resolution by fusing the following features at different scales: the low-level information from skip-connection encoding rich spatial details, and the high-level information encoding more context information. The simple fusion operations like sum or concatenation lack the further processing of local details and neglect the semantic gap between different level features, leading to blurry artefacts in predicted depth maps. The core of predicting sharper edges is to properly handle local details, and it is easy for network to recover accurate depth predictions

if it knows the category and location of features describing object boundaries clearly. Therefore, by using the channel attention mechanism that can make network pay attention to specific channel features, we propose a detail emphasis module to emphasize important details and efficiently fuse features at different scales.

Specifically, we first concatenate the low-level features \mathbf{L} and the high-level features \mathbf{H} , then utilize a convolution layer to obtain \mathbf{U} with the batch normalization [16] to balance the scales of the features

$$U = \sigma(BN(W_1 \otimes f(L, H))), \quad (11)$$

where $f()$ denotes concatenation and \otimes denotes the 3×3 or 1×1 convolution, BN refers to the batch normalization and we use the ReLU as the activation function $\sigma()$.

Next, we squeeze \mathbf{U} to a vector by global average pooling to obtain global context and use two 1×1 convolution layers followed by a sigmoid function to compute a weights vector $\mathbf{V} \in \mathbb{R}^{1 \times 1 \times C}$, to recalibrate channel-wise features and measure the importance of them in the meantime

$$V = \delta(W_2 \otimes \sigma(W_3 \otimes (\frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W U_{i,j}))), \quad (12)$$

where H and W refer to the height and width of \mathbf{U} , and $\delta()$ denotes the sigmoid function. Then we perform the element-wise multiplication between \mathbf{V} and \mathbf{U} to generate re-weight features. As the weights scores in \mathbf{V} indicate the importance of corresponding channels, *i.e.* the channel maps containing critical information will get higher scores, this recalibration operation can adaptively emphasize the crucial details at multiple scales for sharp edges (Fig. 7). Finally, we sum up \mathbf{U} and re-weight features for stability

$$O = V \odot U + U, \quad (13)$$

where \odot denotes element-wise dot product and \mathbf{O} refers to the final outputs. Fig. 4 shows the design of the detail emphasis module, this design recalibrates channel-wise feature responses and produce more precise depth estimation.

4. Experiments

In this section, we show extensive experiments for evaluating the performance of our methods and demonstrate the effectiveness of the proposed approaches.

4.1. Implementation Details

Our model are implemented based on PyTorch [35], trained for 20 epochs on a single Nvidia 3090 with a batch size of 12 and an input/output resolution of 640×192 . We jointly train both depth network and pose network with the Adam Optimizer [20] with $\beta_1 = 0.9$, $\beta_2 = 0.999$. The initial learning rate is set to $1e^{-4}$ and decay to $1e^{-5}$ after 15

epochs. We set the SSIM weight to $\alpha = 0.85$ and smoothness term weight to $\lambda = 1e - 3$.

DepthNet. We implement our depth estimation network as an encoder-decoder architecture. Moreover, we start the ResNet50 encoder with weights pretrained on ImageNet [38] as it has been shown to improve accuracy compared to training from scratch. We set sigmoid follow the output of network and convert result σ to depth with $D = 1/(a\sigma + b)$, where use a and b to constrain D between 0.1 and 100 units.

PoseNet. Our PoseNet is built on ResNet50, modified to accept six channels tensor as input, which allows the adjacent frames to feed into the network. The outputs of PoseNet is the 6-DoF relative pose consist of translation vectors and Euler angles, scaled by a factor of 0.01.

4.2. KITTI Results

We train and evaluate our methods using the KITTI 2015 stereo data set [9]. We adopt the data split of Eigen *et al.* [5] for distance to 80m and use pre-processing to remove static frames before training. Ultimately, this results in 39,810 training monocular triplets, and 4,424 for validation and 697 for evaluation. We report results using the per-image median ground truth scaling during evaluation.

As shown in Table 1, our proposed CADEPTH-Net significantly outperforms the existing SoTA self-supervised approaches in all metrics. In addition, we score dramatically higher in the hardest accuracy metric $\delta < 1.25$, which indicates that our model predicts more accurate and realistically detailed depth estimation than all other competing models. For a fair comparison, we also give the results on various input resolution and training settings, and our model still improves the performance at higher image resolutions. Fig. 5 shows that our model produces sharper depth estimation on thinner structures *e.g.* road signs and poles. Moreover, our model successfully estimates correct depth at the highly reflective car roof (6th row), which are the challenging problems for previous advanced methods. These improvements can be explained by the better perception of scenes and objects afforded by the structure perception module, and the further regularisation provided by detail emphasis module. Additional results can be seen in supplementary material.

4.3. Make3D result

To evaluate the generalization ability of our model on the unseen dataset, we report the quantitative results for the Make3D dataset [39] using our model trained on KITTI 2015 [9]. Following the same evaluation protocol as [10], we test on a center crop of 2×1 ratio and apply median scaling. As shown in Table 3, our approach produces superior results compared with the other SOTA self-supervised methods. Qualitative results can be seen in Fig. 8, which show that our model generates more accurate and sharper depth estimation on the previously unseen Make3D dataset.

