

PYRAMID FEATURE ATTENTION NETWORK FOR MONOCULAR DEPTH PREDICTION

Yifang Xu, Chenglei Peng, Ming Li, Yang Li, and Sidan Du

Nanjing University, Nanjing Institute of Advanced Artificial Intelligence, Nanjing, China
 {mf20230128, dg20230020}@smail.nju.edu.cn, {pcl, yogo, coff128}@nju.edu.cn

ABSTRACT

Deep convolutional neural networks (DCNNs) have achieved great success in monocular depth estimation (MDE). However, few existing works take the contributions for MDE of different levels feature maps into account, leading to inaccurate spatial layout, ambiguous boundaries and discontinuous object surface in the prediction. To better tackle these problems, we propose a Pyramid Feature Attention Network (PFANet) to improve the high-level context features and low-level spatial features. In the proposed PFANet, we design a Dual-scale Channel Attention Module (DCAM) to employ channel attention in different scales, which aggregate global context and local information from the high-level feature maps. To exploit the spatial relationship of visual features, we design a Spatial Pyramid Attention Module (SPAM) which can guide the network attention to multi-scale detailed information in the low-level feature maps. Finally, we introduce scale-invariant gradient loss to increase the penalty on errors in depth-wise discontinuous regions. Experimental results show that our method outperforms state-of-the-art methods on the KITTI dataset.

Index Terms— Depth estimation, channel attention, spatial attention, pyramid feature, deep learning

1. INTRODUCTION

Monocular depth estimation (MDE) is an important task that aims to predict pixel-wise depth from a single RGB image, and has many applications in computer vision, such as 3D reconstruction, scene understanding, autonomous driving and intelligent robots [1]. In the meanwhile, MDE is a technically ill-posed problem as a single image can be projected from an infinite number of different 3D scenes. To solve this inherent ambiguity, one possibility is to leverage prior auxiliary information, such as texture information, occlusion, object locations, perspective, and defocus [2], but it is not easy to effectively extract useful prior information.

More recently, some works on MDE based on encoder-decoder architecture have shown significant improvements in performance by using deep convolutional neural networks (DCNNs) [3]. As backbone for encoder, very powerful deep

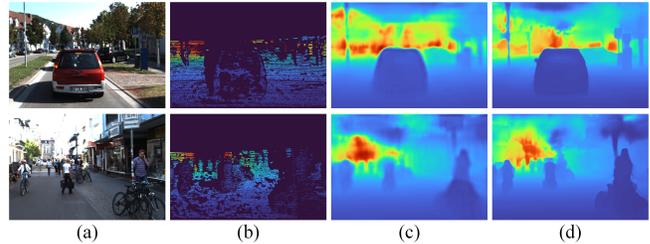


Fig. 1. Depth estimation example. (a) Input RGB image; (b) Ground truth depth; (c) Fu et al. [2]; (d) Ours.

networks such as ResNet [4], DenseNet [5] or ResNext [6] are widely adopted. These networks cascade multiple convolutions and spatial pooling layers to gradually increase the receptive field and generate the high-level depth information. In decoder phase, state-of-the-art methods are based on upsampling layer with global context module [7], skip connection, depth-to-space [8], multi-scale local planar guidance for upsampling operation [3]. These methods directly fuse different scale features without considering their different contributions for MDE, which leads to ambiguous boundaries and discontinuous object surface in predicted depth (see Fig.1 (c)). To tackle these problems, logarithmic discretization for ordinal regression [2] and attention module with structural awareness [9] are introduced to MDE network. However, the high-level and low-level features play different roles in MDE. The existing methods did not consider this aspect, which may affect the effective extraction of depth information.

