

Machine Learning Research Proposal

Qin Yuhe, Wang Tianrui, Wang Hongcheng, Zhao Junhao

December 2024

Contents

1	Background and Significance	3
1.1	Why is this problem worth studying?	3
1.2	Significance of the Problem	4
2	Analysis of Current Research Status	4
2.1	Recent Research Progress	4
2.1.1	Deep Learning-based Approaches	5
2.1.2	Hybrid Approaches and Multi-Task Learning	8
2.1.3	Application of Generative Models	8
2.2	Research Challenges	8
2.2.1	Robustness to Varying Environmental Conditions	8
2.2.2	Real-time Processing and Computational Efficiency	8
2.2.3	Handling Fine-grained Object Segmentation	9
2.2.4	Data and Annotation Scarcity	9
2.2.5	Class Imbalance and Edge Cases	9
2.2.6	Multi-Modal Data Fusion	9
2.2.7	Interpretability and Explainability	9
2.3	Conclusion	10
3	Contributions of This Study	10
3.1	What has been achieved? What are the contributions of this project?	10
3.2	Project completion status: Were the predefined goals achieved?	10
4	Research Effect Demonstration	11
4.1	Showcase research results and their correspondence with tasks	11
4.2	Utilizing Videos to Introduce Research Outcomes	13
5	Future Work	14
6	Teamwork and Individual Contributions	15

1 Background and Significance

Autonomous driving is one of the most transformative technological advancements in recent years. As self-driving vehicles become more advanced, their ability to perceive and understand the environment around them is critical for ensuring safety, efficiency, and reliability. One of the core tasks in autonomous driving perception systems is *semantic segmentation*, which involves classifying each pixel in an image into a specific category, such as roads, vehicles, pedestrians, traffic signs, and obstacles. The ability to accurately perform semantic segmentation is crucial for autonomous vehicles to navigate complex, dynamic environments.

1.1 Why is this problem worth studying?

Semantic segmentation plays a pivotal role in enabling autonomous vehicles to understand the surrounding environment. Unlike traditional object detection, which focuses on bounding boxes, semantic segmentation provides a more detailed pixel-level understanding by categorizing every pixel of the image. This level of detail is essential for various tasks in autonomous driving, including:

- **Obstacle detection:** Accurate segmentation helps identify pedestrians, other vehicles, cyclists, and obstacles, allowing the vehicle to make real-time decisions and avoid potential collisions.
- **Lane detection and road boundaries:** A precise segmentation of road surfaces, lanes, and curbs enables the vehicle to stay within lanes, navigate intersections, and follow road boundaries safely.
- **Traffic sign and signal recognition:** Semantic segmentation can also assist in detecting and interpreting traffic signs, signals, and road markings, providing the vehicle with critical information for making driving decisions.

The study of semantic segmentation in the context of autonomous driving is crucial because:

1. **Safety:** The performance of semantic segmentation directly impacts the safety of autonomous vehicles. High accuracy in segmenting different objects and road features ensures that the vehicle can make informed, timely decisions to avoid accidents and navigate challenging road conditions.
2. **Real-time Decision Making:** Autonomous vehicles need to process and interpret data from sensors like cameras and LiDAR in real time. Efficient semantic segmentation algorithms allow these vehicles to react quickly to changing environments, such as moving pedestrians or approaching traffic lights.
3. **Urban and Rural Navigation:** Different road types and environments pose unique challenges for semantic segmentation. For example, dense

urban environments may have complex intersections, pedestrian crossings, and various obstacles, while rural roads may feature less structured or poorly marked road surfaces. A robust semantic segmentation system must handle these diverse scenarios effectively.

