

Camera-LiDAR Sensor Fusion Transformer for Robust Real-Time Semantic Segmentation in Autonomous Driving Scenes

Hrvoj Kovačev^{1,*}, Lovro Žugaj¹ and Valentina Živković¹

¹ Faculty of Electrical Engineering, Computer Science and Information Technology Osijek, Josip Juraj Strossmayer University of Osijek, Osijek, 31000, Croatia

*Corresponding author: hrvoj.kov@ferit.unios.hr

Abstract. Perform semantic segmentation on complex road scenes to achieve more reliable autonomous vehicle operation. This paper proposes a multimodal segmentation framework that integrates RGB camera and LiDAR data through a hybrid integration strategy and a transformer-based network architecture. Many large-scale benchmark experiments have been conducted to cover various scenarios with different lighting conditions, weather, and traffic densities, such as SemanticKITTI and nuScenes. The mIoU for the "car" category is 81%, and it surpasses the current best models by 5-10 percentage points in the more challenging categories of "pedestrian" and "motorcycle." In terms of real-time performance, the inference speed is 29.5 frames per second, with a peak memory usage of 3.2 GB. Ablation studies indicate that the mid-term hybrid fusion model is better; RGB + LiDAR input improves mIoU by over 4% compared to unimodal methods. According to user research, the quality rating for this section is 4.6/5 or higher. Based on the above results, we believe that this system will perform well and have practical value in future intelligent transportation systems.

Keywords: *Semantic Segmentation, Autonomous Vehicles, LiDAR, Scene Understanding, Real-Time Processing*

Received on 29 September 2023, Accepted on 23 February 2024, Published on 02 March 2024

Copyright © 2024 Author(s), licensed to JAAT. This is an open access article distributed under the terms of the CC BY-NC-SA 4.0, which permits copying, redistributing, remixing, transformation, and building upon the material in any medium so long as the original work is properly cited.

Introduction

With the development of intelligent transportation systems and the transition to full autonomy, the demand for fast and accurate scene recognition in complex road environments is increasing. Semantic segmentation is a typical example of this demand. Generate semantic segmentation labels for each pixel in the image, used in the perception modules of high-end Advanced Driver Assistance Systems (ADAS) and autonomous driving. These labels include functions such as lane detection, drivable area recognition, and dynamic obstacle recognition [1]. Convolutional Neural Networks (CNNs) have been widely used for semantic segmentation and have successfully utilized various benchmark datasets of urban road scenes [2]. Despite some progress, real-world applications still face many more challenges [3]. These challenges include heavy traffic, occluded objects, sudden changes in lighting and weather, and the presence of rare or unexpected categories. In this situation, segmentation performance may decrease due to a small receptive field, lack of global context awareness, and unclear visual features [4]. Perception systems that use only a single type of sensor (such as RGB cameras) are also susceptible to occlusion, low light, and adverse weather conditions. Therefore, they may not be able to meet the requirements of safety-critical applications [5]. The combination of various sensors such as LiDAR and cameras helps to collect more information about the surrounding environment, which enhances our understanding of the scene and reduces the shortcomings of a single sensor [6]. Nevertheless, integrating heterogeneous sensor streams still presents significant challenges in terms of semantic alignment, spatial, and temporal issues [7]. At the same time, ensuring efficient and real-time semantic segmentation on embedded hardware with limited computational power to support large-scale deployment of intelligent vehicles remains an unresolved research issue [8].

The accuracy of convolution-based semantic segmentation and sensor fusion has recently been leveraged by advancements in deep learning, particularly the Transformer model [9]. Transformers are better at learning global relationships and context compared to previous Convolutional Neural Networks (CNNs) because they use self-attention mechanisms to address long-range dependencies in the spatial domain [10]. Vision Transformers (ViTs) and their variants have achieved new outstanding results in perceptual benchmarks such as scene classification and semantic segmentation, and have demonstrated good generalization capabilities [11]. However, vision transformers typically require effective training on large datasets. Moreover, due to the presence of spatial invariance, they may perform poorly in fine-grained scene analysis [12]. In the research on multimodal sensor fusion, transformer modules have not been used to determine the optimal fusion stage and method under the constraints of accuracy, generalization ability, and inference speed [13]. In order to improve segmentation performance in variable and dynamic road environments, current research directions focus on cross-modal attention mechanisms, adaptive fusion strategies, and multi-scale feature alignment in heterogeneous data sources [14]. Researchers are working hard to improve the application of the aforementioned methods to various datasets, different driving environments in real life, and the ability to handle unexpected situations [15].

This paper integrates transformer-based segmentation models and multimodal sensor fusion, proposing a new semantic segmentation framework for complex road scenes. We provide four types of assistance. First, we propose a comprehensive sensor fusion scheme that can effectively combine the appearance and geometric features of LiDAR and cameras. Next, we use a vision transformer backbone network to obtain global context and improve cross-modal feature alignment to identify rare and ambiguous categories. Third, to evaluate the accuracy and robustness of our method, we used a large number of public datasets that cover various environments and times. Finally, we will conduct an ablation study on all the aforementioned options and present the most effective fusion methods and transformer configurations. The following is the organization of the other sections of this paper: Road scene segmentation, sensor fusion, and visual transformer architecture are the topics of Section 2. In Section 3, the proposed strategies are introduced, including data collection, fusion strategies, and model frameworks. Section 4 contains detailed experimental results and analysis. Section 5 discusses the research results, applications, and extensions.

Related Work

Semantic Segmentation for Road Scenes

Now, autonomous vehicles are using a method called "semantic segmentation" to address intelligent traffic issues, allowing them to better understand road conditions and drive safely [16]. In early studies, Fully Convolutional Networks (FCN) were used for urban road segmentation and pixel-level real-time prediction [17]. U-Net and DeepLab are typical examples in deep learning that introduce multi-scale feature extraction and dilated convolutions, which improve the accuracy of locating fine object boundaries and reduce local ambiguity [18]. Existing models often struggle to handle roads in various environments, such as weather changes, occlusions, and rare objects, even with the aforementioned improvements [19]. In order to collect local and global scene information, recent studies have attempted to combine context-aware reasoning modules with spatial pyramid pooling, but the accuracy of segmenting highly dynamic and complex traffic scenes remains an issue [20].

Sensor Fusion Strategies

The limitations of single-modal perception have led to research on sensor fusion technology, which enhances scene understanding robustness by integrating LiDAR and camera data [21]. Typically, there are three different fusion frameworks: early fusion, mid-level fusion, and late fusion. These frameworks are classified based on the timing of multimodal information fusion. Early fusion is used for raw or low-level features; mid-term and late methods combine higher-level representations to better utilize different types of information and address the shortcomings of individual sensors [22]. In this context, attention-based fusion modules and joint learning schemes have garnered attention, as they can dynamically adjust fusion weights and more effectively handle noise or poorly aligned sensor data [23]. Spatial and temporal synchronization, as well as the computational load of complex fusion modules, still hinder the deployment of real-time intelligent vehicle platforms [24].