Method	Train	Resolution	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
SfMLearner [56]†	M	416 × 128	0.183	1.595	6.709	0.270	0.734	0.902	0.959
GeoNet [50]†	M	416 × 128	0.149	1.060	5.567	0.226	0.796	0.935	0.975
DDVO [44]	M	416 × 128	0.151	1.257	5.583	0.228	0.810	0.936	0.974
Struct2depth ‘(M)’ [1]	M	416 × 128	0.141	1.026	5.291	0.215	0.816	0.945	<u>0.979</u>
SIGNet [32]	M	416 × 128	0.133	<u>0.905</u>	5.181	0.208	0.825	0.947	0.981
SGDepth [22]	M	416 × 128	<u>0.128</u>	1.003	<u>5.085</u>	0.206	0.853	0.951	0.978
Monodepth2 [11]	M	416 × 128	<u>0.128</u>	1.087	5.171	<u>0.204</u>	<u>0.855</u>	<u>0.953</u>	0.978
CADepth-Net (Ours)	M	416 × 128	0.116	0.893	4.906	0.192	0.874	0.957	0.981
DF-Net [57]	M	576 × 160	0.150	1.124	5.507	0.223	0.806	0.933	0.973
SGDepth [22]	M	640 × 192	0.117	0.907	4.844	0.196	0.875	0.958	0.980
Monodepth2 [11]	M	640 × 192	0.115	0.903	4.863	0.193	0.877	0.959	0.981
PackNet-SfM [13]	M	640 × 192	0.111	<u>0.785</u>	<u>4.601</u>	0.189	0.878	0.960	<u>0.982</u>
HR-Depth [30]	M	640 × 192	0.109	0.792	4.632	<u>0.185</u>	0.884	<u>0.962</u>	0.983
Johnston <i>et al.</i> [18]	M	640 × 192	<u>0.106</u>	0.861	4.699	<u>0.185</u>	<u>0.889</u>	<u>0.962</u>	<u>0.982</u>
CADepth-Net (Ours)	M	640 × 192	0.105	0.769	4.535	0.181	0.892	0.964	0.983
CC [36]	M	832 × 256	0.140	1.070	5.326	0.217	0.826	0.941	0.975
Monodepth2 [11]	M	1024 × 320	0.115	0.882	4.701	0.190	0.879	0.961	<u>0.982</u>
TrianFlow [54]	M	832 × 256	0.113	0.704	4.581	0.184	0.871	0.961	0.984
HR-Depth [30]	M	1024 × 320	0.106	0.755	<u>4.472</u>	0.181	0.892	0.966	0.984
FeatDepth [40]	M	1024 × 320	0.104	0.729	4.481	<u>0.179</u>	<u>0.893</u>	<u>0.965</u>	0.984
CADepth-Net (Ours)	M	1024 × 320	0.102	<u>0.734</u>	4.407	0.178	0.898	0.966	0.984
DualNet [55]	M	1248 × 384	0.121	0.837	4.945	0.197	0.853	0.955	0.982
SGDepth [22]	M	1280 × 384	0.113	0.880	4.695	0.192	0.884	0.961	0.981
PackNet-SfM [13]	M	1280 × 384	0.107	0.802	4.538	0.186	0.889	0.962	0.981
HR-Depth [30]	M	1280 × 384	0.104	0.727	4.410	0.179	0.894	0.966	0.984
CADepth-Net (Ours)	M	1280 × 384	0.102	0.715	4.312	0.176	0.900	0.968	0.984
EPC++ [29]	MS	832 × 256	0.128	0.935	5.011	0.209	0.831	0.945	0.979
Monodepth2 [11]	MS	640 × 192	0.106	0.818	4.750	0.196	0.874	0.957	0.979
HR-Depth [30]	MS	640 × 192	0.107	0.785	<u>4.612</u>	<u>0.185</u>	<u>0.887</u>	<u>0.962</u>	<u>0.982</u>
DepthHints [47]	MS	640 × 192	<u>0.105</u>	<u>0.769</u>	4.627	0.189	0.875	0.959	<u>0.982</u>
CADepth-Net (Ours)	MS	640 × 192	0.102	0.752	4.504	0.181	0.894	0.964	0.983
Monodepth2 [11]	MS	1024 × 320	0.106	0.806	4.630	0.193	0.876	0.958	0.980
HR-Depth [30]	MS	1024 × 320	0.101	0.716	<u>4.395</u>	<u>0.179</u>	<u>0.899</u>	<u>0.966</u>	<u>0.983</u>
DepthHints [47]	MS	1024 × 320	0.098	0.702	4.398	0.183	0.887	0.963	<u>0.983</u>
FeatDepth [40]	MS	1024 × 320	<u>0.099</u>	<u>0.697</u>	4.427	0.184	0.889	0.963	0.982
CADepth-Net (Ours)	MS	1024 × 320	0.096	0.694	4.264	0.173	0.908	0.968	0.984

Table 1. **Quantitative results on the KITTI Eigen Split.** Comparison of existing methods on KITTI 2015 [9] using the Eigen split for distances up to 80m. All methods in this table are trained on the KITTI dataset without additional datasets or online refinement. Best results are in **bold**, with second-best underlined. For Abs Rel, Sq Rel, RMSE and $RMSE_{log}$ lower is better, and for $\delta < 1.25$, $\delta < 1.25^2$ and $\delta < 1.25^3$ higher is better. In the *Train* column, S: Self-supervised stereo supervision, M: Self-supervised mono supervision. † refers to the newer results from Github. At test time, we scale the estimated depths with median ground-truth LiDAR information.

4.4. Structure Perception Module

To demonstrate the effectiveness of the structure perception module, Fig. 6 presents a comparison of features before and after treatment. For clear visualization, we map the channel maps to the RGB color cube and project them to the original RGB image. Fig. 6 (b) refers to the high-level features produced by encoder, which mainly indicates region-specific responses and various kinds of structural information of 3D scene, *e.g.* vanish point, region with same depth range and the area with same color or texture (*e.g.* sky).

As mentioned earlier, the structure perception module produces the aggregated features (Fig. 6 (c)) by performing the weighted summation of all channel maps. As the attention map describes the feature discrimination and the spatial

relationship of region responses between channels, by aggregating discriminative features at channel dimension, we can make each single channel map get rich scene structure representation and more complete region responses from non-contiguous regions. By doing so, the structure perception module obtains rich contextual information of overall scene geometry perception, which results in better scene understanding and feature representation for depth estimation.

Fig. 6 adequately demonstrates that each channel map obtains more extra depth perception from distant regions *e.g.* foreground (1st row) and midground (2nd row). In addition, it also particularly emphasizes the vanishing point regions, which are naturally a strong cue to understand the geometry of a scene. In the last row, the network originally focuses on foreground objects such as cars and obtains more

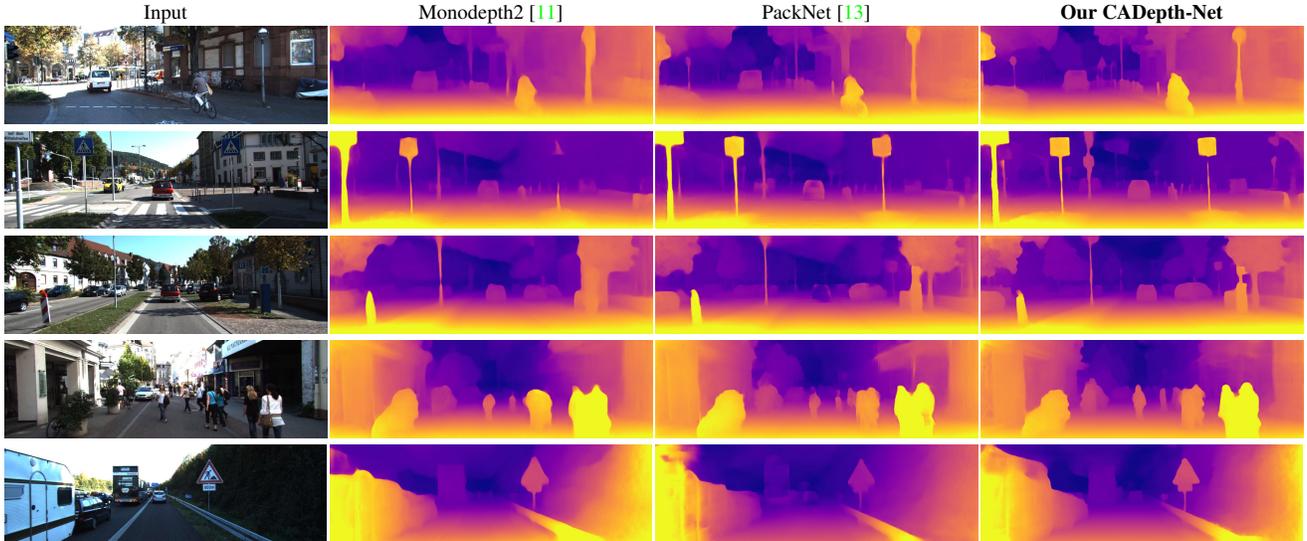


Figure 5. **Qualitative results on the KITTI Eigen split.** Our model consistently predicts sharper boundaries and fine-gained details on thinner objects, *e.g.* trees, pedestrians and signs.

