# Bird's-Eye-View Panoptic Segmentation Using Monocular Frontal View Images

Nikhil Gosala and Abhinav Valada

*Abstract*—**Bird's-Eye-View (BEV) maps have emerged as one of the most powerful representations for scene understanding due to their ability to provide rich spatial context while being easy to interpret and process. Such maps have found use in many real-world tasks that extensively rely on accurate scene segmentation as well as object instance identification in the BEV space for their operation. However, existing segmentation algorithms only predict the semantics in the BEV space, which limits their use in applications where the notion of object instances is also critical. In this work, we present the first BEV panoptic segmentation approach for directly predicting dense panoptic segmentation maps in the BEV, given a single monocular image in the frontal view (FV). Our architecture follows the top-down paradigm and incorporates a novel dense transformer module consisting of two distinct transformers that learn to independently map vertical and flat regions in the input image from the FV to the BEV. Additionally, we derive a mathematical formulation for the sensitivity of the FV-BEV transformation which allows us to intelligently weight pixels in the BEV space to account for the varying descriptiveness across the FV image. Extensive evaluations on the KITTI-360 and nuScenes datasets demonstrate that our approach exceeds the state-of-the-art in the PQ metric by 3.61 pp and 4.93 pp respectively.**

## I. INTRODUCTION

Autonomous vehicles require rich, detailed, and comprehensive understanding of their surroundings for carrying out essential tasks such as collision avoidance and object tracking [1]. Bird's-Eye-View (BEV) maps [2], [3], [4] have gained immense popularity in recent years due to their information-rich and easily interpretable representation of the world. They also capture absolute distances in the metric scale which allow them to be readily deployed in applications such as motion planning and behavior prediction [5]. Many real-world tasks such as path planning and trajectory estimation rely on an accurate semantic scene segmentation as well as object instance identification in the BEV space for their effective operation. However, existing BEV map generation approaches only incorporate semantic information in the BEV maps, which limits their use in many real-world applications that require knowledge of object instances.

In this work, we aim to overcome this limitation by proposing the first BEV panoptic segmentation approach that generates coherent panoptic predictions in the BEV using monocular FV images (Fig. 1). Panoptic segmentation allows for the simultaneous estimation and fusion of both semantic and instance predictions, which enables complete and coherent scene understanding [6]. Existing methods generate semantic BEV maps from monocular images by either (i) using trivial homography such as IPM [7], (ii) unprojecting the 2D image using predicted depth [4], or (iii) using dense transformers to learn the mapping from FV to BEV [3], [8]. The flat-world assumption in

All authors are with the Department of Computer Science, University of Freiburg, Germany.
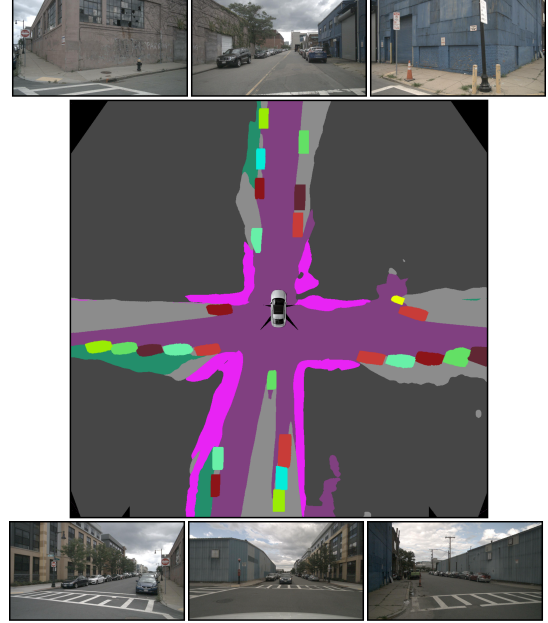


Fig. 1: Full 360° BEV panoptic segmentation output obtained from our PanopticBEV model. Given a monocular image in the frontal view, PanopticBEV directly predicts the panoptic segmentation in the BEV, consisting of both semantic *stuff* classes (road, sidewalk, etc.) and instance-specific *thing* classes (cars, pedestrians, etc.).

IPM-based approaches hinders their performance in regions that lie above the ground plane, while depth unprojection-based methods rely on multi-stage pipelines that fail to reap the benefits of end-to-end learning. In contrast, dense transformer-based approaches have shown immense potential due to their ability to model the complex mapping from FV to BEV without any additional supervision [3], [8]. However, existing methods do not account for the different transformation characteristics of the vertical and flat regions, and thus employ a single transformer across the entire image. This forces these models to learn a common mapping for all the different regions in the image which leads to imprecise BEV predictions.

To address this problem, we propose a dense transformer module that incorporates two distinct transformers to independently map the vertical and flat regions in the input FV image to the BEV. Our proposed PanopticBEV architecture follows the top-down paradigm comprising a modified EfficientDet [9] backbone, the novel dense transformer module, a semantic head, an instance head, and an adaptive panoptic fusion module. We observe that the perspective projection makes mapping far-away objects from FV to BEV extremely challenging. This can be attributed to the fact that a given displacement of far-away objects in the 3D space, maps to a comparatively smaller displacement in their 2D position. To alleviate this problem, we derive a mathematical formulation to quantify the sensitivity of the FV-BEV transformation and employ it to normalize the descriptiveness

across the input image. Moreover, as our proposed approach is the first to tackle the problem of BEV panoptic segmentation, we introduce multiple baselines to facilitate quantitative comparisons. We develop the baselines by combining existing BEV semantic segmentation models with the instance head and fusion modules from state-of-the-art FV panoptic segmentation methods. We perform extensive evaluations of our approach on the KITTI-360 [10] and nuScenes [11] datasets, and demonstrate that it substantially outperforms the state-of-the-art.

Our main contributions can thus be summarized as follows:
1) An end-to-end learning architecture for BEV panoptic segmentation from monocular FV images.
2) A dense spatial transformer module comprising two distinct transformers that independently learn to map vertical and flat regions in the input FV image to the BEV.
3) A mathematical formulation of the FV-BEV transformation sensitivity which we exploit for weighting pixels in the BEV space during the training phase.
4) Several competitive baselines for the novel task of BEV panoptic segmentation.
5) A data processing pipeline to generate panoptic BEV groundtruth labels from annotated LiDAR point clouds.
6) Extensive evaluations along with ablation studies on two standard real-world datasets, KITTI-360 and nuScenes.
7) Publicly available code and pre-trained models at http://rl.uni-freiburg.de/research/panoptic-bev.

## II. RELATED WORK

**FV-BEV Transformation**: Numerous works have been proposed to address the challenging task of estimating the BEV map using monocular images. One common approach is to use Inverse Perspective Mapping (IPM) [12], or variants of it, to project the FV image onto the ground plane using a homography [13], [14], [15]. Several authors address this task as a generative problem and advocate the use of GANs [14], [15], [16]. Other works implicitly transform FV images into the BEV for perception tasks such as 3D object detection [17] and vehicle extent estimation [18].

Very few works, however, address the more specific task of generating BEV segmentation maps using monocular FV images. These works can be broadly classified into two categories: *geometry-agnostic* and *geometry-aware*, based on whether they account for the geometry of the scene while transforming the input FV image into the BEV space. *Geometry-agnostic* approaches do not utilize the scene geometry and fully rely on the representational capacity of the network to learn the transformation. VED [19] and VPN [8] fall under this category of approaches. The former employs a variational encoder-decoder architecture with a fully-connected bottleneck layer, while the latter uses a two-layer multi-layer perceptron to transform the FV features into the BEV space. Discarding scene geometry forces the network to approximate it which makes the output coarser and less accurate. *Geometry-aware* approaches, on the contrary, either exploit the scene geometry explicitly or capture it in the network design implicitly. Cam2BEV [7] and DSM [20] explicitly capture the scene geometry by incorporating IPM into their transformers. However, the use of IPM is limited to pixels on the assumed ground plane and fails for pixels that lie above it, such as those belonging to

buildings and vehicles. Other works [4], [21] incorporate scene geometry by unprojecting 2D color pixels into the 3D space using the estimated monocular depth and then converting them into BEV maps. Nevertheless, multi-stage approaches prevent end-to-end learning which results in sub-optimal BEV map predictions. PON [3] alleviates these problems by proposing an end-to-end approach to estimate the BEV semantic map from a monocular image. However, it does not account for the different transformation characteristics of the vertical and flat regions in the image which limits its performance on certain classes that are inadequately modeled by the transformer. Recently, LSS [2] proposes the estimation of a categorical depth distribution to unproject the FV features into a volumetric lattice, and transform it into the BEV frame. However, it struggles to generalize well across semantic categories resulting in poor segmentation performance for a large number of classes.

**Panoptic Segmentation**: Panoptic segmentation is the task of semantically distinguishing regions in the image at the pixel-level while simultaneously discerning between instances of an object. Existing approaches can be classified into two categories: *proposal-based* and *proposal-free*. *Proposal-based* approaches independently estimate the semantic and instance masks using two separate heads before fusing them to generate the panoptic segmentation output. These approaches typically suffer from the mask-overlapping problem wherein areas around *thing* classes become ambiguous due to the disagreement between the semantic and instance heads. This problem has been mitigated by either (i) weakly supervising *thing* and *stuff* classes using bounding boxes and image-level tags [22], (ii) explicitly constraining the *stuff* and *thing* distributions using a learned binary mask [23], or (iii) performing pixel-wise classification on the combined semantic and instance logits mask [24], [25]. *Proposal-free* approaches, in contrast, yield the panoptic segmentation output by first predicting the semantic label for each pixel, before detecting instances by clustering *thing* pixels together. Panoptic-DeepLab [26] couples bounding box corners and object centers while incorporating a dual-ASPP and dual-decoder structure into each sub-task branch. SSAP [27] proposes grouping pixels using an affinity pyramid with a graph partitioning strategy to detect instances while learning the semantic labels.

Through this work, we address two major limitations of existing approaches, i.e., (i) the inability of the existing transformers to account for the unique transformation characteristics of vertical and flat regions in FV images, and (ii) the lack of object instance information in semantic BEV maps, which hinders using existing methods in many real-world use-cases. To this end, we propose a novel dense transformer module that accounts for the distinct transformation characteristics of the vertical and flat regions in FV images, and present the first BEV panoptic segmentation approach.

## III. TECHNICAL APPROACH

In this section, we first provide a brief overview of the proposed PanopticBEV architecture illustrated in Fig. 2, before diving into the crux of each constituent component. Our network comprises a shared backbone, a dense transformer module, a semantic head, an instance head, and a panoptic fusion
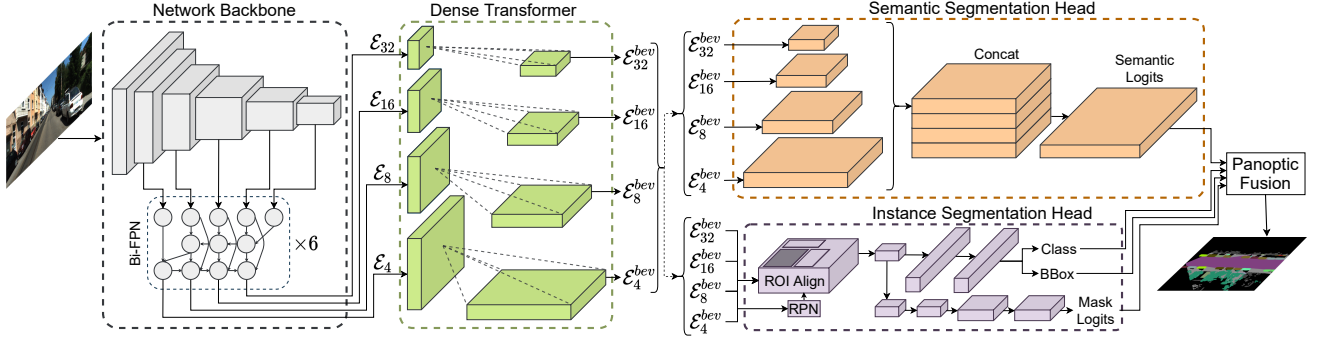
Fig. 2: Topology of the proposed PanopticBEV architecture consisting of a modified EfficientDet backbone (in gray) that generates four feature maps with strides $4, 8, 16,$ and $32$. The resulting multi-scale FV features are independently transformed into the BEV by our novel dense transformer module (in green). The transformed BEV features are then fed into the semantic (in orange) and instance (in violet) heads, followed by the adaptive panoptic fusion module that yields the output. In the figure, the *blocks* represent the shapes of the intermediate tensors obtained after performing a series of mathematical operations.

module. We employ a modified variant of the EfficientDet [9] topology for the backbone which outputs feature maps at four different scales $\mathcal{E}_4, \mathcal{E}_8, \mathcal{E}_{16}$ and $\mathcal{E}_{32}$. The feature maps are then input to the dense transformer module, which consists of two distinct transformers that independently transform the vertical and flat regions in the input FV image to the BEV. The dense transformer then combines the transformed vertical and flat feature maps to yield the corresponding composite BEV features $\mathcal{E}_4^{bev}, \mathcal{E}_8^{bev}, \mathcal{E}_{16}^{bev}$ and $\mathcal{E}_{32}^{bev}$. Subsequently, the transformed feature maps are fed into the semantic and instance heads in parallel, followed by the panoptic fusion module that generates the final BEV panoptic segmentation output.