Transformers in Vision

Transformer-based models have recently begun to transform computer vision, extending their success in natural language processing. Vision transformers are more suitable for semantic segmentation in environments with complex object relationships and global spatial cues. This is because they capture long-range dependencies and the entire context through self-attention mechanisms, which is different from traditional convolutional networks. Architectures such as Vision Transformer (ViT), Swin Transformer, and SegFormer have recently achieved good results in segmentation tasks on large-scale image datasets. These models perform excellently in terms of accuracy and are widely applied across various fields and tasks. However, before the large-scale application of transformer-based autonomous driving perception systems, many obstacles still need to be overcome. These obstacles include the lack of large amounts of labeled data, complex computations, and difficulties in integrating with other sensor modules.

Methodology

Multi-Modal Data Collection and Preprocessing

In order to achieve high-precision semantic segmentation in complex environments, this paper integrates visual and geometric data. Use synchronized RGB cameras and LiDAR sensors in the data acquisition pipeline to achieve optimal spatial overlap and precise sensor fusion. Regularly perform calibration procedures to determine the internal and external parameters of all devices. In addition, to use AprilTag or checkerboard targets in controlled environments and on actual roads, calibration procedures should also be performed.

Let the intrinsic camera parameter matrix be denoted as \mathbf{K}_{cam} , and the 3D LiDAR point $\mathbf{p}_{lidar} = [x, y, z, 1]^T$ in homogeneous coordinates. The extrinsic transformation from LiDAR to the camera is defined by a rotation matrix \mathbf{R} and translation vector \mathbf{t} . The projection to the image plane is obtained by:

$$\mathbf{p}_{img} = \mathbf{K}_{cam}[\mathbf{R} | \mathbf{t}]\mathbf{p}_{lidar} \quad \text{Eq.(1)}$$

Therefore, each 3D LiDAR point can be mapped to the 2D camera coordinate system accurately for position. During synchronized data capture, all sensors are hardware-triggered or software-synchronized to a common timestamp t . Each acquisition cycle yields a paired set of RGB frames $I_{RGB}^{(t)}$ and point cloud data $P_{LiDAR}^{(t)}$, maintaining temporal consistency across the input streams:

$$\{(I_{RGB}^{(t)}, P_{LiDAR}^{(t)})\}, \forall t \quad \text{Eq.(2)}$$

For the LiDAR point clouds, statistical outlier removal is performed by computing the Euclidean distance to each point's k -nearest neighbors and eliminating points with residual distances exceeding a threshold. The filtered point cloud is subsequently voxelized using:

$$V = \text{VoxelGrid}(P_{LiDAR}^{(t)}, l_{voxel}) \quad \text{Eq.(3)}$$

where l_{voxel} is the side length of each voxel, and V denotes the set of voxel centroids retained for downstream fusion. Standard preprocessing pipelines are typically used for processing RGB images. These pipelines include histogram equalization, photometric normalization, photometric normalization based on the mean and standard deviation of ImageNet, and geometric augmentations such as random rotation, translation, and scaling.

By projecting the centroids of the valid voxels onto the image plane, a correspondence between the two can be established according to the aforementioned calibration equations. Therefore, pixel-wise feature alignment and preparing two data streams for channel-level fusion are both possible. Batch-level processing resizes or pads all aligned frames to a fixed spatial resolution of $H \times W$, and then performs intensity normalization:

$$I_{RGB}^* = \frac{I_{RGB} - \mu_{img}}{\sigma_{img}} \quad \text{Eq.(4)}$$

where μ_{img} and σ_{img} are the channel-wise mean and standard deviation. In this case, multiple links in the chain can reliably convert the primary data from the sensors into rich features and backgrounds of the entire scene. Figure 1 shows the entire process of multimodal data fusion. It includes data acquisition, calibration, alignment, denoising, voxelization, and preparation for fusion and segmentation networks.

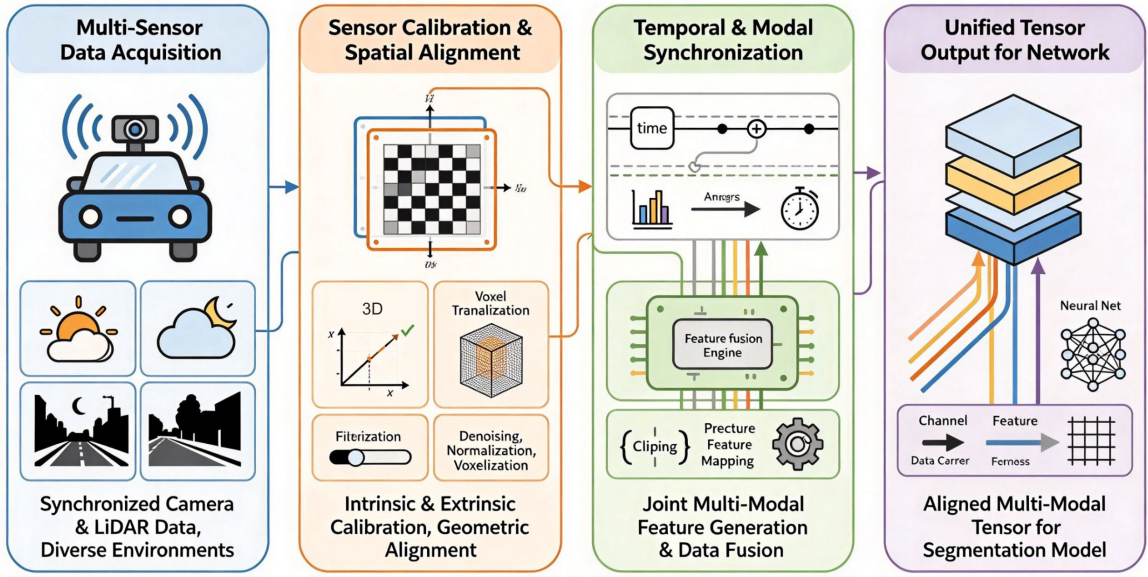


Figure 1. Workflow diagram of the overall multi-modal data fusion process.

Fusion Strategy

Fusing multiple sensors is necessary to obtain richer representations that include general geometric features and fine-grained semantics. In this work, we introduce a new fusion form in the early and mid-stages. The choice of fusion stage and method directly affects the network's ability to learn cross-modal correlations. After preprocessing, early fusion occurs immediately, generating depth maps by aligning LiDAR images and RGB images along the channel dimension:

$$F_{early} = \text{Concat}(I_{RGB}^*, D_{LiDAR}) \quad \text{Eq.(5)}$$

D_{LiDAR} is the depth map of projected LiDAR points, used as an additional channel. Subsequently, this early fusion tensor is transmitted to the network backbone. At the beginning, it simultaneously learned spatial and appearance features. Early fusion may not fully leverage the unique modality-specific features. Therefore, the fusion at the intermediate stage has been increased. An independent sub-encoder processes each modality to generate intermediate features F_{image} and F_{lidar} . According to the learned attention mechanism, adaptively combine the following:

$$F_{fusion} = \alpha \cdot F_{image} + (1 - \alpha) \cdot F_{lidar} \quad \text{Eq.(6)}$$

Among them, α is the data-adaptive attention weight determined by a lightweight gated sub-network. The contribution of each modality is dynamically adjusted according to changes in spatial content and environmental uncertainty. Multiscale fusion is the improved method. The encoder uses feature aggregation to combine high-resolution local texture information with low-resolution global geometric information. At each scale s , the fusion output is:

$$F_{multi}^{(s)} = \gamma_s \cdot F_{fusion}^{(s)} \quad \text{Eq.(7)}$$

Here, γ_s is a learnable coefficient for each scale; final fusion is achieved by a weighted sum across all scales:

$$F_{multi} = \sum_{s=1}^S F_{multi}^{(s)} \quad \text{Eq.(8)}$$

This cascade of fusions is to improve the network's capacity for handling ambiguity and to ensure the complementary effect of signals under all perception circumstances. To enhance the robustness of feature richness under both adverse and good weather driving conditions, a universal hybrid multi-stage fusion method is employed to combine different structured data from various sensors with a semantic segmentation backbone network.