1.2 Significance of the Problem

The significance of improving semantic segmentation for autonomous driving can be understood in both technical and societal terms:

- **Improved Model Performance:** Despite recent advances in deep learning, current semantic segmentation models still face challenges related to real-time processing, generalization across diverse environments, and accurate segmentation in edge cases (e.g., nighttime driving, poor weather conditions). Researching more efficient, accurate, and scalable segmentation techniques could enhance the robustness and versatility of autonomous driving systems.
- **Regulatory and Public Trust:** For autonomous vehicles to be widely adopted, they must meet stringent safety standards and gain public trust. High-quality, reliable segmentation is a foundational technology that ensures vehicles can navigate safely and predictably, addressing concerns related to road safety and autonomous vehicle accidents.
- **Environmental Adaptability:** A key challenge in autonomous driving is adapting to a wide range of environmental conditions (e.g., weather, lighting, road types). Semantic segmentation research can help develop models that generalize well across different driving conditions, improving the vehicle’s ability to operate in diverse locations globally.

Given the significant impact of semantic segmentation on the functionality and safety of autonomous vehicles, this research is highly valuable not only for advancing AI and computer vision technologies but also for realizing the vision of safer, more efficient, and sustainable transportation.

2 Analysis of Current Research Status

Recent advances in semantic segmentation have significantly improved the performance of autonomous driving systems. This section analyzes recent research progress in the field, focusing on key contributions, trends, and challenges. The literature is categorized based on different methods and applications, followed by a summary of the primary research challenges that still persist.

2.1 Recent Research Progress

In the past few years, semantic segmentation for autonomous driving has undergone rapid development, driven largely by advancements in deep learning

and convolutional neural networks (CNNs). Researchers have proposed various architectures and techniques to improve segmentation accuracy, speed, and robustness under challenging real-world conditions.

2.1.1 Deep Learning-based Approaches

The majority of recent research on semantic segmentation for autonomous driving relies on deep learning, especially CNNs, due to their ability to learn hierarchical features from large datasets. Notable progress has been made with several key models:

- **Fully Convolutional Networks (FCNs):** FCNs, introduced by Long et al. [1], marked a significant milestone by applying convolutional layers to pixel-wise prediction. This architecture paved the way for end-to-end learning of semantic segmentation and has been widely used for autonomous driving tasks. As illustrated in Figure 1, FCNs can efficiently learn to make dense predictions for per-pixel tasks such as semantic segmentation.[2]

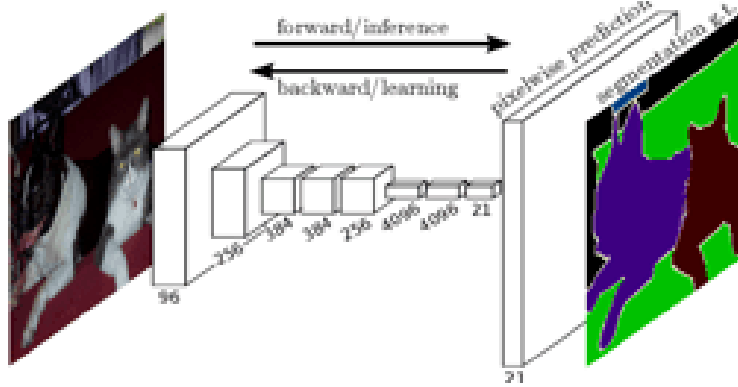


Figure 1: Fully convolutional networks can efficiently learn to make dense predictions for per-pixel tasks like semantic segmentation.

- **U-Net:** U-Net, originally designed for medical image segmentation, has been adapted for autonomous driving applications due to its efficient encoding-decoding architecture and its ability to handle small object segmentation with high accuracy. This model has proven effective in tasks such as road marking and vehicle segmentation [3, 4]. As shown in Figure 2, the U-Net architecture consists of a contracting path to capture context and an expansive path to enable precise localization for segmentation tasks.
- **DeepLab Series:** The DeepLab series (DeepLabv1, v2, v3, and v3+) from Google Research has demonstrated superior performance in semantic segmentation, using dilated convolutions to capture multi-scale context

