# Aerial-PASS: Panoramic Annular Scene Segmentation in Drone Videos

Lei Sun[1], Jia Wang[1], Kailun Yang[2], Kaikai Wu[1], Xiangdong Zhou[1], Kaiwei Wang[1] and Jian Bai[1]

*Abstract*— Aerial pixel-wise scene perception of the surrounding environment is an important task for UAVs (Unmanned Aerial Vehicles). Previous research works mainly adopt conventional pinhole cameras or fisheye cameras as the imaging device. However, these imaging systems cannot achieve large Field of View (FoV), small size, and lightweight at the same time. To this end, we design a UAV system with a Panoramic Annular Lens (PAL), which has the characteristics of small size, low weight, and a $360°$ annular FoV. A lightweight panoramic annular semantic segmentation neural network model is designed to achieve high-accuracy and real-time scene parsing. In addition, we present the first drone-perspective panoramic scene segmentation dataset Aerial-PASS, with annotated labels of track, field, and others. A comprehensive variety of experiments shows that the designed system performs satisfactorily in aerial panoramic scene parsing. In particular, our proposed model strikes an excellent trade-off between segmentation performance and inference speed suitable, validated on both public street-scene and our established aerial-scene datasets.

## I. INTRODUCTION

In the last years, UAV (Unmanned Aerial Vehicle) systems have become relevant for applications in military recognition, civil engineering, environmental surveillance, rice paddy remote sensing, and spraying, etc. [1], [2]. Compared with classic aerial photography and ground photography, the UAV system is more flexible, small in size, low-cost, and suitable for a wider range of application scenarios. Environment perception algorithms based on UAV system needs to be light and efficient enough for the application in mobile computing devices like portable embedded GPU processors.

However, most optical lens of the UAV systems have a small field of view, and often rely on a complex servo structure to control the pitch of the lens, and post-stitch the collected images to obtain a $360°$ panoramic image [2]. The expansion of the field of view is essential for real-time monitoring with UAVs. In traditional methods, the control system of UAV is usually very complicated. Since the drone takes images during flight, the post-stitching algorithm is highly computation-demanding, and the images have parallax and exposure differences, which renders the reliability of image stitching rather low. To address this problem, we have designed a lightweight Panoramic Annular Lens (PAL)

Fig. 1. Overview of our Aerial-PASS system. Panoramic images captured from the UAV with a PAL camera are unfolded, and with the designed semantic segmentation model, field and track classes are predicted at the pixel level.

especially suitable for UAV systems. The system does not require a complicated servo system to control the attitude of the lens. The optical axis of the lens is placed vertically on the ground, and the cylindrical field of view can be horizontally upward by $10°$ and downward by $60°$ (Fig. 1).

To support fast on-board remote sensing, we further propose a lightweight real-time semantic segmentation network for panoramic image segmentation. The network has an efficient up-sampling module and a multi-scale pooling module for learning objects of different scales in the panoramic images. To facilitate credible evaluation, we collect with our PAL-UAV system and present the first drone-perspective panoramic scene segmentation benchmark Aerial-PASS with annotated labels of critical field sensing categories. We find a superior balance between accuracy and inference speed for the network towards efficient aerial image segmentation and it also achieves the state-of-the-art real-time segmentation

performance on the popular *Cityscapes* dataset [3].

The contributions of this paper are summarized as follows:

- The designed PAL lens has the advantages of 4*K* high resolution, large field of view, miniaturization design, and real-time imaging capabilities, etc. It can be applied to UAV survey and identification, autonomous driving, security monitoring, and other fields.
- An efficient real-time semantic segmentation network is proposed for panoramic images and it achieves a state-of-the-art accuracy-speed trade-off on *Cityscapes*. A set of comparison experiments is conducted on the panoramic images collected by our PAL-UAV system.
- An annotated Aerial Panoramic dataset is presented for the first time, which is conducive to the rapid and robust segmentation of target objects in a large field of view from UAV perspectives. Particularly, this work focuses on pixel-wise segmentation of track and field objects.

## II. RELATED WORKS

### A. Panoramic Annular Imaging

The optical lens in the drone provides an indispensable visual aid for target recognition and detection. A single large field of view lens can bring a wider field of view to the drone. The current ultra wide-angle optical systems include fisheye lenses [4], catadioptric optical systems [5], and Panoramic Annular Lens (PAL) imaging systems [6]. The structure of the PAL is simpler, smaller in size, easy to carry, and better in imaging quality. These advantages make the PAL become the research focus of large field of view optical systems.