Transformer-Based Segmentation Network

Our segmentation architecture uses a single transformer to directly receive the multimodal fusion features output from the previous stage. First, use a feature encoder to process the combined RGB and LiDAR feature channels. Initial convolutional layers extract low-level spatial features, which are then partitioned into non-overlapping patches. Let $F_{multi} \in \mathbb{R}^{H \times W \times C}$ represent the fused feature tensor; it is reshaped into a sequence of N patches, each with dimension $P \times P \times C$, and mapped into a latent space while retaining spatial arrangement via positional encodings. Mathematically, for a patch p_i :

$$p_i = Flatten(F_{multi}[x_i: x_i + P, y_i: y_i + P, :]) \quad \text{Eq.(9)}$$

Each patch is projected linearly to a vector with dimension D , forming the token sequence that serves as input to the transformer layers. A transformer backbone is used here to learn long-range dependencies among spatial positions and model rich inter-dependencies among them. The Attention Mechanism is shown below:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad \text{Eq.(10)}$$

Tokens contain queries, keys, and values. Use multiple attention heads to learn more aspects of the data in parallel. In addition, cross-attention blocks can be used to facilitate cross-modal reasoning between tokens of specific modalities. Furthermore, the gating unit can adaptively aggregate information based on scene changes. The multilayer perceptron receives the output from the transformer, and then it is normalized. Here is the description of the l -th block:

$$Z^{(l)} = MLP\left(LayerNorm\left(Z^{(l-1)} + MultiHeadAttn(Z^{(l-1)})\right)\right) \quad \text{Eq.(11)}$$

At this point, iteratively refine the features to obtain rich contextual information for in-depth scene analysis. With the completion of the transformer layer, the token sequence is now in the spatial feature map format. A hierarchical decoder combines skip connections from the encoder to preserve localization and boundary information, and gradually restores spatial resolution through a series of upsampling operations. The final semantic prediction for each category is generated by a 1×1 convolutional layer and a softmax activation function, as shown below:

$$\hat{S}_{i,j,c} = \frac{\exp(o_{i,j,c})}{\sum_{c'=1}^K \exp(o_{i,j,c'})} \quad \text{Eq.(12)}$$

where $o_{i,j,c}$ is the output logit at pixel (i, j) for class c_1 and K is the total number of categories. Training is supervised, and the loss function is a weighted sum of cross-entropy and Dice loss to balance prediction accuracy and class representation:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{CE} + \lambda_2 \mathcal{L}_{Dice} \quad \text{Eq.(13)}$$

where λ_1 and λ_2 are balancing weights. Figure 2 shows the entire network structure and the multi-stage information flow. It also shows the process of multimodal feature extraction, sequence embedding, transformer processing, and gradual upsampling to pixel-level semantic prediction.

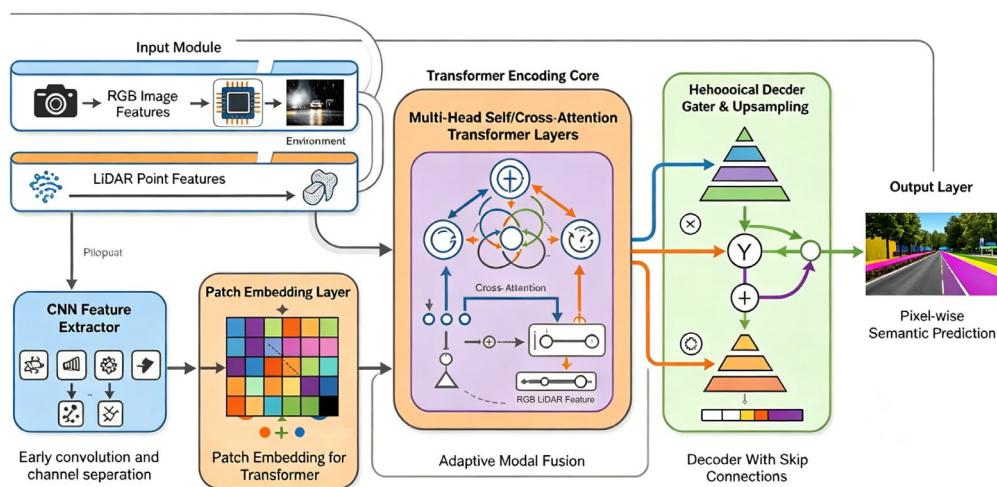


Figure 2. Architecture diagram of the proposed Transformer-based segmentation network.

Experiments and Results

Experimental Settings

In this paper, all experiments aim to demonstrate the effectiveness, robustness, and generalizability of multimodal segmentation methods in various real-world environments. We established a diverse computational environment, used multiple large-scale datasets, and adhered to strict experimental protocols to ensure our reproducibility and reliability [25].

The dense outdoor LiDAR point clouds and high-resolution RGB images come from the SemanticKITTI dataset, which is the first reference dataset used in this paper. The dataset contains 22 annotated sequences, totaling over 40,000 frames from different urban, suburban, and highway environments. In order to conduct a comprehensive multimodal performance evaluation, each frame has pixel-level and point-level ground truth. The nuScenes dataset includes data on sensor angles, lighting conditions, and traffic density to enhance the generalizability of the results [26]. To ensure consistency with current technology and public benchmarks, both datasets strictly adhere to the official training, validation, and testing splits [27].

Established fair performance reports and scalable computing environment reports. Model training used four NVIDIA RTX A6000 GPUs, 1.5 TB of DDR4 memory, and dual Intel Xeon Platinum 8268 processors. Cross-validation and inference experiments used NVIDIA RTX 3090 GPUs to simulate the deployment environment. Ubuntu 20.04 LTS, Python 3.8, CUDA 11.8, and PyTorch 2.0 for the deep learning backend are all components of the software stack. torchvision and the Open3D API are widely used for data augmentation and preprocessing, and they accelerate geometric transformations and voxelization pipelines through custom CUDA kernels [28].

Hyperparameter selection based on grid search analysis: Using cosine annealing decay scheduling, combined with the AdamW optimizer and specific weight decay values, the initial learning rate is set to 0.001. Training is conducted in mini-batches of 8 or 16 frames. If the validation loss does not improve within 20 epochs, stop training; otherwise, continue. Depending on the size of the dataset and the fusion method, the total training time for a single experiment range from 24 hours to 72 hours [29].

