

1 **SEMI-SUPERVISED SEMANTIC SEGMENTATION ON VEHICLE-MOUNTED**  
2 **FISH-EYE CAMERA IMAGES**

3  
4  
5

6 **Sneha Paul, Corresponding Author**

7 Ph.D. Student, Concordia Institute for Information Systems Engineering  
8 Computer Science and Visual Arts Integrated Complex, 1515 St. Catherine W., Montreal, Canada  
9 Email: sneha.paul@mail.concordia.ca

10

11 **Zachary Patterson**

12 Professor, Concordia Institute for Information Systems Engineering  
13 Computer Science and Visual Arts Integrated Complex, 1515 St. Catherine W., Montreal, Canada  
14 Tel: 514-848-2424 ext. 8830, Email: zachary.patterson@concordia.ca

15

16 **Nizar Bouguila**

17 Professor, Concordia Institute for Information Systems Engineering  
18 Computer Science and Visual Arts Integrated Complex, 1515 St. Catherine W., Montreal, Canada  
19 Tel: 514-48-2424 ext. 5663, Email: nizar.bouguila@concordia.ca

20

21

22 Word Count: 4877 words + 7 table(s)  $\times$  250 = 6627 words

23

24

25

26

27

28

29 Submission Date: May 14, 2026

**1 ABSTRACT**

2 The use of large Field of View (FoV) cameras employing fish-eye lenses benefits various computer  
3 vision applications in autonomous driving. Although deep learning has been successful in various  
4 computer vision applications using regular perspective images, its potentiality in fish-eye camera  
5 applications still needs to be explored due to the availability of relatively small datasets for fully-  
6 supervised learning. Semi-supervised learning can be a potential solution to address this challenge.  
7 In this study, we focus on exploring the use of semi-supervised learning for the semantic segmen-  
8 tation of fish-eye images. To the best of our knowledge, this is the first work on semi-supervised  
9 semantic segmentation of fish-eye images. We implement three popular semi-supervised methods  
10 from the perspective image literature for fish-eye image segmentation and conduct extensive anal-  
11 ysis in different settings. Based on the best practices observed in these methods, we propose three  
12 variants of semi-supervised segmentation models for fish-eye images. Our proposed approaches  
13 are evaluated on the WoodScape dataset, which is collected from vehicle-mounted fish-eye cam-  
14 eras. Detailed experiments demonstrate an 8.89% enhancement of the proposed semi-supervised  
15 method compared to fully supervised learning using the same amount of labeled data. We also per-  
16 form an extensive sensitivity analysis to exhibit the effectiveness and advantages of the proposed  
17 method in this study.

18

19 *Keywords:* Semi-supervised learning, Semantic Segmentation, fish-eye images, Autonomous driv-  
20 ing, Autonomous vehicles, AED50 - Committee on Artificial Intelligence and Advanced Comput-  
21 ing.

## 1 INTRODUCTION

2 Currently, various computer vision applications in autonomous driving are benefiting from  
3 the use of large Field of View (FoV) cameras (1). These types of cameras use fish-eye lenses,  
4 also known as wide-angle lenses, to capture a broad view by using extensive non-linear mapping  
5 instead of regular perspective projection. However, this results in a strong radial distortion to the  
6 captured images. One potential solution to the distortion issue is to use transformation techniques  
7 to convert fish-eye images into regular perspective images (2). However, this compromises the rich  
8 feature information contained at the boundary of the images. As a result, the detection process is  
9 hindered. Another possible solution is to use the distorted fish-eye images directly in the model, as  
10 suggested by some existing research (3, 4). However, the current literature on computer vision is  
11 not well-optimized for learning from fish-eye cameras, thus providing a need for further research  
12 in this area.

13 Semantic segmentation, an important computer vision task, focuses on classifying each  
14 pixel in an image and grouping together the pixels that belong to a specific object. Thus it seg-  
15 ments a scene into bounding areas or masks containing different objects. 2D perspective images  
16 are the very first domain to which semantic segmentation was used in the context of autonomous  
17 driving (5). However, the application of segmentation in autonomous driving has been extended to  
18 other domains as well, e.g. panoptic segmentation (6), point cloud (7), depth images (8), thermal  
19 images (9), and so on. Various segmentation datasets on autonomous driving, e.g. Cityscape (10),  
20 KITTI (11), Toronto City (12), ApolloScape (13), WoodScape (1) have leveraged the computer  
21 vision research on autonomous driving. The significance of semantic segmentation in autonomous  
22 driving lies in its ability to enable the vehicle’s sensor systems to precisely comprehend and differ-  
23 entiate various objects and areas in the surroundings. It leads to well-informed decisions based on  
24 specific information about the road, obstacles, pedestrians, traffic signs, and other crucial elements,  
25 ultimately ensuring safer and more authentic autonomous navigation.

26 However, the application of semantic segmentation for autonomous driving on fish-  
27 eye images is not as well-explored as in other domains. Very little research on fish-eye image  
28 segmentation is available in the literature, which has previously explored different approaches  
29 to address the problem of radial distortion in fish-eye images (14–16) due to non-linear mapping.  
30 Most of these researches rely on models trained on synthetic fish-eye images, which are obtained by  
31 applying different fish-eye transformations to the perspective images. However, these approaches  
32 do not perform well when using real-world fish-eye data. Therefore, this research focuses on  
33 learning semantic segmentation from real-world fish-eye data. Moreover, existing literature on  
34 semantic segmentation on fish-eye autonomous driving images focuses largely on multi-modal,  
35 multi-task models (1, 3, 17), gaining an advantage in learning pixel-level representation from dif-  
36 ferent modalities and similar tasks. However, these models are computationally expensive due to  
37 the high number of learnable parameters. To this end, a single-task, uni-modal model will be a po-  
38 tential solution. In addition, the current literature on semantic segmentation focuses on supervised  
39 learning (2, 15, 18), which requires a large amount of labeled data to train deep learning models  
40 effectively.

41 Since semantic segmentation involves classifying each pixel in an image, annotation is  
42 more costly and laborious than other computer vision tasks, such as object classification. Conse-  
43 quently, there is a lack of publicly available semantically annotated fish-eye segmentation datasets.  
44 Semi-supervised learning has emerged as a potential solution to this problem, allowing models  
45 to learn from a small amount of labeled data while leveraging a larger amount of unlabeled data.