The PAL imaging system was proposed by Greguss in 1986 [7]. In 1994, Powell designed the infrared band large field of view PAL imaging system, and showed that the spherical surface type can be used to obtain better imaging quality [8]. A cluster of scientific research teams has also made great research progress, designing PAL optical systems with longer focal length [9], larger field of view [10], higher resolution [11], and better imaging quality [12]. The PAL designed in this work is small in size, easy to carry, and it has a field of view of $(30° − 100°) × 360°$ and a 4*K* high resolution, which can realize real-time large field of view high-resolution single sensor imaging. Thereby, it is perfectly suitable for UAV survey, unmanned driving, security monitoring, and other application fields [13], [14].

### B. Panoramic Scene Segmentation

Beginning with the success of Fully Convolutional Networks (FCNs) [15], semantic segmentation can be achieved in an end-to-end fashion. Subsequent composite architectures like PSPNet [16] and DeepLab [17] have attained remarkable parsing performance. Yet, due to the use of large backbones like ResNet-101 [18], top-accuracy networks come with prohibitively high computational complexity, which are not suitable for mobile applications such as driving- and aerial-scene segmentation in autonomous vehicles or drone videos. Plenty of light-weight networks emerge such as SwiftNet [19], AttaNet [20], and DDRNet [21], both seeking a fast and precise segmentation. In our previous works, we have leveraged efficient networks for road scene parsing applications like nighttime scene segmentation [22] and unexpected obstacle detection [23].

Driving and aerial scene segmentation clearly benefit from expanded FoV, and standard semantic segmentation on pinhole images has been extended to operate on fisheye images [24], omnidirectional images [25], and panoramic annular images [26]. The latest progress include omni-supervised learning [27] and omni-range context modeling [28] to promote efficient 360° driving scene understanding. In particular, Yang *et al.* [26] proposed a generic framework by using a panoramic annular lens system, which intertwines a network adaptation method by re-using models learned on standard semantic segmentation datasets. Compared to driving scene segmentation, aerial image segmentation is still predominantly based on pinhole images [29]–[32]. In this paper, we lift 360° panoramic segmentation to drone videos and propose an Aerial-PASS system to explore the superiority of the ultra-wide FoV for security applications.

## III. PROPOSED SYSTEM

The proposed system consists of the hardware part and the algorithm part. In the hardware part, we have equipped the UAV with our designed PAL camera and use it to collect panoramic image data. To efficiently parse panoramic images, we have designed a lightweight semantic segmentation model with a novel multi-scale receptive field pyramid pooling module for learning a robust feature representation required for 360° image segmentation.

### A. UAV with PAL Camera

The PAL system designed and used in this work has the characteristics of large field of view and lightweight structure, which can realize real-time, single-sensor imaging on drones. The PAL system follows the imaging principle of Flat Cylindrical Perspective (FCP), which can reflect light by the panoramic block from the lateral field of view around the optical axis 360°, and then enter the subsequent lens group and image on the two-dimensional image surface. The field of view of the PAL system is generally greater than 180°, which is no longer applicable to the classic principle of object-image similarity. At this point, we introduce negative distortion in the optical design to control the height of the image surface, and the commonly used characterization method is F-Theta distortion.

The PAL system designed in this work is composed of 10 standard spherical lenses, and the glass materials are all from the CDGM company. The PAL block is composed of two glass materials glued together. By coating the transmission film and reflection film on its surface, the light path can be folded. Its structure diagram is shown in Fig. 2.