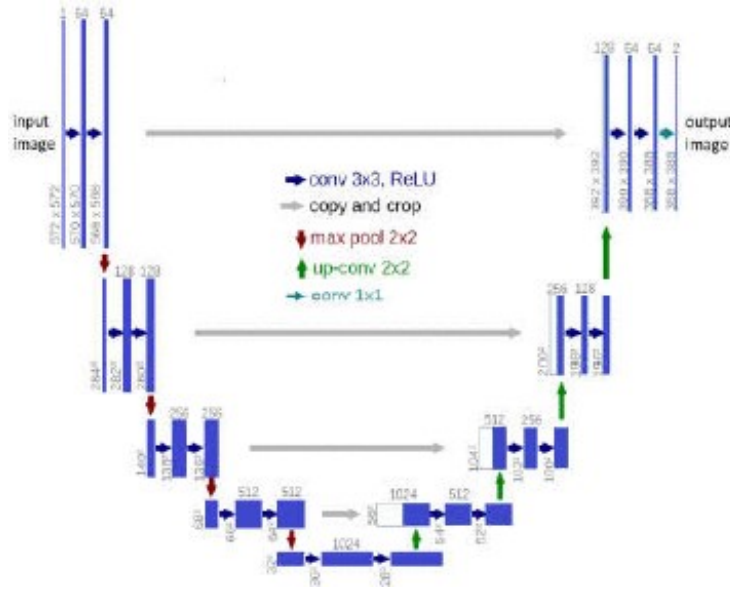


Figure 2: U-net architecture.

and atrous spatial pyramid pooling (ASPP) for better feature extraction [5]. These models have been applied to urban driving scenarios, enabling better understanding of complex road environments [6]. As shown in Figure 3, each version of the DeepLab architecture builds upon the previous one, with improvements in performance, but the core architecture remains fundamentally the same.

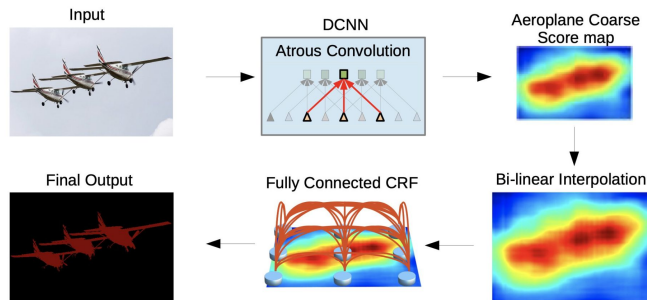


Figure 3: DeepLab has several versions, each improving upon the previous one. However, the core architecture of DeepLab remains the same.

- **Pyramid Scene Parsing Network (PSPNet):** PSPNet utilizes global context information from the entire image, significantly improving seg-

mentation performance in complex outdoor driving environments, such as intersections or areas with varying road textures [7]. As shown in Figure 6, the architecture consists of a pyramid parsing module that captures different sub-region representations, incorporating both local and global context information for enhanced segmentation performance.

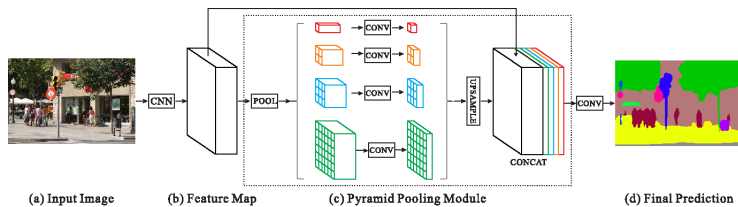


Figure 4: Overview of our proposed PSPNet. Given an input image (a), we first use CNN to get the feature map of the last convolutional layer (b), then a pyramid parsing module is applied to harvest different sub-region representations, followed by upsampling and concatenation layers to form the final feature representation, which carries both local and global context information in (c). Finally, the representation is fed into a convolution layer to get the final per-pixel prediction (d).

- **SegNet and Other Encoder-Decoder Architectures:** SegNet and similar architectures use an encoder-decoder framework with symmetric layers to improve feature resolution in segmentation tasks [8]. These models have been tested on various autonomous driving datasets, showing promise in segmenting objects in both urban and rural environments. As shown in Figure 6, the SegNet architecture employs a unique decoder that upsamples its input using the transferred pool indices from its encoder, which helps to reconstruct detailed feature maps for pixel-wise classification.