1 However, to the best of our knowledge, there is no existing work on semi-supervised semantic  
2 segmentation for fish-eye data. Therefore, in our study, we explore three existing methods for  
3 semi-supervised segmentation (19–21) on fish-eye images and combine the best practices to pro-  
4 pose a new model for fish-eye semantic segmentation. We also perform an in-depth sensitivity  
5 analysis on different model-specific hyperparameters of our proposed methods. We evaluate our  
6 proposed method on the WoodScape dataset (1), which consists of real-world fish-eye images. We  
7 show an 8.89% improvement over fully supervised learning with the same amount of labeled data  
8 by applying the proposed method.

9 The key contribution of this work can be summarized as follows:

- 10 • We present the first study on semi-supervised segmentation in the context of fish-eye  
11 camera images in the context of autonomous vehicles or otherwise.
- 12 • We adopt different popular semi-supervised segmentation methods for this task and com-  
13 pare and benchmark their performance.
- 14 • We also present a detailed sensitivity study on each of the methods in this study for  
15 further understanding of the impact of different components.
- 16 • From the study on existing methods, we combine the best practices and propose a new  
17 method for semi-supervised fish-eye camera segmentation that outperforms all off-the-  
18 shelf methods for this task.

## 19 RELATED WORKS

20 In this section, we will discuss the existing literature on semantic segmentation models for  
21 autonomous driving. Then we will elaborate on the available work in semantic segmentation on  
22 regular perspective and fish-eye images in other domains that are relevant to our work.

### 23 Segmentation for Autonomous Driving

24 Various existing work deals with semantic segmentation for autonomous driving scenarios  
25 from supervised settings. For example, (22) and (23) worked on multi-class semantic segmen-  
26 tation, whereas (24) focused on road segmentation. (25) proposed a fully-convolutional network  
27 (FCN) based architecture, a one-stage pipeline using fully-connected layers to replace the convo-  
28 lutional layers for predicting coarse output maps. These coarse maps are upsampled by deconvolu-  
29 tion to produce dense pixel-level labels. (23) extended this idea by using an encoder-decoder-based  
30 CNN network. However, these networks work in supervised settings, where labeled data is neces-  
31 sary, making semantic segmentation a laborious and costly task.

32 To this end, a great deal of research on unsupervised settings has been proposed, requiring  
33 no (26), or little (27) labeled data. The Asynchronous Teacher-Student Optimization (ATSO)  
34 algorithm (28) used continual learning by partitioning unlabeled training data into two subsets.  
35 One subset of the training data was used for fine-tuning the model that updated the labels of the  
36 other subset. (29) proposed a semi-supervised segmentation network using adversarial learning.  
37 (30) proposed a multi-task, self-supervised segmentation network that is efficient for classification,  
38 segmentation, and detection tasks when labeled data is unavailable. Though these models perform  
39 satisfactorily on perspective images, their performance deteriorates for fish-eye images due to high  
40 radial distortion. Thus in-depth study is necessary on the unsupervised segmentation domain for  
41 fish-eye autonomous driving images.

## 1 **Perspective Image Segmentation**

2 Semantic segmentation is a popular dense prediction task in computer vision that has seen a  
3 significant number of developments in the past few years, including CNN and Vision transformer  
4 (ViTs) based models. Vision Transformers are gaining popularity in Semantic segmentation for the  
5 advantage gained from their special network architecture. These models learn pixel-level feature  
6 representation using a transformer-based encoder along with an attention mechanism and a bottle-  
7 neck layer (31, 32). SegViT (33) stands apart from conventional ViTs due to its use of an Attention-  
8 to-Mask (ATM) module. This module effectively transfers learned similarity maps, derived from  
9 trainable parameters and spatial feature maps, into segmentation masks. They also propose a  
10 shrunk structure of the SegViT model based on a query-based up-sampling (QU) and query-based  
11 down-sampling (QD) module that reduces a model’s computational cost by 40%. SegFormer (34)  
12 is another ViT-based semantic segmentation model that incorporates a transformer-based encoder  
13 with a lightweight full MLP-based decoder. SegFormer does not require positional encoding that  
14 needs to be interpolated, resulting in reduced performance of the model when the resolution of  
15 training and testing data is different. Thus it introduces Mix-FFN, which is a combination of con-  
16 volution block with MLP that can still capture locational information without positional encoding.  
17 None of these SOTA models have been explored in the concept of fish-eye image segmentation.

## 18 **Fish-eye Image Segmentation**

19 There has been a bit of research on fish-eye image segmentation in recent years that ad-  
20 dresses distortion at the edge of the image. One such approach, proposed in (14), uses an Effi-  
21 cient Residual Factorized Network- (ERFNet) based architecture to handle distortion in synthetic  
22 fish-eye data. Another approach, proposed in (15), uses an adaptive deformable CNN-based ar-  
23 chitecture that can adapt to fish-eye images while being trained on pinhole camera images, thus  
24 avoiding the laborious task of data transformation. However, due to the lack of large public fish-  
25 eye datasets, most existing literature converts regular perspective images into fish-eye images. In  
26 (16), a method is proposed to reuse rectilinear images as fish-eye images to address this issue. In  
27 (18), a seven-degree of freedom (DoF) augmentation method is proposed, which uses rectilinear  
28 images during training to transform them into fish-eye images. This method simulates fish-eye  
29 images from rectilinear images captured from different cameras with varying focal lengths, posi-  
30 tions, and orientations. Although the seven DoF augmentation method is effective in increasing  
31 the model’s mean Intersection over Union (mIoU) compared to highly distorted real-world fish-  
32 eye data, it falls short in situations where the distortion is not similar among images, which is a  
33 common criterion in real-world scenarios. As a result, existing models trained on transformed or  
34 synthetic fish-eye data cannot achieve superior performance when using real-world fish-eye data.  
35 Another key observation we find from very recent literature is that multi-modal or multi-task data  
36 is jointly used to learn based on fish-eye images (3, 35). However, such a dataset is costly to collect  
37 and computationally demanding to train. Therefore, the focus of our research is to learn semantic  
38 segmentation from real-world fish-eye data using the image alone.

39 In summary, the current research on semantic segmentation in autonomous driving has pri-  
40 marily focused on different modalities, excluding fish-eye images, particularly in situations where  
41 labeled data is scarce and annotation is expensive and time-consuming. Moreover, there is a notice-  
42 able gap in studies conducted on real-world fish-eye images captured from autonomous vehicles  
43 (AVs). To this end, our work introduces a semi-supervised semantic segmentation approach specifi-  
44 cally designed for real-world fish-eye images within the domain of autonomous driving.