We have chosen the Inspire 2 series drone module developed by DJI to be equipped with a PAL, which is placed vertically downward to cover a wider field of view on the ground and to avoid the stray light problem caused by direct sunlight. The drone system is flying around 100*m* above the track and field to collect image data from multiple directions.
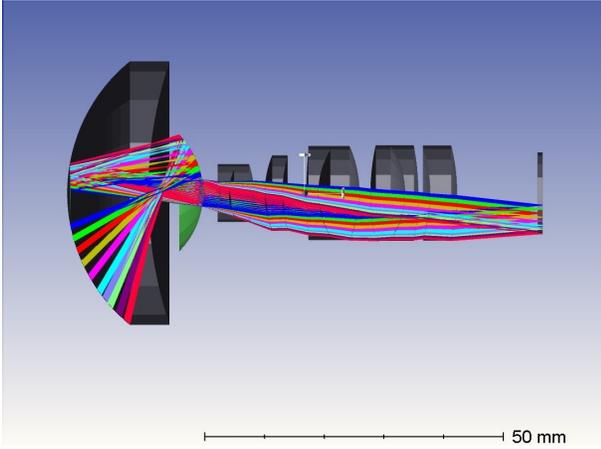
Fig. 2. Shaded model of the designed PAL imaging system. It consists of 6 groups of 10 lenses. Different colors represent light in different field of view.
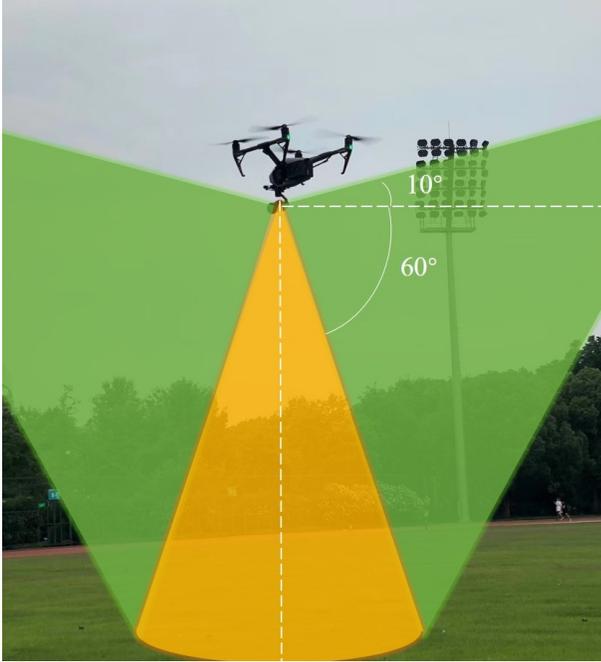


Fig. 3. The schematic diagram of the drone system. The green area represents the imaging area and the FoV is $(30° \sim 100°) \times 360°$. The yellow area represents the blind area.

The schematic diagram of the experiment is shown in Fig. 3. The lateral field of view of the PAL system involved in imaging is $10°$ horizontally upward and $60°$ horizontally downward, and the overall field of view is $360° \times 70°$.

### B. Aerial Scene Segmentation Model

For our segmentation method, there are three requirements. The model must be light enough to meet the real-time inference speed demand for future migration to portable computing devices. In addition, The model needs to parse the image with a high accuracy. Further, the model should have multi-scale receptive fields for panoramic images with a ultra-wide angle. Inspired by SwfitNet [19] and RFNet [23],

we have designed a lightweight novel U-Net like model with multi-scale receptive fields.

*1) Model Architecture:* The proposed network architecture is shown in Fig. 4. We adopt ResNet-18 [18] as our backbone, which is a mainstream light-weight feature extraction network. With ImageNet [33] pre-trained weights, we can benefit from regularization induced by transfer learning and small operation footprint for real-time prediction. The feature maps after each layer of ResNet is fused with the feature maps in the upsampling modules through skip connection with $1 \times 1$ convolution modules. To increase the receptive field, the feature maps are transmitted to Efficient Deep Aggregation Pyramid Pooling module. In the decoder, feature maps are added with lateral features from an earlier layer of the encoder element-wisely, and after convolution, the blended feature maps are upsampled by bilinear interpolation. Through skip connections, high-resolution feature maps full of detail information are blended with low-resolution feature maps with rich semantic information.
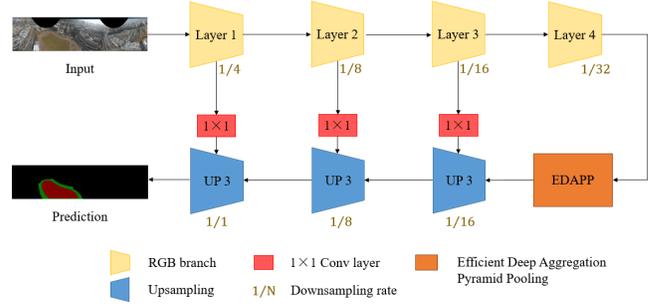


Fig. 4. The architecture of the proposed segmentation model.