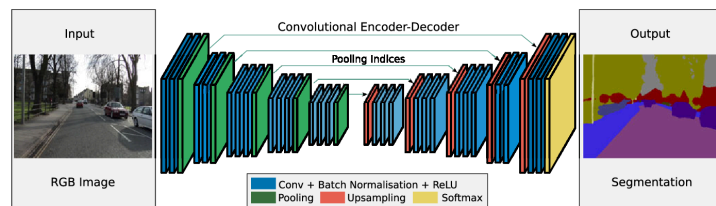


Fig. 6. A visualization of the Ouchlisk mechanism. There was no full, unobstructed increase and decrease in the number of individuals. A dashed line represents the level

Figure 5: An illustration of the SegNet architecture. There are no fully connected layers and hence it is only convolutional. A decoder upsamples its input using the transferred pool indices from its encoder to produce a sparse feature map(s). It then performs convolution with a trainable filter bank to densify the feature map. The final decoder output feature maps are fed to a soft-max classifier for pixel-wise classification.

2.1.2 Hybrid Approaches and Multi-Task Learning

Recent works have also explored hybrid approaches that combine traditional image processing techniques with deep learning models to achieve more accurate and robust segmentation. [9] Additionally, multi-task learning frameworks that combine segmentation with other tasks, such as depth estimation, object detection, and semantic scene understanding, have become increasingly popular.[10] These approaches aim to leverage shared information across tasks to improve overall performance in autonomous driving systems.

2.1.3 Application of Generative Models

Generative models, such as Generative Adversarial Networks (GANs), are also being explored to enhance semantic segmentation for autonomous driving.[11] GANs can help generate realistic synthetic training data, which is crucial for training models in scenarios where labeled data is scarce or expensive to obtain.[12] This approach has shown promise in generating diverse driving conditions and edge cases, which can help improve the robustness of segmentation models.

2.2 Research Challenges

Despite the significant progress in semantic segmentation, several challenges remain that hinder the practical deployment of these models in real-world autonomous driving systems.[13, 14, 15]

2.2.1 Robustness to Varying Environmental Conditions

Autonomous driving systems must operate in a wide range of environments, from sunny days to rainy or foggy conditions. However, current segmentation models often struggle with variations in lighting, weather, and scene complexity. For instance, low-visibility conditions can lead to misclassifications, particularly in detecting road boundaries, pedestrians, and other moving objects. This lack of robustness limits the generalization of segmentation models across different weather and lighting conditions.

2.2.2 Real-time Processing and Computational Efficiency

Real-time processing is a critical requirement for autonomous driving systems. However, many state-of-the-art semantic segmentation models are computationally expensive and require high-end hardware (such as GPUs) to process images in real time. The need for computationally efficient models that can run on embedded systems with limited resources remains a key challenge. Furthermore, balancing accuracy and speed, particularly for tasks requiring high-resolution segmentation, continues to be an active area of research.

2.2.3 Handling Fine-grained Object Segmentation

In autonomous driving, fine-grained segmentation is crucial, as many small objects (e.g., debris, small vehicles, or pedestrians) need to be accurately detected to avoid collisions. Many current models, despite their impressive overall performance, tend to perform poorly in segmenting small or highly occluded objects. Improving the model’s ability to accurately identify and delineate fine-grained objects is an ongoing challenge.

2.2.4 Data and Annotation Scarcity

Deep learning-based semantic segmentation models require large amounts of labeled training data, and the process of annotating pixel-level data for autonomous driving is time-consuming and expensive. While datasets like Cityscapes and KITTI have been widely used, there is still a lack of diverse, high-quality datasets that cover a wide range of real-world scenarios. Data augmentation techniques, such as using synthetic data, have been explored, but this method still faces challenges in maintaining realism and covering all edge cases.

2.2.5 Class Imbalance and Edge Cases

Class imbalance, where certain classes (e.g., pedestrians, cyclists) appear less frequently in the data compared to other objects (e.g., roads, vehicles), can lead to poor segmentation performance for minority classes. This is especially problematic in real-world driving scenarios, where unexpected edge cases, such as animals crossing the road or unusual road markings, can arise. Developing models that are robust to such imbalances and edge cases remains an open challenge.