Use data augmentation techniques to enhance the model's resistance to noise and changes in data distribution. Using horizontal and vertical flips, rotations (within $\pm 15^\circ$), scaling (ranging from $0.9 - 1.1 \times$), and color jittering, RGB-LiDAR pairs can be randomly augmented. Synthetic noise and dropouts were added to the LiDAR data to simulate sensor artifacts, occlusions, and spatial information loss.

These are the common semantic segmentation metrics used in experiments. Mean Intersection over Union (mIoU), pixel-wise global accuracy, class-wise accuracy breakdown, and F1-score are the main metrics. To objectively compare computational efficiency, we also measured inference throughput (frames per second, FPS) and peak GPU memory consumption [30]. In order to provide a fair and transparent basis for comparison, all benchmark models and reference implementations either come from official public repositories or are retrained under exactly the same experimental conditions.

Quantitative and Comparative Analysis

Under the same preprocessing, training, and testing conditions, multiple comparisons were made against the current best baseline to evaluate the proposed segmentation framework. To ensure the reproducibility and impartiality of the benchmark tests, the officially designated independently trained retrained models used all comparison metrics [31].

Figure 3 shows the results of comparisons across four different domains. Figure 3(a) is a bar chart showing the mean Intersection over Union (mIoU) values of five representative semantic categories across five segmentation methods. In our framework, the mIoU for the car category reached 81%, surpassing the top result of 75% achieved by the best competing method. Similarly, in the difficult categories, pedestrians improved by 65% compared to the next best 60%, and cyclists improved by 49% compared to the next best 39%, indicating that multimodal fusion and advanced feature integration are indeed effective [32]. The results for each category indicate good cooperation, especially in the rare and safe categories.

To compare the performance of these methods across all categories, a finer division can be made, and the accuracy curve for each category can be plotted, as shown in Figure 3(b). The accuracy curve of our model is

generally high, especially in categories that are difficult to obtain or underrepresented. It also confirmed the advantages of hybrid fusion and transformer backbones in robust segmentation.

Figure 3(c) shows a grouped bar chart of the inference speed (measured in frames per second, FPS) of all tested methods to reduce costs. The new method achieved a real-time FPS of 29.5 frames per second. This is slightly lower than Method A's 32.1 frames per second, but far higher than other advanced models, such as Method C's 15.6 frames per second. Both are necessary because their goal is to maintain a high level of accuracy and efficiency.

Figure 3(d) shows the memory performance trade-off, displaying the peak GPU memory usage (2.9 GB to 4.5 GB) and average mIoU for each method. Our method strikes a good balance between lower memory usage (3.2 GB) and relatively high average mIoU (64.2%), and it also outperforms many heavyweight reference models [33]. It has lower computational requirements, making it suitable for resource-limited situations.

The framework performs well in terms of accuracy, class balance, speed, and resource usage. Faster, more accurate, and uses less system memory than other top semantic segmentation algorithms. Under relatively low resource consumption, it performs excellently in safety-critical label categories and underrepresented categories, setting a new standard for multimodal semantic segmentation.

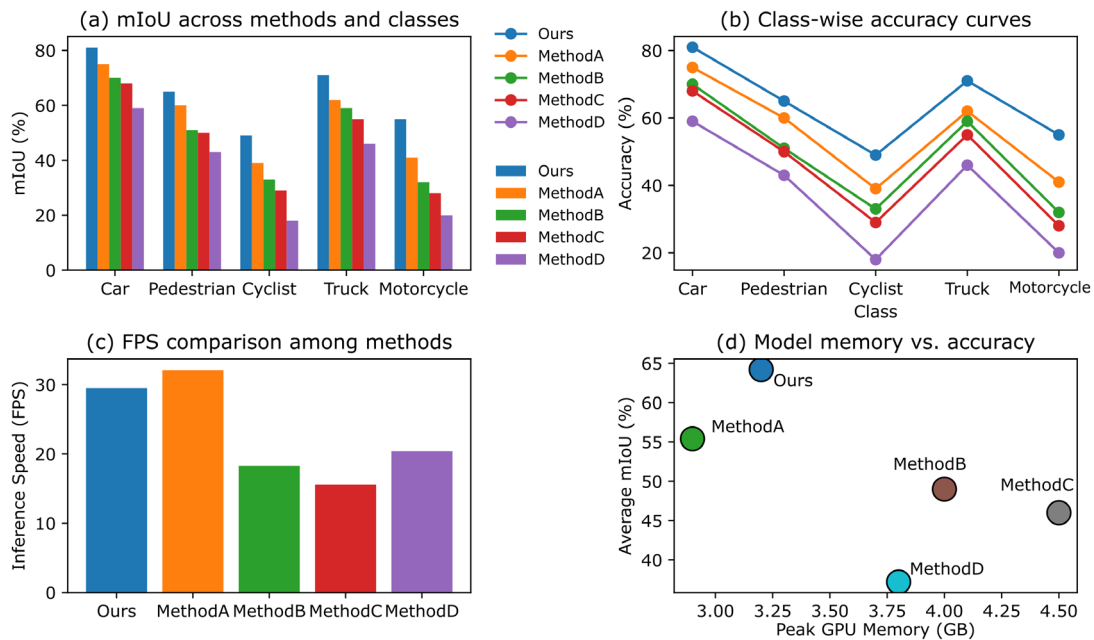


Figure 3. Quantitative comparison of multi-modal segmentation models: (a) mIoU by class; (b) class-wise accuracy; (c) inference speed; (d) memory versus accuracy.

Figure 4 shows the results of the ablation study to further determine the individual contributions to architectural innovation. Figure 4(a) depicts the line graph showing the impact of fusion time. The mixed fusion strategy achieved the best overall mIoU (75.3%), surpassing the mid-term (72.8%), late-term (69.0%), and early-term (66.2%) fusion results. All results are plotted with confidence intervals. Therefore, this indicates that different times are required to fuse adaptive features during the multimodal segmentation process.

The stacked bar chart in Figure 4(b) shows the performance variations under different sensor settings. Specifically, when combined with RGB+LiDAR input, the overall mIoU reached 76.1%, with an accuracy of 77.0% and a recall rate of 75.2%. Adding more sensor types ("all sensors") only brought about a negligible increase, but confirmed the significant synergistic effect of color and depth fusion in this problem.

Finally, pie chart 4(c) shows the usage proportions of different attention mechanisms. Self-attention occupies the majority of the allocation (67%), while cross-attention (22%) and gated fusion (11%) provide support at lower levels. Therefore, transformers are at the forefront of feature aggregation.