In this paper, we propose a novel monocular depth estimation network named Pyramid Feature Attention Network (PFANet). In order to enhance the global structural information in high-level features, we introduced the Dense version of Atrous Spatial Pyramid Pooling (Dense ASPP) [10], which is generally utilized in pixel-level semantic segmentation. Since Dense ASPP applies sparse convolutions with various expansion rates, these convolutions expand receptive field of the high-level features. And then we design Dual-scale Channel Attention Module (DCAM) to aggregate global context and local information at different scales in high-level features. During training process, DCAM assigns larger weight to the channels that play an important role in MDE. Considering the spatial relationship of visual features, we design Spatial Pyramid Attention Module (SPAM) to fuse the attention of multi-scale low-level features. This module improves the detailed

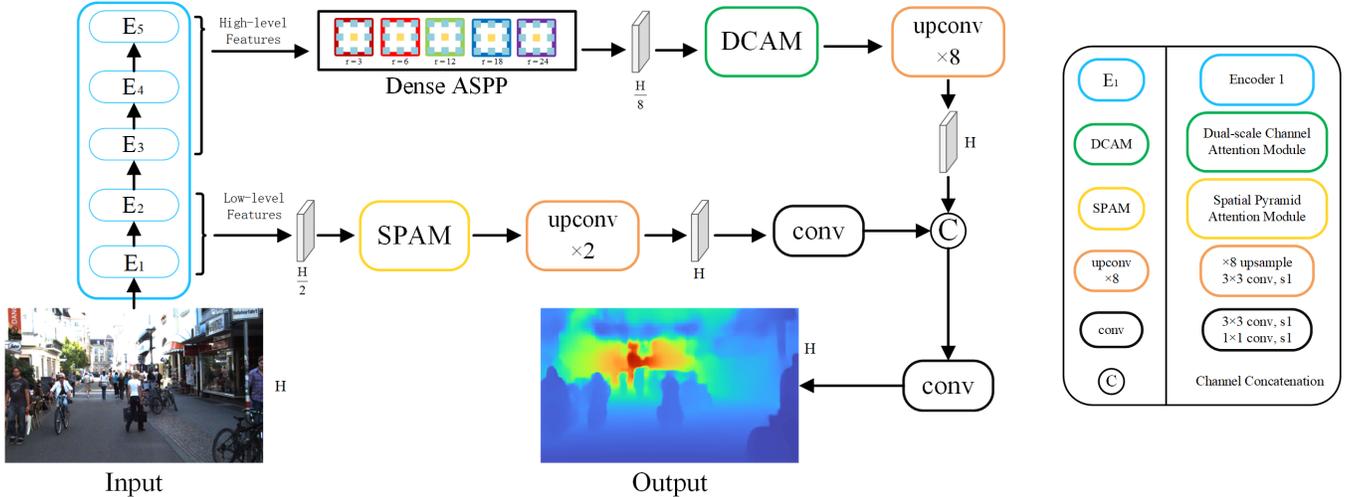


Fig. 2. The overview of Pyramid Feature Attention Network. The network is composed of E_i (the i -th level of encoder), Dense ASPP [10], Dual-scale Channel Attention Module and Spatial Pyramid Attention Module. The high-level features are from E_3 , E_4 and E_5 . The low-level features are from E_1 and E_2 .

local information in the low-level features, which clearer object edge and smoother object surface in prediction depth. Besides, we introduce scale-invariant gradient loss [11] to lead the network to learn more detail of object edges. With the above operations, the proposed PFANet can produce good depth maps (see Fig.1 (d)). In summary, our contributions are as follows:

1) We propose a novel Pyramid Feature Attention Network (PFANet) for MDE. For high-level features, we design Dual-scale Channel Attention Module (DCAM) to aggregate global context and local information. For low-level features, we design Spatial Pyramid Attention Module (SPAM) to capture more detailed information.

2) We introduce scale-invariant gradient loss to emphasize the depth discontinuity at different object boundaries and enhance smoothness in homogeneous regions.

3) The proposed method achieves state-of-the-art results on KITTI dataset.