*2) Efficient Deep Aggregation Pyramid Pooling:* For aerial panoramic images, many objects are rather small and only take up little ratio of pixels. Receptive field is extremely important in the task for a fine-grain segmentation. Inspired by the context extraction module in DDRNet [21], we develop an Efficient Deep Aggregation Pyramid Pooling (EDAPP) module. Fig. 5 shows the architecture of EDAPP. First, we perform average pooling with kernel size of 5, 9, 17, and global pooling respectively. Single $3 \times 3$ or $1 \times 1$ convolutions in Spatial Pyramid Pooling (SPP) [34] is not enough, so after $1 \times 1$ convolution upsampling to the same resolution, to efficiently blend the multi-scale contextual information better, we propose to leverage a combination of $3 \times 1$ convolution and $1 \times 3$ convolution. Another stream consists of only a $1 \times 1$ convolution. Asserting an input $x$, the process can be summarized in the following:

$$y_k = \begin{cases} C_{1\times1}(x), & k=1; \\ C_{1\times3}(C_{3\times1}(UP(C_{1\times1}(P_{2^k+1,2^{k-1}}(x)))+y_{k-1}), & 1<k<n; \\ C_{3\times3}(UP(C_{1\times1}(Pglobal(x)))+y_{k-1}), & k=n. \end{cases} \quad (1)$$

Here, $C$ denotes convolution, $UP$ denotes bilinear upsample, $P$ and $P_{global}$ denote pooling and global pooling respectively. $i$ and $j$ of the $P_{i,j}$ denote the kernel size and stride of the pooling layer. A pair of $3 \times 1$ convolution and
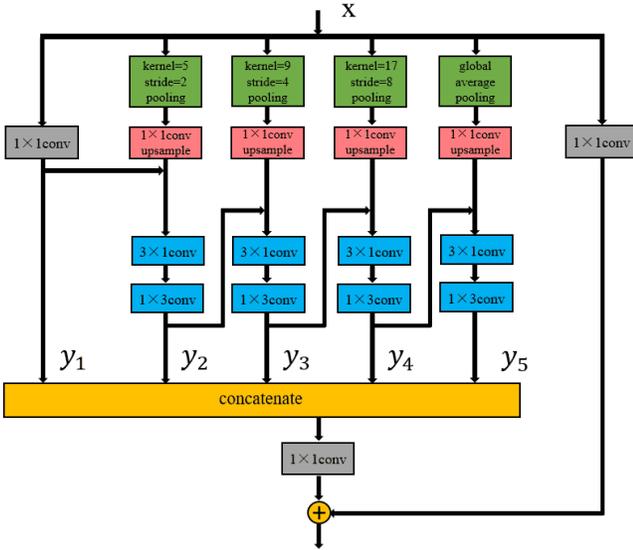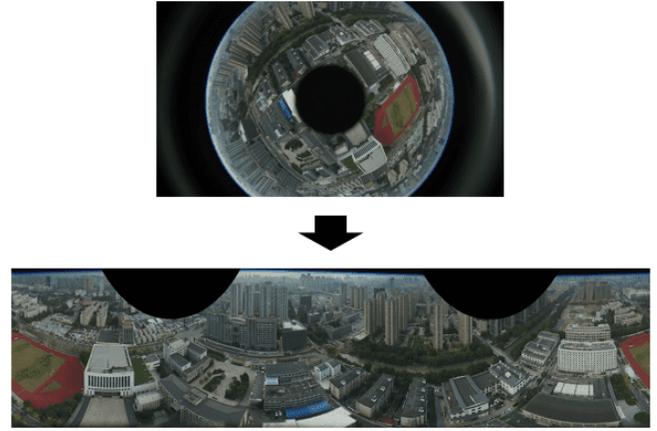
Fig. 5. The architecture of the EDAPP module.



Fig. 6. The unfolding process of the PAL image. Limited to the size of CMOS, top and bottom sides of the imaging plane are blocked in the image, resulting in two scalloped shadows in the unfolded image.