### A. Network Backbone

The backbone of our network is built upon the EfficientDet [9] architecture which has shown tremendous potential on both segmentation and detection tasks while being computationally efficient. Specifically, we employ the EfficientDet-D3 topology in the PanopticBEV architecture to achieve the right trade-off between computational complexity and representational capacity. However, this can readily be replaced with any of the other EfficientDet variants according to the available computational budget. We adapt this backbone to output feature maps with strides 4-32 instead of the conventional 8-128 by replacing the input to the first BiFPN layer with feature maps of strides $4, 8, 16,$ and $32$ from the EfficientNet encoder. This enables the semantic head to use higher resolution features and consequently improves spatial scene understanding as well as reduces the depth ambiguity in the BEV space.

### B. Dense Transformer

Our proposed dense transformer is based on the principle of how different regions in the 3D world are projected onto a perspective 2D image. Specifically, a column belonging to a flat region in the FV image projects onto a perspectively-distorted area in the BEV, whereas a column belonging to a vertical non-flat region maps to an orthographic projection of a volumetric region in the BEV space. Fig. S.1 in the supplementary material illustrates this phenomenon. To tackle this problem, we employ two distinct transformers to independently map features from the vertical and flat regions in the FV to the BEV. Fig. 3 shows an overview of our dense transformer module. Each scale $k$ of the backbone features $\mathcal{E}$ is first fed to a semantic masking module $\mathcal{M}_k$ to generate the vertical and
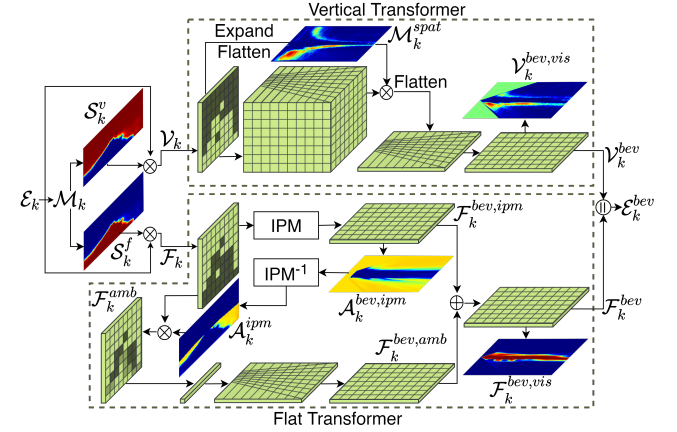


Fig. 3: Illustration of the dense transformer module consisting of a distinct vertical and flat transformer. The vertical transformer uses a volumetric lattice to model the intermediate 3D space which is flattened to generate the vertical BEV features, and the flat transformer uses IPM followed by our Error Correction Module (ECM) to generate the flat BEV features. The vertical and flat BEV features are then merged to generate the composite BEV feature map.

flat semantic masks $\mathcal{S}_k^v$ and $\mathcal{S}_k^f$ respectively. We then compute the Hadamard product between the semantic masks and the backbone features to yield the vertical and flat features $\mathcal{V}_k$ and $\mathcal{F}_k$. Subsequently, we independently transform $\mathcal{V}_k$ and $\mathcal{F}_k$ into their BEV representations $\mathcal{V}_k^{bev}$ and $\mathcal{F}_k^{bev}$ using their respective transformers. We then combine these features in the BEV space to generate the composite BEV feature map $\mathcal{E}_k^{bev}$. A more detailed architectural diagram is depicted in Fig. S.2 of the supplementary material.

*1) Vertical Transformer:* We model the vertical transformer to implicitly capture the intricate relationship between the FV and BEV for vertical regions. To this end, we first expand the 2D encoder features $\mathcal{E}_k$ of shape $C \times H_k \times W_k$ into a perspectively distorted 3D volumetric grid of size $Z_k \times C \times H_k \times W_k$ using 3D convolutional filters. Simultaneously, we generate a spatial occupancy mask $\mathcal{M}_k^{spat}$ from $\mathcal{E}_k$, which estimates the probability of a pixel being occupied by a vertical element in the BEV. We then multiply $\mathcal{M}_k^{spat}$ with the 3D volumetric grid to constrain the spatial extents of vertical regions in the 3D grid. We actively supervise $\mathcal{M}_k^{spat}$ using the BEV groundtruth to guide the transformer during the training phase. We then collapse the spatially-attended 3D grid along the height dimension to generate features of size $C \times Z \times W$ in the BEV space. Finally, we correct the perspective distortion in the BEV feature

map, carried from the perspective input image, by resampling the feature map using the known camera intrinsics as described in [3] to generate the final vertical BEV features $\mathcal{V}_k^{bev}$.

*2) Flat Transformer:* The IPM algorithm reinforced with a learnable error correction module (ECM) forms the basis of our flat transformer. The IPM algorithm estimates a homography matrix $M$ which when multiplied with the FV features transforms them into the BEV space. It is mathematically sound and parameter-free which allows it to be used in a wide range of scenarios. However, due to its flat-world assumption, it is inapplicable for feature points that lie above the defined ground plane. Since $\mathcal{F}_k$, by definition, is largely devoid of vertical elements, IPM provides a good basis to transform $\mathcal{F}_k$ into $\mathcal{F}_k^{bev,ipm}$. However, since the flat regions in the real-world are not perfectly flat, IPM introduces errors into $\mathcal{F}_k^{bev,ipm}$. We resolve these errors using a learnable ECM whose architecture is inspired by our earlier observation.

To this end, we first estimate regions in FV where the IPM transformation is ambiguous, and then minimize the ambiguity by focusing the ECM on these regions. We estimate the FV ambiguity map $\mathcal{A}_k^{ipm}$ by first computing the BEV confidence map $\mathcal{C}_k^{ipm,bev}$, then estimating the BEV ambiguity map from it as $\mathcal{A}_k^{ipm,bev} = 1 - \mathcal{C}_k^{ipm,bev}$, and subsequently projecting it into the FV using $M^{-1}$. We then multiply $\mathcal{A}_k^{ipm}$ with $\mathcal{F}_k$ to obtain the ambiguous FV features $\mathcal{F}_k^{amb}$. We also account for regions ignored by IPM, i.e., flat regions above the principal point, by adding the features from these regions to $\mathcal{F}_k^{amb}$. ECM then processes these FV features to generate the ambiguity-correction features $\mathcal{F}_k^{bev,amb}$ in the BEV. ECM achieves this by first collapsing the FV features along the height dimension into a bottleneck dimension $B$ before horizontally expanding it to obtain the BEV features. Since ECM only corrects for errors made by IPM and does not predict the entire FV-BEV mapping, we use parameter-efficient 2D convolutions instead of parameter-intensive fully-connected layers used in the competing baselines. This significantly reduces the parameters in our model and promotes model efficiency. We then add $\mathcal{F}_k^{bev,amb}$ to $\mathcal{F}_k^{bev,ipm}$, and refine it using a residual block to generate flat BEV features $\mathcal{F}_k^{bev}$. During the training phase, flat regions in the BEV groundtruth actively supervises $\mathcal{F}_k^{bev,vis}$, obtained from $\mathcal{F}_k^{bev}$, to guide the ECM and $\mathcal{C}_k^{ipm}$ estimation.

Finally, $\mathcal{V}_k^{bev}$ and $\mathcal{F}_k^{bev}$ are concatenated in the BEV space and processed using a single 2D convolution layer to generate the composite BEV feature map $\mathcal{E}_k^{bev}$.

*C. Semantic and Instance Heads*

The semantic and instance heads of our PanopticBEV architecture follow the topology proposed in EfficientPS [25]. Both heads process the composite BEV feature maps, $\mathcal{E}_4^{bev}, \mathcal{E}_8^{bev}, \mathcal{E}_{16}^{bev}$ and $\mathcal{E}_{32}^{bev}$, and output the semantic logits and instance logits respectively. Briefly, the semantic head uses DPC and LSFE [25] modules with depthwise separable convolutions to separately process feature maps of different scales before augmenting them using feature alignment connections. These features are then upsampled to stride 4, concatenated along the channel dimension, and processed using a $1 \times 1$ convolution to generate semantic features with $N_{stuff} + N_{thing}$ channels. We further upsample these features to the output resolution and apply the softmax function to obtain the semantic logits.

We employ a modified Mask-RCNN architecture [28] with depthwise separable convolutions for the instance head. The instance head follows a two-stage approach wherein the first stage uses a Region Proposal Network (RPN) to output a set of region proposals and objectness scores for each input level. The second stage processes these proposals to extract region-specific features which are then used to generate bounding box, class, and mask predictions. To generate optimal instance predictions, we use anchors of scales $4, 8, 16$ and ratios $0.5, 1, 2$, and set the RPN NMS threshold to $0.7$. Further, we set region-specific NMS and score thresholds to $0.3$, $0.1$ and $0.2$, $0.3$ for the KITTI-360 and nuScenes datasets respectively.

*D. Panoptic Fusion Module*

Our panoptic fusion module builds upon the approach proposed in EfficientPS [25]. The EfficientPS fusion module first merges the per-pixel logits from the semantic and instance heads to generate panoptic logits having $N_{stuff} + N_{instance}$ channels. It then computes the *stuff* and *thing* class predictions using the argmax operation before copying them onto an empty canvas which results in many pixels being classified as *unknown*. To this end, we discard the argmax operation and copy steps, and introduce a cross entropy-based panoptic loss on the generated panoptic logits. This new fusion module enables end-to-end training of our PanopticBEV model, thereby enabling the model to learn panoptic-specific features and improving the overall PQ score. We evaluate the influence of this approach in the ablation study presented in Sec. IV-E.

*E. Losses*

We train PanopticBEV using six loss functions: a weighted cross entropy loss with hard mining for the semantic head ($\mathcal{L}_{sem}$), the standard Mask-RCNN [28] loss for the instance head ($\mathcal{L}_{inst}$), a cross-entropy loss on the panoptic segmentation output ($\mathcal{L}_{po}$), and binary cross entropy losses on the vertical-flat mask logits ($\mathcal{L}_{vf}$) as well as the vertical ($\mathcal{L}_v$) and flat ($\mathcal{L}_f$) region prediction masks. The final loss $\mathcal{L}$ is thus computed as

$$\mathcal{L} = w_c w_s \mathcal{L}_{sem} + \mathcal{L}_{inst} + \mathcal{L}_{po} + \mathcal{L}_{vf} + \lambda_v \mathcal{L}_v + \lambda_f \mathcal{L}_f, \quad (1)$$

where $w_c$ and $w_s$ refer to the class and sensitivity-based weights described in detail below, and $\lambda_v = \lambda_f = 10$.

*1) Class-based Weighting:* The class-based weight $w_c$ addresses the class imbalance in the dataset by increasing the weights of infrequent classes such *car* and *truck*. We compute the weight of a class as the inverse square root of its relative pixel frequency. However, due to the large difference between class weights, sometimes in the order of a magnitude, the infrequent class overflows into the frequent class resulting in fuzzy class boundaries. We address this problem by gradually decreasing the weight around infrequent classes using a linear combination of the frequent and infrequent class weights. To this end, we employ the L1-distance from the boundary of the infrequent class up to a distance of 20 pixels to compute the weight pertaining to each component. Mathematically, we compute the weight of point $\mathbf{p}$ which is $d$ pixels away from the infrequent class boundary using

$$w_c^{\mathbf{P}} = (20 - d)w_{infreq} + (d)w_{freq}, \quad (2)$$

where $d \leq 20$, $w_{freq}$ and $w_{infreq}$ represent the weights of the frequent and infrequent classes respectively.