2. RELATED WORK

Monocular Depth Estimation. As a pioneering work, Saxena et al. [12] propose to learn depth from visual cues based on Markov Random Field (MRF). Eigen et al. [1] introduce deep learning network that make coarse global prediction and refine it with local information, and extend it to a multi-scale network for depth estimation [13]. Since then, given the success of DCNNs in computer vision, more and more depth estimation networks have been proposed. Laina et al. [14] use a fully convolutional architecture with residual upsampling blocks to tackle the high-dimension regression problem. Jiao et al. [15] apply semantic segmentation network to assist depth estimation, and propose attention-driven loss that address long-tail distribution of depth values. The lat-

est SOTA network DORN [2] models MDE as an ordinal regression problem, to address the increase in error with depth magnitude, via spacing-increasing discretization strategy.

Attention Mechanism. Attention mechanism is derived from human perception, it can selectively focus on the prominent parts to capture useful information in entire scene. Similarly, attention mechanism is also suitable for various computer vision tasks, such as image classification, depth estimation, etc. SENet [16] proposes channel attention module to adaptively recalibrate channel-wise feature responses by explicitly modeling the interdependence between channels. CBAM [17] introduces spatial attention module based on channel attention module, and concatenates two modules for adaptive feature refinement. Wang et al. [18] design pyramid diverse attention (PDA) to learn multi-scale diverse local representations automatically, leading to network focus on different local patches.

3. OUR METHOD

3.1. Overview

In this paper, we propose Pyramid Feature Attention Network based on encoder-decoder architecture. DenseNet-161 [5] pre-trained on ILSVRC as our encoder. Decoder is composed of Dense ASPP [10], Dual-scale Channel Attention Module, and Spatial Pyramid Attention Module.

The proposed network architecture is shown in Fig.2. Input a single RGB image with resolution H . In encoder, the five convolutional blocks $\{E_1, E_2, E_3, E_4, E_5\}$ output feature maps with different resolutions that are $H/2$, $H/4$, $H/8$, $H/16$ and $H/32$ respectively. The high-level features are from E_3 , E_4 and E_5 . The low-level features are from E_1 and E_2 , which upsample to resolution of E_2 . After the backbone net-

