# Panoramic Panoptic Segmentation: Insights Into Surrounding Parsing for Mobile Agents via Unsupervised Contrastive Learning

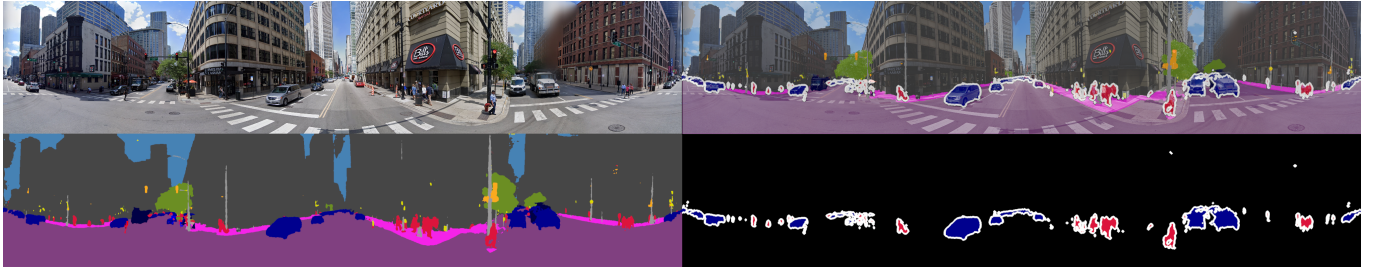Alexander Jaus[1], Kailun Yang[1], and Rainer Stiefelhagen[1]



**Fig. 1:** Within this work, we differentiate between various levels of image understanding: The original image (first row, left) can be interpreted as a panoramic semantic map (second row, left) by assigning a label to each pixel without differentiating between different instances of countable objects. Instances of countable objects are distinguished in the panoramic instance understanding (second row, right). The panoramic panoptic understanding (first row, right), which is the proposed method in this paper, builds on top of the previous understandings by eliminating their shortcomings: If possible different instances are distinguished and we guarantee that a label is assigned to each pixel.

*Abstract*—In this work, we introduce panoramic panoptic segmentation, as the most holistic scene understanding, both in terms of Field of View (FoV) and image-level understanding for standard camera-based input. A complete surrounding understanding provides a maximum of information to a mobile agent. This is essential information for any intelligent vehicle to make informed decisions in a safety-critical dynamic environment such as real-world traffic. In order to overcome the lack of annotated panoramic images, we propose a framework which allows model training on standard pinhole images and transfers the learned features to the panoramic domain in a cost-minimizing way. The domain shift from pinhole to panoramic images is non-trivial as large objects and surfaces are heavily distorted close to the image border regions and look different across the two domains. Using our proposed method with dense contrastive learning, we manage to achieve significant improvements over a non-adapted approach. Depending on the efficient panoptic segmentation architecture, we can improve $3.5-6.5\%$ measured in Panoptic Quality (PQ) over non-adapted models on our established Wild Panoramic Panoptic Segmentation (WildPPS) dataset. Furthermore, our efficient framework does not need access to the images of the target domain, making it a feasible domain generalization approach suitable for a limited hardware setting. As additional contributions, we publish WildPPS: The first panoramic panoptic image dataset to foster progress in surrounding perception and explore a novel training procedure

combining supervised and contrastive training.

*Index Terms*—Panoptic Segmentation, Autonomous Driving, Scene Understanding, Contrastive Learning

## I. INTRODUCTION

PANOPTIC segmentation is the so far most complete segmentation task to describe the context of an image [1]. It seamlessly addresses stuff and thing classes and thus unifies semantic segmentation which does not differentiate between instances and instance segmentation which falls short to segment uncountable objects. Both are critical pieces of information for the task of scene understanding in an autonomous driving setting. Not distinguishing between different instances of cars or pedestrians does not allow to anticipate the dynamics of individuals as they tend to interact with each other.

Typically, in a street-scene context, stuff classes such as *roads* or *sidewalks* can be used to find traversable areas, whereas *cars* or *pedestrians* which are represented by the thing class, can be interacted with by the mobile agent in order to protect vulnerable road users or avoid moving obstacles. A purely semantic segmentation approach [2], [3], which fails to identify the number of instances of the thing class, is insufficient, as the number of pedestrians or cars is an important piece of information in order to plan in accordance with the situation [4].

Despite the high level of information provided by panoptic image understanding on standard pinhole images, as it is the de facto center of research at the moment, these works lack a crucial source of information which comes from the limited Field of View (FoV) of the input image. As real-world traffic is a highly dynamic environment in which a lot of movement is happening around the participants, the problem of people or vehicles moving out of the FoV of the mobile agent is easily
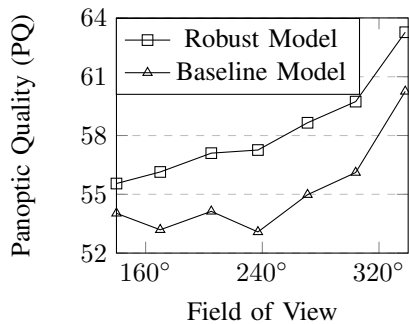
**Fig. 2:** Comparison of a standard Seamless Scene Segmentation model [7] to its robust counterpart depending on the Field of View (FoV) on the WildPPS dataset. Larger FoVs tend to widen the gap between the two models confirming the effectiveness of the proposed Panoramic Robust Feature (PRF) framework.

imaginable [2]. This poses severe problems due to the lack of information containing the entire surrounding and the inability of the agent to make proper decisions which may even lead to accidents [5]. Thus, both pieces of information are equally important: the image should cover the entire surrounding and should be understood holistically [6].

One approach to overcome the aforementioned FoV problem is to use multi-camera or LiDAR systems [8], [9], [10]. These systems, however, introduce a further level of complexity as they either need to stitch images of multiple camera sources or resort to the combination of the RGB image information with LiDAR views in order to get a holistic view of the entire surrounding [11]. Within this work, we aim to introduce a simple, efficient, and comprehensive approach, which is suitable for a mobile agent setup.

To this intent, we introduce Panoramic Panoptic Segmentation, to enable complete image understanding based on standard camera input. There are a multitude of tasks which need to be addressed in order to properly operate mobile agents within a real-world scenario such as localization, motion planning, or robot controls. The one problem which these tasks have in common is that they rely on an interpretation of sensor data in order to collect information about the mobile agent as well as the world [12], [13]. While there are multiple approaches to receive this information, this work focuses on scene understanding via standard camera input to derive information about the current state. Our task combines panoptic image level understanding and the wide FoV of panoramic images [14], tackling both previously discussed problems in a cost minimizing way because besides pure performance, a major design requirement for the derived approach is hardware and energy efficiency.

It is well known that segmentation models require a substantial amount of annotated data in order to be trained in a supervised fashion [15]. However, dense, pixel-wise annotations are notoriously time-consuming and extremely expensive to produce, especially for panoramic images with a large FoV and many small objects co-occurring in the complex surroundings.

Due to this lack of available annotated panoramic images, we are forced to train our segmentation models on standard pinhole images, which introduces a domain shift to the proposed panoramic panoptic image segmentation task. A common problem when training and testing under different image distributions is the performance drop of models in the unseen target domain [2]. The domain shift from pinhole- to panoramic images is no exception. The differences between standard pinhole images and panoramic images, as shown in Fig. 1, are quite apparent. Streets span through the entire image, whereas sidewalks are no longer present only on the side and now form bow-shaped islands in central parts of the image. These properties have not been observed by the model during the training and make their correct segmentation challenging. As they represent the traversable areas within a street scene, it is crucial to segment them in a correct way in order to avoid accidents or navigation failures. In order to overcome these difficulties, we propose the Panoramic Robust Feature (PRF) framework which allows us to generate robust backbones via a contrastive pretext task. Contrastive learning does not only encourage similar features to be represented in a similar manner but more important, it pushes dissimilar features away from each other [16], [17]. This leads to well separated clusters in the latent space of the backbone which proves to mitigate distribution shift performance drops. The robust backbones can be inserted into the target panoptic model and the model can be trained in a standard supervised fashion on pinhole images. We find that our proposed approach significantly outperforms non-adapted models in a variety of segmentation scenarios.

It furthermore comes with the advantage that we can make use of any large-scale dataset, as our method works best in our experiments when using the same domain for pretraining and supervised training. This does not only eliminate the need for labels within the target domain during training but even for images in the target domain making our approach feasible for domain generalization problems. Diverse autonomous transportation application scenarios, in which it may even be difficult to obtain more images from the target domain due to the need for specialized equipment, could benefit from the proposed method as well. By adapting target models via our proposed PRF framework, we achieve significant improvements of over $6\%$ overall and over $7\%$ for the difficult stuff class measured in Panoptic Quality (PQ) over non-adapted baseline models on our extended panoramic panoptic dataset WildPPS [6]. In particular, the aforementioned challenging classes *road* and *sidewalk* on which non-adapted models perform poorly are now well addressed, which can be seen in Fig. 10. PRF consistently brings generalization gains on different FoVs cropped from the panoramas, as illustrated in Fig. 2. The proposed approach, thereby, attains state-of-the-art performance on both the public PASS [2] and our established WildPPS benchmarks.

Finally, we discuss the feasibility of our efficient implementations of various state-of-the-art panoptic models for panoramic panoptic segmentation. Our findings indicate that models using well established baseline models such as the Seamless Scene Segmentation model [7] which uses Mask R-CNN [18] do outperform more advanced transformer [19] based models or fully convolutional panoptic models.

On a glance, we deliver the following contributions:

1) We introduce the novel panoramic panoptic segmentation

task, which aims to yield a $360°$ holistic surrounding understanding, providing dense semantic- and instance information at the pixel level.

2) We propose Panoramic Robust Feature (PRF), a contrastive-learning-driven panoramic panoptic segmentation framework with efficient architectures for mobile agents in road-driving scene parsing scenarios.

3) We create WildPPS, the first panoramic panoptic segmentation dataset with images collected in cities located around the world, to foster robust surrounding perception.

4) The proposed method achieves state-of-the-art accuracy on both public PASS and our established WildPPS benchmarks. For a variety of CNN- and transformer-based panoptic segmentation models, our approach significantly improves their performances in PQ by $3.5-6.5\%$.

**Difference from the Conference Paper:** This paper is an extension of our conference work [6]. Within this publication, we significantly increase the insights into the novel task introduced in our conference paper [6], which received the Best Paper Award (Third Place) at the 2021 IEEE Intelligent Vehicles Symposium, by adding the following contributions:

1) We double the number of annotated images forming the WildPPS dataset which results in more stable and reliable estimates for the performance of models.

2) We add extensive experiments with multiple state-of-the-art panoptic segmentation models using various new adaptation techniques and extend the previously established efficiency setting to the new models.

3) We add a benchmark of the proposed framework against a purely contrastive learning based training of the backbone.

4) We take our proposed robust model to the test against multiple more powerful panoptic baseline models on WildPPS.

5) We provide novel insights into our established idea which combines supervised and contrastive training procedures suitable for a limited hardware setting.

6) A more detailed description of the proposed framework and other enhanced parts such as related works and additional qualitative feature- and segmentation map analyses.