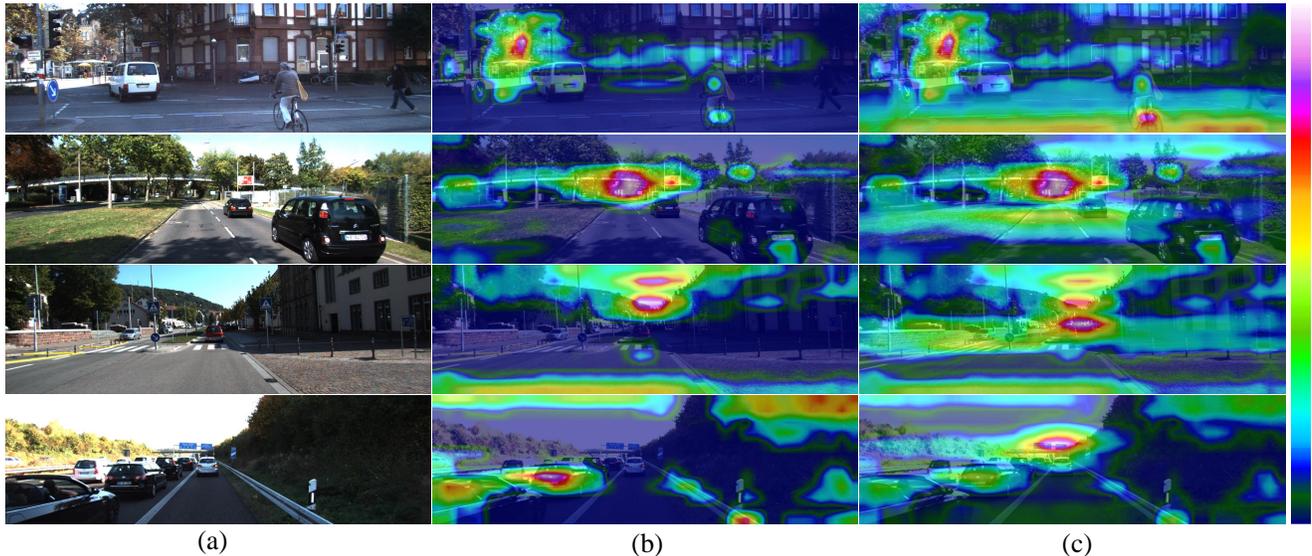


Figure 6. **The visualizations of the structure perception module.** (a) Input image. (b) Input feature maps. (c) Corresponding output of the structure perception module. All feature maps are projected onto RGB images for clear visualization. The structure perception module explicitly enhances the perception of scene structure and feature representation.

relative depth relationships after adding the vanishing point information. More visualization results are described in the supplementary material.

4.5. Detail Emphasis Module

The detail emphasis module generates a weights vector to re-calibrate channel maps and emphasizes the informative features for subsequent transformations. The weights scores produced by the sigmoid function are between 0 and 1, indicating the importance of corresponding features, *i.e.* the more important of the channel, the higher score is

earned. As shown in Fig. 7, we report the top n ($n = 10$) feature maps with the highest scores to show which features the network focuses on. We observe that our model mainly highlights the crucial low-level features that describe the object boundaries precisely. This observation is consistent with the fact that the decoder mainly uses detail information to recover spatial resolution gradually. In general, the detail emphasis module adaptively selects the informative local details and helps network handle and locate the object edges for sharper depth prediction. More visualization results are provided in the supplementary material.

	Backbone	SPM	DEM	Para	Time	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Baseline (MD2 ResNet18 [11])	R18			14.84 M	11.95 ms	0.115	0.903	4.863	0.193	0.877	0.959	0.981
Baseline + SPM	R18	✓		14.84 M	12.13 ms	0.110	0.843	4.739	0.188	0.883	0.961	0.982
Baseline + DEM	R18		✓	18.74 M	15.44 ms	0.111	0.851	4.746	0.189	0.881	0.961	0.982
CADepth-Net ResNet18 (full)	R18	✓	✓	18.74 M	15.77 ms	0.110	0.812	4.686	0.187	0.882	0.962	0.983
Baseline (MD2 ResNet50 [11])	R50			34.57 M	22.52 ms	0.110	0.831	4.642	0.187	0.883	0.962	0.982
Baseline + SPM	R50	✓		34.57 M	22.80 ms	0.107	0.784	4.589	0.185	0.887	0.963	0.982
Baseline + DEM	R50		✓	58.34 M	28.01 ms	0.107	0.759	4.557	0.183	0.884	0.964	0.983
CADepth-Net ResNet50 (full)	R50	✓	✓	58.34 M	28.41 ms	0.105	0.769	4.535	0.181	0.892	0.964	0.983

Table 2. **Ablation Studies.** We evaluate the performance of our structure perception module (SPM) and detail emphasis module (DEM) contributions with Monodepth2 (MD2) [11] as the baseline. R: ResNet, Para: parameters. All models in this table are trained with monocular self-supervised (M) and standard resolution (640×192). The inference time is tested on a single RTX3090 GPU.

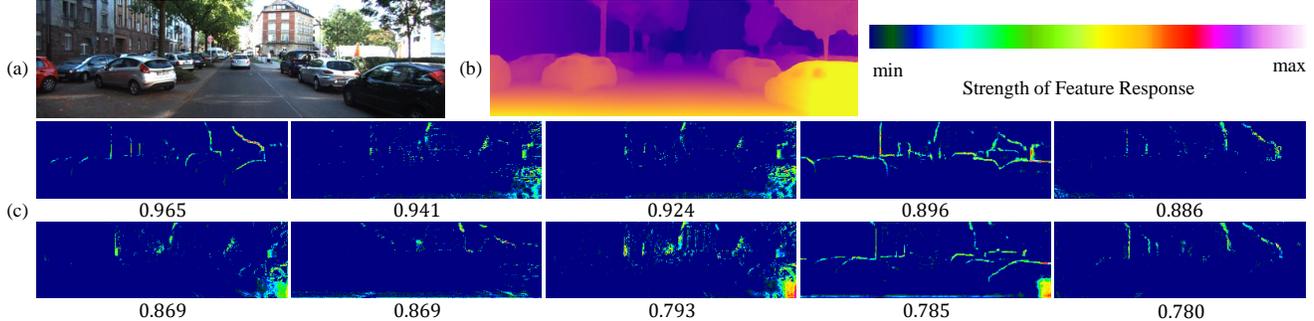


Figure 7. **The visualizations of detail emphasis module.** (a) Input image. (b) Predicted depth map. (c) Top n ($n = 10$) feature maps with highest weights score that range in $[0, 1]$. Detail emphasis module mainly highlights the crucial local details.

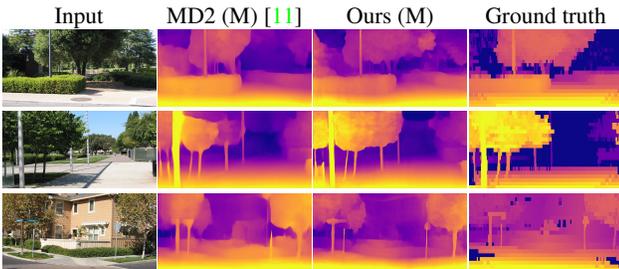


Figure 8. **Qualitative Make3D results.** All methods were trained on KITTI using monocular supervision.

	Type	Abs Rel	Sq Rel	RMSE	\log_{10}
Karsch [19]	D	0.428	5.079	8.389	0.149
Liu [28]	D	0.475	6.562	10.05	0.165
Laina [24]	D	0.204	1.840	5.683	0.084
Monodepth [10]	S	0.544	10.94	11.760	0.193
Zhou [56]	M	0.383	5.321	10.470	0.478
DDVO [44]	M	0.387	4.720	8.090	0.204
Monodepth2 [11]	M	0.322	3.589	7.417	0.163
Ours	M	0.312	3.086	7.066	0.159

Table 3. **Make3D results.** All self-supervised mono (M) results benefit from median scaling.