*2) Sensitivity-based Weighting:* This weighting scheme normalizes the varying descriptiveness across the FV image due to the perspective projection. The perspective projection makes mapping far away objects from the FV into the BEV significantly more challenging, resulting in high uncertainty in distant regions. We address this problem by introducing the concept of FV-BEV sensitivity which we define as the change observed in the FV when a pixel is displaced by a unit value in the BEV. Accordingly, close and distant regions mapped into the BEV have high and low sensitivities respectively. Using the camera projection equation, we obtain

$$u = \frac{f_x x}{z} + c_x, \qquad v = \frac{f_y y}{z} + c_y, \qquad (3)$$

where $u$, $v$ are the image coordinates of a 3D point $\mathbf{p}$ located at coordinates $(x, y, z)$, $f_x$, $f_y$ are the focal lengths of the camera in terms of pixels, and $c_x$, $c_y$ denote the image center offset. We compute the sensitivity by first quantifying the effect of moving a pixel by an infinitesimal amount and then applying the orthographic BEV projection constraints on it. The pixel displacement in the FV can be denoted as

$$\vec{dr} = \vec{du} + \vec{dv}, \qquad (4)$$

where $\vec{du}$ and $\vec{dv}$ are estimated using the gradient of Eq. (3) as

$$\vec{du} = \frac{f_x}{z}\vec{dx} + 0\vec{dy} - \frac{f_x x}{z^2}\vec{dz}, \qquad (5)$$

$$\vec{dv} = 0\vec{dx} + \frac{f_y}{z}\vec{dy} - \frac{f_y y}{z^2}\vec{dz}. \qquad (6)$$

Substituting Eq. (5) and Eq. (6) in Eq. (4), we obtain

$$\vec{dr} = \frac{f_x}{z}\vec{dx} + \frac{f_y}{z}\vec{dy} - \frac{f_x x + f_y y}{z^2}\vec{dz}. \qquad (7)$$

Setting $\vec{dy}$ to 0 to account for the orthographic projection of the 3D points onto a fixed plane, and estimating the sensitivity map $S$ by computing the norm of $\vec{dr}$, we obtain

$$S = \|\vec{dr}\|_2 = \frac{\sqrt{f_x^2 z^2 + (f_x x + f_y y)^2}}{z^2}. \qquad (8)$$

The sensitivity weight, $w_{sens}$ is then computed by weighting $S$ by a constant $\lambda_s = 10$, scaling it using the log function, inverting, and normalizing it to give

$$w_{sens} = 1 + \frac{1}{log(1 + \lambda_s S)}. \qquad (9)$$

An illustration of this weighting function is shown in Fig. S.4 of the supplementary material.

## IV. EXPERIMENTAL EVALUATION

In this section, we present quantitative and qualitative evaluations of our proposed PanopticBEV model, and provide detailed ablation studies that demonstrates the efficacy of our contributions. We primarily use the panoptic quality (PQ) metric as the main evaluation criteria, but we also report the recognition quality (RQ), segmentation quality (SQ), and mean Intersection-over-Union (mIoU) scores for completeness.

### A. Datasets

We evaluate PanopticBEV on the large-scale KITTI-360 [10] and nuScenes [11] datasets. As the datasets themselves do not provide BEV panoptic segmentation groundtruth annotations, we generate the labels using a five-step process as depicted in Fig. S.3 of the supplementary material. First, we accumulate LiDAR points belonging to static objects over multiple frames and store the points belonging to dynamic objects for later use. Second, we transform both the accumulated static and dynamic point clouds into the BEV coordinates using the known ego pose and camera extrinsics. We subsequently project the transformed point cloud onto the XZ-plane using an orthographic projection to generate a sparse BEV image.

Third, we densify the sparse BEV image using a series of morphological dilate and erode operations on each class. We then project the 3D bounding boxes onto the BEV image and fuse them with the densified dynamic points to obtain the object instance masks. To prevent tree canopies from occluding the underlying classes, we add pixels belonging to the *vegetation* class at the end and only to regions that do not contain any other label. Fourth, we introduce a new stuff label called *occlusion*, depicted using light-gray in Fig. S.3, to account for regions occluded by other classes and thus not visible in the FV image. We define a pixel to be occluded if it is lower than a previously seen pixel along a 2D ray cast from the camera across the BEV image. Lastly, we zero-out the pixels lying outside the field-of-view of the camera and crop the resulting image to the required dimensions to obtain the BEV panoptic segmentation labels. Tab. S.1 in the supplementary material summarizes the various parameters used to generate the labels.

For the KITTI-360 dataset, we use sequences 0, 2-9 for training and hold out sequence 10 for validation, and for the nuScenes dataset, we follow the train-val split specified in [3] to obtain 702 train and 148 validation sequences.

### B. Training Protocol

We train PanopticBEV using an image of size $1408 \times 768$ pixels for KITTI-360 and $768 \times 448$ pixels for nuScenes. We augment the dataset using random horizontal flips, and random perturbations of the image brightness, contrast and saturation. We initialize the EfficientDet backbone with weights pre-trained on the COCO dataset while the remaining layers are randomly initialized using Xavier with biases set to zeros. We optimize the network with SGD for a total of 20 epochs on KITTI-360 and 30 epochs on nuScenes. We use a batch size of 8, a momentum of 0.9, and a weight decay of 0.0001. We employ a multi-step training schedule with an initial learning rate of 0.005 and decay it by a factor of 0.5 and 0.2 at epochs 10 and 15 for KITTI-360, and at epochs 15 and 25 for nuScenes.

### C. Quantitative Results

We evaluate the performance of our PanopticBEV model in comparison with IPM [12] and four novel baselines. For the IPM baseline, we apply the IPM algorithm on the panoptic segmentation masks obtained from the state-of-the-art EfficientPS [25] model. Further, we create four baselines by combining two state-of-the-art BEV semantic segmentation models, View Parsing Network (VPN) [8] and Pyramid Occupancy Network (PON) [3], with the instance head and

| Dataset | Method | PQ | SQ | RQ | PQ$^{Th}$ | SQ$^{Th}$ | RQ$^{Th}$ | PQ$^{St}$ | SQ$^{St}$ | RQ$^{St}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| KITTI-360 | IPM [12] + EPS [25] | 3.93 | 25.87 | 6.46 | 0.01 | 12.99 | 0.01 | 6.18 | 33.24 | 10.14 |
| | VPN [8] + EPS [25] | 17.62 | 64.74 | 25.93 | 10.45 | **67.37** | 14.87 | 21.72 | 63.24 | 32.24 |
| | VPN [8] + PDL [26] | 16.41 | **64.27** | 24.58 | 7.48 | 65.51 | 11.26 | 21.52 | **63.55** | 32.19 |
| | PON [3] + EPS [25] | 14.95 | 57.77 | 22.27 | 8.92 | 64.86 | 12.86 | 18.38 | 53.71 | 27.64 |
| | PON [3] + PDL [26] | 14.95 | 62.56 | 21.77 | 9.46 | 63.84 | 13.35 | 18.08 | 61.82 | 26.59 |
| | PanopticBEV (Ours) | **21.23** | 63.89 | **31.23** | **12.97** | 65.59 | **18.60** | **25.96** | 62.92 | **38.46** |
| nuScenes | IPM [12] + EPS [25] | 5.63 | 35.13 | 8.62 | 0.04 | 13.87 | 0.07 | 9.35 | 49.29 | 14.32 |
| | VPN [8] + EPS [25] | 14.35 | 63.67 | 21.16 | 6.35 | 66.16 | 9.52 | 19.69 | 62.00 | 28.92 |
| | VPN [8] + PDL [26] | 14.91 | 64.44 | 22.01 | 7.76 | **68.62** | 11.39 | 19.67 | 61.64 | 29.08 |
| | PON [3] + EPS [25] | 14.52 | 61.91 | 21.06 | 9.28 | 62.69 | 13.50 | 18.01 | 61.39 | 26.11 |
| | PON [3] + PDL [26] | 14.72 | 63.04 | 21.21 | 8.98 | 65.40 | 12.78 | 18.54 | 61.46 | 26.83 |
| | PanopticBEV (Ours) | **19.84** | 64.38 | **28.44** | **14.64** | 66.37 | **20.39** | 23.30 | **63.05** | **33.81** |

TABLE I: Evaluation of BEV panoptic segmentation performance on the KITTI-360 and nuScenes datasets. All scores are in [%].

| Dataset | Method | Road | Side. | Build. | Wall | Manm. | Veg. | Ter. | Occ. | Per. | 2-Wh. | Car | Truck | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| KITTI-360 | IPM [12] | 53.50 | 15.04 | 8.14 | 1.99 | - | 21.97 | 18.93 | 0.00 | 0.03 | 0.80 | 6.79 | 4.21 | 11.95 |
| | VED [19] | 65.37 | 29.94 | **31.65** | 8.96 | - | 38.93 | 28.67 | 38.93 | 0.01 | 0.06 | 27.17 | 9.41 | 25.37 |
| | VPN [8] | 70.98 | 35.58 | 22.56 | 13.46 | - | 37.32 | 31.59 | 43.27 | 3.91 | 4.83 | 38.17 | 10.60 | 28.39 |
| | PON [3] | 73.37 | 33.98 | 27.60 | 9.14 | - | 36.84 | 32.97 | 45.31 | 1.56 | 2.95 | 36.96 | 14.53 | 28.66 |
| | PanopticBEV† (Ours) | **75.50** | **40.08** | 28.68 | **16.41** | - | **40.91** | **35.58** | **48.29** | 4.76 | 8.46 | **42.48** | 15.30 | **32.40** |
| nuScenes | IPM [12] | 50.56 | 8.69 | - | - | 18.79 | 21.42 | 10.44 | 0.00 | 0.12 | 0.09 | 6.00 | 1.21 | 11.73 |
| | VED [19] | 73.68 | 23.20 | - | - | 34.07 | 33.47 | 29.28 | 32.14 | 1.58 | 1.95 | 29.67 | 22.74 | 28.18 |
| | VPN [8] | 73.16 | 23.82 | - | - | 33.03 | 32.27 | 29.47 | 31.01 | 2.54 | 6.25 | 30.72 | 23.55 | 28.58 |
| | PON [3] | 74.07 | 23.25 | - | - | 31.56 | 34.40 | 29.03 | 32.21 | 2.94 | 5.56 | 32.21 | 27.56 | 29.28 |
| | PanopticBEV† (Ours) | **77.32** | **28.55** | - | - | **36.72** | 35.06 | **33.56** | **36.65** | 4.98 | 9.63 | **40.53** | **33.47** | **33.65** |

TABLE II: Evaluation of BEV semantic segmentation performance. All values are in [%] and '-' indicates that the respective class is not present in the dataset.

| Dataset | Method | # Params (M) Trans. | Total | MAC (G) | Runtime (ms) |
|---|---|---|---|---|---|
| KITTI-360 | IPM [12] + EPS [25] | - | 45.0 | 418.1 | 140.6 |
| | VPN [8] + EPS [25] | 152.5 | 192.1 | 559.3 | **61.7** |
| | VPN [8] + PDL [26] | 152.5 | 175.9 | 496.9 | 66.7 |
| | PON [3] + EPS [25] | 58.2 | 97.6 | 719.2 | 254.9 |
| | PON [3] + PDL [26] | 58.2 | 92.2 | 655.1 | 282.9 |
| | PanopticBEV (Ours) | **9.5** | **39.5** | **379.4** | 277.5 |
| nuScenes | IPM [12] + EPS [25] | - | 45.0 | **117.0** | 120.8 |
| | VPN [8] + EPS [25] | 61.9 | 101.5 | 440.0 | **59.5** |
| | VPN [8] + PDL [26] | 61.9 | 85.3 | 306.3 | 67.8 |
| | PON [3] + EPS [25] | 62.0 | 101.4 | 859.0 | 302.4 |
| | PON [3] + PDL [26] | 62.0 | 95.6 | 697.9 | 307.8 |
| | PanopticBEV (Ours) | **9.9** | **39.8** | 377.7 | 238.7 |

TABLE III: Comparison of model efficiency on KITTI-360 and nuScenes.

panoptic fusion module from two FV panoptic segmentation networks EfficientPS (EPS) and Panoptic-DeepLab (PDL) [26]. Tab. I presents the results from this comparison on both the KITTI-360 and nuScenes datasets.

We observe that our proposed PanopticBEV model outperforms all the baselines by a large margin on both the datasets. PanopticBEV achieves an improvement of 3.61 pp over the best performing baseline VPN + EPS on the KITTI-360 dataset, and an improvement of 4.93 pp over the best performing baseline VPN + PDL on the nuScenes dataset in terms of the PQ score. Moreover, we observe a significant improvement in the RQ score as compared to the VPN-based baselines which signifies that our model achieves better detection performance. Furthermore, the consistent improvement in PQ$^{Th}$ and PQ$^{St}$ scores can be attributed to our dense transformer module which independently transforms the vertical and flat regions resulting

in richer BEV features. We also observe that the baselines do not generalize well across both the datasets. For instance, VPN + EPS achieves the best performance on KITTI-360, but performs the worst among learnable-transformer models on nuScenes. Whereas, PanopticBEV consistently outperforms all the baselines by a large margin on both datasets, demonstrating its effective generalization ability.

We also evaluate the performance of PanopticBEV for the task of semantic segmentation, by discarding the instance head and the panoptic fusion module. We denote this model as PanopticBEV† and compare its performance with IPM [12] and three state-of-the-art BEV segmentation methods, namely, Variational Encoder Decoder (VED) [19], VPN [8], and PON [3]. Tab. II presents the results of this comparison on both the datasets. We observe that our PanopticBEV† model once again substantially outperforms the existing methods, thereby achieving state-of-the-art performance. We observe a significant improvement in performance for classes such as *road*, *sidewalk*, *two-wheeler*, *car*, and *truck*. This improvement in both the vertical and flat semantic classes can be attributed to the targeted transformations performed by our dense transformer. The region-specific vertical and flat transformers capture the intricate relationship pertaining to these regions resulting in improved spatial as well as boundary estimates.