$1 \times 3$ convolution have the same receptive field as a single $3 \times 3$ kernel, enhanced directional feature extraction, but less computational complexity and faster inference speed. This architecture helps extracting and integrating deep information with different scales by different pooling kernel sizes. All the features maps are concatenated and blended through a $1 \times 1$ convolution. Finally, a $1 \times 1$ convolution compress all the feature maps and after that we also add a skip connection with a $1 \times 1$ convolution for easier optimization.

## IV. EXPERIMENTS

### A. Aerial-PASS Dataset

We collected data in 4 different places using the UAV with a PAL camera, and all the data were under sufficient illumination conditions. In the height of about 100 meters, we collected videos with the length of about 3 hours. Limited to the size of the CMOS in the camera, the image produced can't show all the imaging plane of the lens. In the following deployment phase, the PAL system was calibrated using the interface provided by the omnidirectional camera toolbox [35]. Before training, the PAL image was unfolded to a familiar rectangle image. The unfolded process is depicted in the following equations:

$$i = \frac{r - r_1}{r_2 - r_1} \times height \qquad (2)$$

$$j = \frac{\theta}{2\pi} \times width \qquad (3)$$

Here, $i$ and $j$ denote the index of x and y axis of the unfold image, respectively. $r_1$ and $r_2$ are the internal and external radii of the PAL image. Width and height are the image size of the unfolded image. In our experiment, we unfolded the PAL image to a $2048 \times 512$ image. Fig. 6 shows the unfolding process.

We annotated all 462 images out of the 3-hour-long video. We created pixel-wise fine labels on the most critical classes relevant to the application of track detection, and we randomly split out 42 images for the test set. As far as we know, This is the first aerial panoramic image dataset with semantic segmentation annotations in the community.

### B. Training Details

All the experiments were implemented with PyTorch 1.3 on a single 2080Ti GPU with CUDA10.0, cuDNN 7.6.0. We chose Adam [36] for optimization with a learning rate of $5 \times 10^{-4}$, where cosine annealing learning rate scheduling policy [37] is adopted to adjust the learning rate with a minimum value of $5 \times 10^{-4}$ in the last epoch and weight decay was set to $1 \times 10^{-4}$. The ResNet-18 [18] backbone was initialized with pre-trained weights from ImageNet [33] and the rest part of the model was initialized with the Kaiming initialization method [38]. We updated the pre-trained parameters with 4 times smaller learning rate and weight decay rate. The data augment operations consist of scaling with random factors between 0.5 and 2, random horizontal flipping, and random cropping with an output resolution of $512 \times 512$. Models were trained for 100 epochs with a batch-size of 6. We used the standard "Intersection over Union (IoU)" metric for the evaluation:

$$IoU = \frac{TP}{TP + FP + FN} \qquad (4)$$

### C. Results and Analysis

Based on the Aerial-PASS dataset, we have created a benchmark to compare our proposed network with other two competitive real-time semantic segmentation networks: the single-scale SwiftNet [19] (similar backbone with our network) and ERF-PSPNet [39] (a lightweight network designed for panoramic segmentation [26]). All the networks were trained on the training set of Aerial-PASS with the same training strategy and tested on the testing set of the dataset. Table I shows the numerical performance comparisons of the three efficient networks. Our proposed model outperforms
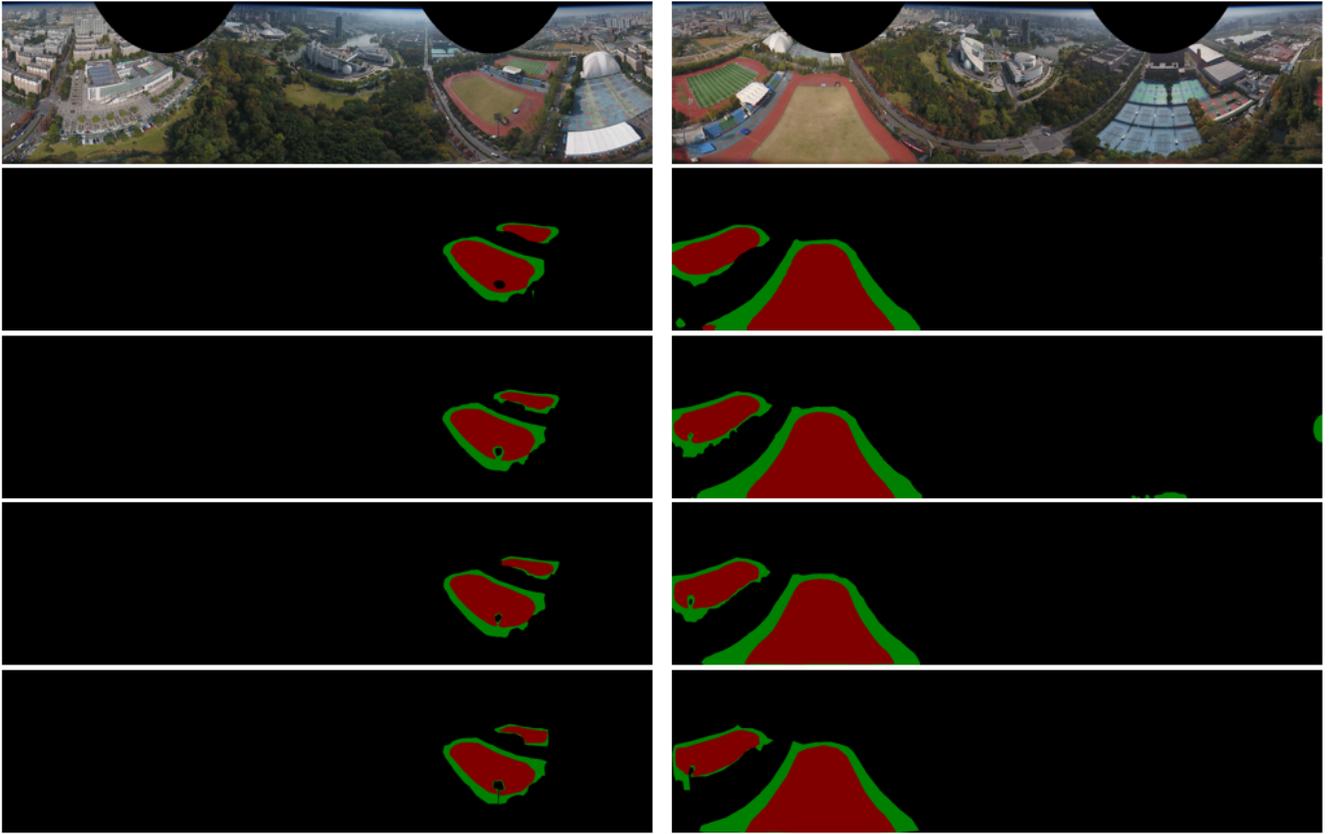
Fig. 7. Qualitative semantic segmentation results. From top to bottom row: RGB input image, ERF-PSPNet, SwiftNet, our method, and ground truth.

| Network | Track | Field | Others | Mean |
|---|---|---|---|---|
| ERF-PSPNet | 64.16% | 97.67% | 92.02% | 84.62% |
| SwiftNet (single scale) | 63.15% | 98.76% | 91.63% | 84.52% |
| **Ours** | 67.67% | 99.06% | 92.16% | 86.30% |

both networks designed for panoramic segmentation (ERF-PSPNet) and semantic segmentation (SwifNet) by clear gaps.

Fig. 7 shows visualizations of inference labels of our proposed method and other two models, in which green denotes the track and red denotes the field. All the input images are the unfolded PAL images. The labels show the qualitative result of the proposed method. As we can find that our method performs well in both large-scale objects like field and small-scales objects like the boundary of track and other objects thanks to our EDAPP module designed for multi-scale feature learning.

### D. Comparison with the State of the Art

To further compare with other state-of-the-art network, we also trained our network on *Cityscapes* [3], which is a large-scale RGB dataset that focuses on semantic understanding of urban street scenes. It contains 2975/500/1525 images in the training/validation/testing subsets, both with finely annotated labels on 19 classes. The images cover 50 different cities with a full resolution of $2048 \times 1024$. We trained our network on the training set of the *Cityscapes* dataset and test our network on the validation set. Table II shows the IoU result of our method and other mainstream real-time semantic segmentation models. Our method has achieved an excellent balance between accuracy and inference speed. Fig. 8 shows some representative inference results of our proposed method. Overall, the qualitative results verify the generalization capacity of our proposed network for both challenging large-FoV aerial image segmentation and high-resolution driving scene segmentation.

| Network | MIoU | Speed (FPS) |
|---|---|---|
| FCN8s [15] | 65.3% | 2.0* |
| DeepLabV2-CRF [17] | 70.4% | n/a |
| ENet [40] | 58.3% | 76.9* |
| ERF-PSPNet [39] | 64.1% | 20.4 |
| SwiftNet [19] | 72% | 41.0 |
| Ours | 72.8% | 39.4 |

* Speed on half resolution images.