## II. RELATED WORK

### A. Semantic, Instance, and Panoptic Segmentation

Pixel-wise image, semantic- and instance-specific segmentation have advanced exponentially, driven by the architectural advances of deep Convolutional Neural Networks (CNNs) [18], [20], [21]. The Fully Convolutional Network (FCN) [21] views semantic segmentation as an end-to-end pixel classification task, followed by encoder-decoder architectures [22], [23], [24] that significantly enhance segmentation performance by aggregating multi-scale contextual features. Researchers further come up with promising approaches, by utilizing boundary cues [25], appending attention blocks [26], [27], and leveraging vision transformers [28], [29], [30] to improve FCN-based dense semantic understanding. Additionally, there have been some attempts for semantic segmentation

using supervised and unsupervised communication [31] or leveraging knowledge distillation [32], [33].

Early instance segmentation models are built upon the "detection followed by segmentation" principle, with Mask R-CNN [18] being a well-known architecture, forming the basis of many contemporary box-based networks [34], [35]. Recent box-free methods like SOLO [36] and CondInst [37] devise FCN-like solutions, which offer simpler architecture with comparable segmentation performance to box-based models. Moreover, additional methods tackle instance-specific image segmentation via clustering [38], [39], class-agnostic mask generation [40], [41], and contour-based techniques [4], [42].

The lately introduced panoptic segmentation task [1] unifies semantic- and instance segmentation in a single system, facilitating to recognize both things and stuff in urban driving scenes, which are of important relevance for autonomous vehicles. A multitude of works [7], [43], [44] implement panoptic segmentation with a universal framework, showing the significance of the new task to driving scene parsing. Among these approaches, some develop based on state-of-the-art instance segmentation methods like Panoptic FPN [45] extending Mask R-CNN with a semantic branch for stuff segmentation [46], [47]. Another cluster of approaches enhances semantic segmentation architectures like Panoptic-DeepLab [48] extending DeepLab [23]. Further, there are works that capture the relations among semantic categories and instances for providing richer contexts to enhance visual understanding in panoptic segmentation [49].

More recently, there are single-path architectures without using separate branches, which are realized via conditional convolution filters [50], [51], learnable kernels [52], [53], and unified panoptic embeddings [54]. Transformer models have also been introduced into panoptic segmentation due to their capability to model long-range dependencies [19], [55], [56]. Inspired by DETR [55], MaX-DeepLab [57] and Mask-Former [58] view unified image segmentation from a mask classification perspective. MaX-DeepLab builds atop [59] by extending the axial-attention backbone with a dual-path framework combining CNNs and transformers in the head network. Building upon MaskFormer, Mask2Former [60] devises masked attention to extract localized features by constraining cross-attention with predicted mask regions. Panoptic Seg-Former [61] leverages an efficient deep supervision mask decoder and a query decoupling strategy to delve deep into panoptic segmentation. In this work, considering that fast responses are critical for autonomous driving, we build on CNN- and transformer-based panoptic segmentation architectures [7], [52], [60]. We make use of our efficient system and address panoptic segmentation in $360°$ panoramic imagery to pursue a unified and complete surrounding understanding.

### B. Panoramic Scene Segmentation

Modern scene segmenters are mostly designed to work with pinhole images on mainstream datasets such as Cityscapes [15] and Mapillary Vistas [62]. To enlarge the Field of View (FoV), many surrounding understanding platforms are based on fisheye images or multiple cameras [63], [64], [65], [66],

[67]. However, this either comes with severe distortions in particular around the fisheye image borders or leads to being cumbersome with a lot of multi-camera calibration work. Motivated by the prospect of attaining $360°$ perception with a single panoramic camera, recent works [68], [69], [70] build directly on this modality, but they rely on synthetic collections that are far less diverse than pinhole databases [15], [62]. A high variety of images is however critical for yielding robust segmentation models for real-world perception.

In contrast, the Panoramic Annular Semantic Segmentation (PASS) framework [2] reuses knowledge in pinhole data to produce robust models suitable for $360°$ images, and it is further augmented by DS-PASS [71] via a detail-sensitive design with lateral attention connections. In [27], context-aware omni-supervised models are taken to the wild, fulfilling panoramic semantic segmentation in a single pass with enhanced generalizability. P2PDA [72] explicitly tackles panoramic segmentation from a domain adaptation perspective by transferring from the label-rich pinhole domain to the label-scarce panoramic domain. In [73], a distortion convolutional module is developed to correct the panoramic image distortion according to the image-forming principle. Trans4PASS [3] architects a distortion-aware transformer with deformable token mixers for adapting to panoramic images. However, existing $360°$ perception systems only render semantic- or instance-specific segmentation, whereas a recent video panoptic segmentation work [74] based on the Waymo open dataset only offers a coverage of $220°$. To achieve a unified and holistic scene understanding, this work advocates and first addresses panoramic panoptic segmentation, extending panoptic segmentation models with dense contrastive learning regimens that intertwine pixel-level consistency propagation for robust segmentation across pinhole- and panoramic imagery.

### C. Unsupervised Dense Contrastive Learning

Currently, the most appealing approaches for learning representations without labels are unsupervised contrastive learning [16], [75], [76] tasks. Contrastive learning methods learn visual representations in a discriminative way by contrasting similar, positive pairs against dissimilar, negative pairs, which is promising to yield generalized features for robust predictions in previously unseen domains. The training pairs are often generated from augmented views of image samples, and thereby the previous methods are mostly designed for image classification tasks, which does not ensure more accurate pixel-wise segmentation [16]. Different from previous works, we aim to develop a contrastive learning regimen for pixel-level tasks. Taking the wide FoV of omnidirectional data into consideration, in this paper, we put forward a learning framework for fine-grained panoptic segmentation operating on panoramic images towards a holistic understanding.

A few latest contrastive training methods [17], [77], concurrent to our work, also address dense prediction tasks. DenseCL [17] implements self-supervision by optimizing a pairwise contrastive (dis)similarity loss between two views of input images, whereas pixel-level pretext tasks are introduced for learning dense feature representations in [77]. FisheyePix-Pro [78] attempts to pretrain a contrastive learning based

model directly on fisheye images. Cross-image pixel contrast has been leveraged for semantic segmentation by looking beyond single images [79], [80], [81] and enforcing pixel embeddings belonging to the same semantic class to be more similar than embeddings from different classes. Some methods also explore pixel-to-region contrast [79], [80], [82], [83] as a complementary to the pixel-to-pixel contrast strategy. More recently, MaskCo [84] introduces contrastive mask prediction for visual representation learning. ORL [85] leverages image-level self-supervised pretraining as the prior to discover object-level semantic correspondence. DSC [86] models semantic category decision boundaries at a dense level with multi-granularity contrastive learning. In this work, we step further to open $360°$ panoramic scenes and explore generalization effects from dense contrastive learning in a cost-minimizing way.

## III. PROPOSED FRAMEWORK

As the target of this work is to establish the most holistic surrounding understanding based on standard images, we are dealing with downstream tasks operating on a pixel level going beyond image level classification tasks. This observation is a key design principle we considered when mitigating the performance drop between the pinhole image source distribution and the panoramic target domain. Within this section, we show how to overcome the lack of annotated panoramic images by proposing the Panoramic Robust Feature (PRF) framework which generates robust features from pinhole images in a cost minimizing way, suitable for a mobile agent setting.

The proposed procedure consists of two steps in order to adapt panoptic target models. The first step focuses on the backbone of the model which is pretrained in a pixel-level contrastive fashion, the second step is the standard supervised training of the panoptic model.

During the first step, the feature space, which was learned via the standard supervised ImageNet [87] classification task, serves as a starting point. Despite the proven ability of the commonly-used ImageNet weights for transfer learning tasks, we find that facing domain shift problems such as transferring from standard pinhole- to panoramic images causes severe performance drops. This behaviour may be caused by the fact that features learned via a final linear classification task are not encouraged to be well separated beyond what is necessary for a linear classifier, as shown in Fig. 3. Bearing in mind that intelligent vehicles often entail flexible model development and deployment, we aim to propose an approach which allows model modifications or retraining on a hardware setting which can run on a mobile agent.

Due to the limited availability of purely contrastive model backbones and the required access to at least 8 GPUs [16], [77] to train them from scratch in order to achieve comparable results to the widely available ImageNet [87] pretrained feature extractors, we unite the two approaches in our mixed training approach and take the best of both worlds. We use the widely available supervised ImageNet [87] pretrained feature extractors and modify these weights using the contrastive task as described in Section III-A. After adapting the backbone, a standard supervised training can be applied which will be introduced in Section III-B. This makes the proposed Panoramic
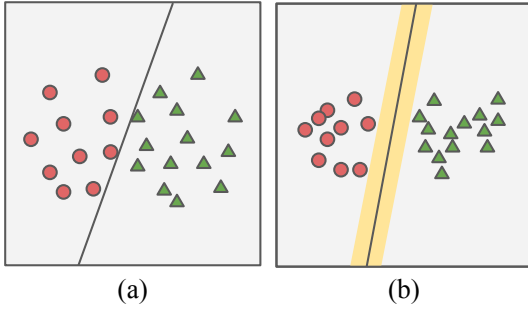
**Fig. 3:** Simplified representation of the feature embeddings obtained by a supervised training task via a linear classification (a) and a contrastive pretraining task (b). Supervised training features are only separated from each other to the extend which is necessary to maximize the performance of a linear classifier. It can be easily seen that the learned feature representation in (a) leads to the same classification result as the feature representation in (b) if the distribution of the data does not change by much. Via contrastive training, we do not only encourage similar features to be represented in a similar way but also dissimilar features to be represented as dissimilar as possible. As shown in (b), this leads to more robust features which can facilitate the downstream task in the face of distribution changes.

Robust Feature (PRF) framework a drop-in, model-agnostic procedure, which does not require changes in model architectures between pretraining and application-oriented fine-tuning. Furthermore, the required pretraining phase runs on a single GPU and finishes within very reasonable times, making our mixed training approach ideal for a limited hardware setting.

### A. Contrastive Pretraining

The contrastive pretraining phase starts with the commonly-used supervised training weights $\theta_0^{enc}$ obtained by the 1000 classes ImageNet [87] classification task for a ResNet18 backbone [20].

From a given image, we crop two random views, augment them, and feed one of them through the encoder and the other through its momentum encoder [16] initialized as $\theta_0 = \theta_0^{enc}$. The augmentation techniques include standard procedures such as color jittering, flipping, gray-scale conversion, or solarization. We use a standard ResNet18 encoder with a projection head on top. Following the work of [77], the projection head is a simple network consisting of a $1 \times 1$ convolutional layer, projecting the features into a 2048-dimensional latent space, a batch normalization layer, a ReLU nonlinearity layer, and a final $1 \times 1$ convolutional layer projecting the features to a target dimensionality of 256. We obtain two feature maps containing the features from the respective crops. The basic idea of the approach is that if the views overlap at a certain point, the feature maps of the corresponding pixels should be very similar. This idea builds upon the inherent translation equivariance of CNNs. After each step, we update the momentum encoder according to the following rule:

$$\theta_t = \theta_{t-1} \times \beta + (1 - \beta) \times \theta_t^{enc} \qquad (1)$$

This update rule leads to an exponentially weighted moving average encoder of the regular encoder. Following [88], we set $\beta = 0.99$ which leads to both of the networks having similar weights, hence they should produce similar results for similar inputs which enables us to exploit two expected consistencies across the two feature maps. We briefly explain them and their effect on panoramic panoptic image segmentation in this section and refer the interested reader to [77]. Panoramic images are characterized by their wide FoV which requires the kernels of the backbone to analyze a plethora of different objects and pay attention to details as well as global contexts. **Spatial Contrastive Loss:** This loss encourages the network to focus on the details of feature computations as it enforces the backbone to compute features that are similar for spatially close pixels and dissimilar for more distant pixels. This is achieved by comparing two feature vectors $x_i$ and $x_i'$ computed from pixel $p_i$ by the regular encoder and the momentum encoder respectively. If features originate from the same or spatially close pixels (indicated by the green frame in Fig. 4) within the original image space, the features should be similar, despite their different locations within the two crops. Consider the first cropped view $F$ and pixel $p_i$ located in $F$. We define $\tilde{F}_i^P$ as the pixels located in the second view, which are spatially close to $p_i$ measured within the original image space. The rest of the pixels in the second view are denoted by $\tilde{F}_i^N$ as their distance from $p_i$ exceeds a certain threshold. $\tau$ is a normalization parameter set to $0.3$ as suggested by [77]. The overall loss, as shown in Equation 2, is averaged across all pixels in each of the views, and the final loss is obtained as the average of the losses of the two views.

$$L_s(p_i) = -\log \frac{\sum_{j \in \tilde{F}_i^P} e^{\frac{cos(x_i, x_j')}{\tau}}}{\sum_{j \in \tilde{F}_i^P} e^{\frac{cos(x_i, x_j')}{\tau}} + \sum_{k \in \tilde{F}_i^N} e^{\frac{cos(x_i, x_k')}{\tau}}} \qquad (2)$$

**Global Propagation Loss:** This loss focuses on global consistency among similar pixels by propagating globally similar pixels onto the current pixel as shown by the green arrows in Fig. 4. This guarantees that the model assigns the same features to similar pixels beyond the FoV of the current kernel location, making globally similar pixels semantically similar, which is in particular important for large-FoV panoramic images. An additional effect of this feature propagation is that it induces a smoothing effect, which prefers smooth outputs over fragmented ones as it averages out small differences.