### D. Evaluation of Model Efficiency

In this section, we evaluate the efficiency of our PanopticBEV model on both the datasets. From Tab. III, we observe that our model is more than two-times more parameter efficient than the baselines and uses significantly fewer Multiply-Accumulate (MAC) operations. A large chunk of the efficiency can be attributed to our dense transformer module that consumes

| Model | $\mathcal{T}_v$ | $\mathcal{T}_f^{IPM}$ | ECM | $w_c$ | $w_s$ | Scales | Fusion | PQ | SQ | RQ | PQ$^{\text{Th}}$ | SQ$^{\text{Th}}$ | RQ$^{\text{Th}}$ | PQ$^{\text{St}}$ | SQ$^{\text{St}}$ | RQ$^{\text{St}}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| M1 | - | - | - | - | - | 4-32 | Ours | 13.12 | 50.65 | 20.06 | 2.45 | 33.03 | 3.74 | 19.21 | 60.72 | 29.39 |
| M2 | ✓ | - | - | - | - | 4-32 | Ours | 10.95 | 49.79 | 17.18 | 2.36 | 30.75 | 3.56 | 15.86 | 60.67 | 24.97 |
| M3 | - | ✓ | - | - | - | 4-32 | Ours | 12.35 | 44.46 | 19.00 | 2.59 | 16.52 | 3.91 | 17.92 | 60.43 | 27.61 |
| M4 | ✓ | ✓ | - | - | - | 4-32 | Ours | 20.20 | 62.89 | 29.81 | 10.51 | 63.61 | 15.62 | 25.74 | 62.47 | 37.92 |
| M5 | ✓ | ✓ | ✓ | - | - | 4-32 | Ours | 20.33 | 63.94 | 29.94 | 11.99 | 65.99 | 17.20 | 25.09 | 62.77 | 37.23 |
| M6 | ✓ | ✓ | ✓ | ✓ | - | 4-32 | Ours | 20.38 | **64.73** | 29.73 | 12.33 | **67.90** | 17.38 | 24.99 | 62.91 | 36.79 |
| M7 | ✓ | ✓ | ✓ | ✓ | ✓ | 4-32 | Ours | **21.23** | 63.89 | **31.23** | **12.97** | 65.59 | **18.60** | **25.96** | 62.92 | **38.46** |
| M8 | ✓ | ✓ | ✓ | ✓ | ✓ | 8-64 | Ours | 20.55 | 63.68 | 30.37 | 10.49 | 64.26 | 15.73 | 26.30 | 63.35 | 38.74 |
| M9 | ✓ | ✓ | ✓ | ✓ | ✓ | 4-32 | EPS [25] | 20.53 | 64.85 | 29.72 | 11.72 | 65.55 | 16.39 | 25.56 | 64.46 | 37.33 |

TABLE IV: Ablation study on the various architectural components proposed in our PanopticBEV model. The results are reported on the KITTI-360 dataset.
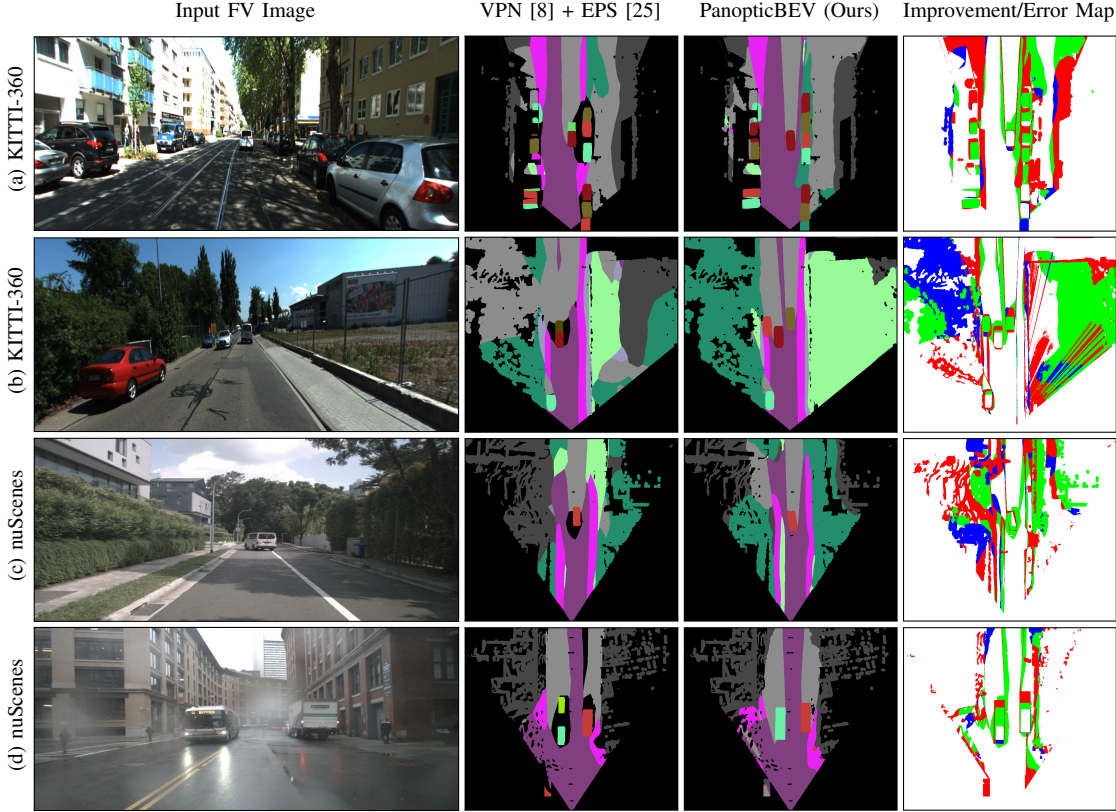


Fig. 4: Qualitative results of BEV panoptic segmentation in comparison with the best performing baseline from Tab. I on the KITTI-360 and nuScenes datasets. We show the Improvement/Error map which depicts pixels misclassified by the baseline and correctly predicted by the PanopticBEV model in green, pixels misclassified by PanopticBEV and correctly by the baseline in blue, and pixels misclassified by both models in red.

more than six-times fewer parameters as compared to its counterparts. This can be primarily attributed to the use of 2D and 3D convolutions instead of the fully-connected layers in the dense transformer module. Our PanopticBEV model has an inference time of nearly $280\,\text{ms}$ on KITTI-360 and $240\,\text{ms}$ on nuScenes, most of which is due to the expensive 3D convolution operations in $\mathcal{T}_v$. Given that our model uses significantly fewer MAC operations, we believe that optimizing the implementation of 3D convolutions would result in a substantial decrease in the inference time and enable real-time performance.

### E. Ablation Study

In this section, we study the influence of various architectural components proposed in this work. Tab. IV presents the results of this study on the KITTI-360 dataset. We begin with model M1 comprising of a bare-bones variant of PanopticBEV that maps the input FV image to the BEV space without any transformer or associated loss functions, and add the various

components to it. Upon adding only the vertical transformer ($\mathcal{T}_v$) in model M2, or only the horizontal transformer ($\mathcal{T}_f$) in model M3, we observe a significant decrease in model performance. This is due to the fact that $\mathcal{T}_v$, by itself, is not powerful enough and corrupts the BEV feature space, while $\mathcal{T}_f^{IPM}$, by itself, distorts the features of objects above the ground plane resulting in a significant $14.23\,\text{pp}$ drop in the SQ$^{\text{Th}}$ score. However, when both $\mathcal{T}_v$ and $\mathcal{T}_f^{IPM}$ are used to independently transform the vertical and flat regions into the BEV (model M4), we observe a notable improvement of $7.08\,\text{pp}$ in the PQ score, which demonstrates the utility of our region-specific transformers. Note that the model M4 already outperforms all the BEV panoptic segmentation baselines reported in Tab. I. In model M5, we incorporate ECM to account for irregularities in the flat regions, which further increases the PQ score to $20.33\%$. We then use our spatial class-based weighting scheme in model M6 to prevent overflowing of infrequent *thing* classes resulting in a noticeable improvement

in the SQ and SQ^Th scores. Finally, we employ our novel sensitivity-based weighting function in model M7 which leads to an improvement of $0.98\,\mathrm{pp}$ in the PQ score and $1.5\,\mathrm{pp}$ in the RQ score. This large improvement in the RQ score demonstrates that our sensitivity-based weighting scheme enables the model to achieve a good balance between precision and recall of the matched segments. We denote model M7 as our proposed PanopticBEV architecture. Additionally, model M8 shows the improvement in performance due to using features of strides 4-32 in our backbone instead of 8-64 used in the standard EfficientDet architecture, and model M9 shows the improvement due to our panoptic fusion scheme as compared to that used in the EfficientPS architecture. We also perform an additional ablation study to analyze the impact of multi-scale features on the performance of our model, which we present in Sec. S.3 of the supplementary material.

*F. Qualitative Evaluation*

We qualitatively evaluate the performance our PanopticBEV model in comparison to the best performing baseline, VPN + EPS, in Fig. 4. We observe from the Improvement/Error map that our model accurately segments all the object instances in the scene despite being only partially visible in the FV image. In Fig. 4(a), we observe that our model segments the white car in front of the ego-vehicle as a single instance and also accurately segments the parked vehicles on the right, while the baseline fails to do so. This observation also extends to Fig. 4(b) in which all the three cars in the distance are accurately segmented by our model. Fig. 4(c) demonstrates the ability of our model to segment instances of objects with the right orientation. We observe that the baseline incorrectly segments objects having an orientation that is not parallel to the optical axis in the FV image, e.g., the white van angled towards the left. Lastly, Fig. 4(d) demonstrates that our model effectively estimates the BEV panoptic segmentation predictions even in challenging weather conditions. We provide additional qualitative results of the BEV panoptic and BEV semantic segmentation in Fig. S.5, Fig. S.6, Fig. S.7, Fig. S.8 and Fig. S.9 of the supplementary material.

## V. CONCLUSION

In this paper, we present the first end-to-end trainable BEV panoptic segmentation architecture that takes monocular images in the FV as input and predicts coherent panoptic segmentation masks in the BEV. Our PanopticBEV architecture incorporates the proposed dense transformer module which uses two distinct transformers to independently transform features belonging to vertical and flat regions in the input FV image to the BEV. We also introduce a sensitivity-based weighting scheme to account for the varying levels of descriptiveness across the FV image by intelligently weighting pixels in the BEV space. Using extensive evaluations on the KITTI-360 and nuScenes datasets, we demonstrate that our model outperforms both the BEV panoptic and semantic segmentation baselines, thereby setting the new state-of-the-art for both these tasks.