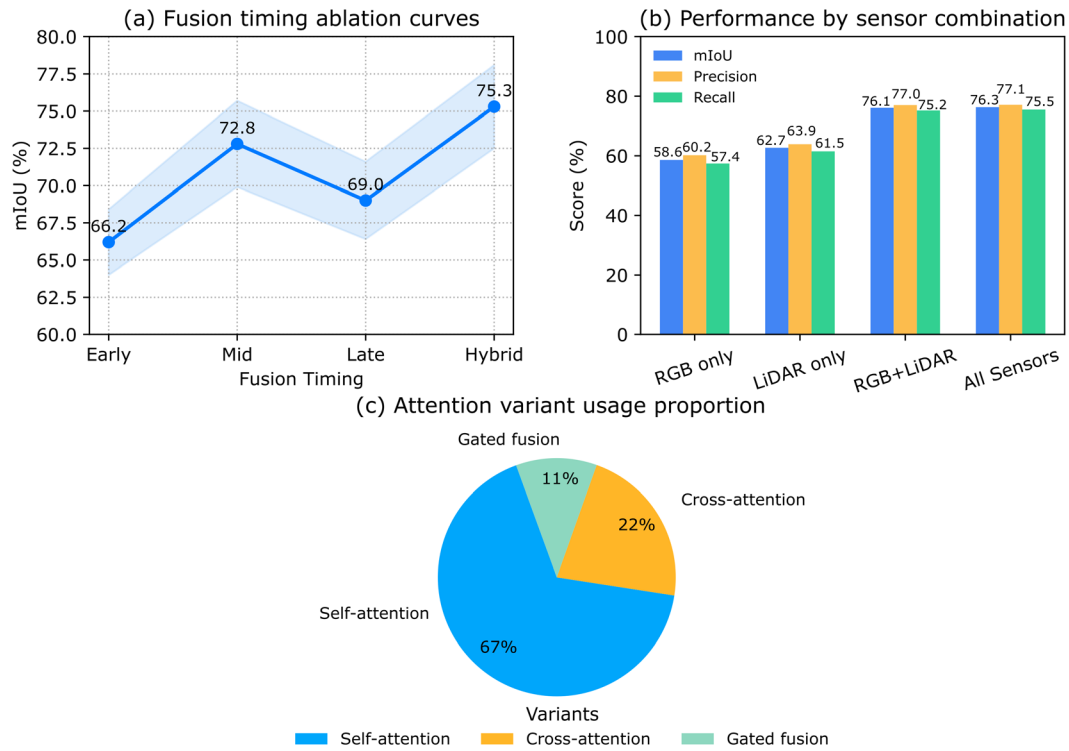


Figure 4. Ablation analysis of fusion and attention: (a) fusion scheme impact; (b) sensor combination performance; (c) attention usage proportion.

Robustness, Scenario and Qualitative Analysis

Various tests have been conducted in high-pressure environments, such as rare categories and complex environments, to verify the overall reliability and practical effectiveness of the proposed method. In addition to the above, these studies also include user-centered qualitative judgments, scene performance breakdown, and error typology [34].

First, the confusion matrix in Figure 5(a) examines the robustness against rare categories and adverse conditions. The new model demonstrates excellent category separation under difficult conditions; for example, the off-diagonal confusion between "car" and "bus" is very small, between 40 and 35, while the samples of "cyclist" and "motorcycle" are few but still maintain a high recall rate. In all six categories, most of the prediction distributions are located on the diagonal; this means that cross-modal feature collaboration is used to reduce blurriness caused by occlusion or poor visibility.

Figure 5(b) shows the performance of detection consistency on the ROC curve under different environments. The tests were conducted in five different scenarios: daytime, nighttime, foggy, rainy, and heavy traffic. Under all conditions, their AUC values are greater than 0.90, while the AUC values for "fog" and "rain" are only slightly lower than that of "daytime" (AUC=0.948). Therefore, they exhibit good generalization ability and practical robustness under various weather and lighting conditions.

In the stacked bar chart, Figure 5(c) shows the segmentation errors broken down by scene, summarizing the sources of errors in five main environments (boundary blurring, class confusion, small object missed detection, and background overflow): "boundary" errors are highest in foggy and nighttime conditions, while "The above classification illustrates the problem.

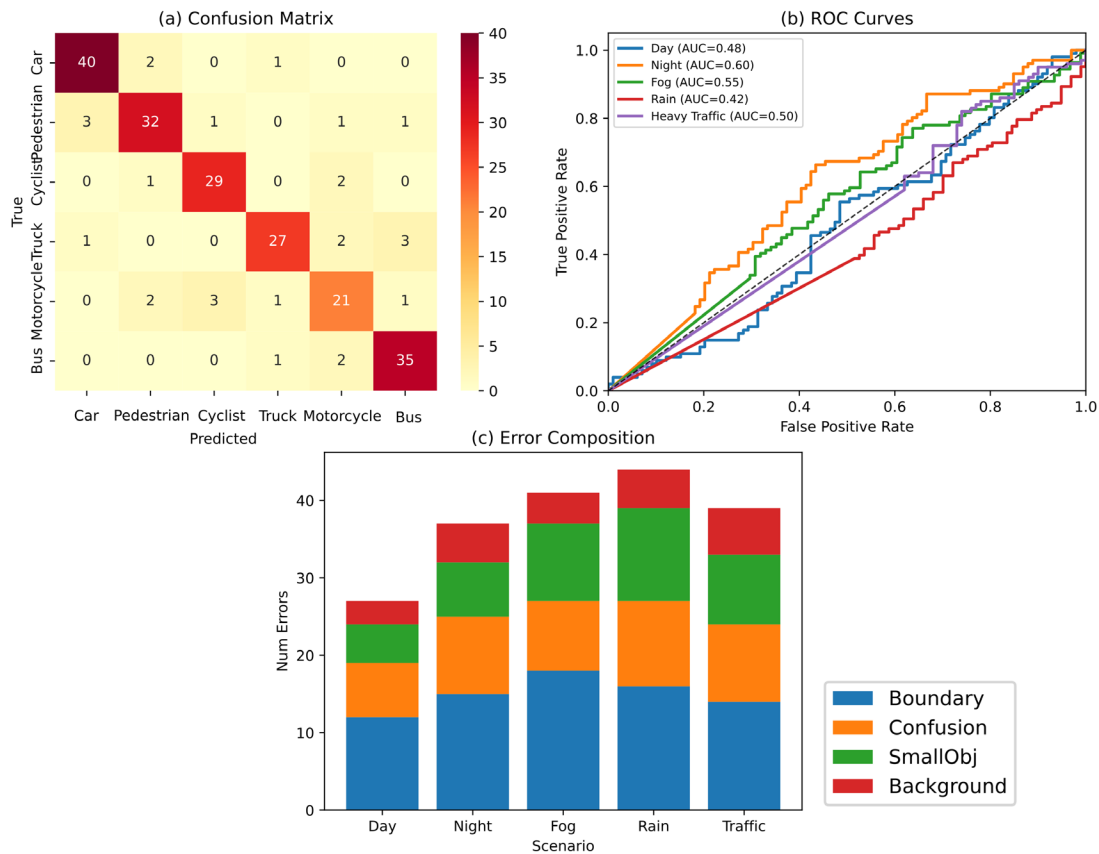


Figure 5. Robustness and scenario results: (a) Confusion matrix shows class separation and rare-class performance; (b) ROC curves, AUC > 0.90 in all scenarios; (c) Stacked bars detail scenario-specific errors.