A smoothed feature vector $x_i^{smooth}$ can be calculated according to Equation 3 as a weighted sum consisting of all the projected features $g(x)$ in the entire view. The weights are determined by the similarity of the two respective feature vectors $x_i$ and $x_j$.

$$x_i^{smooth} = \sum_{j \in F} \max(cos(x_i, x_j), 0)^2 * g(x_j) \qquad (3)$$

Here, the transformation $g(x)$ can be computed via a $1 \times 1$ convolution keeping the number of channels constant.

Finally, the feature maps obtained by the encoder and the momentum encoder should produce consistent results for spatially close pixels $p_i$ and $p_j$ in the image space. This behavior is encouraged by the global propagation loss shown in Equation 4.

$$L_{GloPro} = -cos(x_i^{smooth}, x_j') - cos(x_j^{smooth}, x_i') \qquad (4)$$
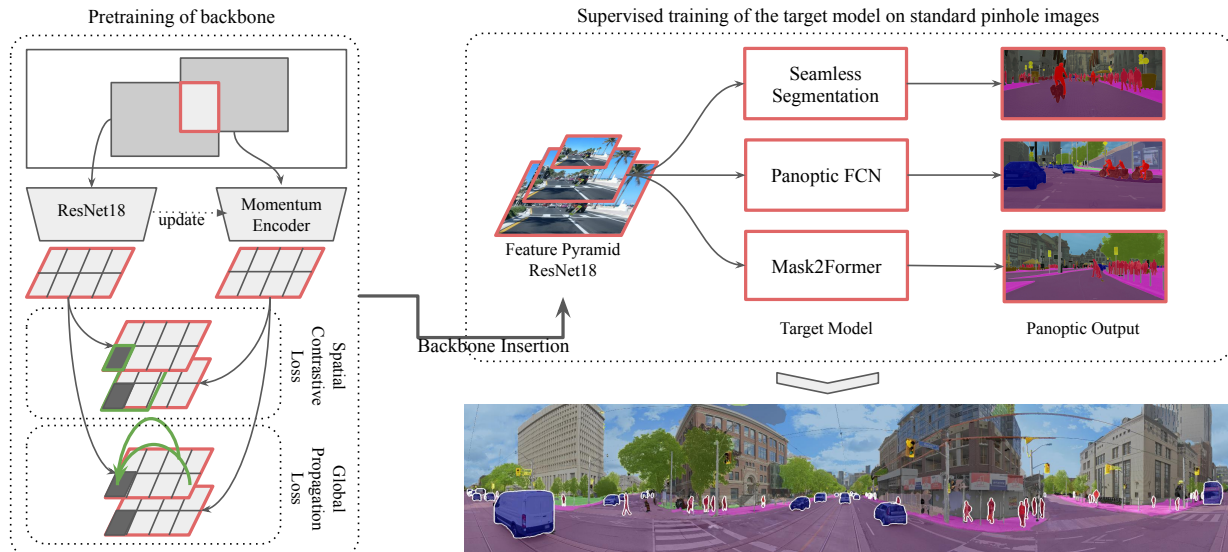
**Fig. 4:** The Panoramic Robust Feature (PRF) framework consists of two steps: A contrastive pretraining step as shown on the left side and the standard supervised training of the target model. During the pretraining, the features of the backbone are made robust by contrasting spatially close pixels against more distant pixels (green box), and globally-similar pixels are propagated onto each other. The robust backbone is inserted into the target model and the next step is to use the standard supervised training procedure on a pinhole image dataset such as Cityscapes [15]. After training the robust model, it can produce high-quality panoptic annotations on our panoramic image dataset WildPPS. Within this work, we experiment with different target models in order to prove the effectiveness of the proposed approach and to identify models which are most suitable for panoramic panoptic segmentation.

We find in our experiments that the additive combination of the Spatial Contrastive loss and the Global Propagation loss according to Equation 5 is very beneficial for our desired application and produces the expected smooth results. Unless stated otherwise, we set $\alpha=1$. More details about the experiments, ablation studies as well as qualitative results are shown in Sections IV-C and IV-F.

$$L_{Pretrain} = L_s + \alpha \times L_{GloPro} \qquad (5)$$

### B. Supervised Training

One advantage of the proposed framework is its simple feasibility since it is a drop-in method that does not require any changes in the model architecture or the model training procedure. The only required changes may be due to the chosen mobile agent efficiency setting which will be briefly discussed in Section IV-B. The panoptic interpretation of the panoramic image can be obtained in a single forward pass and does not require any post-processing steps such as image stitching or fusion of different information sources.

The main design criteria of our approach are driven by the need for flexibility and efficiency. We want to establish a drop-in model-agnostic procedure for the pretraining as well as the supervised training phase. The proposed methods and models should work as efficiently as possible since energy is a valuable good in any mobile setting. We thus restrict the entire setup to a maximum number of two GPUs which we find a reasonable upper bound for a mobile agent.

The proposed efficient modifications of current state-of-the-art panoptic segmentation models follow our conference paper [6] and meet our established criteria, since we replace the standard ResNet50 [20] with a ResNet18 [20] backbone. This reduces the number of parameters in the backbone by

54% from 25.6 million to 11.7 million but still remains a model-agnostic drop-in method. A more detailed efficiency comparison with the two different backbones plugged into the respective models as well as a comparison in inference speed are shown in Table IV.

Our selection of panoptic target models is motivated by our goal to determine if certain types of models perform better than others in panoramic panoptic segmentation and to identify the most promising architectures to advance in this task. In order to provide a fair comparison, we make sure that all of the selected models are withing the same order of magnitude regarding the number of parameters, required Flops and processing time. We verify this in Table IV and Table V. We differentiate between two-branch models combining the outputs of well-established instance and semantic segmentation models via a heuristic merging principle, single-branch models, and transformer-based models. The following models serve as representatives for the aforementioned categories:

- **Seamless-Scene-Segmentation** [7]**:** The seamless scene segmentation model combines two well-established baseline models: the Mask R-CNN [18] performing instance segmentation on the input and a DeepLabV3 [89] inspired second branch computes the semantic segmentation maps. Finally, a fusion step similar to [1] combines the output of the two branches to the panoptic image. The model serves as an example of a two-branch model, which inherently treats things and stuff differently.
- **Panoptic-FCN** [52]**:** Whereas previous models inherently treat things and stuff as different branches, the approach of this model is to treat them in a unified single-branch approach. It uses a kernel generator, which generates a dedicated convolutional kernel for each individual in-

stance and each stuff region. Aside from the kernels, a high-resolution feature map is computed. The authors use the semantic branch of the Panoptic-FPN model [45] as their feature encoder. The final panoptic results can be obtained via a convolution operation of the learned kernels with the encoded features. This eliminates the need for the final merging operation of the stuff and thing results as it was necessary in the two-branch models.

- **Mask2Former** [60]**:** Finally, we want to explore the capability of transformer-based models on panoramic panoptic image segmentation. We pick the Mask2Former model as a well-performing model with a general architecture which allows to perform semantic, instance, and panoptic segmentation without architectural changes. A transformer decoder in combination with an MLP computes the mask embeddings and class predictions for the masks based on the image features. The output of the model consists of a set of predicted masks, which are obtained by calculating the dot product between the mask embeddings and the encoded image features. Depending on the target task, the set of masks is post-processed to match the expected output format.

The quality of the panoramic output is calculated according to the Panoptic Quality (PQ) measure [1] as shown in Equation 6.

$$PQ = \frac{\sum_{(p,g)\in TP} IoU(p,g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|} \qquad (6)$$

Each of the efficient models generates panoramic panoptic segmentation in a single shot and is ready to be deployed in intelligent vehicle systems, delivering a complete and robust surrounding understanding.

## IV. EXPERIMENTS AND DISCUSSIONS

### A. Experiment Datasets

**Cityscapes dataset.** By doing pioneering work in the novel task of panoramic panoptic image segmentation, we are faced with two data-related challenges which result from the lack of panoramic images with panoptic annotations. As we cannot generate a dataset of sufficient size which can serve as a panoramic panoptic training dataset, we are forced to train on publicly available pinhole image datasets. Due to the wide acceptance in the field, we rely on Cityscapes [15] for the supervised training procedure of the target models. Cityscapes is a dataset of street scene images captured under similar weather and lighting conditions in 50 different locations in Germany and Switzerland. It consists of 2975 training, 500 validation, and 1525 test images with fine panoptic annotations.

**Mapillary Vistas dataset.** Mapillary Vistas [62] is a diverse street scene dataset which consists of $18k$ training, $2k$ validation, and $5k$ test images which are captured from all around the globe. The images have been collected under various weather and lighting conditions. In this work, we use the Mapillary Vistas dataset to benchmark our most promising models against previous state-of-the-art semantic segmentation models on the panoramic PASS dataset [2], as it uses the labeling scheme of Mapillary Vistas.

**WildPPS: Wild Panoramic Panoptic Segmentation dataset.** In order to evaluate the models on the panoramic panoptic segmentation task, we introduce WildPPS [6], which is the first publicly available panoramic panoptic image dataset. In the latest version of WildPPS, we double the number of annotations, which results in 80 panoramic images with fine panoptic annotations for the most important street scene classes. We provide annotations for the stuff classes *street* and *sidewalk* which are necessary to find traversable areas and the thing classes *car* and *person* that are essential in order to interact with other road users. WildPPS builds on top of WildPASS [90] and contains panoramic images of 40 cities from all around the globe. It is a very diverse dataset containing street scenes not only from European or American cities but also from historically underrepresented regions such as Southeast Asia or Africa. The images are captured under various weather, lightning, and environmental conditions, which makes the dataset a challenging target domain. We adopt the Cityscapes annotation style in order to minimize the threshold to evaluate panoptic models on WildPPS.