## 1 METHOD

2 In this work, we explore three well-known semi-supervised segmentation methods that were  
 3 originally developed for regular images: CPS (19), MeanTeacher (20), and CPS with CutMix (36).  
 4 An overview of these methods is presented in Figure 1. However, none of these methods are  
 5 specifically tailored for fish-eye image segmentation. To address this limitation, we propose a  
 6 novel semi-supervised method that draws inspiration from existing methods in the literature and  
 7 incorporates best practices discovered during our preliminary studies. Here, we will first provide  
 8 an overview of existing semi-supervised segmentation methods on perspective images that are used  
 9 in this study. Next, we will discuss the proposed method.

### 10 Overview of Existing Semi-supervised Segmentation Methods

11 In (19), a new approach for semi-supervised learning of semantic segmentation was pro-  
 12 posed that learns from the labeled and unlabeled perspective images simultaneously. The approach  
 13 is based on consistency regularization and is called Cross Pseudo Supervision (CPS). CPS aims  
 14 to increase the consistency between two different segmentation networks for similar input im-  
 15 ages. Both labeled and unlabeled data are fed into the two different networks, and for each unlabeled  
 16 image, the model generates a one-hot pseudo-prediction map. This pseudo-prediction map  
 17 is then used to supervise the other segmentation network and vice versa. This way, the pseudo-  
 18 segmentation map is used as an additional signal to provide supervision for the unlabeled data.

19 Formally, CPS generates the prediction as follows:

$$20 P_1 = f(X; \theta_1), \quad (1)$$

$$21 P_2 = f(X; \theta_2). \quad (2)$$

22 Here,  $P_1$  and  $P_2$  are the segmentation prediction obtained from the two segmentation  
 23 networks after applying the softmax function. Note that the two segmentation networks are similar  
 24 in structure but initialized and updated differently. The predictions are then used to generate the  
 25 pseudo-labels. The whole process can be described as follow:

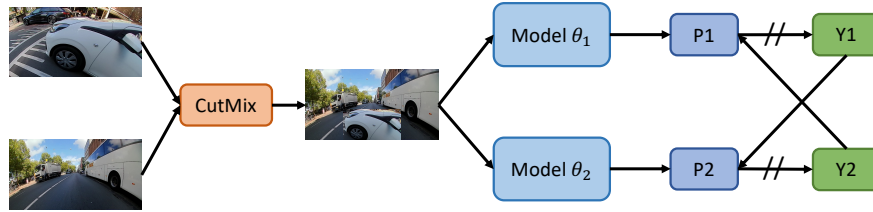
$$26 X \rightarrow X \rightarrow \begin{matrix} f(\theta_1) \rightarrow P_1 \rightarrow Y_1 \\ \searrow f(\theta_2) \rightarrow P_2 \rightarrow Y_2 \end{matrix}, \quad (3)$$

28 where,  $Y_1$  and  $Y_2$  are one-hot pseudo-labels.

29 The CPS method imposes two loss functions, a supervised loss  $\mathcal{L}_s$  and pseudo su-  
 30 pervised loss (unlabeled loss)  $\mathcal{L}_{cps}$ . The supervised loss,  $\mathcal{L}_s$ , is a standard cross-entropy loss  
 31 applied at pixel level.

$$32 \mathcal{L}_s = \frac{1}{|\mathcal{D}^l|} \sum_{X \in \mathcal{D}^l} \frac{1}{W \times H} \sum_{i=0}^{W \times H} (\ell_{ce}(\mathbf{p}_{1i}, \mathbf{y}_{1i}^*) \\ 33 + \ell_{ce}(\mathbf{p}_{2i}, \mathbf{y}_{2i}^*)), \quad (4)$$

34 The pseudo supervised loss (unsupervised loss),  $\mathcal{L}_{cps}$  is also a standard cross en-  
 35 tropy loss applied bi-directionally on the unlabeled data  $\mathcal{D}_u$  and predicts pseudo-labels. The CPS  
 36 loss can be represented as:



**Figure 1 Overview of the semi-supervised segmentation model**

$$\begin{aligned}
 1 \quad \mathcal{L}_{cps}^u &= \frac{1}{|\mathcal{D}_u|} \sum_{X \in \mathcal{D}_u} \frac{1}{W \times H} \sum_{i=0}^{W \times H} (\ell_{ce}(\mathbf{p}_{1i}, \mathbf{y}_{2i}) \\
 2 \quad &\quad + \ell_{ce}(\mathbf{p}_{2i}, \mathbf{y}_{1i})). \tag{5}
 \end{aligned}$$

3 The overall cross pseudo supervised loss is the sum of the supervised loss  $\mathcal{L}_s$  on the labeled data,  
 4  $\mathcal{D}_l$  and the pseudo supervised loss  $\mathcal{L}_{cps}$  on the unlabeled data  $\mathcal{D}_u$ .

$$5 \quad \mathcal{L} = \mathcal{L}_s + \lambda \mathcal{L}_{cps}, \tag{6}$$

6 where,  $\lambda$  is a hyper-parameter that trade-off the importance of the two losses.

7 MeanTeacher (20) is another semi-supervised learning method that does not train both  
 8 encoders, which is expensive and less intuitive. Rather, inspired by the semi-supervised literature  
 9 of object recognition (37), it updates the teacher encoder as an exponential moving average (EMA)  
 10 of the online encoder (student encoder), denoted as  $\bar{\theta}$ . The loss function of MeanTeacher can be  
 11 represented as:

$$\begin{aligned}
 12 \quad X &\rightarrow X_1 \rightarrow f(\theta) \rightarrow P_1 \\
 13 \quad \searrow X &\rightarrow X_2 \rightarrow f(\bar{\theta}) \rightarrow P_2 \tag{7}
 \end{aligned}$$

14 where,  $X_1$  and  $X_2$  are two different augmentations of images an input image  $X$ . The consistency  
 15 regularization works in a way to align the consistency map  $P_1$  predicted by the student network  
 16 to the consistency map  $P_2$  of the teacher network. At training time, the  $P_1$  is supervised by  $P_2$ .  
 17 However, the teacher network does not back-propagate the loss but is updated as the exponential  
 18 moving average (EMA) of the student network. As a regularization method, EMA helps to stabi-  
 19 lize the student model's updates, smoothing out fluctuations and facilitating knowledge transfer to  
 20 the student network. This technique enhances the model's capacity to handle unlabeled data and  
 21 improves overall performance on both labeled and unlabeled datasets.  $f(\theta)$ . It reduces the redun-  
 22 dant calculation and cost of training while keeping the weights of the teacher network similar to  
 23 the student network.