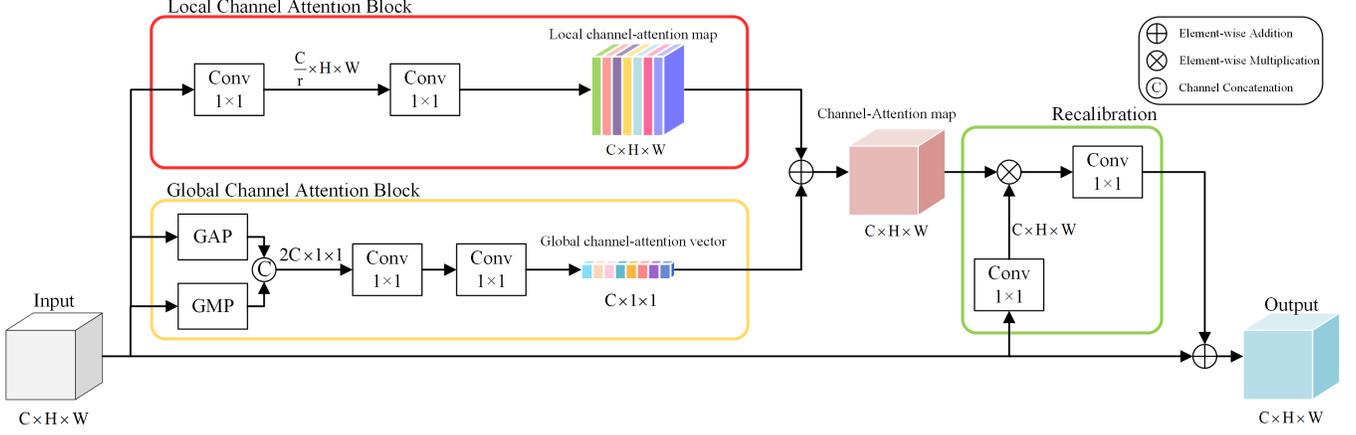


Fig. 3. The architecture of Dual-scale Channel Attention Module (DCAM). It consists of two blocks: local channel attention block and global channel attention block. The outputs of two blocks are fused to generate the channel attention map. Recalibration block is utilized to calibrate the channel attention map and further extract useful information for MDE. GAP denotes global average pooling layer. GMP denotes global max pooling layer.

work, for high-level features, we apply Dense ASPP module to fuse high-level features and expand the receptive field. This module produces an $H/8$ feature map via various dilated convolutional operations. The dilation rates r are 3, 6, 12, 18 and 24 respectively. Following Dense ASPP, we place DCAM to extract global context and local information from high-level features. Then, we apply SPAM to capture spatial information at multi-scale from the low-level feature maps. To get the high-level and low-level features with same resolution H , we utilize up-convolutional layer, which consists of upsampling operations and a 3×3 convolutional layer. Finally, they are concatenated and fed into the final convolutional layer to get the depth estimation \tilde{d} .

3.2. Dual-scale Channel Attention Module

The previous channel attention methods based on the squeeze and excitation [16] capture global context from the feature maps. However, this way ignores local information in features. To aggregate global context and local information simultaneously, we propose Dual-scale Channel Attention Module, as shown in Fig.3. DCAM consists of global channel attention block, local channel attention block and recalibration block. Its core idea is to implement channel attention on different scales.

In global channel attention block, average pooling layer and max pooling layer apply to reduce computational cost. To aggregate global context across channels, the dimension of the input feature maps are fused and reduced to $2C/r$ by 1×1 convolution layer, where C is the dimension of input $x \in \mathbb{R}^{C \times H \times W}$, r is reduction rate. And then this block produces global channel-attention vector $g(x) \in \mathbb{R}^{C \times 1 \times 1}$ via another 1×1 convolution layer. Similarly, in local channel attention block, we place two 1×1 convolution layers to extract local information across channels, which generate lo-

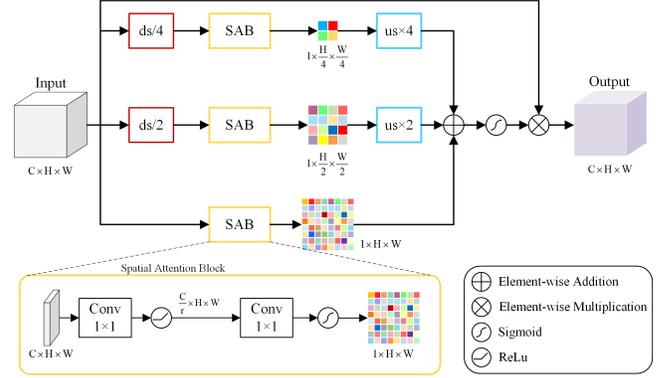


Fig. 4. The architecture of Spatial Pyramid Attention Module (SPAM). Ds/4 refers to /4 downsampling operation. Usx4 refers to $\times 4$ upsampling operation. Spatial attention blocks learn the spatial attention map, these three maps form a pyramid structure.

cal channel-attention map $l(x) \in \mathbb{R}^{C \times H \times W}$. Local channel-attention map and global channel-attention vector are fused to channel-attention map $A_c(x) \in \mathbb{R}^{C \times H \times W}$, before calibration (see Eq.(1)). Thus, we can effectively employ channel attention information and avoid introducing interference. Finally, we calibrate the original channel-attention map to improve feature representation, and new channel-attention map $\tilde{A}_c(x)$ as shown in Eq.(2).

$$A_c(x) = l(x) \otimes g(x) \quad (1)$$

$$\tilde{A}_c(x) = h[f(x) \otimes A_c(x)] \oplus x \quad (2)$$

where f and h are 1×1 convolutional layers in recalibration block. \oplus denotes element-wise addition. And \otimes denotes element-wise multiplication. Note that each convolutional layer is followed by an activation function ReLU.