### B. Experimental Setup

**Baseline setups.** As mentioned before, we use the Cityscapes dataset [15] due to its wide acceptance within the community as the dataset of choice to perform the supervised training. All the experiments performed within this work were designed and executed with a mobile agent setting in mind. For the supervised training procedure, we restrict ourselves to the usage of two GPUs and otherwise follow the proposed training procedure of the authors of the respective models as closely as possible. The only noteworthy differences in the supervised training procedures result from adapting the original training procedures which typically rely on at least 8 GPUs [7], [52], [60] to a 2 GPU setting, which is achieved by reducing the batch size and adapting the learning rate to the lower batch size. We experiment with the three different target models mentioned in Section III-B and train them according to the previously discussed procedure. In order to assess the performance of the Panoramic Robust Feature (PRF) framework, a baseline model for each of the target models serves as a benchmark. The baseline model is simply the target model only trained in a supervised way on Cityscapes without a prior backbone adaption. The baseline models are evaluated on the panoramic WildPPS dataset and compared against multiple robust models.

**Pretraining settings.** In order to generate robust models, we experiment with three different pretraining settings. All of the following pretraining experiments are performed on a single GPU and use the ImageNet pretrained weights as a starting point as described in Section III-A. Reusing the knowledge in these weights reduces the adaptation time by two orders of magnitude compared to a purely contrastive training, which will be described later.

We achieve the best improvements of more than $6\%$ measured in PQ over the different non-adapted baseline models by pretraining on the same Cityscapes dataset on which we perform the supervised training. This eliminates the need for

**TABLE I:** Overview of the main results of the proposed method. The respective target model is shown on the left whereas the respective robust adaptation technique is listed in the header. Bold numbers indicate the best performance for the respective target model.

| Panoptic Model | Baseline | Panoramic Robust Feature (PRF) Approaches | | | |
|---|---|---|---|---|---|
| | | SGD | LARS | LARS Large | Pure |
| Seamless Segmentation [7] | 59.43% | 63.24% | 62.76% | 61.54% | **63.65%** |
| Seamless Segmentation PQ Stuff | 61.55% | 68.10% | 59.93% | 55.84% | **69.36%** |
| Seamless Segmentation PQ Things | 57.30% | 58.38% | 65.60% | **67.25%** | 57.93% |
| Panoptic FCN [52] | 49.29% | 50.77% | 44.82% | 41.26% | **52.93%** |
| Panoptic FCN PQ Stuff | 60.05% | 62.39% | 55.35% | 48.71% | **66.60%** |
| Panoptic FCN PQ Things | 38.53% | 39.16% | 34.28% | 33.80% | **39.26%** |
| Mask2Former [60] | 49.47% | 54.93% | **55.95%** | 52.79% | 51.91% |
| Mask2Former Stuff | 63.83% | 69.60% | **70.91%** | 69.15% | 69.69% |
| Mask2Former Things | 35.11% | 40.26% | **40.99%** | 36.43% | 34.13% |

any other data than the one already available for the supervised training procedure. We train for 90 epochs using a batch size of 100 images. This procedure can be achieved on a single GPU by resizing the Cityscapes images to $256{\times}512$ pixels before feeding them into the framework.

- Following [77], we use a LARS optimizer [91], a base learning rate of $0.4$ and a cosine learning rate schedule with restarts every 30 epochs. This serves as a regularization since the final model can be interpreted as an average across the restarts. We furthermore regularize the training by decaying the weights using a $1e^{-5}$ penalty.
- The LARS optimizer was designed to work well with large batch sizes, hence we add a line of experiments in which the batch size is doubled to 200. This setting will be referred to as LARS large.
- In order to compare the LARS optimizer setting to a more standard one, we also make use of a SGD optimizer with momentum and a step learning rate scheduler in which we decrease the learning rate from $1e^{-3}$ by a factor of 10 every 30 epochs. The weight decay remains at $1e^{-5}$.

**Purely contrastive counterpart.** Finally, we want to benchmark our novel training approach which combines supervised and unsupervised training against a purely contrastive trained backbone. In order to train a ResNet18 backbone [20] from scratch, we perform our contrastive pretraining procedure starting with randomly initialized weights. In order to compare the contrastive backbone to the supervised backbones, we pretrain it on ImageNet [87] for 100 epochs with a batch size of 256 images on 4 GPUs using the LARS optimizer and a cosine learning rate schedule. Ideally, we would want to use 8 GPUs, however the performance of the target models on the Cityscapes validation set is on par with that of their supervised counterpart, which makes us reasonably confident regarding the contrastive training. One exception is the Seamless Scene Segmentation model for which the pure backbone scores lower. We suppose that the supervised training procedure should be adapted to the contrastive backbone, but for a fair comparison, we keep the model. The total training time of this contrastive backbone was about $3.5$ days on 4 GPUs compared to around $2h$ for the Panoramic Robust Feature (PRF) adaptation on a single GPU.

*C. Experimental Results*

The main goal of this section is twofold: We want to verify the effectiveness of the proposed PRF framework. This is achieved by comparing the results of the robust versions of the target models to the non-adapted baseline models. The second goal is to identify the most promising models to perform panoramic panoptic segmentation and answer the question if certain models are by design more suitable to excel at this novel and challenging task.

**On the effectiveness of the PRF framework.**

The effectiveness of the PRF framework can be confirmed by comparing the WildPPS Panoptic Quality Scores of the Baseline models with those of the adapted models as shown in Table I. Regarding the Seamless Scene Segmentation model [7], the effectiveness clearly can be confirmed as all of the different pretraining settings outperform the baseline model. The best performing adaptation is achieved by training with the purely trained backbone followed by the SGD setting.

Regarding the Panoptic FCN model [52], we find the same order of models outperforming the non-adapted target model. The LARS-adapted backbones do not work well and even worsen the performance. We conjecture, that the supervised training procedure of this model must be changed in order to profit from the LARS-adapted backbone features. The SGD-adapted backbone is, in contrast to the LARS-adapted backbone, still very similar to the original backbone as it is shown in Fig. 9, so the model training which has been optimized for the supervised baseline backbone still works quite well for the SGD backbone. The same line of arguments accounts for the purely trained backbone which produces more well-behaved embeddings without extreme feature isolation. We also investigate this behaviour in more detail within the qualitative analysis section. It seems that the FCN model tends to confuse zebra crossings with sidewalks, as shown in Fig. 10.

Finally, the transformer-based Mask2Former model [60] can be improved reliably using the PRF framework. The best performing configuration is the LARS-based approach confirming that models that can capture the more heavily separated latent space modification benefit from it.

Summarizing the presented results, we are confident regarding the effectiveness of the PRF framework. We suggest the SGD-based backbone as the default setting since we find it consistently outperforming the baseline and generating the second-best results in all settings using the training procedures
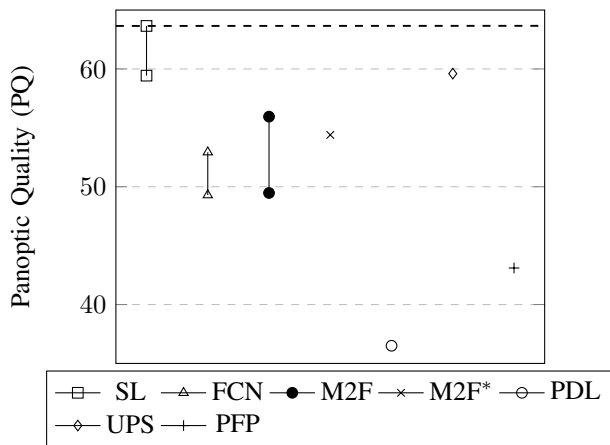
**Fig. 5:** Comparison of our proposed efficient robust models with multiple panoptic baseline models. SL, FCN, M2F represent the Seamless-Scence-Segmentation [7], the Panoptic-FCN [52] and the Mask2Former [60] model. These are our efficient baseline models along with their robust adaptations. M2F*, PDL, UPS and PFP represent the Mask2Former, Panoptic-DeepLab [48], the UPSNet [43], and the Panoptic-FPN [45] model respectively. All of the latter models use a ResNet50 backbone [20] and are trained in a standard fashion.

proposed by the authors of the target models. In order to maximize the performance of a given target model, it can be beneficial to open up to the LARS or the Pure settings. We furthermore want to point out that in comparison to the purely contrastive backbone, we achieve comparable or better results using the orders of magnitude cheaper mixed training approaches, namely SGD, LARS, or LARS Large. This confirms the suitability of our framework for the desired mobile application setting.

**On the maximization of segmentation performance.** The second posed question was on maximizing the performance of Panoramic Panoptic Segmentation. This question can be easily addressed, as the Seamless Scene Segmentation model [7] outperforms any other model by a large margin, as it can be seen in Table I, and provides the best score for PQ Things across all tested scenarios. The baseline model not only outperforms all the other baseline models but also the best-adapted versions of those models. We argue that this is largely due to the fact that the Seamless model uses two separate state-of-the-art models to compute the semantic segmentation map and the instance segmentation results, which are then combined into the final panoptic parsing result. These models have proven themselves in countless scenarios and have been updated and refined over the years, which makes them very robust by design.

Regarding the remaining models, we cannot recommend the Panoptic FCN model [52] due to the low performance and the inability to cope with the modified LARS backbones.

The Mask2Former model [60], despite not performing well in its baseline scenario, can be improved to a solid $55.95\%$ PQ score and yields the best result for PQ Stuff among all the models, which might very well be attributed to the used attention mechanisms allowing global information exchange particularly relevant for the heavily distorted stuff classes.

Nevertheless, our clear recommendation at this point is to use our robust Seamless Scene Segmentation model for
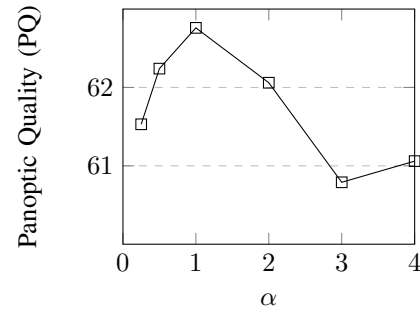


**Fig. 6:** Ablation study on the $\alpha$ hyper-parameter for the LARS setting. The x-axis indicates the weight between the Spatial Contrastive and the Global Propagation Loss as shown in Equation 5.

domain shift problems.

**Comparison to panoptic baseline models.** In order to prove the effectiveness of our robust efficient models against other panoptic models, we test the performance of standard panoptic models on WildPPS. We report the results in Fig. 5. We trained the standard models as closely as possible to the settings proposed by the authors. For Panoptic-DeepLab, we replace the average pooling in the ASPP module with global average pooling in order to be able to use the model with changing image resolutions during inference. The performances of the baseline models on the Cityscapes validation dataset are comparable with each delivering about the same results as in the original papers. Their varying performances on WildPPS highlights once more the difficulty of the panoramic domain shift which not every model can cope with equally well. Not only is it visible that the proposed PRF framework yields consistent improvements over the non-adapted models, it furthermore can be seen that the proposed robust Seamless-Scene-Segmentation model outperforms a large variety of different models. This holds despite the fact, that the other baseline models were trained with a multiple of GPUs and use the more powerful ResNet50 backbone.

**Parameter study on $\alpha$.** As discussed in Section III-A, both the Spatial Contrastive loss and the Global Propagation loss play important roles in capturing both local as well as global image information. In order to determine the relative importance of the two losses, we perform an ablation study on $\alpha$ as introduced in Equation 5. Following [6], [77], we use the LARS setting of the Seamless Scene Segmentation model [7] to perform this study on the novel extended WildPPS dataset. The results are shown in Fig. 6 and indicate that within reasonable bounds around $\alpha=1$, the results are quite similar. Deviating further from 1 tends to decrease the performance. We obtain the best performance for exactly $\alpha=1$ which deviates from the results of the conference paper in which $\alpha=2$ performs best, however, due to the larger dataset used in the current work and the comparable results within $0.5\leq\alpha\leq2$, we believe that the performance within the stated interval is comparable and a more fine-grained hyper-parameter search could be restricted to the this range in order to maximize the performance.