## V. CONCLUSION

In this study, we propose a lightweight UAV system Aerial-PASS with a designed Panoramic Annular Lens (PAL) camera and a real-time semantic segmentation network for
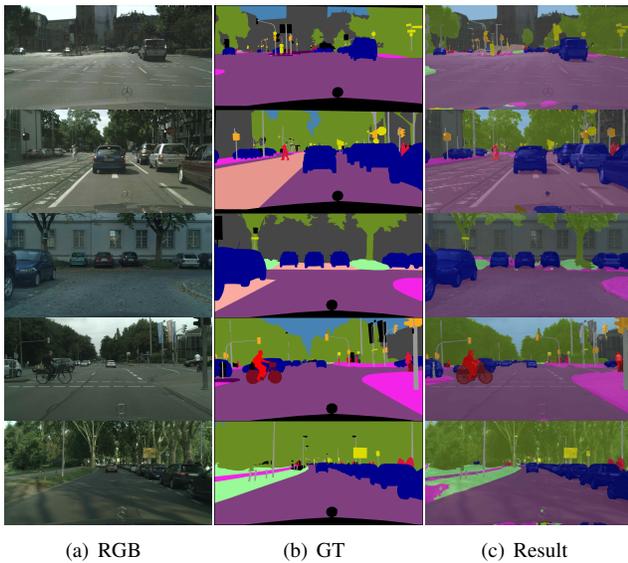
(a) RGB      (b) GT      (c) Result

Fig. 8. Qualitative results of our method on the *Cityscapes* dataset.

aerial panoramic image collection and segmentation. The minimization-dedicated PAL camera equipped in the UAV can be used for collecting annular panoramic images without requiring a complicated servo system to control the attitude of the lens. To classify the track and field in the images at the pixel wise, we collect and annotate 462 images and propose an efficient semantic segmentation network. The proposed network has multi-scale reception fields and an efficient backbone, which outperforms other competitive networks on our Aerial-PASS dataset and also has reached the state-of-the-art performance on the *Cityscapes* dataset with 39.4 Hz in full resolution on a single 2080Ti GPU processor. In the future, we aim to transplant the algorithm to the portable embedded GPU on the UAV for more field tests.

## REFERENCES

[1] S. M. Adams and C. J. Friedland, "A survey of unmanned aerial vehicle (UAV) usage for imagery collection in disaster research and management," in *International Workshop on Remote Sensing for Disaster Response*, 2011.

[2] S. D'Oleire-Oltmanns, I. Marzolff, K. D. Peter, and J. B. Ries, "Unmanned aerial vehicle (UAV) for monitoring soil erosion in morocco," *Remote Sensing*, 2012.

[3] M. Cordts *et al.*, "The cityscapes dataset for semantic urban scene understanding," in *CVPR*, 2016.

[4] Y. Xiong and K. Turkowski, "Creating image-based VR using a self-calibrating fisheye lens," in *CVPR*, 1997.

[5] S. Baker and S. K. Nayar, "A theory of catadioptric image formation," in *ICCV*, 1998.

[6] D. Lehner, A. Richter, D. Matthys, and J. Gilbert, "Characterization of the panoramic annular lens," *Experimental Mechanics*, 1996.

[7] P. Greguss, "Panoramic imaging block for three-dimensional space," Jan. 28 1986. US Patent 4,566,763.

[8] I. Powell, "Panoramic lens," *Applied Optics*, 1994.

[9] S. Niu, J. Bai, X.-y. Hou, and G.-g. Yang, "Design of a panoramic annular lens with a long focal length," *Applied Optics*, 2007.

[10] X. Zhou, J. Bai, C. Wang, X. Hou, and K. Wang, "Comparison of two panoramic front unit arrangements in design of a super wide angle panoramic annular lens," *Applied Optics*, 2016.

[11] J. Wang, X. Huang, J. Bai, K. Wang, and Y. Hua, "Design of high resolution panoramic annular lens system," in *Optical Sensing and Imaging Technology*, 2019.