## ACKNOWLEDGMENT

## REFERENCES

[1] J. V. Hurtado, R. Mohan, W. Burgard, and A. Valada, "Mopt: Multi-object panoptic tracking," *arXiv preprint arXiv:2004.08189*, 2020.
[2] J. Philion and S. Fidler, "Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d," in *European Conf. on Computer Vision*, 2020.
[3] T. Roddick and R. Cipolla, "Predicting semantic map representations from images using pyramid occupancy networks," in *IEEE Conf. on Computer Vision and Pattern Recognition*, June 2020.
[4] M. H. Ng, K. Radia, J. Chen, D. Wang, I. Gog, and J. E. Gonzalez, "Bev-seg: Bird's eye view semantic segmentation using geometry and semantic point cloud," *arXiv preprint arXiv:2006.11436*, 2020.
[5] N. Radwan, W. Burgard, and A. Valada, "Multimodal interaction-aware motion prediction for autonomous street crossing," *The International Journal of Robotics Research*, vol. 39, no. 13, pp. 1567–1598, 2020.
[6] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollar, "Panoptic segmentation," in *IEEE Conf. on Computer Vision and Pattern Recognition*, June 2019.
[7] L. Reiher, B. Lampe, and L. Eckstein, "A sim2real deep learning approach for the transformation of images from multiple vehicle-mounted cameras to a semantically segmented image in bird's eye view," in *Int. Conf. on Intelligent Transportation Systems*, 2020.
[8] B. Pan, J. Sun, H. Y. T. Leung, A. Andonian, and B. Zhou, "Cross-view semantic segmentation for sensing surroundings," *IEEE Robotics & Automation Letters*, vol. 5, no. 3, pp. 4867–4873, 2020.
[9] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Conf. on Computer Vision and Pattern Recognition*, 2020.
[10] J. Xie, M. Kiefel, M.-T. Sun, and A. Geiger, "Semantic instance annotation of street scenes by 3d to 2d label transfer," in *IEEE Conf. on Computer Vision and Pattern Recognition*, 2016.
[11] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," *arXiv preprint arXiv:1903.11027*, 2019.
[12] H. A. Mallot, H. H. Bülthoff, J. Little, and S. Bohrer, "Inverse perspective mapping simplifies optical flow computation and obstacle detection," *Biological cybernetics*, vol. 64, no. 3, pp. 177–185, 1991.
[13] S. Ammar Abbas and A. Zisserman, "A geometric approach to obtain a bird's eye view from an image," in *IEEE/CVF Int. Conf. on Computer Vision Workshops*, 2019.
[14] X. Zhu, Z. Yin, J. Shi, H. Li, and D. Lin, "Generative adversarial frontal view to bird view synthesis," in *Int. Conf. on 3D Vision*, 2018.
[15] T. Bruls, H. Porav, L. Kunze, and P. Newman, "The right (angled) perspective: Improving the understanding of road scenes using boosted inverse perspective mapping," in *IEEE Intelligent Vehicles Symp.*, 2019.
[16] K. Mani, S. Daga, S. Garg, S. S. Narasimhan, M. Krishna, and K. M. Jatavallabhula, "Monolayout: Amodal scene layout from a single image," in *IEEE Wint. Conf. on Appl. of Computer Vision*, 2020, pp. 1689–1697.
[17] T. Roddick, A. Kendall, and R. Cipolla, "Orthographic feature transform for monocular 3d object detection," *British Mac. Vision Conf.*, 2019.
[18] A. Palazzi, G. Borghi, D. Abati, S. Calderara, and R. Cucchiara, "Learning to map vehicles into bird's eye view," in *Int. Conf. on Image Analysis and Processing.* Springer, 2017, pp. 233–243.
[19] C. Lu, M. J. G. van de Molengraft, and G. Dubbelman, "Monocular semantic occupancy grid mapping with convolutional variational encoder–decoder networks," *IEEE Robotics & Automation Letters*, 2019.
[20] S. Sengupta, P. Sturgess, L. Ladický, and P. H. S. Torr, "Automatic dense visual semantic mapping from street-level imagery," in *Int. Conf. on Intelligent Robots and Systems*, 2012, pp. 857–862.
[21] S. Schulter, M. Zhai, N. Jacobs, and M. Chandraker, "Learning to look around objects for top-view representations of outdoor scenes," in *European Conf. on Computer Vision*, 2018, pp. 787–802.
[22] Q. Li, A. Arnab, and P. H. Torr, "Weakly- and semi-supervised panoptic segmentation," in *European Conf. on Computer Vision*, 2018.
[23] J. Li, A. Raventos, A. Bhargava, T. Tagawa, and A. Gaidon, "Learning to fuse things and stuff," *arXiv preprint arXiv:1812.01192*, 2018.
[24] Y. Xiong, R. Liao, H. Zhao, R. Hu, M. Bai, E. Yumer, and R. Urtasun, "Upsnet: A unified panoptic segmentation network," in *IEEE Conf. on Computer Vision and Pattern Recognition*, 2019, pp. 8810–8818.
[25] R. Mohan and A. Valada, "Efficientps: Efficient panoptic segmentation," *Int. Journal of Computer Vision*, 2021.
[26] B. Cheng, M. D. Collins, Y. Zhu, T. Liu, T. S. Huang, H. Adam, and L.-C. Chen, "Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation," *arXiv preprint arXiv:1911.10194*, 2020.
[27] N. Gao, Y. Shan, Y. Wang, X. Zhao, Y. Yu, M. Yang, and K. Huang, "Ssap: Single-shot instance segmentation with affinity pyramid," in *Int. Conf. on Computer Vision*, 2019, pp. 642–651.
[28] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Int. Conf. on Computer Vision*, 2017, pp. 2980–2988.

# Bird's-Eye-View Panoptic Segmentation Using Monocular Frontal View Images

## - Supplementary Material -

Nikhil Gosala and Abhinav Valada

In this supplementary material, we present additional details about our novel dense transformer module and additional figures to illustrate the principle behind it. We then describe our methodology for generating BEV panoptic segmentation groundtruth labels for the KITTI-360 and nuScenes datasets. Furthermore, we provide additional qualitative results for both the BEV panoptic segmentation and BEV semantic segmentation.

## S.1. TECHNICAL APPROACH

In this section, we illustrate the principle governing our dense transformer module and also provide more details pertaining to the topology of the dense transformer. We then present an illustration of our sensitivity-based weighting function and describe its formulation.

### A. Dense Transformer Principle

We design our novel dense transformer module based on the principle of how different regions of the 3D world are projected into a 2D image, as illustrated in Fig. S.1. A column belonging to flat regions in the world maps to a perspectively-distorted area in the BEV space. Since flat regions are fully observable unless occluded by another object, the transformation of the flat regions into the BEV involves correcting the perspective distortion and inferring the missing information in distant regions using the learned model.

Conversely, a column belonging to a vertical region maps to an orthographic projection of a volumetric region in the BEV space. Being projections of 3D volumetric objects such as vehicles and humans, vertical regions are not fully observable and often completely lack a dimension. For instance, a car is not fully observable because of the absence of information pertaining to its spatial extents. Furthermore, their depth in the world, as captured from a monocular camera, is also ambiguous which further makes the problem even more challenging. Transforming a vertical non-flat object into the BEV thus requires the prediction of both its spatial location and extents which the model learns using a data-driven paradigm in our setting.

### B. Dense Transformer Architecture

Our proposed dense transformer module transforms the intermediate features of the network backbone, using two distinct transformers that independently transform features belonging to vertical and flat regions in the input FV image to the BEV coordinates. Fig. S.2 presents the detailed topology of our dense transformer module. The semantic masking module
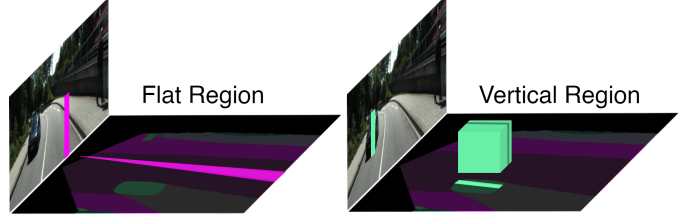


Fig. S.1: Illustration of the contrasting transformation patterns observed for the vertical and flat regions when transforming a monocular FV image into the BEV space.

$\mathcal{M}_k$, first processes each scale $\mathcal{E}_k$ from the network backbone to predict the vertical and flat semantic masks $\mathcal{S}_k^v$ and $\mathcal{S}_k^f$. $M_k$ achieves this using a sequence of three 2D convolutional layers with $3 \times 3$ kernels. We then obtain the vertical and flat FV features $\mathcal{V}_k$ and $\mathcal{F}_k$, by computing the Hadamard product between $\mathcal{E}_k$ and the corresponding semantic mask. We actively supervise $\mathcal{S}_k^v$ and $\mathcal{S}_k^f$ using the FV groundtruth vertical-flat masks to guide $\mathcal{M}_k$ during the training phase.

For the KITTI-360 dataset, we create the FV vertical-flat semantic groundtruth by grouping corresponding classes in the FV panoptic segmentation groundtruth. On the contrary, due to the lack of such groundtruth labels in the nuScenes dataset, we generate pseudo-labels for the FV vertical-flat masks using the approach described in Sec. S.2.

The vertical transformer processes the vertical FV feature map $\mathcal{V}_k$, to generate the vertical BEV features $\mathcal{V}_k^{bev}$. We first expand $\mathcal{V}_k$ into a 3D volumetric lattice using a single 3D convolutional layer with a $3 \times 3$ kernel. Simultaneously, we generate a spatial occupancy mask $\mathcal{M}_k$ for vertical regions to estimate the probability of a pixel being occupied by a vertical element in the BEV. We generate this spatial occupancy mask by first expanding the number of channels in $\mathcal{V}_k$ to the depth dimension ($Z$) using a sequence of 2D convolutions, and then flattening it along the height dimension. We then broadcast this spatial occupancy mask over the volumetric lattice to constrain the spatial extents of the vertical regions in the 3D grid. Subsequently, we reshape the 3D volumetric lattice and collapse it along the height dimension using a 3D convolutional layer with a $3 \times 3$ kernel to generate the perspectively-distorted vertical features in the BEV space. We correct the perspective distortion in the vertical BEV feature map, carried forward from the perspective projection of the input FV image, by resampling the feature map using the known camera intrinsics and BEV projection resolution using the approach proposed in [3]. We further process the output of the resampling step using a 2D convolutional layer with a 3x3 kernel to generate the final vertical BEV features $\mathcal{V}_k^{bev}$.
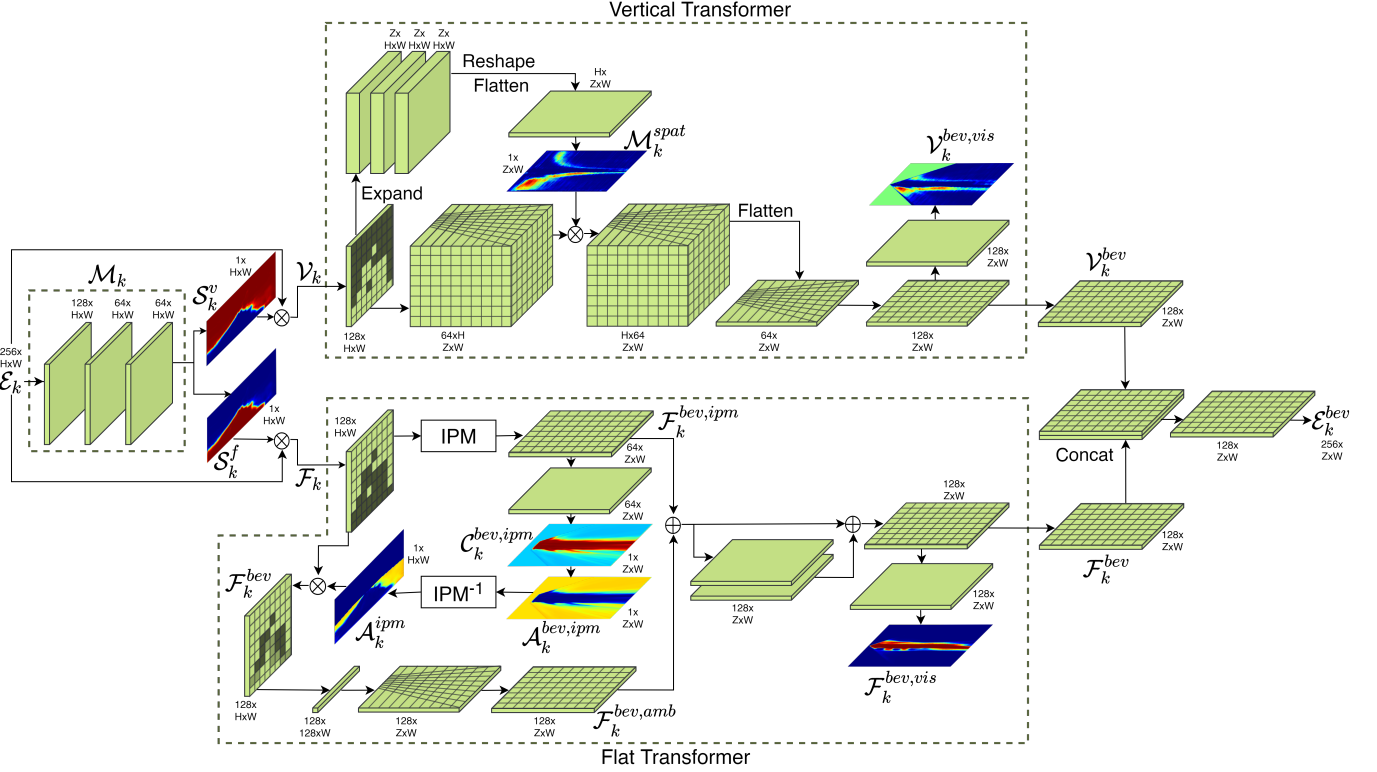
Fig. S.2: Detailed architectural diagram of our novel dense transformer module. Every convolution is followed by an in-place activated batch norm layer that applies both the batch norm and the non-linearity using a single function. All intermediate convolutions use a kernel of size $3 \times 3$ while all the channel-mapping and output convolutions use a $1 \times 1$ kernel. The perspective distortion in the feature maps is represented using slanted lines while the feature maps unwarped using a resampling operation are represented using a regular checkerboard-like grid. The outputs of the semantic masking module $\mathcal{S}_k^f$ and $\mathcal{S}_k^v$ are actively supervised by a FV vertical-flat mask during the training phase. The vertical spatial attention mask $\mathcal{M}_k^{spat}$ and the flat region estimation mask $\mathcal{F}_k^{bev,vis}$ are supervised using the BEV groundtruth during training. $H$, $W$ in the figure represent the height and width of the input feature map, and $Z$ represents the depth of the BEV prediction.

We transform the flat FV feature map $\mathcal{F}_k$ into the flat BEV feature map $\mathcal{F}_k^{ipm}$, using our flat transformer module. The flat transformer is based on the IPM algorithm reinforced with an error correction module (ECM) to account for the errors introduced by IPM. The IPM algorithm generates a non-learnable homography $M$, which when multiplied with the FV features generates features in the BEV. Due to its flat-world assumption, the IPM algorithm is applicable only to feature points that lie on the defined ground plane. Since $\mathcal{F}_k$, by definition, contains only flat regions, using the IPM algorithm to transform the flat FV features into the BEV generates acceptable results. However, since flat regions in the 3D world as not perfectly flat, using only the IPM algorithm introduces errors into the BEV prediction.