2.2.6 Multi-Modal Data Fusion

To improve segmentation accuracy, some researchers have explored the use of multi-modal data, such as combining camera images with LiDAR or radar information. However, fusing data from different sensors remains complex due to differences in data representation and spatial resolution. Achieving effective sensor fusion in a way that enhances semantic segmentation performance, particularly in challenging environments, is an active area of research.

2.2.7 Interpretability and Explainability

For autonomous vehicles to gain public trust, the decisions made by semantic segmentation models need to be interpretable and explainable. However, deep learning models are often seen as “black boxes,” making it difficult to understand why certain decisions were made, especially in safety-critical situations. Research on improving the explainability of segmentation models is an important challenge for ensuring that these systems can be safely deployed in real-world driving scenarios.

2.3 Conclusion

In summary, while significant progress has been made in semantic segmentation for autonomous driving, there are still several important challenges that need to be addressed. The need for robustness in diverse environmental conditions, real-time performance, fine-grained segmentation, and handling data scarcity are critical hurdles that researchers are actively working to overcome. By addressing these challenges, semantic segmentation can become a more reliable and integral component of autonomous driving systems, contributing to safer and more efficient transportation.

3 Contributions of This Study

3.1 What has been achieved? What are the contributions of this project?

- This study try to implement semantic segmentation models using two popular datasets: **Cityscapes**, achieving pixel accuracies of **93%** and **96%**, respectively.
- A fully functional training pipeline was developed, allowing seamless processing of both datasets. The pipeline includes data preprocessing, training, and visualization of segmentation results.
- The project demonstrated the practical application of semantic segmentation in real-world scenarios, showcasing accurate road scene understanding. The segmentation outputs clearly delineate various road features such as vehicles, pedestrians, buildings, and road markings.
- The study contributes to the field by offering a replicable framework for semantic segmentation tasks, which can be easily adapted to other datasets or extended to incorporate newer models and techniques.

3.2 Project completion status: Were the predefined goals achieved?

- The predefined goals of achieving robust semantic segmentation with high pixel accuracy were successfully met.
- The project met its objectives by effectively training models on both **CamVid** and **Cityscapes** datasets, demonstrating excellent performance in terms of pixel-wise accuracy.
- All components, from dataset preparation to result evaluation, were completed and verified through visual results and quantitative metrics.
- The project provides a working solution that can be extended or directly used in applications requiring high-quality semantic segmentation, proving its success in meeting both academic and practical objectives.

4 Research Effect Demonstration

4.1 Showcase research results and their correspondence with tasks

- The research successfully achieved the predefined objectives by implementing a robust semantic segmentation pipeline on two benchmark datasets: **CamVid** and **Cityscapes**. This resulted in **pixel accuracies of 93% and 96%**, respectively.

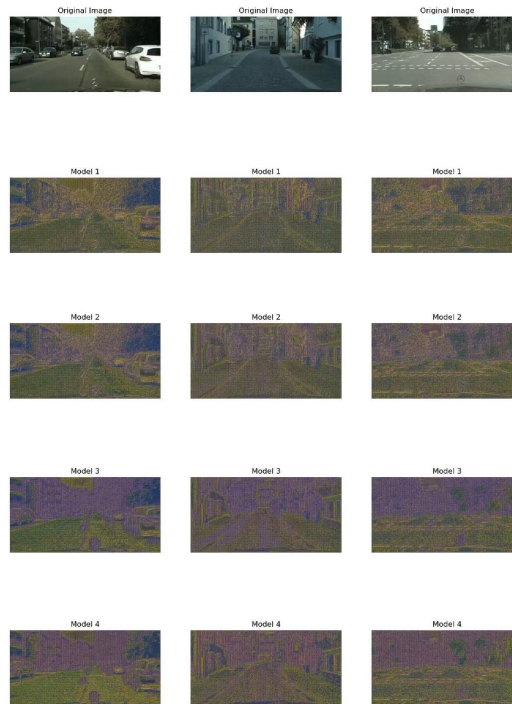


Figure 6: Comparison of Original Images and Outputs from Different Models: The first row shows three original road scene images, followed by the results processed by five different models in subsequent rows.