4.6. Ablation Study

For further demonstrating the performance improvements of our provided methods, Table 2 and supplemental material show the ablation study of our various components with Monodepth2 [11] as the baseline. It shows that all of our contributions achieve a steady improvement in almost all evaluation measures and obtain a consistent performance gain on different backbones, showing the robustness of our CADepth-Net to the backbone architecture capacity. Note that the structure perception module shows superior performance on metric $\delta < 1.25$, with little time cost and no additional parameters, demonstrating the improvements benefit from the better scene understanding rather than an increase in network complexity. Finally, the combination of all our modules with ResNet50 achieves the best results, with 59M parameters and an inference time of 28ms on an RTX3090

GPU, meets the requirements of real-time applications.

5. Conclusion

In this paper, we introduce a novel architecture named channel-wise attention-based depth estimation network, with two effective components, the structure perception module and the detail emphasis module. The structure perception module aggregates the discriminative features by capturing the long-range dependencies to obtain the context of scene structure and rich feature representation. Additionally, the detail emphasis module employs the channel attention mechanism to highlight objects' boundaries information and efficiently fuse different level features. Furthermore, the experiments demonstrate that our CADepth-Net produces sharper depth estimation and achieves the state-of-the-art results on KITTI datasets.

References

- [1] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8001–8008, 2019. [6](#)
- [2] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015. [12](#)
- [3] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. [13](#)
- [4] Guilherme N DeSouza and Avinash C Kak. Vision for mobile robot navigation: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 24(2):237–267, 2002. [1](#)
- [5] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015. [1](#), [2](#), [5](#), [12](#)
- [6] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *arXiv preprint arXiv:1406.2283*, 2014. [1](#), [2](#)
- [7] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2002–2011, 2018. [1](#), [2](#), [13](#)
- [8] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3146–3154, 2019. [2](#)
- [9] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. [1](#), [5](#), [6](#), [13](#)
- [10] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 270–279, 2017. [1](#), [5](#), [8](#), [13](#)
- [11] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3828–3838, 2019. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [13](#), [15](#)
- [12] Matan Goldman, Tal Hassner, and Shai Avidan. Learn stereo, infer mono: Siamese networks for self-supervised, monocular, depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. [13](#)
- [13] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raveentos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2485–2494, 2020. [1](#), [2](#), [6](#), [7](#), [13](#), [15](#)
- [14] Vitor Guizilini, Rui Hou, Jie Li, Rares Ambrus, and Adrien Gaidon. Semantically-guided representation learning for self-supervised monocular depth. *arXiv preprint arXiv:2002.12319*, 2020. [2](#)
- [15] Xiaoyang Guo, Hongsheng Li, Shuai Yi, Jimmy Ren, and Xiaogang Wang. Learning monocular depth by distilling cross-domain stereo networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 484–500, 2018. [1](#)
- [16] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. [5](#)
- [17] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances in neural information processing systems*, 28:2017–2025, 2015. [3](#)
- [18] Adrian Johnston and Gustavo Carneiro. Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4756–4765, 2020. [2](#), [6](#), [13](#)
- [19] Kevin Karsch, Ce Liu, and Sing Bing Kang. Depth transfer: Depth extraction from video using non-parametric sampling. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2144–2158, 2014. [8](#)
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [21] KITTI Single Depth Evaluation Server. http://www.cvlibs.net/datasets/kitti/eval_depth.php?benchmark=depth_prediction. 2017. [12](#)
- [22] Marvin Klingner, Jan-Aike Termöhlen, Jonas Mikolajczyk, and Tim Fingscheidt. Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance. In *European Conference on Computer Vision*, pages 582–600. Springer, 2020. [2](#), [6](#), [13](#)
- [23] Shu Kong and Charless Fowlkes. Pixel-wise attentional gating for parsimonious pixel labeling. *arXiv preprint arXiv:1805.01556*, 2018. [13](#)
- [24] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth international conference on 3D vision (3DV)*, pages 239–248. IEEE, 2016. [2](#), [8](#)
- [25] Bo Li, Yuchao Dai, and Mingyi He. Monocular depth estimation with hierarchical fusion of dilated cnns and soft-weighted-sum inference. *Pattern Recognition*, 83:328–339, 2018. [13](#)
- [26] Ruibo Li, Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, and Lingxiao Hang. Deep attention-based classification network for robust depth prediction. In *Asian Conference on Computer Vision*, pages 663–678. Springer, 2018. [13](#)

- [27] Fayao Liu, Chunhua Shen, and Guosheng Lin. Deep convolutional neural fields for depth estimation from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5162–5170, 2015. 2
- [28] Miaomiao Liu, Mathieu Salzmann, and Xuming He. Discrete-continuous depth estimation from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 716–723, 2014. 8
- [29] Chenxu Luo, Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, Ram Nevatia, and Alan Yuille. Every pixel counts++: Joint learning of geometry and motion with 3d holistic understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2624–2641, 2019. 6, 13
- [30] 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. *arXiv preprint arXiv:2012.07356*, 2020. 6
- [31] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5667–5675, 2018. 2, 13
- [32] Yue Meng, Yongxi Lu, Aman Raj, Samuel Sunarjo, Rui Guo, Tara Javidi, Gaurav Bansal, and Dinesh Bharadia. Signet: Semantic instance aided unsupervised 3d geometry perception. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9810–9820, 2019. 6
- [33] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3061–3070, 2015. 1
- [34] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison. Dtam: Dense tracking and mapping in real-time. In *2011 international conference on computer vision*, pages 2320–2327. IEEE, 2011. 1
- [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. In *NeurIPS-W*, 2017. 5
- [36] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12240–12249, 2019. 2, 6, 13
- [37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 1
- [38] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 5
- [39] Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE transactions on pattern analysis and machine intelligence*, 31(5):824–840, 2008. 5, 13
- [40] Chang Shu, Kun Yu, Zhixiang Duan, and Kuiyuan Yang. Feature-metric loss for self-supervised learning of depth and egomotion. In *European Conference on Computer Vision*, pages 572–588. Springer, 2020. 6
- [41] Jaime Spencer, Richard Bowden, and Simon Hadfield. Defeat-net: General monocular depth via simultaneous unsupervised representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14402–14413, 2020. 2
- [42] Jonas Uhrig, Nick Schneider, Lukas Schneider, Uwe Franke, Thomas Brox, and Andreas Geiger. Sparsity invariant cnns. In *2017 international conference on 3D Vision (3DV)*, pages 11–20. IEEE, 2017. 12, 13
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 2
- [44] Chaoyang Wang, José Miguel Buenaposada, Rui Zhu, and Simon Lucey. Learning depth from monocular videos using direct methods. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2022–2030, 2018. 2, 3, 6, 8, 13
- [45] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 2
- [46] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 3
- [47] Jamie Watson, Michael Firman, Gabriel J Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2162–2171, 2019. 6
- [48] Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, and Ram Nevatia. Lego: Learning edge with geometry all at once by watching videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 225–234, 2018. 2
- [49] Zhenheng Yang, Peng Wang, Wei Xu, Liang Zhao, and Ramakant Nevatia. Unsupervised learning of geometry with edge-aware depth-normal consistency. *arXiv preprint arXiv:1711.03665*, 2017. 2
- [50] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1983–1992, 2018. 2, 6, 13
- [51] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 340–349, 2018. 13
- [52] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In

International conference on machine learning, pages 7354–7363. PMLR, 2019. [2](#)

- [53] Zhenyu Zhang, Chunyan Xu, Jian Yang, Ying Tai, and Liang Chen. Deep hierarchical guidance and regularization learning for end-to-end depth estimation. *Pattern Recognition*, 83:430–442, 2018. [13](#)
- [54] Wang Zhao, Shaohui Liu, Yezhi Shu, and Yong-Jin Liu. Towards better generalization: Joint depth-pose learning without posenet. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9151–9161, 2020. [6](#)
- [55] Junsheng Zhou, Yuwang Wang, Kaihuai Qin, and Wenjun Zeng. Unsupervised high-resolution depth learning from videos with dual networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6872–6881, 2019. [6](#)
- [56] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1851–1858, 2017. [1](#), [2](#), [6](#), [8](#), [13](#)
- [57] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *Proceedings of the European conference on computer vision (ECCV)*, pages 36–53, 2018. [6](#)

Supplementary Material

This document provides additional details and results concerned with paper "Channel-Wise Attention-Based Network for Self-Supervised Monocular Depth Estimation". The supplementary material is organized as follows: Section 1 provides the details of depth estimation network, Section 2 and Section 3 reports the quantitative results on the KITTI improved ground truth and online evaluation server, Section 4 collects additional ablation experiments, Section 5 provides the additional qualitative comparisons, and Section 6 shows visualization results of our methods.