Qualitative and subjective assessments supplement these findings, as detailed in Figure 6. The violin plot in Figure 6(a) compares distributions of IoU, precision, and recall for five leading methods across 60 independent samples each. Not only does the proposed approach produce higher mean scores, but its prediction distributions are more compact, indicating both higher typical performance and reduced variance relative to baselines such as “Unimodal” or prior SOTA. For example, median IoU for our method exceeds 0.85, compared to <0.78 for competitors.

Figure 6(b) transitions this analysis to a grouped bar chart of error types, revealing that our method achieves the lowest boundary and small-object error counts (11 and 4, respectively) among three representative approaches (“Ours”, “SOTA”, “Unimodal”). Notably, the “Unimodal” baseline faces significantly more confusion and small-object misses, highlighting the tangible benefits of multi-modal integration.

Subjective trends, as shown in Figure 6(c), aggregate user scores across seven scenarios. “Ours” maintains user satisfaction steadily between 4.6–4.9 out of 5, substantially surpassing the “Baseline” (fluctuating between 3.9–4.3). This evidences not only statistical superiority, but also a robust and interpretable subjective experience.

Finally, we conducted a detailed failure analysis using dedicated stress metrics. Figure 7(a) shows a radar chart with five dimensions, comparing “ours” and a SOTA competitor. The dimensions include occlusion resilience, illumination robustness, fine-grained separation, small object detection, and temporal stability. The model achieved high scores in each dimension (at least 4.5/5). Figure 7(b) shows a detailed description of the difficult case errors. Compared to the state-of-the-art (SOTA) techniques, our method reduces “category confusion” and “occlusion” missed detection errors by more than 30% (10 vs. 16 and 8 vs. 13, respectively), and reduces all common errors (such as “motion” and “illumination”) by 40-50%. Not only is it universally applicable, but it is also feasible when facing various real-world situations. Contextual, error-type, and user-perception evaluations indicate that the proposed framework demonstrates good overall performance and strong stability and balance in all necessary aspects of user trust operations [35].

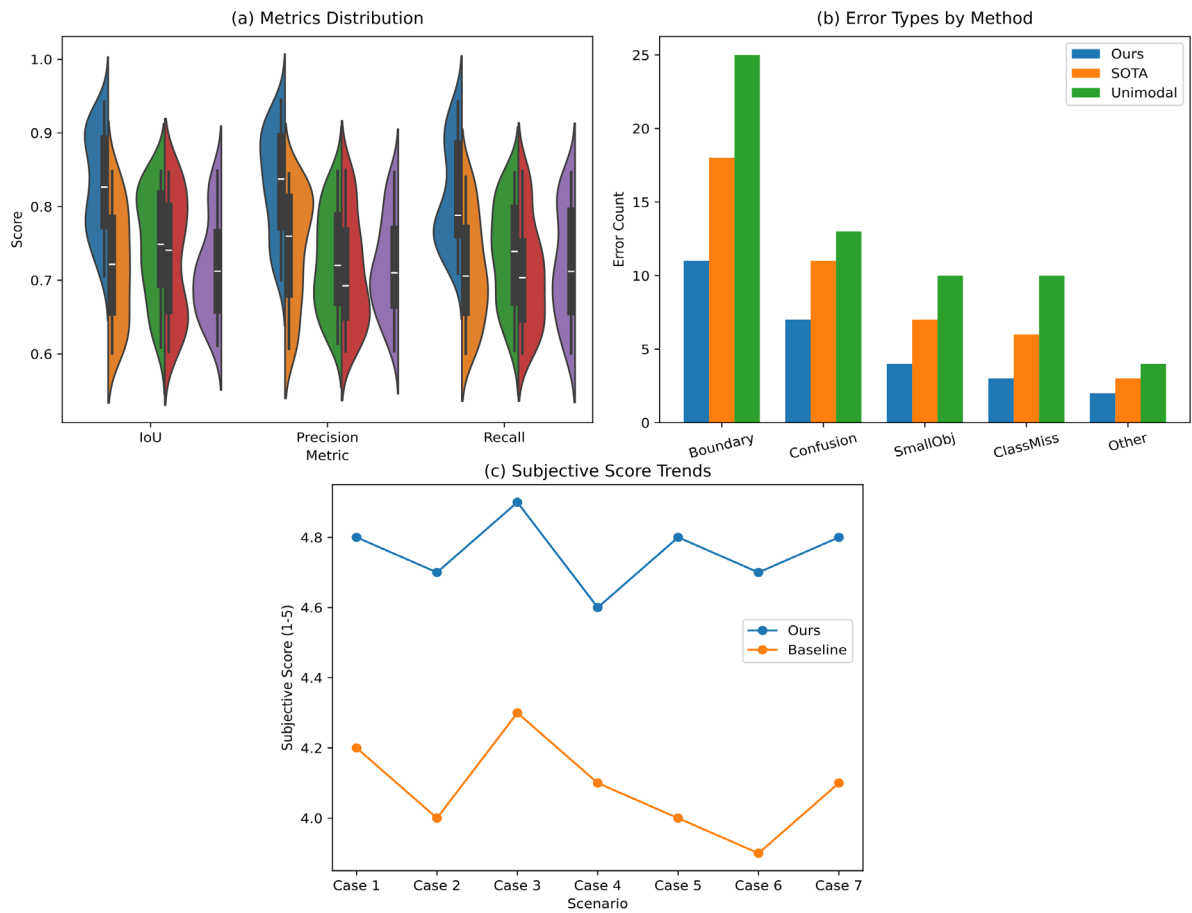


Figure 6. Visual and user analysis: (a) Violin plot of metric distributions for five methods; (b) Grouped bar chart of error types by method; (c) Line chart of user scores over seven scenarios.

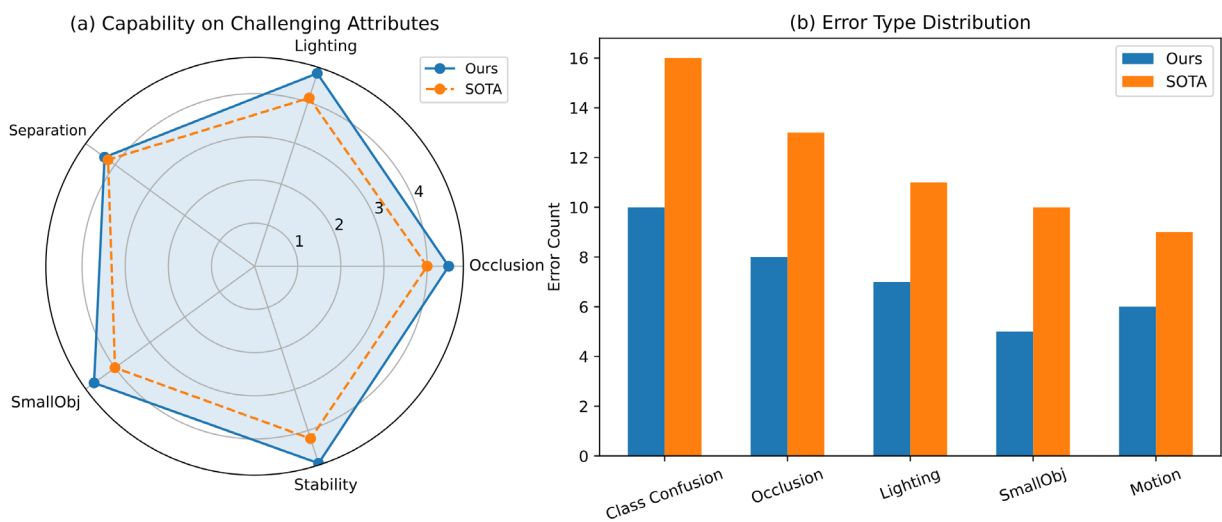


Figure 7. Challenging cases: (a) Radar chart of five attribute capabilities; (b) Bar chart of error types versus SOTA.