**Generalization to different FoVs.** Generally speaking, a difficulty when transferring from pinhole- to panoramic images is the increased Field of View (FoV) and the distortion resulting from mapping spherical image data to 2D planes. If
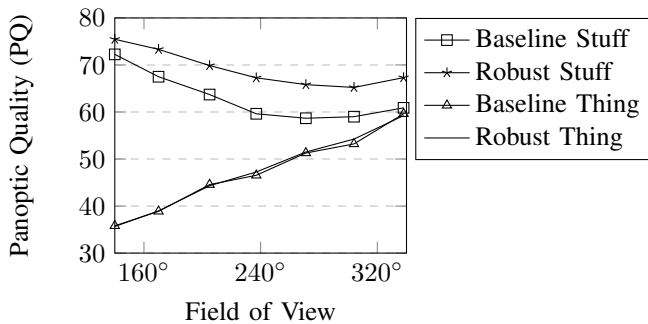
**Fig. 7:** Detailed comparison of a standard Seamless Scene Segmentation model [7] to its robust counterpart depending on the Field of View (FoV) on the WildPPS dataset.

the proposed robust modifications help to mitigate the targeted distribution shift from pinhole to panoramic images, we expect two behaviors: We should consistently outperform the baseline model across different FoVs and the performance gains should generally increase with increasing FoV. Both behaviors can be seen in Fig. 2. As the increase in PQ with increasing FoVs is counter-intuitive, we explore this behavior in more detail by separating the Panotic Quality into the parts PQ Stuff and PQ Things. The results are shown in Fig. 7. As we crop starting from the center of the panorama, the Stuff areas tend to look similar to the Cityscapes dataset. By widening the FoV, the stuff areas get stretched across the entire panoramic FoV which is unusual, especially for the sidewalk class. The performance drop in PQ Stuff is mitigated by the robust backbone features. The counter-intuitive increase in PQ is caused by the increase in PQ Things which can easily be explained as cropping from the center in panoramic images tends to crop out small and distant objects, whereas the larger and easier to segment objects are on the side of the image. The difference in PQ Things of the two models is negligible, as cars and persons usually are present at different locations during training thus not leading to positional priors. On top, they are on average less affected by panoramic distortions due to their smaller surface area.

The presented results indicate that our model would also be potentially applicable and beneficial in other wide-FoV scenarios such as fisheye image segmentation and surround-camera perception, *e.g.,* the $220°$ sensing of the Waymo open dataset [74], which are left for future exploration.

**Generalization on Mapillary Vistas.** Besides domain shifts from standard pinhole images to panoramic images, a related common problem is to have pinhole images obtained from a different distribution. As this is not the main focus of this work, we just want to address this problem briefly within this section. We extend the results of our conference version [6] and add the novel LARS Large and Pure setting. Since Mapillary Vistas is a much more diverse dataset compared to Cityscapes, we map the learned Cityscapes labels to the corresponding Vistas label and ignore all the classes for which the model does not predict labels. The mapping is straightforward with two exceptions: We map the Cityscapes classes *Rider* and *Train* to the Mapillary classes *Motorcyclist* and *Other Vehicle* respectively. The results are shown in Table II and indicate that not only the Seamless Scene Segmentation model is a very

**TABLE II:** Comparison of results of the Seamless Scene Segmentation model [7] on the Mapillary Vistas dataset [62].

| Pretrain Setting | PQ | PQ Stuff | PQ Things |
|---|---|---|---|
| Baseline Model | 30.2% | 38.6% | 20.7% |
| **Pretrain Cityscapes SGD** | **31.9%** | **41.7%** | 20.7% |
| Pretrain Cityscapes LARS | 31.3% | 40.2% | **21.2%** |
| Pretrain Cityscapes LARS Large | 30.14% | 38.31% | 20.81% |
| Pretrain Cityscapes Pure | 27.89% | 34.26% | 20.60% |

suitable model due to the inherent robustness for domain shift problems, but the results once more confirm the effectiveness of the PRF framework. We want to point out that the reported numbers in Table II cannot be compared to the performance reports in the of the original Seamless-Scene-Segmentation paper. Our model was trained on Cityscapes and only performs inference on the classes available in Cityscapes which is only a fraction of the much more diverse Mapillary Vistas dataset.

### D. Comparison to Previous Models on PASS

We compare our best model identified in Section IV-C which is the Seamless-Scene-Segmentation model [7] to previous works. As this work is the first to achieve panoptic segmentation on panoramic images, we can only compare it to semantic segmentation results on the publicly available Panoramic Annular Semantic Segmentation (PASS) benchmark [2]. PASS includes $400$ unfolded panoramic annular images with pixel-level semantic annotations on 6 navigation-relevant classes for evaluation with a resolution of $692{\times}2048$.

These compared networks, experimented by [90], cover state-of-the-art accuracy-oriented segmentation models Seg-Net [22], PSPNet50 [24], DenseASPP [92], and DANet [26]. Efficiency-oriented models are also included for a comprehensive comparative analysis, spanning ENet [93], CGNet [94], ERFNet [95], PSPNet18 [24], ERF-PSPNet [2], SwifNet [96], and SwaftNet [71]. They all view the wide-FoV panoramic image as a single input without separating it into multiple segments. Under a fair comparison, we measure the Intersection over Union (IoU) of the different classes. We also measure the complexity in MACs on an input resolution of $512{\times}1024$ and the number of parameters (#Params) of our method by following the setup of [97].

Besides restricting ourselves to the semantic output of our model, we point out that our efficient modification is also among the most efficient models within the comparison. We apply our proposed PRF framework to train a ResNet18 backbone followed by supervised training on Mapillary Vistas [62], as PASS is designed with the same labeling strategy as Vistas. We follow the supervised training procedure described in Section III-B for a total of 23 epochs. Staring with an original learning rate of $1e^{-2}$, we reduce by a factor of 10 after $95k$ and $125k$ iterations due to the observed loss saturation. The results shown in Table III illustrate that compared to the previous state-of-the-art segmentation networks, we surpass them by a large margin by using our proposed framework. In particular, the safety-critical class *person* is well addressed by our model, outstripping the second-best model by $17.3\%$ in IoU. While our model is among the larger models, there are two important observations to make. Model size only seems to play an inferior role as our model does outperform larger models

**TABLE III:** Per-class accuracy analysis in Intersection over Union (IoU) and mean IoU (mIoU) on the public Panoramic Annular Semantic Segmentation (PASS) dataset [2].

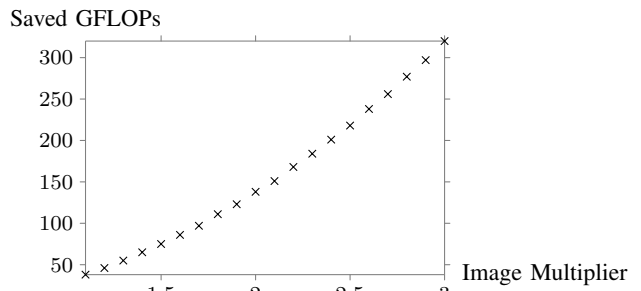| Network | Car | Road | Sidewalk | Crosswalk | Curb | Person | mIoU | MACs | #PARAMs |
|---|---|---|---|---|---|---|---|---|---|
| SegNet [22] | 57.5% | 52.6% | 17.9% | 11.3% | 11.6% | 3.5% | 25.7% | 398.3G | 28.4M |
| PSPNet (ResNet50) [24] | 76.2% | 67.9% | 34.7% | 19.7% | 27.3% | 22.6% | 41.4% | 403.0G | 53.3M |
| DenseASPP (DenseNet121) [92] | 65.8% | 62.9% | 30.5% | 8.7% | 23.0% | 8.7% | 33.3% | 78.3G | 8.3M |
| DANet (ResNet50) [26] | 70.0% | 67.8% | 35.9% | 21.3% | 12.6% | 25.9% | 38.9% | 114.1G | 47.4M |
| ENet [93] | 59.4% | 59.6% | 27.1% | 16.3% | 15.4% | 8.2% | 31.0% | 4.9G | 0.4M |
| CGNet [94] | 65.2% | 56.9% | 23.7% | 3.8% | 11.2% | 21.4% | 30.4% | 7.0G | 0.5M |
| ERFNet [95] | 70.0% | 57.3% | 25.4% | 22.9% | 15.8% | 15.3% | 34.3% | 30.3G | 2.1M |
| PSPNet (ResNet18) [24] | 64.1% | 67.7% | 31.2% | 15.1% | 17.5% | 12.8% | 34.8% | 235.0G | 17.5M |
| ERF-PSPNet [2] | 71.8% | 65.7% | 32.9% | 29.2% | 19.7% | 15.8% | 39.2% | 26.6G | 2.5M |
| ERF-PSPNet (Omni-Supervised) [90] | 81.4% | **71.9%** | **39.1%** | 24.6% | 26.4% | 44.1% | 47.9% | 26.6G | 2.5M |
| SwiftNet [96] | 67.5% | 70.0% | 30.0% | 21.4% | 21.9% | 13.7% | 37.4% | 41.7G | 11.8M |
| SwaftNet [71] | 76.4% | 64.1% | 33.8% | 9.6% | 26.9% | 18.5% | 38.2% | 41.8G | 11.9M |
| **Ours** | **84.7%** | 68.8% | 37.5% | **50.2%** | **27.4%** | **61.4%** | **55.0%** | 115.2G | 34.8M |

**TABLE IV:** Efficiency comparison between the standard models and our proposed efficient modifications. As an efficiency benchmark, we use the number of parameters (in million parameters) and the number of GFLOPs on the training dataset Cityscapes [15] and the evaluation dataset WildPPS. Images from both datasets are not resized and are fed at resolutions of $1024 \times 2048$ and $400 \times 2048$ for Cityscapes and WildPPS respectively.

| Target Model | Standard | Efficient | Savings |
|---|---|---|---|
| Parameter Comparison in Millions | | | |
| Seamless Segmentation [7] | 51.5 | 34.8 | 32.0% |
| Panoptic FCN [52] | 36.8 | 23.7 | 35.6% |
| Mask2Former [60] | 44.0 | 30.9 | 13.1% |
| FLOPs Comparison on Cityscapes in GFLOPs | | | |
| Seamless Segmentation [7] | 594.8 | 485.9 | 22.4% |
| Panoptic FCN [52] | 574.7 | 476.2 | 20.7% |
| Mask2Former [60] | 527.4 | 419.6 | 25.7% |
| FLOPs Comparison on WildPPS in GFLOPs | | | |
| Seamless Segmentation [7] | 282.3 | 239.5 | 17.9% |
| Panoptic FCN [52] | 233.5 | 193.5 | 20.7% |
| Mask2Former [60] | 215.1 | 171.3 | 25.6% |



**Fig. 8:** Difference in GFLOPs between the standard ResNet50 backbone [20] and the efficient ResNet18 backbone [20] with increasing image size. The x-axis denotes the scalar with which the height and the width of the original WildPPS [6] resolution are multiplied. On the very left, the difference is below 50 GFLOPs at around $400 \times 2048$ pixels and increases to over 300 GFLOPs for images three times the resolution of WildPPS.

by a large margin and there are many small models which perform excellent. Secondly, model adaptation techniques are considerably more important than pure sizes. This can be seen as the second best model was adapted according to the "Omni-Supervised" scheme [90] and is rather small in size. The two adapted models do perform considerably better in heavily-distorted classes such as *sidewalk* or *curb*, proving the effectiveness of adaptation techniques for the panoramic domain shift. Our proposed model combines the advantages of a light-weight distortion-aware model with the performance of the latest object detection models such as Mask R-CNN [18] which explains the good performance on all of the classes.