3.3. Spatial Pyramid Attention Module

The high-level feature maps always processed by channel attention module, since it can capture channel’s dependency. However, this ignores structural information of the feature maps. To extract more local detailed information from the low-level feature maps, we proposed the spatial pyramid attention module, which utilizes the spatial pyramid structure. Fig.4 depicts the paradigm of SPAM. This module contains downsampling operation, spatial attention block and upsampling operation. Suppose the input low-level features $y \in \mathbb{R}^{C \times H \times W}$ via down-sampling operation, get down-sampling feature maps y_i ($i=1, 2, 3$), the resolution is 1, 1/2, 1/4 of the input, respectively. Then spatial attention block learns the spatial attention map $s(y_i) \in \mathbb{R}^{C \times H \times W}$, as shown in Eq. (3), these three blocks form a pyramid structure. Finally, through upsampling operation, we fuse the multi-scale spatial attention map. The output of SPAM $A_s(y)$ can be presented as Eq. (4).

$$s(y_i) = \sigma(\text{Conv}_2(\delta(\text{Conv}_1(y_i)))) \quad (3)$$

$$A_s(y) = \sigma[s(y_1) \oplus s(y_2) \oplus s(y_3)] \otimes y \quad (4)$$

where Conv_1 and Conv_2 refer to convolutional layers in spatial attention block. δ refers to ReLU function, σ refers to Sigmoid function.

3.4. Training Loss

The loss function to constraint our network contains two terms, i.e., scale-invariant loss (in log space) L_d and scale-invariant gradient loss L_s . We describe in detail each loss item as follows.

Scale-invariant loss. Scale invariant loss is proposed in [1] by Eigen et al., as shown in Eq. (5).

$$L_d(e) = \frac{1}{T} \sum_i e_i^2 - \frac{\lambda}{T^2} \left(\sum_i e_i \right)^2 \quad (5)$$

where $e_i = \log(\tilde{d}_i) - \log(d_i)$. \tilde{d}_i denotes ground truth of depth. d_i denotes predicted depth. $\lambda = 0.5$. T refers to the number of pixels with valid ground truth depth value. By rewritig Eq. (5):

$$L_d(e) = \frac{1}{T} \sum_i e_i^2 - \frac{1}{T^2} \left(\sum_i e_i \right)^2 + \frac{(1-\lambda)}{T^2} \left(\sum_i e_i \right)^2 \quad (6)$$

we can think of Eq. (6) as the sum of variances and weighted mean square error in log space. Setting $\lambda = 0$ is L2-norm, and setting $\lambda = 1$ is scale invariant error. Inspired by [3], we set $\lambda = 0.85$ in this work to accelerate minimizing the variance of error.

Scale-invariant gradient loss. Scale-invariant gradient loss (see Eq. (7)) is defined by [11] by Ummenhofer et al., which based on gradient loss. This loss function emphasizes sharpness on object boundaries and increases smoothness with