We account for these errors using a learnable ECM that is optimized alongside the IPM algorithm during the training phase. The ECM works by estimating regions where the IPM could potentially be erroneous and applies the ECM to these regions. To this end, we first compute the confidence in the IPM transformation $C_k^{bev,ipm}$, by a applying a single 2D convolution with a $3 \times 3$ kernel followed by a channel mapping layer with a $1 \times 1$ kernel to the output of IPM $F_k^{bev,ipm}$. The IPM ambiguity is then computed using the equation $\mathcal{A}^{bev,ipm} = 1 - \mathcal{C}^{bev,ipm}$. Examples of the IPM confidence and ambiguity maps are shown in Fig. S.2.

We transform the IPM ambiguity map from the BEV back into the BEV using the inverse of the estimated IPM homography, i.e., $M^{-1}$. We then multiply the FV ambiguity map with the flat FV features to mask out regions of high confidence while retaining only the ambiguous regions. We also estimate flat features ignored by IPM, i.e., features above the principal point, and add them to the FV ambiguity map to allow the ECM to operate on such areas as well. We then collapse these ambiguous flat features along the height dimension to a bottleneck dimension of size $B$ using a $3 \times 3$ convolution, followed by a 2D convolution in the bottleneck dimension to further refine the collapsed features. Subsequently, we expand the bottleneck features along the depth dimension using a 2D convolution and further refine the expanded feature map using another 2D convolution. Here, we account for the perspective distortion in the BEV feature map by resampling it using the known camera intrinsics using the approach described in [3], to generate the ambiguity correction features in the BEV $\mathcal{F}_k^{bev,amb}$. We then add $\mathcal{F}_k^{bev,amb}$ to the output of IPM $\mathcal{F}_k^{bev,ipm}$, and refine it using a residual block consisting of two 2D convolutional layers with a $3 \times 3$ kernel, followed by another 2D convolutional layer to generate the final flat BEV feature map $\mathcal{F}_k^{bev}$.

The vertical and flat BEV feature maps $\mathcal{V}_k^{bev}$ and $\mathcal{F}_k^{bev}$ are subsequently concatenated along the channel dimension and processed using a 2D convolutional layer to generate the final composite feature map $\mathcal{E}_k^{bev}$ in the BEV coordinates.
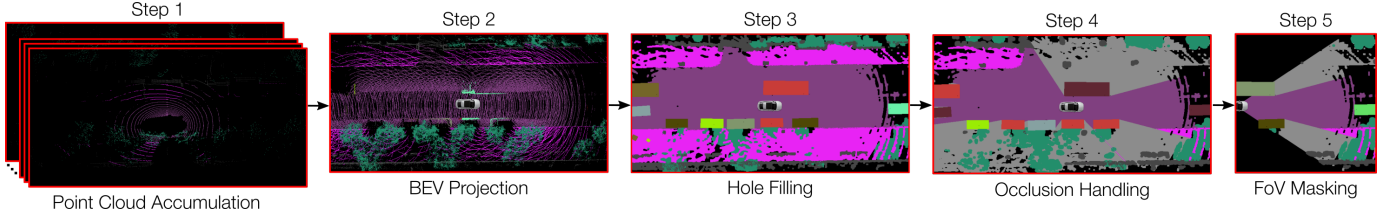
Fig. S.3: Illustration of outputs from each stage of the BEV panoptic segmentation ground truth generation pipeline. In the first stage, point clouds are motion-compensated and accumulated over multiple time steps to generate a relatively dense point cloud. The accumulated point cloud is then orthographically projected into the BEV using the ego pose. The third stage densifies the projected BEV image using a series of morphological dilate and erode operations on each class. Simultaneously, 3D bounding boxes are used to densify regions belonging to *thing* classes. In the fourth stage, an occlusion mask representing regions occluded by other classes is generated, and the last stage masks the regions outside the field-of-view of the camera.

| Dataset | Image Size | Resolution | Dilation | | | | | Erosion | | | | |
|---------|-----------|------------|----------|--------|------|---------|---------|----------|--------|------|---------|---------|
| | | | $St_T$ | $St_S$ | Veg. | $Th_V$ | $Th_P$ | $St_T$ | $St_S$ | Veg. | $Th_V$ | $Th_P$ |
| KITTI-360 | $768 \times 704$ | $0.074\,\mathrm{m/px}$ | 3 | 9 | 9 | 9 | 7 | 3 | 5 | 3 | 5 | 5 |
| nuScenes | $896 \times 768$ | $0.077\,\mathrm{m/px}$ | 3 | 9 | 9 | 9 | 7 | 3 | 5 | 3 | 5 | 5 |

TABLE S.1: The parameters used to generate the panoptic BEV ground truth from annotated LiDAR point clouds. In our setting, the BEV image has a resolution of $7.4\,\mathrm{cm/px}$ and $7.6\,\mathrm{cm/px}$ for KITTI-360 and nuScenes dataset respectively. In the table, $St_T$ refers to tall classes consisting of *wall*, *pole*, *traffic light* and *traffic sign*, while $St_S$ refer to short stuff classes which comprises of classes *ground*, *road*, *sidewalk*, *parking*, *fence* and *terrain*. *Veg.* refers to the *vegetation* class, and $Th_V$ and $Th_P$ refers to all vehicle and person thing classes respectively.
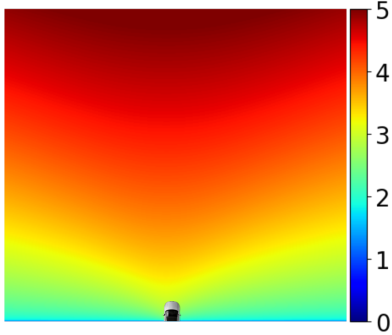


Fig. S.4: Illustration of the sensitivity-based weighting function across the BEV space. The ego-vehicle, depicted by the car, is located at the bottom of the plot. Close regions are highly sensitive in the FV image and can be easily mapped to the BEV, resulting in the sensitivity weight described by Eq. (9) being low for such regions. On the contrary, mapping far-away regions is much more difficult which results in their low sensitivity and accordingly a high sensitivity weight.

### C. Sensitivity-Based Weighting

The sensitivity-based weighting function accounts for the varying descriptiveness across the FV image by intelligently weighting pixels in the BEV space. Owing to the perspective projection of 2D cameras, the apparent motion observed in the FV image when a 3D point close to the camera is moved by a unit value is significantly larger than the motion observed when a far-away pixel is moved by a similar amount. In other words, close regions have a high sensitivity while far regions have low sensitivity. This disparity makes it extremely difficult to differentiate between small changes in distance for the far-away regions. To address this disparity, we introduce a sensitivity-based weighting function as described in Eq. (9) to up-weight pixels belonging to far-away regions. This up-weighting focuses the network on farther regions which helps in improving the performance of the model. Fig. S.4 shows a plot of the sensitivity-based weighting function across the BEV space.

## S.2. DATASETS

### A. KITTI-360 and nuScenes Dataset Preparation

We generate the dense panoptic BEV groundtruth annotations for both the KITTI-360 and nuScenes datasets using a five stage pipeline as depicted in Fig. S.3. The pipeline takes as input, the annotated LiDAR point clouds, 3D bounding boxes, ego vehicle pose and camera extrinsics, and outputs a dense panoptic groundtruth image in the BEV coordinates. In the first stage, static LiDAR points, i.e., points belonging to stationary objects, are accumulated over multiple frames to generate a dense static point cloud. Simultaneously, dynamic points, i.e., points belonging to movable objects, are also stored for use in the downstream stages. In the second stage, both the accumulated static point cloud and the dynamic point cloud are transformed into the BEV coordinate system of the $k^{th}$ frame using the camera extrinsics $M$ and the ego pose for the $k^{th}$ frame $e_k$. This transformed point cloud is then projected onto the XZ-plane using an orthographic projection to generate the initial BEV image. However, this BEV image is extremely sparse with the dynamic objects being invisible and having a sparsity factor greater than $60\%$.

The third stage, thus, densifies this image using a sequence of morphological dilate and erode operations independently on each class. To address the lack of LiDAR points on the far side of dynamic objects, 3D bounding boxes are projected into the BEV image and are intelligently fused with the existing dynamic points to obtain realistic looking instances. To ensure that the tree canopies do not occlude the underlying classes, they are added to the BEV image in the end and only in regions that do not contain any other label. Since it is extremely difficult for the network to hallucinate labels behind occlusions (ex: regions behind cars), the fourth stage generates an occlusion mask using the height map obtained from the second stage. A new stuff label *occlusion*, depicted using light-grey in Fig. S.3, is introduced to incorporate the occlusion mask into densified panoptic BEV image. Lastly, the pixels that lie outside the

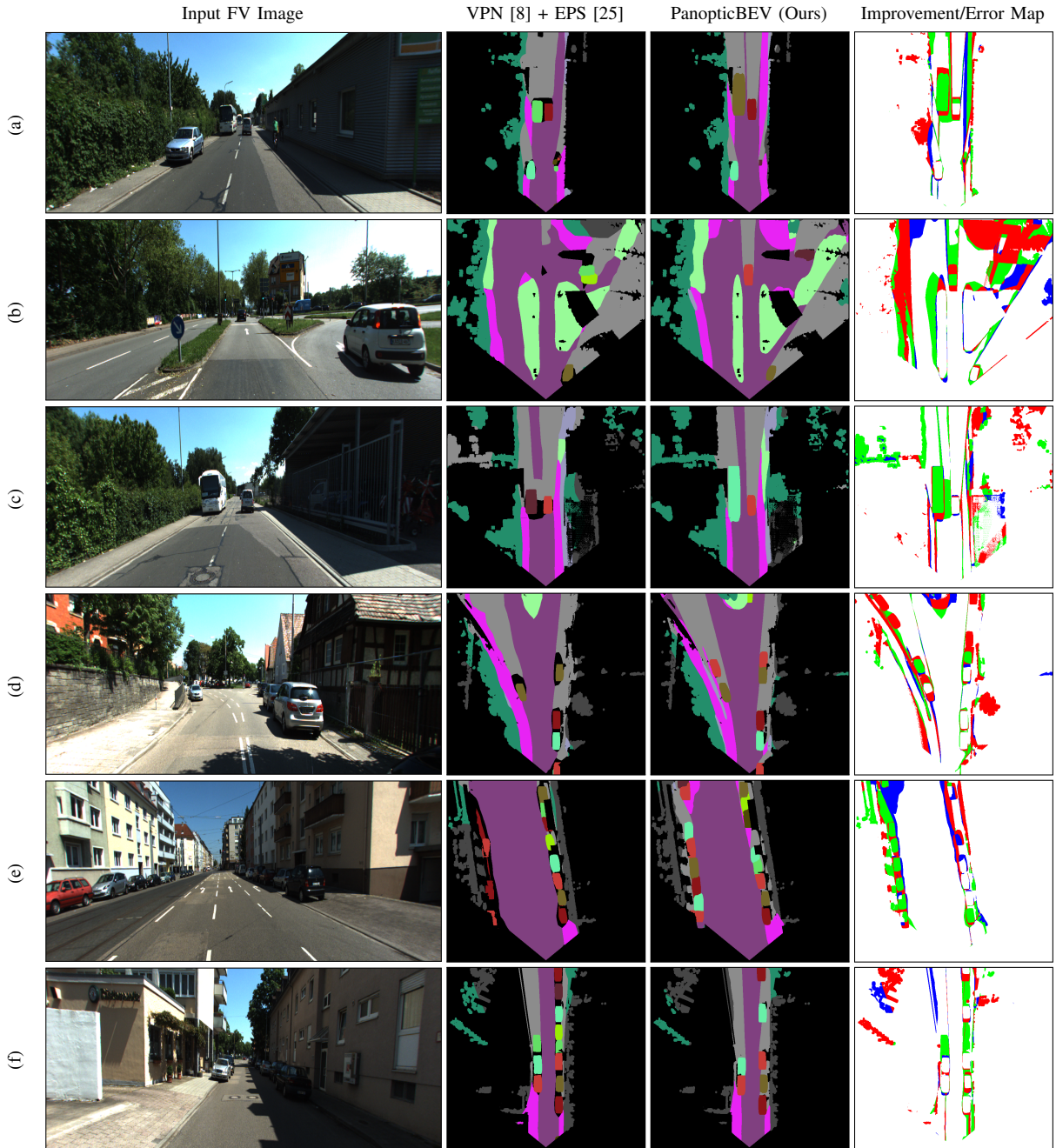|  | Input FV Image | VPN [8] + EPS [25] | PanopticBEV (Ours) | Improvement/Error Map |

Fig. S.5: Additional qualitative results comparing the performance of our PanopticBEV model with the best performing baseline on the KITTI-360 dataset. The rightmost column shows the Improvement/Error map which depicts the pixels misclassified by the baseline but correctly predicted by the PanopticBEV model in green, pixels misclassified by the PanopticBEV model but correctly predicted by the baseline in blue, and pixels misclassified by both models in red.

field-of-view (FoV) of the camera are zeroed-out, and the image is cropped to the required dimensions to generate the final panoptic BEV labels. The parameters that we use to generate the BEV panoptic segmentation labels are summarized in Tab. S.1.