*E. Efficiency Analysis*

**Computation savings.** Regarding the targeted mobile agent application, it is essential to work with hardware- and energy-efficient models. Our proposed efficient model modification approach reduces both, the number of parameters and the required FLOPs, which serve as a proxy for hardware- and energy efficiency respectively, while it still remains a model-agnostic drop-in method which does not require any architectural changes. This allows a straightforward adaption of arbitrary target models, such that even future models should be easily adaptable and deployable. In Table IV, we compare the standard models as proposed by the respective authors with our efficient modifications. Depending on the setting and the target model, we can reduce the number of parameters up to $36\%$ and decrease the required FLOPs to calculate the panoptic predictions up to $26\%$.

As the required FLOPs depend on the size of the input image which can only be expected to increase in the future, we want to point out that the difference in FLOPs between the standard models and our efficient modifications will increase in line with the number of pixels. This relationship is visualized in Fig. 8, where we show the difference in GFLOPs between a ResNet50 [20] and a ResNet18 [20] depending on the number of input pixels. We start from the original WildPPS resolution of $400 \times 2048$ pixels and increase the number of pixels threefold to $1200 \times 6144$ pixels while keeping the typical panoramic aspect ratio constant. This is achieved by multiplying the height and width with the same scalar. This scalar is shown on the x-axis of the graph. Despite the fact that there is more to the difference in FLOPs than the difference in the backbones of the models, this is a realistic proxy, indicating that our efficient modifications become even more important when feeding images of higher resolutions into the model.

**Inference speed.** Besides hardware and energy efficiency, inference speed is a further necessary criterion in order to assess the suitability of an image parsing system for mobile agents. This section addresses this issue by measuring the necessary
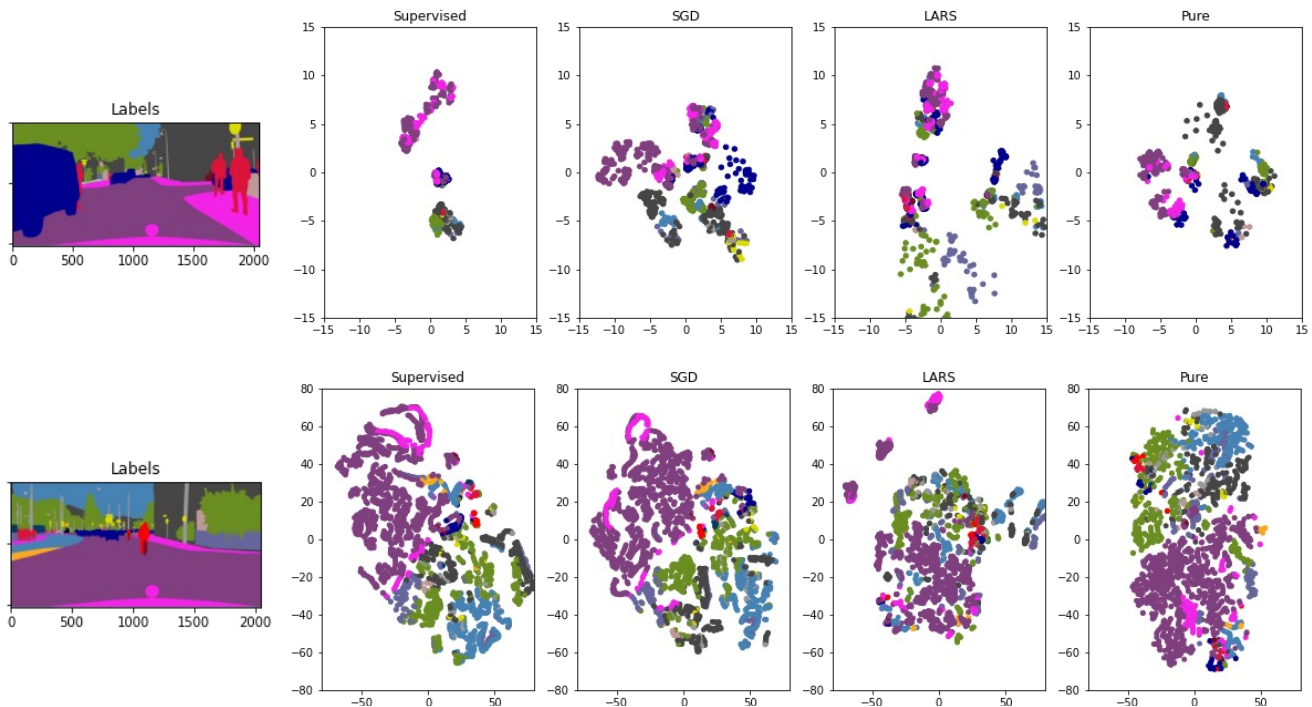
**Fig. 9:** Comparison of the latent space representations of the standard supervised backbone as shown in the left column and the different forms of backbone adaptions in the following columns. SGD, LARS, and Pure denote the mixed trainings according to the SGD and LARS setting as described in Section IV-B and the purely contrastive backbone respectively. In the first row we use the novel h-NNE visualization technique [98], whereas in the bottom row we rely on the established t-SNE dimensionality reduction [99]. For a better comparison among the different adaptation techniques, we set the axis to a reasonable interval. Best viewed on a screen and with zoom.

time to calculate a panoptic street scene prediction for a given single image at full resolution. The reported results, as shown in Table V, assume the model has been built and is located on the GPU. The reported times include the necessary time to ship a given image to the GPU and the execution time to generate the panoptic result including possible necessary post-processing steps such as the heuristic merging procedure used by the Seamless Scene Segmentation model. The experiments have been conducted on a single GeForce RTX 2080Ti GPU and the numbers denote the average inference time of the entire validation dataset.

**TABLE V:** Inference speed of the proposed efficient target models measured as the average Frames Per Second (FPS) on the two relevant image datasets using a single GeForce RTX 2080Ti GPU.

| Model | FPS on Cityscapes | FPS on WildPPS |
|---|---|---|
| Seamless Segmentation [7] | 4.4 | 8.0 |
| Panoptic FCN [52] | 5.5 | 9.1 |
| Mask2Former [60] | 4.5 | 9.0 |

Comparing the inference time between Cityscapes and WildPPS, the Cityscapes images are slower as they contain about $2.5$ times as many pixels. Regarding the inference times on WildPPS, which reflects the desired scenario of a mobile agent operating on panoramic images, our efficient model modifications are able to process between $8$ to $9$ Frames Per Second (FPS). This is below the common 24 FPS threshold, thus the models only offer near real-time performance. There are several possibilities to address this issue. Besides reducing the input resolution which would increase the inference speed but likely decrease the performance, a further step towards real-time processing would be in line with our experimental setup in which we used two GPUs. Using the training setup

for inference could double the FPS by asynchronous image processing. In a two GPU scenario, we are confident that with minor inference modifications or by using slightly more powerful GPUs, a 24 FPS real-time performance is possible for all the proposed models.

### F. Qualitative Analysis

**Feature space analysis.** As stated in Section III, an intuitive explanation regarding the functionality of the Panoramic Robust Feature (PRF) framework is the adaptation and the restructuring of the latent space obtained by the supervised ImageNet classification tasks. In this section, we visualize the introduced adaptations by comparing the backbones used in the baseline model to the robust backbones and the purely contrastive backbone. We use the output of the last pooling layer of the ResNet18 backbone [20], which result in a $32 \times 64$ pixel resolution feature map with $512$ filters which we map to a 2-dimensional visualization space using the novel h-NNE [98] and the well-established t-SNE [99] reduction technique.

The results are shown in Fig. 9 and support the stated intuition. Regarding the h-NNE visualization plots which are shown in the first row, the different adaptation techniques pull apart the features of different labels while features representing the same labels still remain close to each other. Comparing the LARS setting to the SGD setting, we observe that the label separation is performed more aggressively by the LARS optimizer. The same effect can also be seen in the t-SNE plot in the second row. We hypothesize that this larger deformation of the latent space has the potential to improve the results even further compared to the SGD setting, if the model can
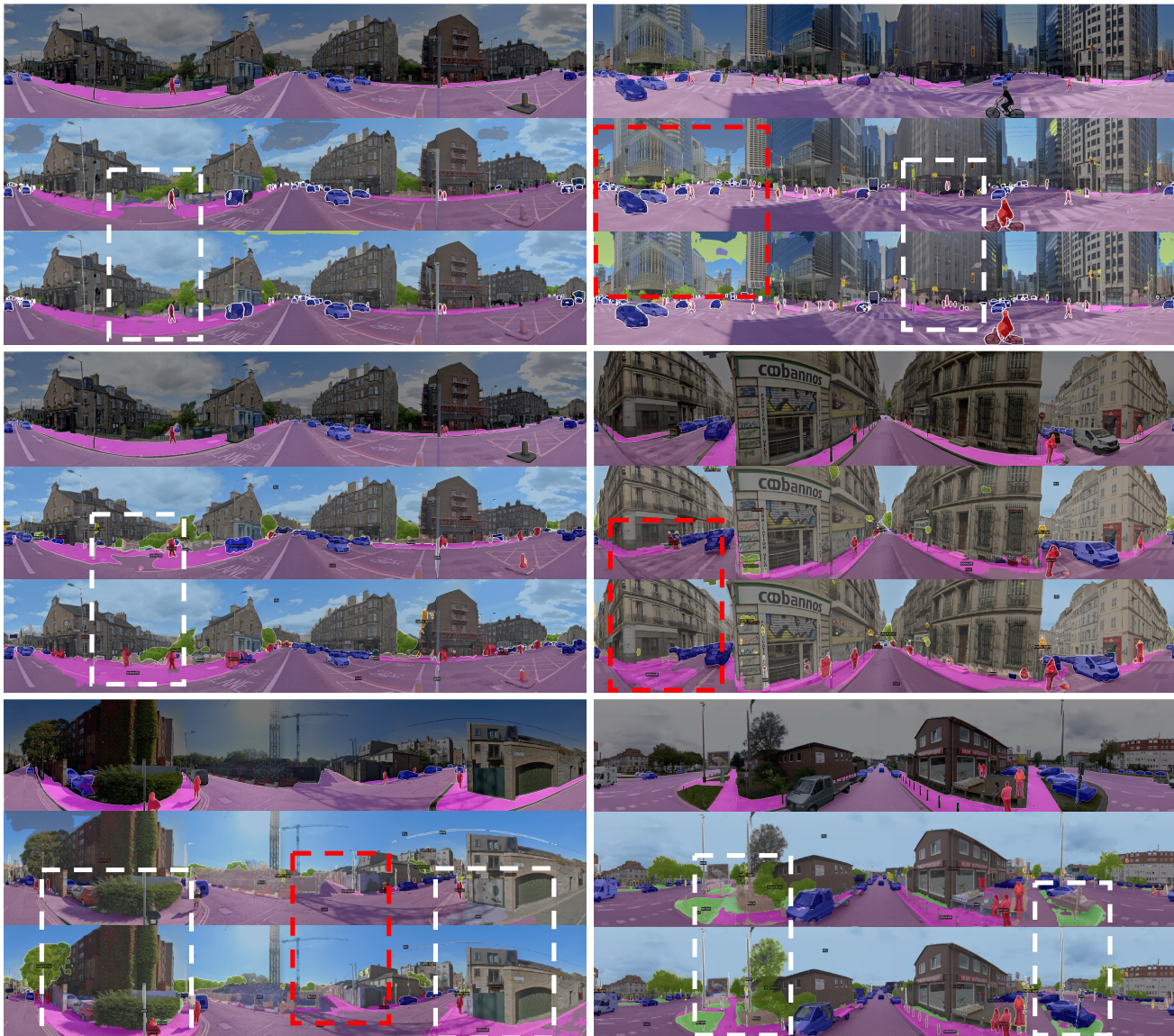
**Fig. 10:** Qualitative comparison of the baseline models with their robust counterparts. In each of the image triplets, the first image shows the manually annotated ground truth, followed by the baseline model, and finally the robust model. The first, second, and third row show examples of the Seamless Scene Segmentation [7], the Panoptic FCN [52], and the Mask2Former [60] model, respectively. Triplets in the first column show results adapted according to the SGD setting, whereas images in the second column show results of the target models adapted via the LARS setting. White boxes show improvements of the robust models. In red we highlight undesired behaviors. Best viewed on a screen and with zoom.

cope with these changes. For the chosen target models, the Mask2Former model seems to be able to capitalize on this, whereas the Panoptic FCN model is not.