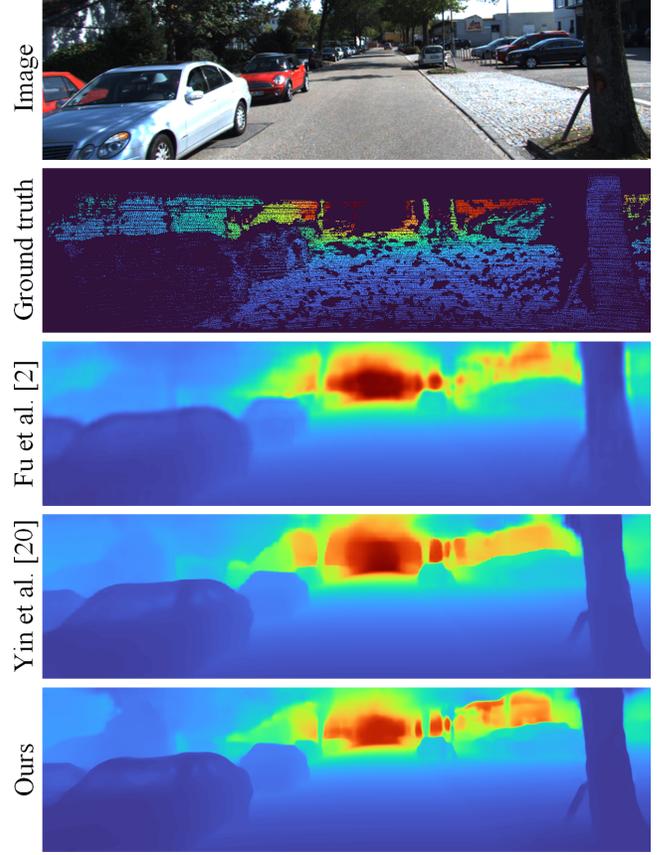


Fig. 5. Visualization of the different methods and our proposed method on KITTI dataset.

in similar fields. To cover gradients at multi-scale, we utilize 5 different spacings $s \in \{1, 2, 4, 8, 16\}$.

$$L_g(g) = \frac{1}{T} \sum_{s \in \{1,2,4,8,16\}} \sum_{i,j} \|\tilde{g}_s(i,j) - g_s(i,j)\|_2 \quad (7)$$

$$g_s(i,j) = \left(\frac{d_{i+s,j} - d_{i,j}}{|d_{i+s,j} + d_{i,j}|}, \frac{d_{i,j+s} - d_{i,j}}{|d_{i,j+s} + d_{i,j}|} \right)^T \quad (8)$$

where \tilde{g}_s and g_s refer to gradient pixel (i,j) in predicted depth map and ground truth respectively. $d_{i,j}$ is the depth value of pixel (i,j).

Total loss. We find that appropriately scaling the range of loss function can accelerate convergence and improve the final predicted result. Our total loss function is defined as follows:

$$L_{total} = \alpha \sqrt{L_d} + \beta \sqrt{L_g} \quad (9)$$

where α and β are constants we set to 10 and 2 for all experiments.

4. EXPERIMENTS

4.1. KITTI Dataset

The KITTI dataset [20] consists of 61 outdoor scene images, each with a resolution of 375×1241. Since previous work is

Table 1. Quantitative results on KITTI using Eigen split.

Method	Accuracy Metric(higher is better)			Error Metric(lower is better)			
	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$	Abs Rel	Sq Rel	RMSE	RMSE(log)
Saxena et al. [12]	0.601	0.820	0.926	0.280	3.012	8.734	0.361
Eigen et al. [1]	0.692	0.899	0.967	0.190	1.515	7.156	0.270
Eigen et al. [13]	0.769	0.950	0.988	0.158	1.210	6.410	0.214
Alhashim et al. [7]	0.886	0.965	0.986	0.093	0.589	4.170	0.171
Fu et al. [2]	0.932	0.984	0.994	0.072	0.307	2.727	0.120
Yin et al. [19]	0.938	0.990	0.998	0.072	-	3.258	0.117
Ours	0.957	0.994	0.999	0.061	0.236	2.699	0.096

Table 2. Experimental results using KITTI Eigen split with various backbone networks.