A. Network Architecture

Except where note, for all experiments, we use a ResNet50 encoder with pretrained ImageNet weights for both depth and pose networks. For depth model, the structure perception module takes the features from encoder as input. In addition, we successively adopt the detail emphasis module at different scales in decoder stage. Note that although there is no skip-connection at the highest resolution, we still use the detail emphasis module for further regularization. For high resolution input, *e.g.* 1024×320 and 1280×384 , we employ a lightweight setup, ResNet18 and 640×192 , for pose encoder at training for memory savings. The depth network architecture is shown in Table 4.

B. KITTI Improved Ground Truth

The evaluation method proposed by Eigen *et al.* [5] for KITTI uses the reprojected LIDAR points to generate the ground truth depth, but does not handle moving objects and occlusions. [42] created a set of high-quality depth maps for the KITTI dataset using five consecutive frames and stereo pairs to handle moving objects better. This improved ground truth depth is provided for 652 (or 93%) of the 697 original test frames contained in the Eigen test split *et al.* [5]. We utilize the same error metrics and evaluation strategy as the main paper, and evaluate our model on these 652 improved ground truth frames without having to retrain each method. As shown in Table 5, we observe that our CADepth-Net still outperforms all existing advanced methods on all metrics.

C. KITTI Evaluation Server Benchmark

In Table 6 we report the results of our models on the KITTI single image depth prediction benchmark [21] which were computed on the KITTI online evaluation server. We train a new model on the new split consisting of 72,084 training examples, 6,060 validation, and 500 test with the same training protocols mentioned in the main paper. As we cannot use median scaling for evaluation, we calculate

Depth network						
layer	k	s	chns	res	input	activation
conv1	7	2	64	2	image	RELU
maxpool	3	2	64	4	conv1	-
econv1	3	1	256	4	maxpool	RELU
econv2	3	2	512	8	econv1	RELU
econv3	3	2	1024	16	econv2	RELU
econv4	3	2	2048	32	econv3	RELU
spm	-	-	2048	32	econv4	-
upconv5	3	1	256	32	spm	ELU [2]
dem5	3	1	1280	16	↑upconv5, econv3	RELU
iconv5	3	1	256	16	dem5	ELU
upconv4	3	1	128	16	iconv5	ELU
dem4	3	1	640	8	↑upconv4, econv2	RELU
iconv4	3	1	128	8	dem4	ELU
disp4	3	1	1	1	iconv4	Sigmoid
upconv3	3	1	64	8	iconv4	ELU
dem3	3	1	320	4	↑upconv3, econv1	RELU
iconv3	3	1	64	4	dem3	ELU
disp3	3	1	1	1	iconv3	Sigmoid
upconv2	3	1	32	4	iconv3	ELU
dem2	3	1	96	2	↑upconv2, conv1	RELU
iconv2	3	1	32	2	dem2	ELU
disp2	3	1	1	1	iconv2	Sigmoid
upconv1	3	1	16	2	iconv2	ELU
dem1	3	1	16	1	↑upconv1	RELU
iconv1	3	1	16	1	dem1	ELU
disp1	3	1	1	1	iconv1	Sigmoid

Table 4. **Depth network architecture.** For symbols in this table, **k**: kernel size, **s**: stride, **chns**: the number of output channels, **res**: the downscaling factor relative to the input image, **input**: corresponds to the input of each layer, **activation**: activation function. ↑ refers to a $2 \times$ nearest-neighbor upsampling. Encoder blocks are denoted by *econv** naming convention. The *spm* and *dem* represent structure perception module and detail emphasis module.

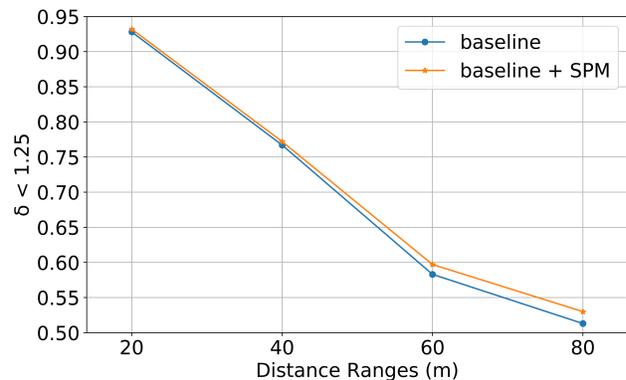


Figure 9. **Depth Evaluation on KITTI binned at different intervals**, calculated independently by only considering ground-truth depth pixels in that range (0-20m, 20-40m, ...).

the scale factor on the 2,000 KITTI training samples which have ground truth depths available. Table 6 shows that our CADepth-Net outperforms the existing self-supervised approaches and significantly reduces the performance gap to

Method	Train	Res	Dataset	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
SfMLearner [56]†	M	416 × 128	CS+K	0.176	1.532	6.129	0.244	0.758	0.921	0.971
Vid2Depth [31]	M	416 × 128	CS+K	0.134	0.983	5.501	0.203	0.827	0.944	0.981
GeoNet [50]	M	416 × 128	CS+K	0.132	0.994	5.240	0.193	0.833	0.953	0.985
DDVO [44]	M	416 × 128	CS+K	0.126	0.866	4.932	0.185	0.851	0.958	0.986
CC [36]	M	832 × 256	K	0.123	0.881	4.834	0.181	0.860	0.959	0.985
EPC++ [29]	M	640 × 192	K	0.120	0.789	4.755	0.177	0.856	0.961	0.987
Monodepth2 [11]	M	640 × 192	K	0.090	0.545	3.942	0.137	0.914	0.983	0.995
Johnston <i>et al.</i> [18]	M	640 × 192	K	0.081	0.484	3.716	0.126	0.927	0.985	0.996
CADepth-Net (Ours)	M	640 × 192	K	<u>0.080</u>	<u>0.442</u>	<u>3.639</u>	<u>0.124</u>	<u>0.927</u>	<u>0.986</u>	<u>0.996</u>
CADepth-Net (Ours)	M	1280 × 384	K	0.076	0.374	3.280	0.115	0.937	0.990	0.997
Zhan FullNYU [51]	D*MS	608 × 160	K	0.130	1.520	5.184	0.205	0.859	0.955	0.981
EPC++ [29]	MS	640 × 192	K	0.123	0.754	4.453	0.172	0.863	0.964	0.989
Monodepth2 [11]	MS	640 × 192	K	0.080	0.466	3.681	0.127	0.926	0.985	0.995
CADepth-Net (Ours)	MS	640 × 192	K	<u>0.076</u>	<u>0.417</u>	<u>3.488</u>	<u>0.120</u>	<u>0.933</u>	<u>0.987</u>	<u>0.996</u>
CADepth-Net (Ours)	MS	1024 × 320	K	0.070	0.346	3.168	0.109	0.945	0.991	0.997

Table 5. **Quantitative results on the KITTI improved ground truth.** Comparison of the existing methods to our CADepth-Net on KITTI 2015 [9] using annotated depth maps from [42]. Best results are in **bold**, with second best underlined. For Abs Rel, Sq Rel, RMSE and $RMSE_{log}$ lower is better, and for $\delta < 1.25$, $\delta < 1.25^2$ and $\delta < 1.25^3$ higher is better. In the *Train* column, S: Self-supervised stereo supervision, M: Self-supervised mono supervision, D*: Auxiliary depth supervision. † refers to the newer results from github. In *Dataset* column, CS: Cityscapes dataset [3], K: KITTI datasets [9].