24 Another semi-supervised method combines CPS with the CutMix augmentation  
 25 (36). CutMix augmented images are fed as input in both of the networks  $f(\theta_1)$  and  $f(\theta_2)$ . The  
 26 rest of the method is similar to (37). Finally, each generated pseudo-label is used to supervise the  
 27 other segmentation network.

## 28 Proposed Method

29 In this paper, we propose a few variants of the existing methods that are inspired by the existing  
 30 literature and aim to improve semi-supervised segmentation on fish-eye images for autonomous  
 31 vehicles. These methods are described in the following subsections. We also attempt to improve  
 32 the existing methods by tuning important parameters of individual methods for fish-eye images.

### 1 *CPS with Confidence Threshold*

2 The first variant we propose is inspired by FixMatch (38), which involves using confidence thresh-  
 3 olding to improve the quality of predicted pseudo-labels on unlabeled samples. FixMatch uses a  
 4 weakly-augmented image as input to the model, which generates the predicted pseudo-labels. We  
 5 transform this prediction into a one-hot pseudo-label. The same image is then subjected to strong  
 6 augmentation and fed into the same model. The one-hot pseudo-label from the weakly-augmented  
 7 image is used to supervise the prediction of the strongly-augmented image, and this process works  
 8 as a consistency regularization by ignoring low-confidence pseudo-labels for supervision. The  
 9 FixMatch objective can be represented as follows:

$$10 \quad X \rightarrow X_s \rightarrow f(\theta) \rightarrow P_s \swarrow \searrow \\
 11 \quad \searrow X_w \rightarrow f(\theta) \rightarrow P_w \rightarrow Y_w, \quad (8)$$

12 where,  $X_w$  and  $X_s$  are the weak and strong augmentations of image  $X$ .

13 We then combine the pseudo-label thresholding idea with CPS. In CPS, one encoder  
 14 generates pseudo-labels that are used as ground truth for the other branch. By incorporating the  
 15 thresholding concept of FixMatch (38), we only consider pseudo-labels that exceed a pre-defined  
 16 threshold. Thus providing more confidence in the model’s performance.

### 17 *CPS with CutMix and Confidence Threshold*

18 As mentioned above, using CutMix augmentation with CPS improves the performance of  
 19 CPS by making the dataset more challenging for the model to learn representation. So, we propose  
 20 a new variant by introducing Cutmix augmentation into CPS with Confidence Thresholding.

### 21 *CPS with CutMix, Confidence Threshold, and Hard Augmentations*

22 In this variant, we introduce the concept of hard augmentations with the previously pro-  
 23 posed variant. Hard argumentation regularizes a model to prevent overfitting on small labeled  
 24 samples. To this end, we combine the following augmentations to form a hard augmentation mod-  
 25 ule: RandomCrop; ColorJitter with brightness, contrast, saturation, hue, and Gaussian blur. The  
 26 model applies Cutmix and hard augmentations on the input and applies threshold on the predicted  
 27 pseudo-labels.

### 28 *Tunning Existing Methods*

29 From our preliminary studies, we identify a few key parameters that are significant to the  
 30 existing methods. Since the existing methods are trained and optimized for perspective images,  
 31 they do not perform well on fish-eye images. To this end, we propose to change some of these  
 32 model-specific parameters to tune them for the fish-eye segmentation task. For this study, we focus  
 33 on parameters of individual methods that we find important for fish-eye segmentation. CPS, with  
 34 the Cutmix augmentation method, contains two important model-specific parameters: the number  
 35 of cutmixed image boxes and the range of these boxes. In our study, we tune both of the parameters  
 36 to perform better in the fish-eye image segmentation. Similarly, for CPS, we tune the CPS weight  
 37 and the learning rate. For MeanTeacher, we tune the momentum and unsupervised loss weight.  
 38 For the other three proposed methods, e.g. CPS with confidence thresholding, CPS with cutmix,  
 39 confidence thresholding, and CPS with cutmix, confidence thresholding, and hard augmentation,  
 40 the threshold value is a crucial hyperparameter that represents a pre-defined threshold to which  
 41 only the exceeded pseudo labels are considered.

## 1 **Dataset Description**

2 WoodScape (*I*) is a dataset specifically designed for autonomous driving tasks captured  
3 using multi-camera fish-eye lenses. It comes with extensive ground-truth annotations for nine  
4 essential tasks, including semantic segmentation, which covers classes like roads, pedestrians,  
5 and traffic signs. The main purpose behind WoodScape is to offer a comprehensive platform for  
6 evaluating computer vision algorithms on fish-eye images. The dataset comprises four surround-  
7 view cameras and nine different tasks, such as depth estimation, 3D bounding box detection, and  
8 soiling detection. Unlike relying on rectification methods, WoodScape encourages the adoption of  
9 computer vision techniques tailored explicitly for fish-eye cameras.

10 This dataset stands out as the first labeled dataset specifically designed for seman-  
11 tic segmentation. The data was collected from three distinct geographical regions: the US, Eu-  
12 rope, and China. To ensure diverse sensor mechanical configurations, the majority of the data was  
13 gathered from saloon vehicles, while a significant subset comes from sports utility vehicles. The  
14 driving scenarios encompass various use cases, including highway driving, urban environments,  
15 and parking scenarios. The dataset provides both intrinsic and extrinsic calibrations for all sen-  
16 sors, along with timestamp files for data synchronization. Quality checks were conducted at every  
17 stage of data collection, and reviewers performed quality assurance for the annotation data. The  
18 sensors used in recording this dataset include four 1MPx RGB fish-eye cameras with a wide 190°  
19 horizontal field of view each, one LiDAR rotating at 20Hz (Velodyne HDL-64E), one GNSS/IMU  
20 setup (NovAtel Propak6 and SPAN-IGM-A1), one GNSS Positioning with SPS (Garmin 18x), and  
21 Odometry signals from the vehicle bus. Overall, WoodScape promises high-quality data and opens  
22 avenues for developing unified multi-task and multi-camera models for advancing autonomous  
23 driving technology.