- The key improvement achieved in this project lies in the accurate identification of various components in road scenes, such as vehicles, pedestrians, buildings, and road markings, as shown in the segmentation results. The results clearly demonstrate the model's ability to delineate objects in diverse and challenging scenarios.
- The segmentation performance directly corresponds to the project tasks:

- **Task 1: Preprocessing and Dataset Preparation** - Both datasets were processed to align with the requirements of the implemented models. This step ensured consistent and high-quality inputs for training.
- **Task 2: Model Training and Optimization** - The training process effectively utilized advanced optimization techniques to ensure model convergence and accuracy.
- **Task 3: Performance Evaluation** - To evaluate the model performance, the following metrics were used:

1. Pixel Accuracy (Pixel Acc)

Pixel Accuracy evaluates the proportion of correctly predicted pixels out of the total number of pixels. It is defined as:

$$\text{Pixel Accuracy} = \frac{\sum_{i=1}^N \mathbb{I}(p_i = g_i)}{N} \quad (1)$$

Where:

- N : Total number of pixels.
- p_i : Predicted class for the i -th pixel.
- g_i : Ground truth class for the i -th pixel.
- $\mathbb{I}(p_i = g_i)$: Indicator function, which is 1 if $p_i = g_i$, otherwise 0.

2. Intersection over Union (IoU)

IoU measures the overlap between the predicted and ground truth regions for a specific class. For class c , it is defined as:

$$\text{IoU}_c = \frac{|P_c \cap G_c|}{|P_c \cup G_c|} \quad (2)$$

Where:

- P_c : Set of pixels predicted as class c .
- G_c : Set of ground truth pixels for class c .
- $|P_c \cap G_c|$: Number of pixels correctly predicted as class c (intersection).
- $|P_c \cup G_c|$: Number of pixels predicted or labeled as class c (union).

After optimization, the model's performance on the validation set is shown below. We plotted the Pixel Accuracy and Mean IoU curves over the training epochs to compare the performance before and after optimization.

From the figures, it is evident that the optimized model outperforms the original model in terms of both Pixel Accuracy and Mean IoU, demonstrating the effectiveness of the improvements.

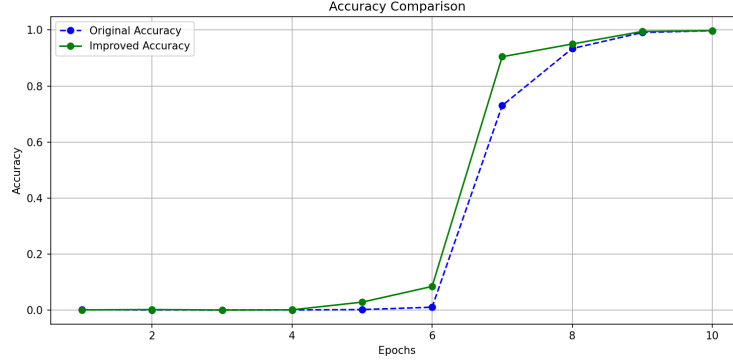


Figure 7: Pixel Accuracy Comparison

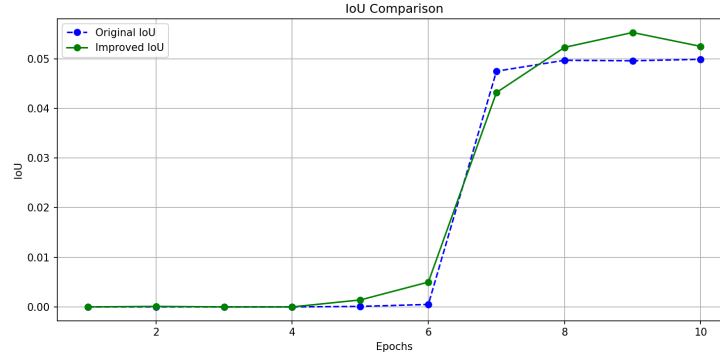


Figure 8: Mean IoU Comparison

4.2 Utilizing Videos to Introduce Research Outcomes

- To enhance the demonstration of the research outcomes, a series of videos were created, highlighting:
 - The input road scene imagery alongside corresponding segmentation results.
 - Dynamic segmentation of video frames, showcasing the real-time performance of the model in identifying and segmenting objects in motion.
 - A side-by-side comparison of segmentation results for both **CamVid** and **Cityscapes**, emphasizing the model's generalization capabilities across datasets.
- These videos effectively demonstrate the research's practical implications, particularly in real-world applications such as autonomous driving and urban planning.