[12] Q. Zhou, Y. Tian, J. Wang, and M. Xu, "Design and implementation of a high-performance panoramic annular lens," *Applied Optics*, 2020.

[13] Y. Fang, K. Wang, R. Cheng, and K. Yang, "CFVL: A coarse-to-fine vehicle localizer with omnidirectional perception across severe appearance variations," in *IEEE IV*, 2021.

[14] H. Chen, W. Hu, K. Yang, J. Bai, and K. Wang, "Panoramic annular SLAM with loop closure and global optimization," *arXiv preprint arXiv:2102.13400*, 2021.

[15] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015.

[16] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017.

[17] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

[18] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[19] M. Oršic, I. Krešo, P. Bevandic, and S. Šegvic, "In defense of pre-trained ImageNet architectures for real-time semantic segmentation of road-driving images," in *CVPR*, 2019.

[20] Q. Song, K. Mei, and R. Huang, "AttaNet: Attention-augmented network for fast and accurate scene parsing," in *AAAI*, 2021.

[21] Y. Hong, H. Pan, W. Sun, and Y. Jia, "Deep dual-resolution networks for real-time and accurate semantic segmentation of road scenes," *arXiv preprint arXiv:2101.06085*, 2021.

[22] L. Sun, K. Wang, K. Yang, and K. Xiang, "See clearer at night: Towards robust nighttime semantic segmentation through day-night image conversion," in *SPIE*, 2019.

[23] L. Sun, K. Yang, X. Hu, W. Hu, and K. Wang, "Real-time fusion network for RGB-D semantic segmentation incorporating unexpected obstacle detection for road-driving images," *IEEE RA-L*, 2020.

[24] Y. Ye, K. Yang, K. Xiang, J. Wang, and K. Wang, "Universal semantic segmentation for fisheye urban driving images," in *IEEE SMC*, 2020.

[25] A. R. Sekkat, Y. Dupuis, P. Vasseur, and P. Honeine, "The omniscape dataset," in *IEEE ICRA*, 2020.

[26] K. Yang, X. Hu, L. M. Bergasa, E. Romera, and K. Wang, "PASS: Panoramic annular semantic segmentation," *IEEE T-ITS*, 2020.

[27] K. Yang, X. Hu, Y. Fang, K. Wang, and R. Stiefelhagen, "Omnisupervised omnidirectional semantic segmentation," *IEEE T-ITS*, 2020.

[28] K. Yang, J. Zhang, S. Reiß, X. Hu, and R. Stiefelhagen, "Capturing omni-range context for omnidirectional segmentation," in *CVPR*, 2021.

[29] M. T. Chiu *et al.*, "Agriculture-vision: A large aerial image database for agricultural pattern analysis," in *CVPR*, 2020.

[30] L. Mou, Y. Hua, and X. X. Zhu, "A relation-augmented fully convolutional network for semantic segmentation in aerial scenes," in *CVPR*, 2019.

[31] S. M. Azimi, P. Fischer, M. Körner, and P. Reinartz, "Aerial LaneNet: Lane-marking semantic segmentation in aerial imagery using wavelet-enhanced cost-sensitive symmetric fully convolutional neural networks," *IEEE Transactions on Geoscience and Remote Sensing*, 2019.

[32] X. Li *et al.*, "PointFlow: Flowing semantics through points for aerial image segmentation," in *CVPR*, 2021.

[33] O. Russakovsky *et al.*, "ImageNet large scale visual recognition challenge," *International Journal of Computer Vision*, 2015.

[34] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," in *ECCV*, 2014.

[35] D. Scaramuzza, A. Martinelli, and R. Siegwart, "A toolbox for easily calibrating omnidirectional cameras," in *IEEE IROS*, 2006.

[36] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.

[37] I. Loshchilov and F. Hutter, "SGDR: Stochastic gradient descent with warm restarts," in *ICLR*, 2017.

[38] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification," in *ICCV*, 2015.

[39] K. Yang, L. M. Bergasa, E. Romera, R. Cheng, T. Chen, and K. Wang, "Unifying terrain awareness through real-time semantic segmentation," in *IEEE IV*, 2018.

[40] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, "ENet: A deep neural network architecture for real-time semantic segmentation," *arXiv preprint arXiv:1606.02147*, 2016.