## 24 **Implementation Details**

25 We follow OmniDet and SemiSeg for the training setup and default hyper-parameters. We  
26 adopt the Omnidet (35) model as a backbone encoder-decoder for the segmentation model. The  
27 visual encoder backbone of OmniDet is a ResNet50 network. For semi-supervised training, we  
28 first divide the total training samples from WoodScape (*I*) into supervised and unsupervised splits.  
29 More specifically, we used 20% of the 8,029 samples as our labeled set and the rest of the samples  
30 as unlabeled set. We train the model for 125 epochs with an Adam optimizer and a learning rate  
31 of 0.0001, and a batch size of 16. A multi-step scheduler with a decay factor of 0.1 and a step size  
32 of [100, 110] is used to adjust the learning rate. The model is trained on an NVIDIA RTX A6000  
33 GPU.

## 34 **RESULTS**

35 This section presents the results of this study. First, we will discuss the performance of  
36 existing methods in fish-eye image segmentation, which we consider the baseline of the study.  
37 Then we will present a discussion of the results of our proposed methods.

### 38 **Baseline**

39 For this and all subsequent experiments, we consider supervised learning with an equal  
40 amount of labeled data (20% to total data) used in the semi-supervised methods as the baseline.  
41 So, the improvement with the semi-supervised method is purely from using unlabeled data (as the  
42 number of labeled samples is the same in both settings). The first row of Table 1 shows the baseline

**Table 1 Comparison between fully-supervised learning and semi-supervised learning with unconstrained unlabeled data.**

Method	mIoU
Fully-supervised (baseline)	54.32
CPS (from SemiSeg (19))	60.31
MeanTeacher (20)	54.20
CPS+CutMix (36)	62.47
Ours (CPS+Conf. thr.)	60.66
Ours (CPS+CutMix+Conf. thr.)	62.38
Ours (CPS+Conf. thr.+Hard. Aug.)	59.43
Ours CPS+CutMix (optimized)	<b>63.21</b>

1 mIoU for our work, which is 54.32. For the rest of the paper, we will discuss the improvement of  
 2 different semi-supervised methods over this baseline.

### 3 Main Results

4 Table 1 summarizes the results for all the existing methods as well as our proposed variants.  
 5 For evaluation, we follow the protocol described in (19). More precisely, the evaluation process is  
 6 unimodal; single-scale testing. The results are presented using the mean intersection-over-union  
 7 (mIoU) metric, which evaluates model accuracy by computing the intersection-to-union ratio be-  
 8 tween predicted and ground truth masks for each class. The mIoU offers an average assessment  
 9 of segmentation performance across all classes, with superior values denoting higher accuracy.  
 10 As mentioned in the Method section, we re-implement three popular semi-supervised segmenta-  
 11 tion methods from the regular image literature; namely, CPS (19), MeanTeacher (20), CPS with  
 12 CutMix (36). From the existing methods, we find the best result for CPS+CutMix. It obtains an  
 13 accuracy of 62.47, which is an 8.15 improvement of the supervised learning with a sample amount  
 14 of labeled samples. The next best performance is shown by the original CPS with a mIoU of 60.31.  
 15 However, MeanTeacher does not perform well, and the result is worse than learning from super-  
 16 vised loss only. Based on these three semi-supervised segmentation methods, we propose three  
 17 new variants, namely CPS+Confidence Thresholding, CPS+Cutmix+Confidence Thresholding,  
 18 and CPS+Cutmix+Confidence Thresholding+Hard Augmentation. Among these three proposed  
 19 methods, the highest performance is obtained from the CPS+Cutmix+Confidence Thresholding  
 20 variant, which is 62.38%. We do not observe comparable results for the CPS+Cutmix+Confidence  
 21 Thresholding+Hard Augmentation method. We suspect that imposing hard augmentation to the  
 22 dataset alters the data distribution and makes it challenging for the model to learn representation.  
 23 However, the highest mIoU of the model in this study, which is **63.21%**, is achieved by our opti-  
 24 mized CPS+Cutmix method.

**Table 2 Performance of CPS with Confidence Thresholding for different threshold values**

Method setting	mIoU
threshold= 0.5	60.56
threshold= 0.75	60.66
threshold= 0.9	59.96

**Table 3 Performance of CPS with CutMix and Confidence Thresholding for different threshold values**

Method setting	mIoU
threshold= 0.5	62.38
threshold= 0.75	61.98
threshold= 0.9	62.12
threshold= 0.95	55.60

## 1 Analysis of the Proposed Method

2 In this section, we report the performance of our proposed semi-supervised learning meth-  
3 ods on the WoodScape dataset.

### 4 *CPS with Confidence Thersholding*

5 Table 2 shows the performance of CPS with Confidence Thresholding for different threshold  
6 values. The results indicate that the highest mIoU value of 60.66 is obtained with a threshold of  
7 0.75, whereas the lowest value of 59.96 is obtained with a threshold of 0.9. Since the threshold  
8 value works as a consistency regularizer by ignoring low-confidence pseudo-labels, the higher  
9 threshold value induces a higher regularization effect on the model and thus affects the model’s  
10 performance by ignoring a larger amount of predicted pseudo-labels for supervision.

### 11 *CPS with CutMix and Confidence Thersholding*

12 Table 3 demonstrates the performance of CPS with CutMix and Confidence Thresholding at  
13 various threshold values. We consider the threshold values of 0.5, 0.75, 0.9, and 0.95 in this study.  
14 The results indicate that the highest mIoU is achieved with a threshold of 0.5, while the lowest was  
15 observed with a threshold of 0.95. We believe the reason behind the following performance of the  
16 model is similar to that explained in the previous subsection.

### 17 *CPS with CutMix, Confidence Thresholding, and Hard Augmentations*

18 Table 4 represents the performance of our proposed final variant, CPS with CutMix, Confi-  
19 dence Thresholding, and Hard Augmentations. The table shows the mIoU scores for the method  
20 with different thresholds: 0.5, 0.75, 0.9, and 0.95. Here, we see an interesting observation as the  
21 model’s performance becomes random with the different threshold values. We suspect that the

**Table 4 Performance of CPS with CutMix, Confidence Thresholding, and Hard Augmentations for different threshold values**

Method setting	mIoU
threshold= 0.5	58.37
threshold= 0.75	57.41
threshold= 0.9	59.43
threshold= 0.95	46.01

**Table 5 Comparison among various combinations of CPS weight and learning rate values for CPS method.**

LR	CPS Weight		
	1.5	1	0.5
0.01	47.36	46.46	46.67
0.001	57.28	57.19	56.10
0.0001	60.31	59.86	60.40

1 application of hard augmentation has made the dataset challenging enough for the model to learn.