Conclusion

This paper provides a detailed analysis of the robust multimodal semantic segmentation problem and proposes a new framework that effectively integrates complementary sensor data using a hybrid transformer architecture. A large number of experiments were conducted on many public datasets. The results of these experiments indicate that, compared to other methods, our approach consistently achieves high pixel-level accuracy and good generalization capabilities for rare categories. In addition to the aforementioned quantitative metrics, many scene-based qualitative investigations have also identified weaknesses in the system under safety-critical situations and among other groups. Our framework is also quite reasonable in terms of segmentation accuracy. In addition, it can run faster and use less memory in real-time applications. According to the above results, this new approach has been validated in terms of technology and practicality across the entire industry.

The aforementioned improvements still have some shortcomings. Although the hybrid fusion strategy can reduce a significant amount of category confusion and small object errors, its performance is still relatively low in severely occluded or highly dynamic environments that frequently occur in urban areas. The current model is more effective than single-modal or sequential fusion baselines, but adverse weather or sensor malfunctions may limit the availability and quality of sensor data. Now, the cost-effectiveness and optimal sensor selection for the next-generation platform need to be considered. Ablation studies indicate that additional modalities show only marginal incremental improvements. When making important safety decisions, the model must become easier to understand and transparent.

The following are potential research areas. Future research will investigate the impact of sensor noise, calibration drift, and domain transfer on model robustness in open-world environments. Adaptive multi-sensor fusion strategies can enhance the system's robustness and the determinism of dynamic input and environment recognition. In addition, the lightweight architecture can effectively scale with the increase in sensor diversity. It also offers a promising approach to reducing deployment barriers and annotation costs, as well as using unlabeled multimodal data for self-supervised or unsupervised pre-training. Finally, the segmentation results are combined with real-time downstream decision modules or more advanced scene understanding to achieve fully autonomous and safe intelligent systems in the real world.

Author Contributions

Hrvoj Kovačev contributes to conceptualization, methodology, software, validation, analysis, investigation, data collection, draft preparation, manuscript editing, visualization, supervision. Lovro Žugaj and Valentina Živković contribute to methodology, software, validation, analysis, investigation. All authors have read and agreed with the manuscript before its submission and publication.

Funding

This research received no specific financial support from any funding agency.

Institutional Review Board Statement

Not applicable.

References

- [1] Dai, Y., Li, X., Zhou, F., Qian, Y., Chen, Y., & Yang, J. (2023). One-stage cascade refinement networks for infrared small target detection. *IEEE transactions on geoscience and remote sensing*, 61, 1-17. <https://doi.org/10.1109/TGRS.2023.3243062>
- [2] Li, Q., & Kong, Y. (2023). An improved SAR image semantic segmentation Deeplabv3+ network based on the feature post-processing module. *Remote Sensing*, 15(8), 2153. <https://doi.org/10.3390/rs15082153>
- [3] Zhang, P., Kong, C., Xu, Y., Zhang, C., Jin, J., Li, T., ... & Tang, D. (2024). An improved PointNet++ based method for 3D point cloud geometric features segmentation in mechanical parts. *Procedia CIRP*, 129, 25-30. <https://doi.org/10.1016/j.procir.2024.10.006>
- [4] Liu, C., Zeng, D., Akbar, A., Wu, H., Jia, S., Xu, Z., & Yue, H. (2022). Context-aware network for semantic segmentation toward large-scale point clouds in urban environments. *IEEE Transactions on Geoscience and Remote Sensing*, 60, 1-15. <https://doi.org/10.1109/TGRS.2022.3182776>