Method	Train	SILog	sqErrorRel	absErrorRel	iRMSE
DHGRL [53]	D	15.47	4.04	12.52	15.72
CSWS [25]	D	14.85	3.48	11.84	16.38
APMoE [23]	D	14.74	3.88	11.74	15.63
DABC [26]	D	14.49	4.08	12.72	15.53
DORN [7]	D	11.77	2.23	8.78	12.98
Monodepth [10]	S	22.02	20.58	17.79	21.84
LSIM [12]	S	17.92	6.88	14.04	17.62
Monodepth2 [11]	M	15.57	4.52	12.98	16.70
SGDepth [22]	M	15.30	5.00	13.29	15.80
Monodepth2 [11]	MS	15.07	4.16	11.64	15.27
Ours	MS	13.34	3.33	10.67	13.61

Table 6. **Results on KITTI depth prediction benchmark.** D refers to ground truth depth supervision, while M and S are monocular and stereo self-supervision respectively.

	Backbone	Abs Rel	Sq Rel	RMSE	$\delta < 1.25$
SGDepth [22]	R18	0.117	0.907	4.844	0.875
monodepth2 [11]	R18	0.115	0.903	4.863	0.877
Ours	R18	0.110	0.812	4.686	0.882
monodepth2 [11]	R50	0.110	0.831	4.642	0.883
Johnston <i>et al.</i> [18]	R101	0.106	0.861	4.699	0.889
Ours	R50	0.105	0.769	4.535	0.892

Table 7. **Comparisons of methods with the same backbone on the KITTI Eigen Split.** R: ResNet. All models are trained with the same settings.

the supervised methods.

D. Additional Ablation Experiments

To further demonstrate the effectiveness of our structure perception module, we evaluate the performance of distant objects with the absolute relative error, Fig. 11 shows that we can effectively reduce the error produced by distant objects. Besides, Fig. 9 reports that our model improves the accuracy at all depth intervals, and the performance gap consistently increases when larger distances are considered, thanks to the better 3D scene geometric perception intro-

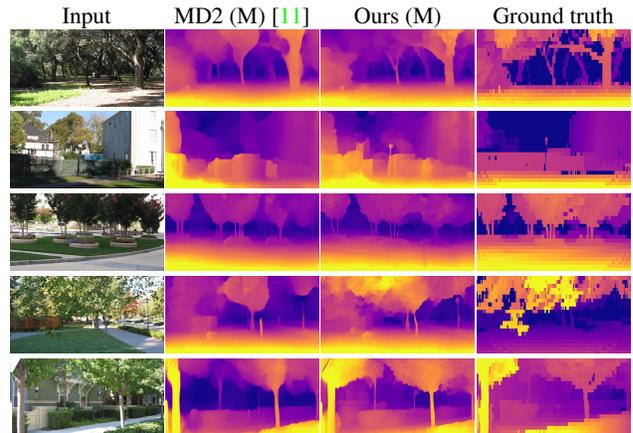


Figure 10. **Qualitative Make3D results.** Our CADepth-Net generates more fine-gained details compared to other method.

duced by the structure perception module. Fig. 12 shows that our method generates more fine-gained details and accurate object boundaries *e.g.* pedestrians and thin road signs by using detail emphasis module individually. For a fair comparison, Table 7 reports that our model outperforms other methods with the same backbone, which means that the gain in performance shown in our experiments is mainly due to the effectiveness of our proposed methods.

E. Additional Qualitative Comparisons

We provide additional qualitative results from the KITTI datasets in Fig. 13. We can observe that compared to existing baselines *e.g.* [11, 13], our CADepth-Net produces higher quality outputs and possesses the clearest border overall. We also show additional results from Make3D datasets [39] in Fig. 10, our methods preserve sharp discon-

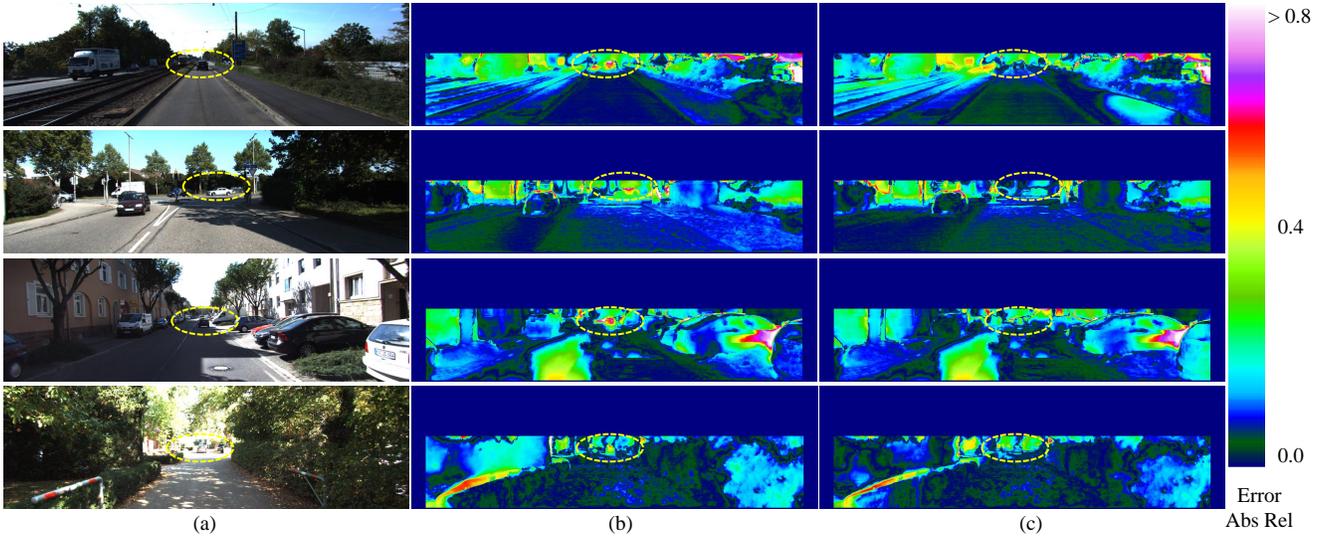


Figure 11. **Qualitative ablation study of Structure Perception Module (SPM).** (a) Input Image. (b) Error maps without SPM. (c) Error maps with SPM. The error induced by distance objects (yellow circles) is improved by the structure perception module.