## 2 Sensitivity Study on Tuning Existing Method

3 In this section, we discuss the results of the sensitivity study on existing methods that we  
4 find important for the performance in fish-eye segmentation.

5 Table 5 shows different combinations of CPS weight and learning rate values for the CPS  
6 method. The table includes the performance of the CPS method of different combinations of CPS  
7 weights of 1.5, 1, and 0.5. We can see that with the change in the CPS weights with a fixed  
8 learning rate, the model’s performance does not change significantly. However, by decreasing the  
9 learning rates from 0.01 to 0.0001 with a constant CPS weight, the model’s performance increases  
10 significantly. The highest mIoU, 60.40% is achieved with a 0.0001 learning rate with 0.5 CPS  
11 weight. So, it can be concluded that the model’s performance is more sensitive to the learning rate  
12 compared to CPS weight.

13 In Table 6, we can see the performance of EMA as a teacher when the unsupervised weight  
14 and momentum values were adjusted. The table includes the method settings and their correspond-  
15 ing mIoU scores. The default setting, with  $W=100$  and  $m=0.99$ , scored 54.20. We also tested other  
16 settings with different values of  $W$  and  $m$ , resulting in different scores. From the analysis, it can  
17 be seen that as the unsupervised weights ( $W$ ) get lower, the model shows better performance. The  
18 highest accuracy in this setting, 57.15%, is achieved from an unsupervised weight of 1.0. For  
19 the momentum values, with 0.999 momentum, the model achieves better performance, 55.62%,  
20 compared with other momentum values. Since momentum adjusts the size of the next update in-  
21 troducing a penalty term with a higher or lower momentum value affects the model’s performance.

**Table 6 Performance of EMA as a teacher for different values of unsupervised weight and momentum.**

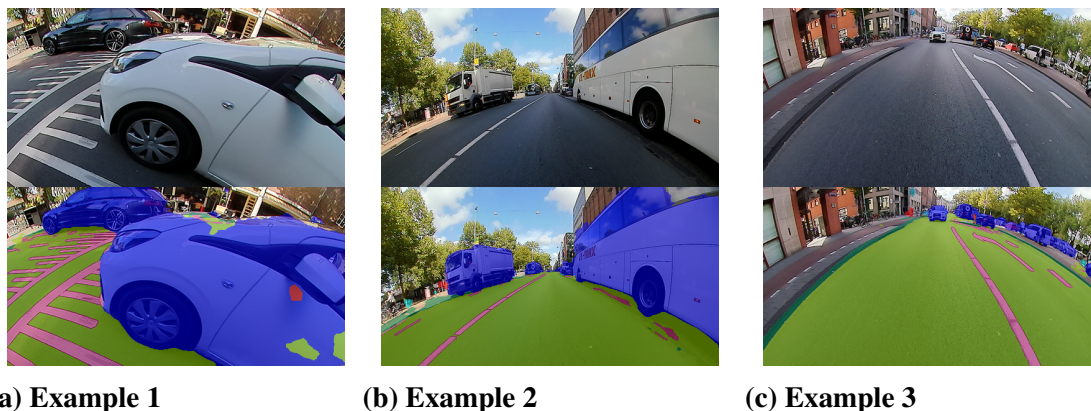
Method setting	mIoU
Default (W=100 & m=0.99)	54.20
W=10.0	56.38
W=1.0	57.15
m=0.999	55.62
m=0.9999	46.21

**Table 7 Performance of CPS with CutMix augmentation with different combinations of the number of cutmixed boxes and the range of boxes of CutMix augmentation**

Method setting	mIoU
Default (box number=3; range= 0.25-0.5)	62.47
nb. of. box=2 and range=0.1 to 0.25	63.13
nb. of. box=2 and range=0.25 to 0.5	62.60
nb. of. box=2 and range=0.4 to 0.75	62.29
nb. of. box=3 and range=0.1 to 0.25	62.02
nb. of. box=3 and range=0.4 to 0.75	62.59
nb. of. box=4 and range=0.1 to 0.25	62.00
nb. of. box=4 and range=0.25 to 0.5	61.91
nb. of. box=4 and range=0.4 to 0.75	<b>63.21</b>

1 To understand the performance of CPS with CutMix augmentation, in table 7, we conduct  
2 experiments using different combinations of the number of cutmixed boxes and the range of these  
3 boxes. The table demonstrates the method setting and the corresponding mIoU value. This exper-  
4 iment aims to determine the optimal combination of box number and range that would yield the  
5 highest mIoU value. The default setting for the experiment had a box number of 3 and a range  
6 of 0.25-0.5, with a mIoU value of 62.47. The highest performance, **63.21%**, is achieved by 4 cut  
7 mix boxes with a range of 0.4 to 0.75. Also, note that, with the increment of the cutmix range, the  
8 performance of the model increases as the overlapping area of the two images increases. Thus the  
9 model gets a wider view of the images to capture the underlying information. The highest mIoU  
10 of the proposed methods, **63.21%**, is achieved from 4 boxes with a range of 0.4 to 0.75, which is  
11 a **8.89%** improvement over the supervised baseline.

12 Overall, the results of the experiment are promising, as they showed that CPS with  
13 CutMix augmentation can significantly improve the performance of the semi-supervised segmen-  
14 tation model for fish-eye images.



**Figure 2 Example of segmentation generated by semi-supervised segmentation model.**

## 1 Qualitative Analysis

2 Figure 2 show some images generated by our proposed semi-supervised segmentation method  
 3 (CPS+cutmix optimized). The sample images show perfect detection of most of the objects and  
 4 their areas.