Variant	#Params	Accuracy Metric(higher is better)			Error Metric(lower is better)			
		$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$	Abs Rel	Sq Rel	RMSE	RMSE(log)
DenseNet-161 [5]	46.6M	0.955	0.993	0.998	0.065	0.251	2.788	0.096
ResNet-101 [4]	68.0M	0.956	0.993	0.999	0.063	0.242	2.721	0.097
ResNext-101 [6]	112.3M	0.957	0.994	0.999	0.061	0.236	2.699	0.096

based on the training set and test set divided by Eigen et al. [1], we also follow it to compare with those works. The training set contains 23488 images from 32 different scenes, and the test set contains 697 images from 29 scenes. The maximum depth of the image in the KITTI dataset is 80.

4.2. Implementation Details

We implement the proposed network based on public deep learning framework PyTorch. In training phase, we use Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\varepsilon = 10^{-6}$. The learning strategy applies polynomial decay with initial learning rate $lr = 10^{-4}$ and power $p = 0.9$. We train our network on two NVIDIA TITAN RTX GPU with 24GB memory. Epoch is set to 50 with batch size 32, which applies to all experiments of this work. As the backbone network for encoder, we use ResNet [4], ResNext [6] and DenseNet [5] pre-trained in ILSVRC dataset. Upconvolution in decoder uses the bilinear neighbor upsampling followed by convolutional layer. Downsampling operation and upsampling operation in spatial pyramid attention module utilizes the nearest neighbor method. We set reduction ration r is 16. To improve training results, we augment images before input to the network using random rotating, random horizontal flipping, random brightness, contrast and color changing. We randomly crop the image size 352×704 to train the network.

4.3. Evaluation Metrics

We quantitatively compare our network with state-of-the-art methods both using the following commonly used metrics:

– Accuracy with Threshold t : $\delta = \max(\frac{d_i}{d_i}, \frac{d_i}{d_i}) < t$, for $t \in \{1.25, 1.25^2, 1.25^3\}$

– Absolute Relative Error (Abs Rel): $\frac{1}{N} \sum_{i=1}^N \frac{|\tilde{d}_i - d_i|}{d_i}$

– Squared Relative Error (Sq Rel): $\frac{1}{N} \sum_{i=1}^N \frac{\|\tilde{d}_i - d_i\|^2}{d_i}$

– Root Mean Squared Error (RMSE): $\sqrt{\frac{1}{N} \sum_{i=1}^N \|\tilde{d}_i - d_i\|^2}$

– Root Mean Squared Error in log space (RMSElog): $\sqrt{\frac{1}{N} \sum_{i=1}^N \|\log(\tilde{d}_i) - \log(d_i)\|^2}$

where N is total number of pixels that the ground truth values are available. \tilde{d}_i and d_i are predicted depth values and ground truth for pixel i .

4.4. Evaluation Results

Table 1 shows quantitative results compared with the state-of-the-art methods. Our network far outperforms all existing methods. As shown in Fig. 5, our method shows much more precise object boundaries and much more continuous object surfaces. To prove the effectiveness of our proposed method, we utilize various backbone network as encoder, and keep other settings. Table 2 provides experimental results. And the results show that ResNext-101 achieve state-of-the-art result.

4.5. Ablation Study

To investigate the importance of different modules in our method, we conduct the ablation study. It can be seen from Table 3 that the proposed model contains all modules (i.e. DCAM, SPAM, scale-invariant gradient loss) to achieve the best performance, which demonstrates that all modules are necessary to get the best monocular depth estimation result.

5. CONCLUSION

In this paper, we present a novel monocular depth estimation network named Pyramid Feature Attention Network to exploit depth information from different levels and address am-

Table 3. Result from the ablation study using KITTI dataset. Baseline: a network composed of only encoder, Dense ASPP, convolution layer and upconvolution layer; L: the network introduce scale-invariant gradient loss. All methods use DenseNet-161 as encoder.