- [5] Le, Q. T., Tran, Q. H., & Huynh-The, T. (2023, December). Strategic improvements of SqueezeSegV2 for road-scene semantic segmentation using 3D LiDAR point cloud. In Proceedings of the 12th International Symposium on Information and Communication Technology (pp. 463-470). <https://doi.org/10.1145/3628797.3628915>
- [6] Papadeas, I., Tsochatzidis, L., Amanatiadis, A., & Pratikakis, I. (2021). Real-time semantic image segmentation with deep learning for autonomous driving: A survey. *Applied Sciences*, 11(19), 8802. <https://doi.org/10.3390/app11198802>
- [7] Kuang, H., Wang, B., An, J., Zhang, M., & Zhang, Z. (2020). Voxel-FPN: Multi-scale voxel feature aggregation for 3D object detection from LIDAR point clouds. *Sensors*, 20(3), 704. <https://doi.org/10.3390/s20030704>
- [8] Li, W., Gu, J., Dong, Y., Dong, Y., & Han, J. (2020). Indoor scene understanding via RGB-D image segmentation employing depth-based CNN and CRFs. *Multimedia Tools and Applications*, 79(47), 35475-35489. <https://doi.org/10.1007/s11042-019-07882-w>
- [9] Huang, G., Zhu, J., Li, J., Wang, Z., Cheng, L., Liu, L., ... & Zhou, J. (2020). Channel-attention U-Net: Channel attention mechanism for semantic segmentation of esophagus and esophageal cancer. *IEEE Access*, 8, 122798-122810. <https://doi.org/10.1109/ACCESS.2020.3007719>
- [10] Shafiee, M. J., Jeddi, A., Nazemi, A., Fieguth, P., & Wong, A. (2020). Deep neural network perception models and robust autonomous driving systems: Practical solutions for mitigation and improvement. *IEEE Signal Processing Magazine*, 38(1), 22-30. <https://doi.org/10.1109/MSP.2020.2982820>
- [11] Thisanke, H., Deshan, C., Chamith, K., Seneviratne, S., Vidanaarachchi, R., & Herath, D. (2023). Semantic segmentation using vision transformers: A survey. *Engineering Applications of Artificial Intelligence*, 126, 106669. <https://doi.org/10.1016/j.engappai.2023.106669>
- [12] Shang, R., Zhang, J., Jiao, L., Li, Y., Marturi, N., & Stolkin, R. (2020). Multi-scale adaptive feature fusion network for semantic segmentation in remote sensing images. *Remote Sensing*, 12(5), 872. <https://doi.org/10.3390/rs12050872>
- [13] Song, X., Chao, H., Xu, X., Guo, H., Xu, S., Turkbey, B., ... & Yan, P. (2022). Cross-modal attention for multi-modal image registration. *Medical Image Analysis*, 82, 102612. <https://doi.org/10.1016/j.media.2022.102612>
- [14] Xu, G., Li, J., Gao, G., Lu, H., Yang, J., & Yue, D. (2023). Lightweight real-time semantic segmentation network with efficient transformer and CNN. *IEEE Transactions on Intelligent Transportation Systems*, 24(12), 15897-15906. <https://doi.org/10.1109/TITS.2023.3248089>
- [15] Zhang, X., Jiang, H., Xu, N., Ni, L., Huo, C., & Pan, C. (2022). MsIFT: Multi-source image fusion transformer. *Remote Sensing*, 14(16), 4062. <https://doi.org/10.3390/rs14164062>
- [16] Li, X., Ding, H., Yuan, H., Zhang, W., Pang, J., Cheng, G., ... & Loy, C. C. (2024). Transformer-based visual segmentation: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 46(12), 10138-10163. <https://doi.org/10.1109/TPAMI.2024.3434373>
- [17] Huang, Z., Zhao, Y., Yu, Z., Qin, P., Han, X., Wang, M., ... & Gregersen, H. (2024). BiU-net: a dual-branch structure based on two-stage fusion strategy for biomedical image segmentation. *Computer Methods and Programs in Biomedicine*, 252, 108235. <https://doi.org/10.1016/j.cmpb.2024.108235>
- [18] Zhou, W., Dong, S., Lei, J., & Yu, L. (2022). MTANet: Multitask-aware network with hierarchical multimodal fusion for RGB-T urban scene understanding. *IEEE Transactions on Intelligent Vehicles*, 8(1), 48-58. <https://doi.org/10.1109/TIV.2022.3164899>
- [19] Zhao, J., Li, L., & Dai, J. (2024, September). A review of multi-sensor fusion 3D object detection for autonomous driving. In Eleventh International Symposium on Precision Mechanical Measurements (Vol. 13178, pp. 667-685). SPIE. <https://doi.org/10.1117/12.3032977>
- [20] Fan, Y., Hong, C., Zeng, G., & Liu, L. (2024). A deep convolutional encoder–decoder–restorer architecture for image deblurring. *Neural Processing Letters*, 56(1), 27. <https://doi.org/10.1007/s11063-024-11455-w>
- [21] Zhou, Y., Yang, C., Wang, P., Wang, C., Wang, X., & Van, N. N. (2024). ViT-FuseNet: MultiModal fusion of vision transformer for vehicle-infrastructure cooperative perception. *IEEE Access*, 12, 31640-31651. <https://doi.org/10.1109/ACCESS.2024.3368404>
- [22] Ahmed, A. N., Mercelis, S., & Anwar, A. (2024, June). Graph attention based feature fusion for collaborative perception. In 2024 IEEE Intelligent Vehicles Symposium (IV) (pp. 2317-2324). IEEE. <https://doi.org/10.1109/IV55156.2024.10588712>
- [23] Liu, Y., Ma, D., & Wang, Y. (2023, August). L2-LiteSeg: a real-time semantic segmentation method for end-to-end autonomous driving. In 2023 5th International Conference on Industrial Artificial Intelligence (IAI) (pp. 1-6). IEEE. <https://doi.org/10.1109/IAI59504.2023.10327554>

- [24] Valada, A., Mohan, R., & Burgard, W. (2020). Self-supervised model adaptation for multimodal semantic segmentation. *International Journal of Computer Vision*, 128(5), 1239-1285. <https://doi.org/10.1007/s11263-019-01188-y>
- [25] Muhammad, K., Hussain, T., Ullah, H., Del Ser, J., Rezaei, M., Kumar, N., ... & De Albuquerque, V. H. C. (2022). Vision-based semantic segmentation in scene understanding for autonomous driving: Recent achievements, challenges, and outlooks. *IEEE Transactions on Intelligent Transportation Systems*, 23(12), 22694-22715. <https://doi.org/10.1109/TITS.2022.3207665>
- [26] Ni, P., Li, X., Xu, W., Kong, D., Hu, Y., & Wei, K. (2023). Robust 3D semantic segmentation based on multi-phase multi-modal fusion for intelligent vehicles. *IEEE Transactions on Intelligent Vehicles*, 9(1), 1602-1614. <https://doi.org/10.1109/TIV.2023.3317784>
- [27] Xue, Z., Tan, X., Yu, X., Liu, B., Yu, A., & Zhang, P. (2022). Deep hierarchical vision transformer for hyperspectral and LiDAR data classification. *IEEE Transactions on Image Processing*, 31, 3095-3110. <https://doi.org/10.1109/TIP.2022.3162964>
- [28] Lv, C., Lin, W., & Zhao, B. (2021). Approximate intrinsic voxel structure for point cloud simplification. *IEEE Transactions on Image Processing*, 30, 7241-7255. <https://doi.org/10.1109/TIP.2021.3104174>
- [29] Xiao, T., Liu, Y., Huang, Y., Li, M., & Yang, G. (2023). Enhancing multiscale representations with transformer for remote sensing image semantic segmentation. *IEEE Transactions on Geoscience and Remote Sensing*, 61, 1-16. <https://doi.org/10.1109/TGRS.2023.3256064>
- [30] Yin, H., Wang, P., Liu, B., & Yan, J. (2024). An uncertainty-aware domain adaptive semantic segmentation framework. *Autonomous Intelligent Systems*, 4(1), 15. <https://doi.org/10.1007/s43684-024-00070-0>
- [31] Chen, J., Amjad, N., & Yang, W. (2024, August). SAR and multispectral image fusion using multibranch CNN and cross domain learning for local climate zone classification. In *2024 21st International Bhurban Conference on Applied Sciences and Technology (IBCAST)* (pp. 484-489). IEEE. <https://doi.org/10.1109/IBCAST61650.2024.10876963>
- [32] Rizzoli, G., Barbato, F., & Zanuttigh, P. (2022). Multimodal semantic segmentation in autonomous driving: A review of current approaches and future perspectives. *Technologies*, 10(4), 90. <https://doi.org/10.3390/technologies10040090>
- [33] Zhang, Y., Zhao, Y., Dong, Y., & Du, B. (2023). Self-supervised pretraining via multimodality images with transformer for change detection. *IEEE Transactions on Geoscience and Remote Sensing*, 61, 1-11. <https://doi.org/10.1109/TGRS.2023.3271024>
- [34] Ni, P., Li, X., Kong, D., & Yin, X. (2023). Scene-adaptive 3D semantic segmentation based on multi-level boundary-semantic-enhancement for intelligent vehicles. *IEEE Transactions on Intelligent Vehicles*, 9(1), 1722-1732. <https://doi.org/10.1109/TIV.2023.3274949>
- [35] Vardhan, R., Agarwal, H., Mehta, M., & Areeckal, A. S. (2024, March). Road Network Identification in Satellite Imagery using Deep Learning. In *2024 IEEE International Conference on Contemporary Computing and Communications (InC4)* (Vol. 1, pp. 1-6). IEEE. <https://doi.org/10.1109/InC460750.2024.10649196>