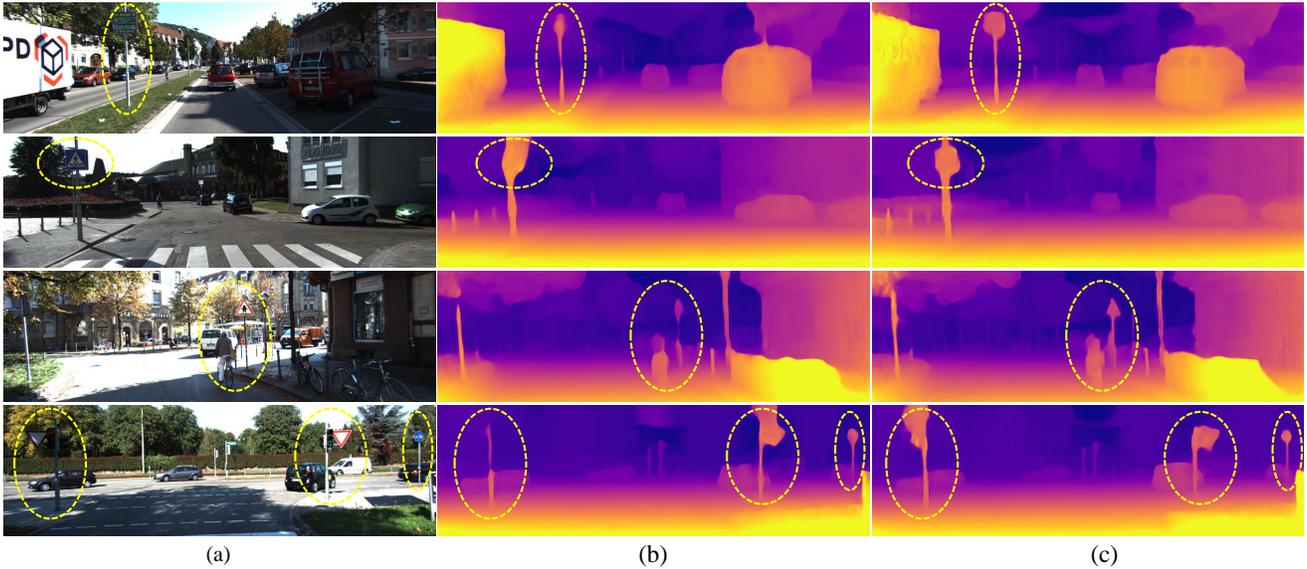


Figure 12. **Qualitative ablation study of Detail Emphasis Module (DEM).** (a) Input image. (b) Predicted depth maps without DEM. (c) Predicted depth maps with DEM. Thanks to the detail emphasis module, we could obtain the more precise and sharper depth estimation.

tinuities in depth prediction results.

F. Additional Visualization Results

To better understand our main contributions, Fig. 14 and Fig. 15 introduce additional visualization results of intermediate features. As shown in Fig. 14, each feature map obtains more region responses from the distant regions and aggregates relative depth relationships over 3D scene. By doing so, our model produces better scene understanding and rich feature representation. Fig. 15 lists the top n ($n = 8$) feature maps with the highest scores in the detail

emphasis module, and we can see that our model highlights critical local details at multiple scales, by assigning higher scores to them.

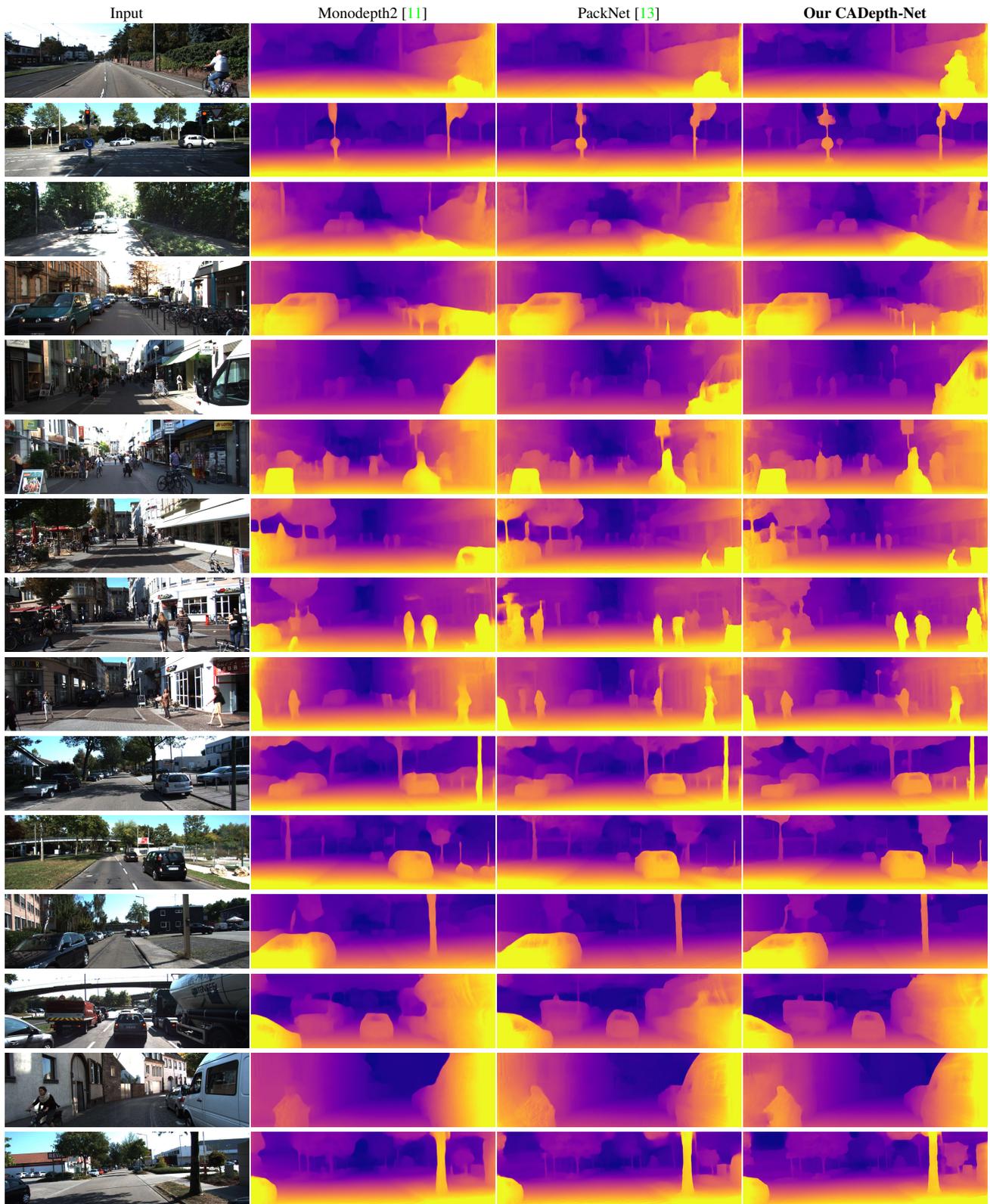


Figure 13. Qualitative results on the KITTI Eigen split. Our model produces higher quality outputs and possesses the most clear border.