The t-SNE visualization technique as shown in the second row provides more insights into the structure within the different labels. Comparing the supervised to the SGD setting, the latent space structure is very similar and this setting only slightly changes the feature encoder.

It is very plausible that the supervised training procedures of the used target models are heavily optimized to maximize the performance on the commonly used supervised backbone features. We observe that the SGD setting being the closest to the original latent space is a reliable choice that outperforms the baseline consistently as shown in Table I. The LARS setting on the other hand has the potential to outperform the

SGD backbone if the model is able to cope with this setting or if the training procedure is fine-tuned to this backbone.

**Segmentation map analysis.** Finally, we compare the segmentation maps of the baseline models to those of the robust modifications, in order to qualitatively assess the effect of the Panoramic Robust Feature framework. The segmentation maps are shown in Fig. 10. It can be easily seen that PRF enables both CNN- and transformer-based panoptic models to produce complete and more reliable panoramic surrounding understandings. The most present and apparent effect is that the robust features help in mitigating positional priors, which are most common for the *sidewalk* class that is typically only present at very similar locations during the supervised training phase. The robust features help to overcome these positional priors and help to improve the sidewalk segmen-

tation. Sometimes this behavior however does go beyond the desired effect, marking zebra crossings or unusual roads as sidewalks. From an application point of view, this behavior, although not desirable, poses less dangers to the environment compared to not being able to capture a sidewalk at all. Some of these behaviors remain unpunished at the moment such as the reported behavior in the second column of the first row, leaving room for further improvements.

## V. CONCLUSION

In this work, we have introduced the novel task of panoramic panoptic image segmentation which is the most holistic scene understanding task based on standard camera input. Only by combining a panoptic image-level understanding with a panoramic field of view, a mobile agent operating in a real-world scenario such as traffic scenes is able to make informed decisions. Semantic Segmentation maps are well suited to identify traversable areas such as roads, whereas instance masks can capture the different traffic participants, hence panoptic image segmentation perfectly combines the advantages of both approaches. We examine different architectures of Panoptic Segmentation models in order to identify the most promising types for the proposed task. As these models are complex and require a substantial amount of labeled data to be trained, we train them on publicly available large-scale image datasets such as Cityscapes. We address the difficult distribution change from pinhole- to panoramic images which is mitigated by the proposed Panoramic Robust Feature framework that leverages a pixel-level contrastive pretext task. The framework allows the training of robust models in a cost-minimizing way, suitable for the targeted mobile agent application. We confirm the effectiveness of the proposed approach by evaluating the target models on the WildPPS dataset, the first panoramic panoptic dataset which will be published in order to foster progress in panoramic panoptic scene parsing. The Seamless Scene Segmentation model proves to be a reliable choice for this new task achieving state-of-the-art performance on WildPPS and the PASS dataset. Finally, we confirm the efficiency of both the proposed panoramic models and the adaptation framework in the face of hardware and energy scarcity to meet the mobile agent setting.

Since the proposed method does not only generalize well to panoramic images but also generates improvements when train and test dataset differ as shown in Table II, we believe that the proposed PRF framework has the capabilities to produce robust models in particular suitable for domain shift problems. Due to the desired real world applications of mobile agents, we think our approach has the potential to perform well in other mobile agent vision related sensing methods such as LiDAR or RADAR and their corresponding tasks such as segmentation of point clouds. As the pinhole to panoramic domain shift is a known problem for the aforementioned sensing modalities, we believe that robust PRF models could help in overcoming these domain shift as well.

We hope that this work sparks interest in the community regarding Panoramic Panoptic Segmentation. In the future, we want to examine different distribution shifts not only

from pinhole- to panoramic images but also from artificially generated panoramic images to real-world panoramic images. Furthermore, we are curious about the effects of the framework when including more annotated classes as some of the undesired behaviors are not yet punished due to the lack of labels for the classes. This may open up additional research questions as different classes are affected differently by the panoramic distortions and may require special treatment.

## REFERENCES

[1] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, "Panoptic segmentation," in *Proc. CVPR*, 2019, pp. 9396–9405.
[2] K. Yang, X. Hu, L. M. Bergasa, E. Romera, and K. Wang, "PASS: Panoramic annular semantic segmentation," *T-ITS*, vol. 21, no. 10, pp. 4171–4185, 2020.
[3] J. Zhang, K. Yang, C. Ma, S. Reiß, K. Peng, and R. Stiefelhagen, "Bending reality: Distortion-aware transformers for adapting to panoramic semantic segmentation," in *Proc. CVPR*, 2022, pp. 16 917–16 927.
[4] S. Peng, W. Jiang, H. Pi, X. Li, H. Bao, and X. Zhou, "Deep snake for real-time instance segmentation," in *Proc. CVPR*, 2020, pp. 8533–8542.
[5] J. Zhang, K. Yang, and R. Stiefelhagen, "Exploring event-driven dynamic context for accident scene segmentation," *T-ITS*, vol. 23, no. 3, pp. 2606–2622, 2022.
[6] A. Jaus, K. Yang, and R. Stiefelhagen, "Panoramic panoptic segmentation: Towards complete surrounding understanding via unsupervised contrastive learning," *Proc. IV*, pp. 1421–1427, 2021.
[7] L. Porzi, S. R. Bulò, A. Colovic, and P. Kontschieder, "Seamless scene segmentation," in *Proc. CVPR*, 2019, pp. 8269–8278.
[8] J. S. Berrio, M. Shan, S. Worrall, and E. Nebot, "Camera-LIDAR integration: Probabilistic sensor fusion for semantic mapping," *T-ITS*, vol. 23, no. 7, pp. 7637–7652, 2022.
[9] K. Peng *et al.*, "MASS: Multi-attentional semantic segmentation of LiDAR data for dense top-view understanding," *T-ITS*, vol. 23, no. 9, pp. 15 824–15 840, 2022.
[10] P. Testolina, F. Barbato, U. Michieli, M. Giordani, P. Zanuttigh, and M. Zorzi, "SELMA: Semantic large-scale multimodal acquisitions in variable weather, daytime and viewpoints," *arXiv preprint arXiv:2204.09788*, 2022.
[11] A. Petrovai and S. Nedevschi, "Semantic cameras for 360-degree environment perception in automated urban driving," *T-ITS*, vol. 23, no. 10, pp. 17 271–17 283, 2022.
[12] W. Ye *et al.*, "PVO: Panoptic visual odometry," *arXiv preprint arXiv:2207.01610*, 2022.
[13] N. Gosala and A. Valada, "Bird's-eye-view panoptic segmentation using monocular frontal view images," *RA-L*, vol. 7, no. 2, pp. 1968–1975, 2022.
[14] S. Gao, K. Yang, H. Shi, K. Wang, and J. Bai, "Review on panoramic imaging and its applications in scene understanding," *TIM*, vol. 71, pp. 1–34, 2022.
[15] M. Cordts *et al.*, "The cityscapes dataset for semantic urban scene understanding," in *Proc. CVPR*, 2016, pp. 3213–3223.
[16] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. CVPR*, 2020, pp. 9726–9735.
[17] X. Wang, R. Zhang, C. Shen, T. Kong, and L. Li, "Dense contrastive learning for self-supervised visual pre-training," in *Proc. CVPR*, 2021, pp. 3024–3033.
[18] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proc. ICCV*, 2017, pp. 2980–2988.
[19] A. Vaswani *et al.*, "Attention is all you need," in *NeurIPS*, vol. 30, 2017, pp. 5998–6008.
[20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770–778.
[21] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. CVPR*, 2015, pp. 3431–3440.
[22] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *TPAMI*, vol. 39, no. 12, pp. 2481–2495, 2017.
[23] 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," *TPAMI*, vol. 40, no. 4, pp. 834–848, 2018.

[24] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proc. CVPR*, 2017, pp. 6230–6239.

[25] T. Takikawa, D. Acuna, V. Jampani, and S. Fidler, "Gated-SCNN: Gated shape CNNs for semantic segmentation," in *Proc. ICCV*, 2019, pp. 5229–5238.

[26] J. Fu *et al.*, "Dual attention network for scene segmentation," in *Proc. CVPR*, 2019, pp. 3141–3149.

[27] K. Yang, J. Zhang, S. Reiß, X. Hu, and R. Stiefelhagen, "Capturing omni-range context for omnidirectional segmentation," in *Proc. CVPR*, 2021, pp. 1376–1386.

[28] S. Zheng *et al.*, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Proc. CVPR*, 2021, pp. 6881–6890.

[29] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "SegFormer: Simple and efficient design for semantic segmentation with transformers," in *NeurIPS*, vol. 34, 2021, pp. 12 077–12 090.

[30] J. Zhang, K. Yang, A. Constantinescu, K. Peng, K. Müller, and R. Stiefelhagen, "Trans4Trans: Efficient transformer for transparent object and semantic scene segmentation in real-world navigation assistance," *T-ITS*, vol. 23, no. 10, pp. 19 173–19 186, 2022.

[31] D. Liang, B. Kang, X. Liu, P. Gao, X. Tan, and S. Kaneko, "Cross-scene foreground segmentation with supervised and unsupervised model communication," *PR*, vol. 117, p. 107995, 2021.

[32] D. Liang, Y. Du, H. Sun, L. Zhang, N. Liu, and M. Wei, "NLKD: Using coarse annotations for semantic segmentation based on knowledge distillation," in *Proc. ICASSP*, 2021, pp. 2335–2339.

[33] R. Liu, K. Yang, H. Liu, J. Zhang, K. Peng, and R. Stiefelhagen, "Transformer-based knowledge distillation for efficient semantic segmentation of road-driving scenes," *arXiv preprint arXiv:2202.13393*, 2022.

[34] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proc. CVPR*, 2018, pp. 8759–8768.

[35] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "YOLACT: Real-time instance segmentation," in *Proc. ICCV*, 2019, pp. 9157–9166.

[36] X. Wang, T. Kong, C. Shen, Y. Jiang, and L. Li, "SOLO: Segmenting objects by locations," in *Proc. ECCV*, 2020, pp. 649–665.