## 5 CONCLUSION

6 Fish-eye camera images differ from perspective camera images in that they capture a Field  
 7 of View (FOV) of  $180^\circ$  to  $360^\circ$ . Because of this characteristic, fish-eye images are potentially  
 8 important for use in autonomous driving. However, the larger FOV comes at a price of higher radial  
 9 distortion and requires specialized models to handle this challenge. Despite the high potentiality  
 10 of this area of fish-eye images, it is not well-explored due to the lack of large publicly annotated  
 11 datasets. Semantic segmentation, which classifies an image at the pixel level, is a popular domain  
 12 of computer vision. However, data annotation for semantic segmentation tasks is more expensive  
 13 and laborious compared to object classification. Semi-supervised learning provides a solution to  
 14 this problem, but it has not been studied in the context of semi-supervised segmentation from fish-  
 15 eye cameras. In this paper, we proposed a semi-supervised semantic segmentation method for  
 16 fish-eye images. We explored various existing semi-supervised methods and benchmarked their  
 17 performances. Based on our findings and best practices from these methods, we proposed three  
 18 new variants for semi-supervised segmentation. An analysis of these methods shows that semi-  
 19 supervised methods can improve up to 8.89% over fully-supervised methods for the same amount  
 20 of labeled data. However, one of the limitations of our work is that we have focused on a subset  
 21 of the available literature concerning semi-supervised semantic segmentation methods designed  
 22 for perspective images. As a result, we recognize the potential for enhancing the performance  
 23 (mIoU) in semantic segmentation tasks for fish-eye images used in autonomous driving by delving  
 24 into other semi-supervised methods already established in the field. We hope that this work will  
 25 encourage the development of more specialized techniques for learning semantic segmentation and  
 26 other computer vision tasks related to fish-eye images, especially when labeled data is scarce in  
 27 the context of autonomous driving.

## 28 ACKNOWLEDGEMENTS

29 We would like to thank Mitacs Accelerate and BusPas Inc. for funding this research.

1 **AUTHOR CONTRIBUTION STATEMENT**

2 The authors acknowledge their involvement in the paper in the following capacities: Sneha Paul,  
3 Zachary Patterson, and Nizar Bouguila contributed to the study's conception and design. Sneha  
4 Paul, Zachary Patterson, and Nizar Bouguila conducted the analysis and interpretation of the re-  
5 sults. All authors thoroughly reviewed the results and granted approval for the final version of the  
6 manuscript.

## 1 REFERENCES

- 2 1. Yogamani, S., C. Hughes, J. Horgan, G. Sistu, P. Varley, D. O’Dea, M. Uricár, S. Milz,  
3 M. Simon, K. Amende, et al., Woodscape: A multi-task, multi-camera fisheye dataset  
4 for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on*  
5 *Computer Vision*, 2019, pp. 9308–9318.
- 6 2. Deng, L., M. Yang, Y. Qian, C. Wang, and B. Wang, CNN based semantic segmentation  
7 for urban traffic scenes using fisheye camera. In *2017 IEEE Intelligent Vehicles Symposium*  
8 *(IV)*, IEEE, 2017, pp. 231–236.
- 9 3. Kumar, V. R., M. Klingner, S. Yogamani, S. Milz, T. Fingscheidt, and P. Mader, Syndistnet:  
10 Self-supervised monocular fisheye camera distance estimation synergized with semantic  
11 segmentation for autonomous driving. In *Proceedings of the IEEE/CVF winter conference*  
12 *on applications of computer vision*, 2021, pp. 61–71.
- 13 4. Shi, H., Y. Li, K. Yang, J. Zhang, K. Peng, A. Roitberg, Y. Ye, H. Ni, K. Wang, and  
14 R. Stiefelhagen, FishDreamer: Towards Fisheye Semantic Completion via Unified Image  
15 Outpainting and Segmentation. *arXiv preprint arXiv:2303.13842*, 2023.
- 16 5. Feng, D., C. Haase-Schütz, L. Rosenbaum, H. Hertlein, C. Glaeser, F. Timm, W. Wiesbeck,  
17 and K. Dietmayer, Deep multi-modal object detection and semantic segmentation for au-  
18 tonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent*  
19 *Transportation Systems*, Vol. 22, No. 3, 2020, pp. 1341–1360.
- 20 6. Kirillov, A., K. He, R. Girshick, C. Rother, and P. Dollár, Panoptic segmentation. In *Pro-*  
21 *ceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019,  
22 pp. 9404–9413.
- 23 7. Wu, B., A. Wan, X. Yue, and K. Keutzer, SqueezeSeg: Convolutional neural nets with  
24 recurrent crf for real-time road-object segmentation from 3d lidar point cloud. In *2018*  
25 *IEEE international conference on robotics and automation (ICRA)*, IEEE, 2018, pp. 1887–  
26 1893.
- 27 8. Valada, A., G. L. Oliveira, T. Brox, and W. Burgard, Deep multispectral semantic scene  
28 understanding of forested environments using multimodal fusion. In *2016 International*  
29 *Symposium on Experimental Robotics*, Springer, 2017, pp. 465–477.
- 30 9. Sun, Y., W. Zuo, and M. Liu, RTFNet: RGB-thermal fusion network for semantic segmen-  
31 tation of urban scenes. *IEEE Robotics and Automation Letters*, Vol. 4, No. 3, 2019, pp.  
32 2576–2583.
- 33 10. Cordts, M., M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke,  
34 S. Roth, and B. Schiele, The cityscapes dataset for semantic urban scene understanding.  
35 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016,  
36 pp. 3213–3223.
- 37 11. Geiger, A., P. Lenz, and R. Urtasun, Are we ready for autonomous driving? the kitti vision  
38 benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*,  
39 IEEE, 2012, pp. 3354–3361.
- 40 12. Wang, S., M. Bai, G. Mattyus, H. Chu, W. Luo, B. Yang, J. Liang, J. Cheverie, S. Fi-  
41 dler, and R. Urtasun, Torontocity: Seeing the world with a million eyes. *arXiv preprint*  
42 *arXiv:1612.00423*, 2016.
- 43 13. Huang, X., X. Cheng, Q. Geng, B. Cao, D. Zhou, P. Wang, Y. Lin, and R. Yang, The  
44 apolloscape dataset for autonomous driving. In *Proceedings of the IEEE conference on*  
45 *computer vision and pattern recognition workshops*, 2018, pp. 954–960.