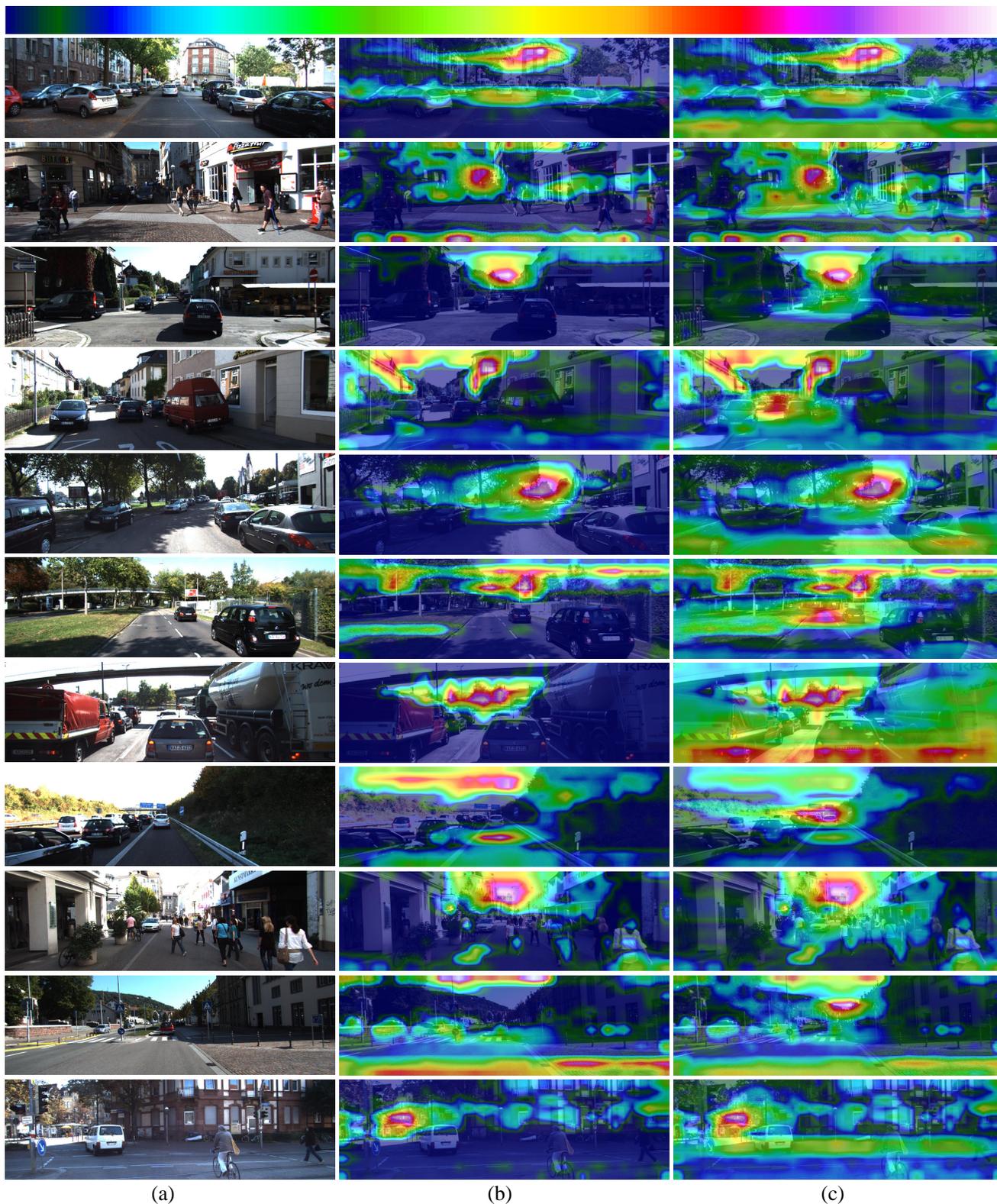


Figure 14. **The visualizations of structure perception module.** (a) Input Image. (b) Input feature maps. (c) Output feature maps. All feature representation are projected onto original image for clear visualizations. Our structure perception module explicitly enhances the scene understanding and feature representation.

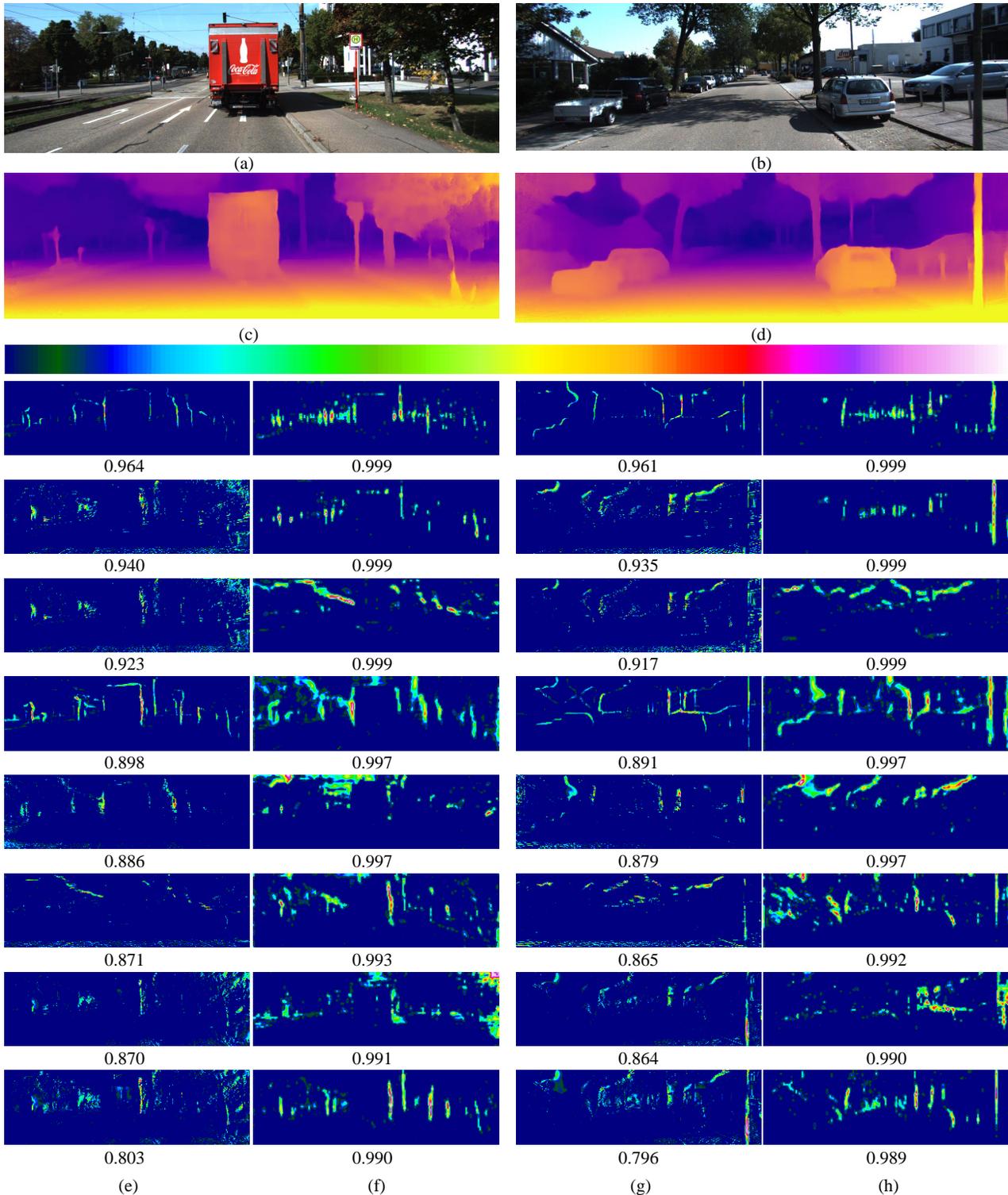


Figure 15. **The visualizations of detail emphasis module.** (a)(b) Input Image. (c)(d) Predicted depth maps. (e) ~ (h) Top n feature maps with highest weights score. Here n is 8 and the weights value range is between 0 and 1. Specifically, (e)(g) are produced by *dem2* (see Table 4) and (f)(h) are generated from *dem3*. The detail emphasis module mainly highlights the crucial local details.