Method	Accuracy Metric(higher is better)			Error Metric(lower is better)			
	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$	Abs Rel	Sq Rel	RMSE	RMSE(log)
baseline	0.928	0.981	0.992	0.086	0.338	3.437	0.158
baseline + DCAM	0.942	0.989	0.996	0.070	0.293	3.015	0.112
baseline + SPAM	0.945	0.990	0.997	0.068	0.253	2.841	0.106
baseline + DCAM + SPAM	0.951	0.992	0.998	0.065	0.251	2.810	0.098
baseline + DCAM + SPAM + L	0.955	0.993	0.998	0.063	0.251	2.788	0.096

biguous object boundaries and discontinuous object surface issues. This network includes two critical modules : Dual-scale Channel Attention Module and Spatial Pyramid Attention Module, which are utilized to improve high-level context features and low-level spatial features, respectively. We also introduce scale-invariant gradient loss for better results. Extensive experimental results on KITTI dataset show that our method outperforms state-of-the-art methods.

6. REFERENCES

- [1] David Eigen, Christian Puhrsch, and Rob Fergus, “Depth map prediction from a single image using a multi-scale deep network,” in *NIPS*, 2014, pp. 2366–2374.
- [2] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao, “Deep ordinal regression network for monocular depth estimation,” in *CVPR*, 2018, pp. 2002–2011.
- [3] 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,” *CoRR*, vol. abs/1907.10326, 2019.
- [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [5] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger, “Densely connected convolutional networks,” *CVPR*, pp. 2261–2269, 2017.
- [6] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He, “Aggregated residual transformations for deep neural networks,” in *CVPR*, 2017, pp. 5987–5995.
- [7] Ibraheem Alhashim and Peter Wonka, “High quality monocular depth estimation via transfer learning,” *CoRR*, vol. abs/1812.11941, 2018.
- [8] Shubhra Aich, Jean Marie Uwabeza Vianney, Md Amirul Islam, Mannat Kaur, and Bingbing Liu, “Bidirectional attention network for monocular depth estimation,” *CoRR*, vol. abs/2009.00743, 2020.
- [9] Tian Chen, Shijie An, Yuan Zhang, Chongyang Ma, Huayan Wang, Xiaoyan Guo, and Wen Zheng, “Improving monocular depth estimation by leveraging structural awareness and complementary datasets,” *ECCV*, 2020.
- [10] Maoke Yang, Kun Yu, Chi Zhang, Zhiwei Li, and Kuiyuan Yang, “Denseaspp for semantic segmentation in street scenes,” in *CVPR*, 2018, pp. 3684–3692.
- [11] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox, “Demon: Depth and motion network for learning monocular stereo,” in *CVPR*, 2017, pp. 5622–5631.
- [12] A. Saxena, Sung H. Chung, and A. Ng, “Learning depth from single monocular images,” in *NIPS*, 2005.
- [13] David Eigen and Rob Fergus, “Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture,” in *ICCV*, 2015, pp. 2650–2658.
- [14] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab, “Deeper depth prediction with fully convolutional residual networks,” in *3DV*, 2016, pp. 239–248.
- [15] Jianbo Jiao, Y. Cao, Yibing Song, and R. Lau, “Look deeper into depth: Monocular depth estimation with semantic booster and attention-driven loss,” in *ECCV*, 2018.
- [16] Jie Hu, Li Shen, and Gang Sun, “Squeeze-and-excitation networks,” in *CVPR*, 2018, pp. 7132–7141.
- [17] S. Woo, Jongchan Park, Joon-Young Lee, and In-So Kweon, “Cbam: Convolutional block attention module,” in *ECCV*, 2018.
- [18] Qiangchang Wang, Tianyi Wu, He Zheng, and Guodong Guo, “Hierarchical pyramid diverse attention networks for face recognition,” in *CVPR*, 2020, pp. 8323–8332.
- [19] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan, “Enforcing geometric constraints of virtual normal for depth prediction,” in *ICCV*. 2019, pp. 5683–5692, IEEE.
- [20] Andreas Geiger, Philip Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The kitti dataset,” *The International Journal of Robotics Research*, vol. 32, pp. 1231–1237, 2013.