- 1 14. Sáez, Á., L. M. Bergasa, E. López-Guillén, E. Romera, M. Tradacete, C. Gómez-Huélamo,  
2 and J. Del Egido, Real-time semantic segmentation for fisheye urban driving images based  
3 on ERFNet. *Sensors*, Vol. 19, No. 3, 2019, p. 503.
- 4 15. Playout, C., O. Ahmad, F. Lecue, and F. Cheriet, Adaptable deformable convolutions for  
5 semantic segmentation of fisheye images in autonomous driving systems. *arXiv preprint*  
6 *arXiv:2102.10191*, 2021.
- 7 16. Blott, G., M. Takami, and C. Heipke, Semantic segmentation of fisheye images. In *Pro-*  
8 *ceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018, pp.  
9 0–0.
- 10 17. Sistu, G., I. Leang, and S. Yogamani, Real-time joint object detection and semantic seg-  
11 mentation network for automated driving. *arXiv preprint arXiv:1901.03912*, 2019.
- 12 18. Ye, Y., K. Yang, K. Xiang, J. Wang, and K. Wang, Universal semantic segmentation for  
13 fisheye urban driving images. In *2020 IEEE International Conference on Systems, Man,*  
14 *and Cybernetics (SMC)*, IEEE, 2020, pp. 648–655.
- 15 19. Chen, X., Y. Yuan, G. Zeng, and J. Wang, Semi-supervised semantic segmentation with  
16 cross pseudo supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vi-*  
17 *sion and Pattern Recognition*, 2021, pp. 2613–2622.
- 18 20. Tarvainen, A. and H. Valpola, Mean teachers are better role models: Weight-averaged  
19 consistency targets improve semi-supervised deep learning results. *Advances in neural*  
20 *information processing systems*, Vol. 30, 2017.
- 21 21. Yun, S., D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, Cutmix: Regularization strat-  
22 egy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF*  
23 *international conference on computer vision*, 2019, pp. 6023–6032.
- 24 22. Schneider, L., M. Jasch, B. Fröhlich, T. Weber, U. Franke, M. Pollefeys, and M. Räscht,  
25 Multimodal neural networks: RGB-D for semantic segmentation and object detection. In  
26 *Image Analysis: 20th Scandinavian Conference, SCIA 2017, Tromsø, Norway, June 12–14,*  
27 *2017, Proceedings, Part I 20*, Springer, 2017, pp. 98–109.
- 28 23. Badrinarayanan, V., A. Kendall, and R. Cipolla, Segnet: A deep convolutional encoder-  
29 decoder architecture for image segmentation. *IEEE transactions on pattern analysis and*  
30 *machine intelligence*, Vol. 39, No. 12, 2017, pp. 2481–2495.
- 31 24. Teichmann, M., M. Weber, M. Zoellner, R. Cipolla, and R. Urtasun, Multinet: Real-time  
32 joint semantic reasoning for autonomous driving. In *2018 IEEE intelligent vehicles sym-*  
33 *posium (IV)*, IEEE, 2018, pp. 1013–1020.
- 34 25. Long, J., E. Shelhamer, and T. Darrell, Fully convolutional networks for semantic segmen-  
35 tation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
36 2015, pp. 3431–3440.
- 37 26. Erkent, Ö. and C. Laugier, Semantic segmentation with unsupervised domain adaptation  
38 under varying weather conditions for autonomous vehicles. *IEEE Robotics and Automa-*  
39 *tion Letters*, Vol. 5, No. 2, 2020, pp. 3580–3587.
- 40 27. Kalluri, T., G. Varma, M. Chandraker, and C. Jawahar, Universal semi-supervised semantic  
41 segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer*  
42 *Vision*, 2019, pp. 5259–5270.
- 43 28. Huo, X., L. Xie, J. He, Z. Yang, W. Zhou, H. Li, and Q. Tian, ATSO: Asynchronous  
44 teacher-student optimization for semi-supervised image segmentation. In *Proceedings of*

- 1        *the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 1235–  
2        1244.
- 3 29.    Hung, W.-C., Y.-H. Tsai, Y.-T. Liou, Y.-Y. Lin, and M.-H. Yang, Adversarial learning for  
4        semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934*, 2018.
- 5 30.    Novosel, J., P. Viswanath, and B. Arsenali, Boosting semantic segmentation with multi-  
6        task self-supervised learning for autonomous driving applications. In *Proc. of NeurIPS-  
7        Workshops*, 2019, Vol. 3.
- 8 31.    Cui, L., X. Jing, Y. Wang, Y. Huan, Y. Xu, and Q. Zhang, Improved Swin Transformer-  
9        Based Semantic Segmentation of Postearthquake Dense Buildings in Urban Areas Using  
10        Remote Sensing Images. *IEEE Journal of Selected Topics in Applied Earth Observations  
11        and Remote Sensing*, Vol. 16, 2022, pp. 369–385.
- 12 32.    Zheng, S., J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr,  
13        et al., Rethinking semantic segmentation from a sequence-to-sequence perspective with  
14        transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
15        recognition*, 2021, pp. 6881–6890.
- 16 33.    Zhang, B., Z. Tian, Q. Tang, X. Chu, X. Wei, C. Shen, et al., Segvit: Semantic segmenta-  
17        tion with plain vision transformers. *Advances in Neural Information Processing Systems*,  
18        Vol. 35, 2022, pp. 4971–4982.
- 19 34.    Xie, E., W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, SegFormer: Sim-  
20        ple and efficient design for semantic segmentation with transformers. *Advances in Neural  
21        Information Processing Systems*, Vol. 34, 2021, pp. 12077–12090.
- 22 35.    Kumar, V. R., S. Yogamani, H. Rashed, G. Sitsu, C. Witt, I. Leang, S. Milz, and P. Mäder,  
23        Omnidet: Surround view cameras based multi-task visual perception network for au-  
24        tonomous driving. *IEEE Robotics and Automation Letters*, Vol. 6, No. 2, 2021, pp. 2830–  
25        2837.
- 26 36.    Zou, Y., Z. Zhang, H. Zhang, C.-L. Li, X. Bian, J.-B. Huang, and T. Pfister, Pseudoseg: De-  
27        signing pseudo labels for semantic segmentation. *arXiv preprint arXiv:2010.09713*, 2020.
- 28 37.    French, G., S. Laine, T. Aila, M. Mackiewicz, and G. Finlayson, Semi-supervised semantic  
29        segmentation needs strong, varied perturbations. *arXiv preprint arXiv:1906.01916*, 2019.
- 30 38.    Sohn, K., D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Ku-  
31        rakin, and C.-L. Li, Fixmatch: Simplifying semi-supervised learning with consistency and  
32        confidence. *Advances in neural information processing systems*, Vol. 33, 2020, pp. 596–  
33        608.