### B. Frontal View Annotations for nuScenes

The nuScenes dataset does not provide dense semantic segmentation or panoptic segmentation groundtruth labels for FV images. This poses a challenge to our training procedure which relies on the vertical-flat groundtruth labels to supervise the semantic masking module in our transformer during the

training phase. We address this challenge by generating pseudo-labels for the vertical and flat regions in the FV image using the EfficientPS model. We train the network using a set of manually annotated vertical-flat masks containing 478 images, and use this network to generate pseudo-labels for all images in the training set.

### S.3. ADDITIONAL ABLATION STUDIES

### A. Multi-scale Features

Multi-scale features are typically employed for the tasks for object detection and instance segmentation, wherein they play

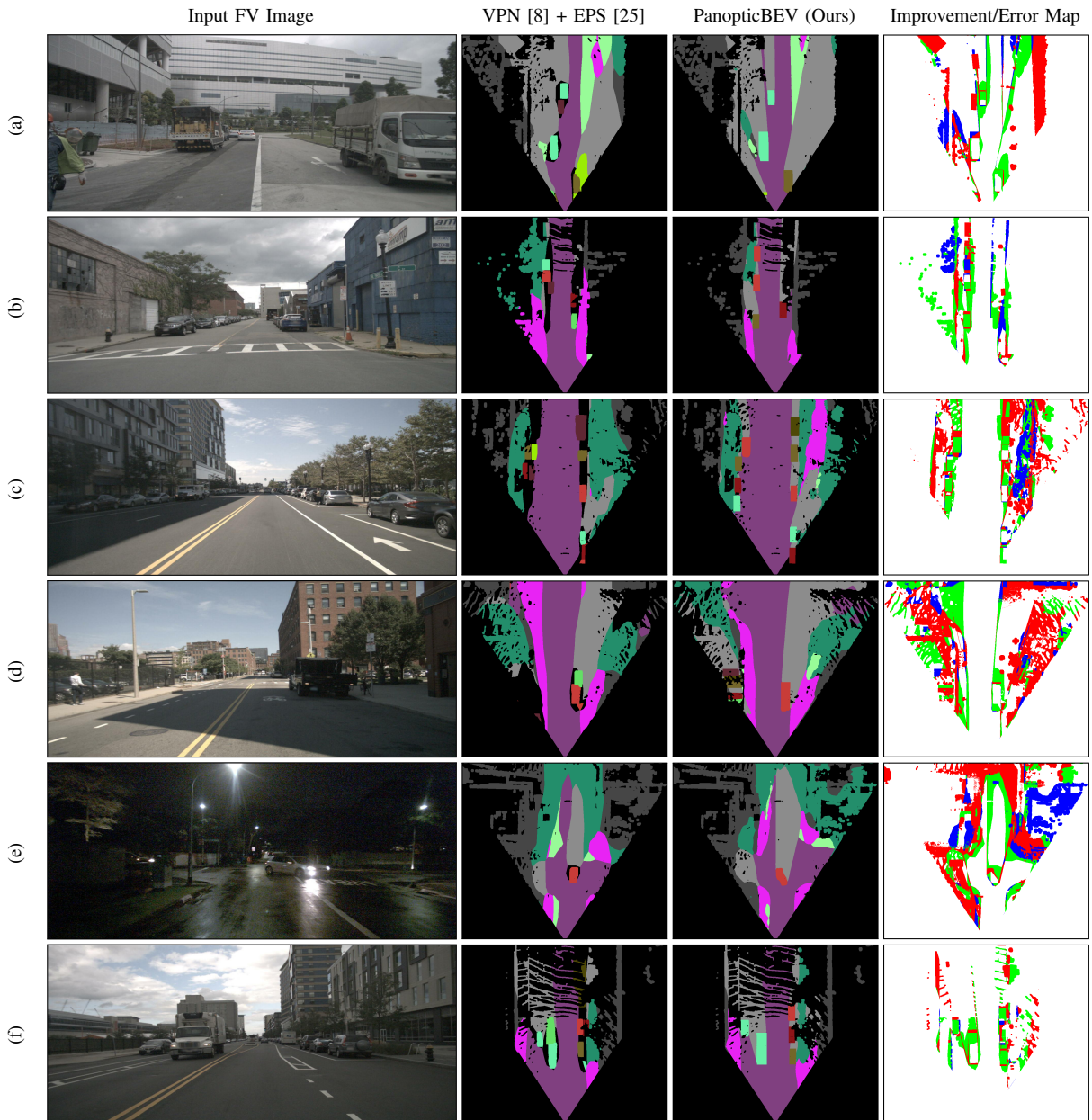| | Input FV Image | VPN [8] + EPS [25] | PanopticBEV (Ours) | Improvement/Error Map |



Fig. S.6: Additional qualitative results comparing the performance of our PanopticBEV model with the best performing baseline on the nuScenes dataset. The rightmost column shows the Improvement/Error map which depicts the pixels misclassified by the baseline but correctly predicted by the PanopticBEV model in green, pixels misclassified by the PanopticBEV model but correctly predicted by the baseline in blue, and pixels misclassified by both models in red.

| Model | $\mathcal{E}_{32}$ | $\mathcal{E}_{16}$ | $\mathcal{E}_8$ | $\mathcal{E}_4$ | PQ | SQ | RQ | mIoU$_{\text{sem}}$ |
|-------|------|------|-----|-----|-------|-------|-------|--------|
| M1 | ✓ | - | - | - | 17.70 | 60.00 | 26.24 | 30.57 |
| M2 | ✓ | ✓ | - | - | 18.88 | **65.36** | 27.86 | 30.47 |
| M3 | ✓ | ✓ | ✓ | - | 20.82 | 63.78 | 30.44 | 31.53 |
| M4 | ✓ | ✓ | ✓ | ✓ | **21.23** | 63.89 | **31.23** | **32.14** |

TABLE S.2: Ablation study on using features from different scales in our PanopticBEV model. The results are reported on the KITTI-360 dataset.

a crucial role in detecting and segmenting objects of different sizes in the image. Since panoptic segmentation encompasses instance segmentation, which in turn encompasses object detection, we hypothesize that multi-scale feature maps are crucial for achieving good panoptic segmentation performance. We validate our hypothesis by performing an ablation study on

the influence of multi-scale feature maps on the performance of our model. We begin with a base model consisting of only the smallest feature scale and iteratively add the larger feature scales to it. We report the Panoptic Quality (PQ), Segmentation Quality (SQ), Recognition Quality (RQ) as well as the semantic mIoU (mIoU$_{\text{sem}}$) for each model. Tab. S.2 presents the results of this ablation study. We observe that model M1 consisting of only the smallest feature scale $\mathcal{E}_{32}$ performs the worst in terms of the PQ metric, achieving a score of only 17.70%. The low PQ score is a consequence of the low RQ score which can be attributed to the poor object detection and instance segmentation performance of this single-scale model. Upon adding $\mathcal{E}_{16}$, we observe a notable $1.18\,\text{pp}$ improvement in the PQ score, most of which can be attributed to a similar increase
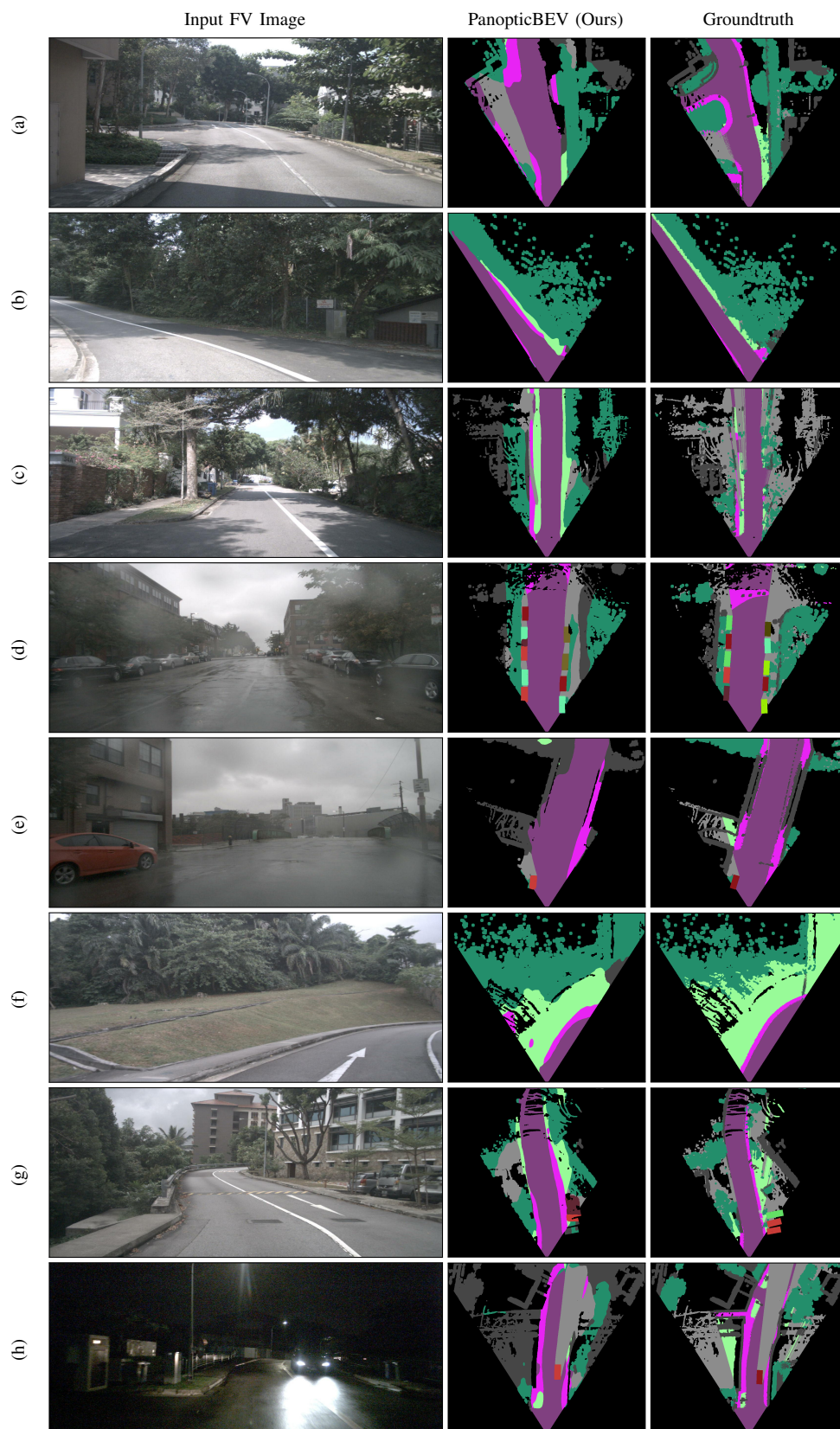
Fig. S.7: Qualitative results showing the performance of our PanopticBEV model in regions with bumpy roads and sudden inclination changes.
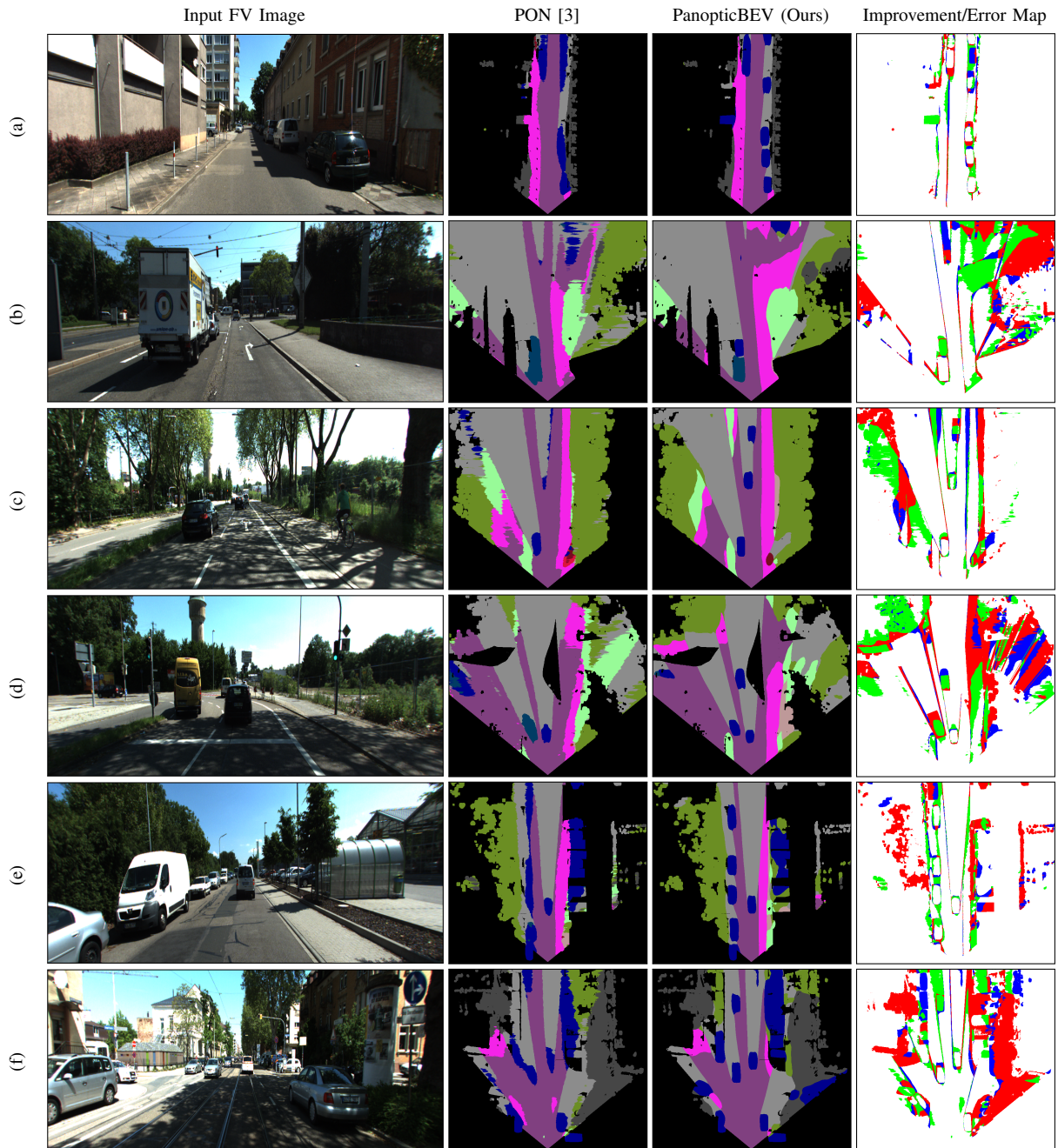
Fig. S.8: Qualitative comparison of BEV semantic segmentation with the best performing previous state-of-the-art model on the KITTI-360 dataset. The rightmost column shows the Improvement/Error map which depicts the pixels misclassified by the previous state-of-the-art but correctly predicted by the PanopticBEV model in green, pixels misclassified by the PanopticBEV model but correctly predicted by the previous state-of-the-art in blue, and pixels misclassified by both models in red.