[37] Z. Tian, C. Shen, and H. Chen, "Conditional convolutions for instance segmentation," in *Proc. ECCV*, 2020, pp. 282–298.

[38] A. Arnab and P. H. S. Torr, "Pixelwise instance segmentation with a dynamically instantiated network," in *Proc. CVPR*, 2017, pp. 441–450.

[39] M. Bai and R. Urtasun, "Deep watershed transform for instance segmentation," in *Proc. CVPR*, 2017, pp. 5221–5229.

[40] P. H. O. Pinheiro, R. Collobert, and P. Dollár, "Learning to segment object candidates," in *Proc. NeurIPS*, vol. 28, 2015, pp. 1990–1998.

[41] J. Dai, K. He, Y. Li, S. Ren, and J. Sun, "Instance-sensitive fully convolutional networks," in *Proc. ECCV*, 2016, pp. 534–549.

[42] E. Xie *et al.*, "PolarMask: Single shot instance segmentation with polar representation," in *Proc. CVPR*, 2020, pp. 12 193–12 202.

[43] Y. Xiong *et al.*, "UPSNet: A unified panoptic segmentation network," in *Proc. CVPR*, 2019, pp. 8818–8826.

[44] Q. Li, X. Qi, and P. H. Torr, "Unifying training and inference for panoptic segmentation," in *Proc. CVPR*, 2020, pp. 13 320–13 328.

[45] A. Kirillov, R. Girshick, K. He, and P. Dollár, "Panoptic feature pyramid networks," in *Proc. CVPR*, 2019, pp. 6392–6401.

[46] H. Liu *et al.*, "An end-to-end network for panoptic segmentation," in *Proc. CVPR*, 2019, pp. 6172–6181.

[47] Y. Li *et al.*, "Attention-guided unified network for panoptic segmentation," in *Proc. CVPR*, 2019, pp. 7026–7035.

[48] B. Cheng *et al.*, "Panoptic-DeepLab: A simple, strong, and fast baseline for bottom-up panoptic segmentation," in *Proc. CVPR*, 2020, pp. 12 472–12 482.

[49] S. Borse, H. Park, H. Cai, D. Das, R. Garrepalli, and F. Porikli, "Panoptic, instance and semantic relations: A relational context encoder to enhance panoptic segmentation," in *Proc. CVPR*, 2022, pp. 1269–1279.

[50] S. Hwang, S. W. Oh, and S. J. Kim, "Single-shot path integrated panoptic segmentation," in *Proc. WACV*, 2022, pp. 3328–3337.

[51] Z. Tian, B. Zhang, H. Chen, and C. Shen, "Instance and panoptic segmentation using conditional convolutions," *TPAMI*, vol. 45, no. 1, pp. 669–680, 2023.

[52] Y. Li *et al.*, "Fully convolutional networks for panoptic segmentation," in *Proc. CVPR*, 2021, pp. 214–223.

[53] W. Zhang, J. Pang, K. Chen, and C. C. Loy, "K-net: Towards unified image segmentation," in *Proc. NeurIPS*, vol. 34, 2021, pp. 10 326–10 338.

[54] T. Kerola, J. Li, A. Kanehira, Y. Kudo, A. Vallet, and A. Gaidon, "Hierarchical lovász embeddings for proposal-free panoptic segmentation," in *Proc. CVPR*, 2021, pp. 14 413–14 423.

[55] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Proc. ECCV*, 2020, pp. 213–229.

[56] Q. Yu *et al.*, "CMT-DeepLab: Clustering mask transformers for panoptic segmentation," in *Proc. CVPR*, 2022, pp. 2560–2570.

[57] H. Wang, Y. Zhu, H. Adam, A. Yuille, and L.-C. Chen, "MaX-DeepLab: End-to-end panoptic segmentation with mask transformers," in *Proc. CVPR*, 2021, pp. 5463–5474.

[58] B. Cheng, A. Schwing, and A. Kirillov, "Per-pixel classification is not all you need for semantic segmentation," in *Proc. NeurIPS*, vol. 34, 2021.

[59] H. Wang, Y. Zhu, B. Green, H. Adam, A. Yuille, and L.-C. Chen, "Axial-DeepLab: Stand-alone axial-attention for panoptic segmentation," in *Proc. ECCV*, 2020, pp. 108–126.

[60] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, "Masked-attention mask transformer for universal image segmentation," in *Proc. CVPR*, 2022, pp. 1290–1299.

[61] Z. Li *et al.*, "Panoptic SegFormer: Delving deeper into panoptic segmentation with transformers," in *Proc. CVPR*, 2022, pp. 1280–1289.

[62] G. Neuhold, T. Ollmann, S. R. Bulò, and P. Kontschieder, "The mapillary vistas dataset for semantic understanding of street scenes," in *Proc. ICCV*, 2017, pp. 5000–5009.

[63] L. Deng, M. Yang, H. Li, T. Li, B. Hu, and C. Wang, "Restricted deformable convolution-based road scene semantic segmentation using surround view cameras," *T-ITS*, vol. 21, no. 10, pp. 4350–4362, 2020.

[64] S. K. Yogamani *et al.*, "WoodScape: A multi-task, multi-camera fisheye dataset for autonomous driving," in *Proc. ICCV*, 2019, pp. 9307–9317.

[65] C. Eising, J. Horgan, and S. Yogamani, "Near-field perception for low-speed vehicle automation using surround-view fisheye cameras," *T-ITS*, 2021.

[66] V. R. Kumar *et al.*, "OmniDet: Surround view cameras based multi-task visual perception network for autonomous driving," *RA-L*, vol. 6, no. 2, pp. 2830–2837, 2021.

[67] Y. Liao, J. Xie, and A. Geiger, "KITTI-360: A novel dataset and benchmarks for urban scene understanding in 2D and 3D," *TPAMI*, 2022.

[68] C. Zhang, S. Liwicki, W. Smith, and R. Cipolla, "Orientation-aware semantic segmentation on icosahedron spheres," in *Proc. ICCV*, 2019, pp. 3532–3540.

[69] A. R. Sekkat, Y. Dupuis, P. Vasseur, and P. Honeine, "The OmniScape dataset," in *Proc. ICRA*, 2020, pp. 1603–1608.

[70] S. Orhan and Y. Bastanlar, "Semantic segmentation of outdoor panoramic images," *SIVP*, vol. 16, no. 3, pp. 643–650, 2022.

[71] K. Yang, X. Hu, H. Chen, K. Xiang, K. Wang, and R. Stiefelhagen, "DS-PASS: Detail-sensitive panoramic annular semantic segmentation through SwaftNet for surrounding sensing," in *Proc. IV*, 2020, pp. 457–464.

[72] J. Zhang, C. Ma, K. Yang, A. Roitberg, K. Peng, and R. Stiefelhagen, "Transfer beyond the field of view: Dense panoramic semantic segmentation via unsupervised domain adaptation," *T-ITS*, vol. 23, no. 7, pp. 9478–9491, 2022.

[73] X. Hu, Y. An, C. Shao, and H. Hu, "Distortion convolution module for semantic segmentation of panoramic images based on the image-forming principle," *TIM*, vol. 71, pp. 1–12, 2022.

[74] J. Mei *et al.*, "Waymo open dataset: Panoramic video panoptic segmentation," in *Proc. ECCV*, 2022.

[75] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *Proc. CVPR*, 2018, pp. 3733–3742.

[76] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proc. ICML*, 2020, pp. 1597–1607.

[77] Z. Xie, Y. Lin, Z. Zhang, Y. Cao, S. Lin, and H. Hu, "Propagate yourself: Exploring pixel-level consistency for unsupervised visual representation learning," in *Proc. CVPR*, 2021, pp. 16 684–16 693.

[78] R. Cheke, G. Sistu, C. Eising, P. van de Ven, V. R. Kumar, and S. Yogamani, "FisheyePixPro: Self-supervised pretraining using fisheye images for semantic segmentation," in *Proc. AVM*, 2022.

[79] W. Wang, T. Zhou, F. Yu, J. Dai, E. Konukoglu, and L. Van Gool, "Exploring cross-image pixel contrast for semantic segmentation," in *Proc. ICCV*, 2021, pp. 7303–7313.

[80] H. Hu, J. Cui, and L. Wang, "Region-aware contrastive learning for semantic segmentation," in *Proc. ICCV*, 2021, pp. 16 291–16 301.

[81] F. Zhang, P. Torr, R. Ranftl, and S. Richter, "Looking beyond single images for contrastive semantic segmentation learning," in *Proc. NeurIPS*, vol. 34, 2021, pp. 3285–3297.

[82] T. Xiao, C. J. Reed, X. Wang, K. Keutzer, and T. Darrell, "Region similarity representation learning," in *Proc. ICCV*, 2021, pp. 10 539–10 548.

[83] S. Liu, S. Zhi, E. Johns, and A. J. Davison, "Bootstrapping semantic segmentation with regional contrast," in *Proc. ICLR*, 2022.

[84] Y. Zhao, G. Wang, C. Luo, W. Zeng, and Z.-J. Zha, "Self-supervised visual representations learning by contrastive mask prediction," in *Proc. ICCV*, 2021, pp. 10 160–10 169.

[85] J. Xie, X. Zhan, Z. Liu, Y. Ong, and C. C. Loy, "Unsupervised object-level representation learning from scene images," in *Proc. NeurIPS*, vol. 34, 2021.

[86] X. Li *et al.*, "Dense semantic contrast for self-supervised visual representation learning," in *Proc. MM*, 2021, pp. 1368–1376.

[87] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. CVPR*, 2009, pp. 248–255.

[88] J.-B. Grill *et al.*, "Bootstrap your own latent-A new approach to self-supervised learning," in *NeurIPS*, vol. 33, 2020, pp. 21 271–21 284.

[89] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017.

[90] K. Yang, X. Hu, and R. Stiefelhagen, "Is context-aware CNN ready for the surroundings? Panoramic semantic segmentation in the wild," *TIP*, vol. 30, pp. 1866–1881, 2021.

[91] Y. You, I. Gitman, and B. Ginsburg, "Large batch training of convolutional networks," *arXiv preprint arXiv:1708.03888*, 2017.

[92] M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang, "DenseASPP for semantic segmentation in street scenes," in *Proc. CVPR*, 2018, pp. 3684–3692.

[93] 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.

[94] T. Wu, S. Tang, R. Zhang, J. Cao, and Y. Zhang, "CGNet: A light-weight context guided network for semantic segmentation," *TIP*, vol. 30, pp. 1169–1179, 2021.

[95] E. Romera, J. M. Alvarez, L. M. Bergasa, and R. Arroyo, "ERFNet: Efficient residual factorized ConvNet for real-time semantic segmentation," *T-ITS*, vol. 19, no. 1, pp. 263–272, 2018.

[96] M. Orsic, I. Kreso, P. Bevandic, and S. Segvic, "In defense of pre-trained ImageNet architectures for real-time semantic segmentation of road-driving images," in *Proc. CVPR*, 2019, pp. 12 607–12 616.

[97] K. Yang, X. Hu, Y. Fang, K. Wang, and R. Stiefelhagen, "Omnisupervised omnidirectional semantic segmentation," *T-ITS*, vol. 23, no. 2, pp. 1184–1199, 2022.

[98] M. S. Sarfraz, M. Koulakis, C. Seibold, and R. Stiefelhagen, "Hierarchical nearest neighbor graph embedding for efficient dimensionality reduction," in *Proc. CVPR*, 2022, pp. 336–345.

[99] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE," *JMLR*, vol. 9, no. 11, 2008.