in the RQ score. We observe a similar trend upon adding the $\mathcal{E}_8$ and $\mathcal{E}_4$ feature scales with the PQ value increasing by $1.94\,\mathrm{pp}$ and $0.41\,\mathrm{pp}$ respectively. This consistent increase in the PQ score upon adding the different feature scales, largely driven by a corresponding increase in the RQ score, indicates that multi-scale features play a crucial role in the object detection performance in the BEV space. Furthermore, we observe from Tab. S.2 that the semantic segmentation performance of the model follows a similar trend and increases from $30.57\%$ when using only $\mathcal{E}_{32}$ to $32.14\%$ when using all the four feature scales. This increase in performance can be attributed to the

presence of both semantically rich small-scale features as well as contextually-rich large-scale feature maps in the model M4. We can thus conclude that the use of multi-scale features improves the performance of our model both in terms of the PQ metric as well as the semantic mIoU score.

### S.4. ADDITIONAL QUALITATIVE RESULTS

#### A. Panoptic Segmentation

We qualitatively evaluate the performance of our proposed PanopticBEV model in comparison to the best performing baseline VPN [8] + EPS [25] on both the KITTI-360 and
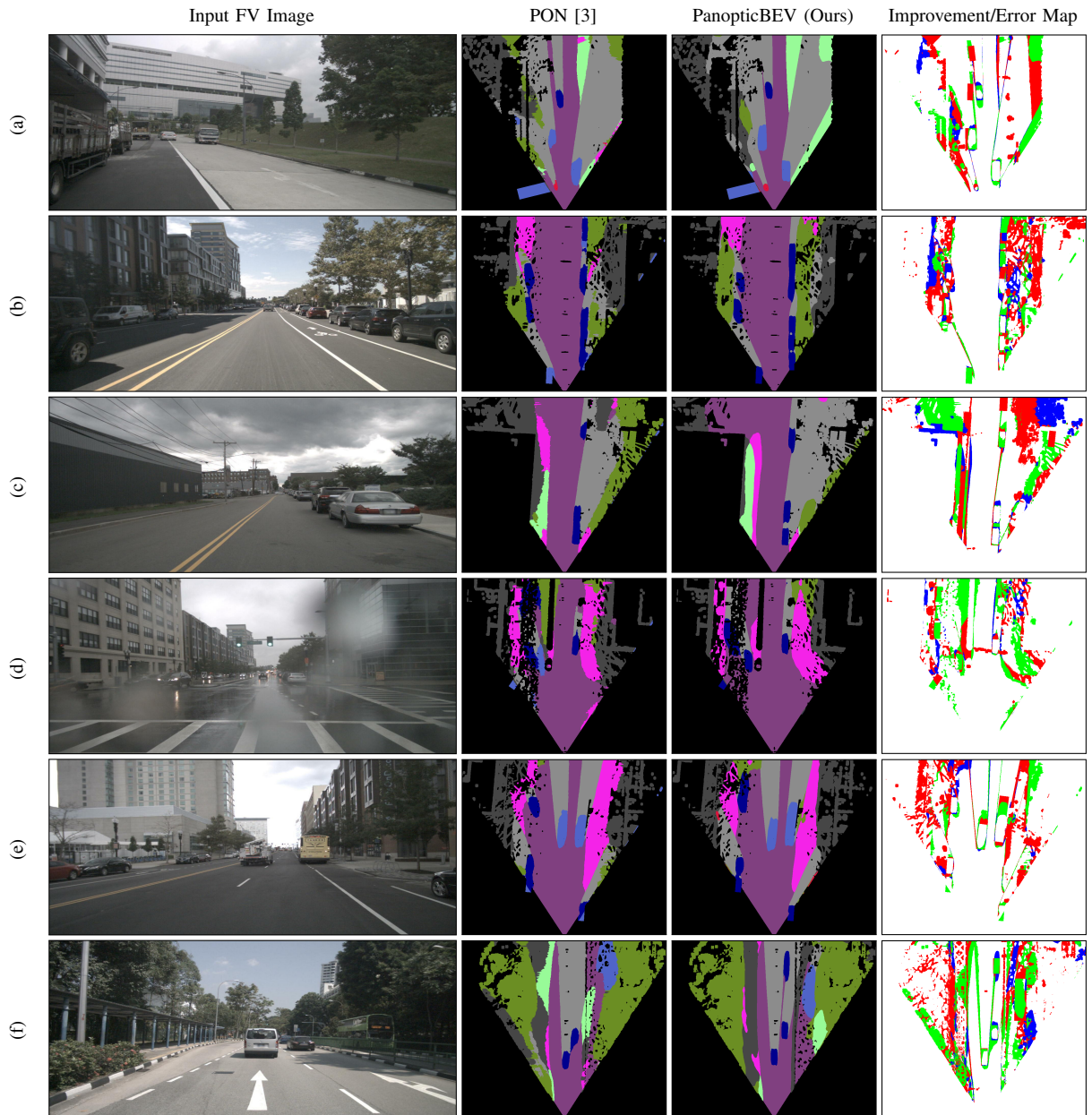
Fig. S.9: Qualitative comparison of BEV semantic segmentation with the best performing previous state-of-the-art model on the nuScenes dataset. The rightmost column shows the Improvement/Error map which depicts the pixels misclassified by the previous state-of-the-art but correctly predicted by the PanopticBEV model in green, pixels misclassified by the PanopticBEV model but correctly predicted by the previous state-of-the-art in blue, and pixels misclassified by both models in red.

nuScenes datasets. Fig. S.5 and Fig. S.6 present the qualitative comparison for the KITTI-360 and nuScenes datasets respectively. We observe in both Fig. S.5 and Fig. S.6 that our model consistent performs better than the VPN + EPS baseline over a wide range of traffic scenarios and environmental conditions. In Fig. S.5(a) and Fig. S.5(c), we observe that our PanopticBEV model accurately estimates the spatial location as well as the semantic category of the bus while the baseline fails to do so in both scenarios. Furthermore, we observe in Fig. S.5(e) and Fig. S.5(f) that our PanopticBEV model accurately estimates the number of traffic participants even in crowded scenarios, while VPN + EPS often fails to detect the instances in the scene and also often hallucinates objects. Furthermore, the predictions of

*stuff* classes such as road and sidewalk is often consistent in the outputs of both the approaches, with PanopticBEV having slightly better and sharper results compared to the VPN + EPS model. We make similar observations in Fig. S.6 in which our PanopticBEV model consistently outperforms the VPN + EPS baseline for both the *stuff* and *thing* classes. For example, in Fig. S.6(a), Fig. S.6(b), Fig. S.6(c), Fig. S.6(d), and Fig. S.6(f) our PanopticBEV model reliably segments most of the vehicles in the scene, even when they are at large distances and are occluded. Similar to the observation made earlier, the VPN + EPS baseline often splits an instance into multiple instances segments and also fails to detect many object instances in the scene. In Fig. S.6(e), we observe that

our PanopticBEV model generalizes effectively to night-time scenes which can be seen in the accurate segmentation of the white car. Whereas, the VPN + EPS model fails to do so and predicts an incorrect orientation for the vehicle. This general trend of our model to accurately detect the number, position, and the orientation of *thing* classes can be specifically attributed to our dense transformer module. Independently processing the vertical and flat regions allows the distinct region-specific transformers to learn vertical and flat cues which helps the network make accurate predictions and significantly improves the overall performance of our model.

### B. Regions with Non-Flat Ground

In this section, we qualitatively evaluate the impact of a non-flat ground on the performance of our approach. Such a situation usually occurs when the road is either bumpy or has sudden changes in elevation (e.g. climbing up / descending down a hill). In such scenarios, using only the IPM algorithm results in an incomplete and distorted projection of flat regions in the BEV space due to the change in the extrinsic transformation between the ground and the camera. We account for this disparity by applying our learnable Error Correction Module (ECM) on regions where the IPM transformation is ambiguous, as well as on regions neglected by the IPM algorithm (e.g. flat regions above the principal point of the camera which contains road segments when the road is bumpy or inclined). We present multiple examples in Fig. S.7 to convey this idea as well as to qualitatively evaluate the performance of our model in such scenarios. We observe that our model accurately predicts the flat regions even when they are bumpy and inclined. For instance, in Fig. S.7(d), the road has a significant upward inclination, but our model is able to accurately predict the road surface. However, due to angle of the road surface, the sidewalk in the distance is extremely difficult to observe and is incorrectly predicted by our model as *road*. In Fig. S.7(f), the road slopes downwards while the grass on the edges slopes upwards. Our model is also able to account for such dynamic changes in the flat regions and accurately predicts both the *road* and *terrain* classes. These results show that our PanopticBEV

model generates accurate panoptic segmentation maps even in the presence of bumpy and hilly terrain, thus enabling its deployment over a wide range of challenging terrains.

### C. Semantic Segmentation

We qualitatively evaluate the performance of BEV semantic segmentation by comparing with the previous state-of-the-art model PON [3]. Fig. S.8 and Fig. S.9 present the results of this comparison for the KITTI-360 and nuScenes datasets respectively. We observe from Fig. S.8(a), Fig. S.8(b), Fig. S.8(e), and Fig. S.8 (f) that our PanopticBEV model efficiently captures the spatial extents of the *thing* classes resulting in well-defined segmentation of object boundaries. Whereas, the predictions from the PON model shows multiple instances of *thing* class fused together in a single blob instead of multiple distinct objects. In Fig. S.8(b) and Fig. S.8(d), we observe that PON incorrectly classifies the truck and the car respectively, while our PanopticBEV model accurately segments them. From the Improvement/Error map shown in the third column of Fig. S.8, we see that the number of green pixels significantly exceed the number of blue pixels indicating that our PanopticBEV model generates more accurate semantic predictions compared to the previous state-of-the-art PON model on the KITTI-360 dataset.

By analyzing the results on the nuScenes dataset shown in Fig. S.9, we observe that our PanopticBEV model consistently outperforms the previous state-of-the-art PON model both in terms of distinguishable vehicle boundaries as well as the accuracy of the semantic class predictions. Fig. S.9(d) shows an example of an extremely complex image wherein most regions of input image are occluded by rain drops. Even in such challenging conditions, our PanopticBEV model yields superior BEV semantic segmentation maps which are more accurate compared the output from the PON model. This can be primarily attributed to our dense transformer module which enables the model to learn consistent and coherent features for both the vertical and flat regions in the image. Moreover, these results demonstrate that our PanopticBEV model is extremely versatile allowing it to be used in a wide range of challenging environmental conditions and traffic situations.