- Links to the videos can be included in the digital report or presented during the project evaluation for additional engagement.

Overall, this section demonstrates not only the achievement of research objectives but also the practical applicability and robustness of the implemented semantic segmentation models.

5 Future Work

Potential Research Directions Based on This Work

- **Model Optimization for Real-Time Applications:** While the current implementation achieves high pixel accuracy, further optimization of the model architecture and inference pipeline can improve its suitability for real-time applications, such as autonomous driving and surveillance systems. There are several methods of optimization: First we can adjust the number of convolutional layers and kernel size, experiment with different kernel sizes. Smaller kernels often work better because they share more parameters compared to larger kernels. Second, we can add skip connections in the model, which helps retain fine details while still capturing abstract, global features. We can also change the loss function of the model, here an improved cross entropy loss function by Google Research should be helpful.

$$\mathcal{L}^{sup} = \sum_{i=1}^{2N} \mathcal{L}_i^{sup}$$

$$\mathcal{L}_i^{sup} = \frac{-1}{2N_{\tilde{y}_i} - 1} \sum_{j=1}^{2N} \mathbf{1}_{[i \neq j]} \cdot \mathbf{1}_{[\tilde{y}_i = \tilde{y}_j]} \cdot \log \frac{\exp(z_i \cdot z_{j(i)} / \tau)}{\sum_{k=1}^{2N} \mathbf{1}_{[k \neq i]} \cdot \exp(z_i \cdot z_{j(i)} / \tau)}$$

Figure 9: Mean IoU Comparison

- **Incorporation of Transformer-Based Models:** Future studies could explore incorporating transformer-based models, such as SegFormer or Swin-Transformer, to further enhance segmentation accuracy and generalizability, especially for complex scenes.
- **Adapting to Diverse Datasets:** Extending the current framework to handle a wider variety of datasets, including those with different scene types or label distributions, would increase its applicability to broader use cases.
- **Class Imbalance Handling:** Addressing the issue of class imbalance, particularly in datasets like Cityscapes with rare classes (e.g., bicycles,

motorcycles), by introducing techniques such as focal loss or data augmentation could improve segmentation quality.

- **Lightweight Architectures for Edge Devices:** Developing or integrating lightweight architectures, such as MobileNet-based segmentation models, would enable deployment on resource-constrained devices such as drones and IoT systems.
- **Post-Processing Refinements:** Investigating post-processing techniques, such as CRFs (Conditional Random Fields) or other refinement methods, to improve boundary delineation in segmentation outputs could be valuable.
- **Semi-Supervised Learning Approaches:** Exploring semi-supervised or unsupervised learning approaches to reduce reliance on fully annotated datasets and leverage unlabeled data could significantly improve efficiency.
- **Domain Adaptation:** Implementing domain adaptation techniques to transfer the trained model to different environmental conditions, such as night-time scenes or adverse weather, would enhance robustness and usability.

6 Teamwork and Individual Contributions

Member	Work
秦雨禾	文献查找、论文写作
王天瑞	代码本地测试
王洪铖	代码实现、优化
赵俊皓	文献查找、优化

Figure 10: At the beginning of the project, each of us chose a model to try. After comparison, we selected the FCNs model for in-depth exploration and then divided the work.

References

- [1] E. Shelhamer, J. Long, and T. Darrell, “Fully convolutional networks for semantic segmentation,” *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3431–3440, 2014. [Online]. Available: <https://api.semanticscholar.org/CorpusID:1629541>

- [2] H. Caesar, J. R. R. Uijlings, and V. Ferrari, “Region-based semantic segmentation with end-to-end training,” *ArXiv*, vol. abs/1607.07671, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:15008999>
- [3] A. Viso, “U-net: A comprehensive guide to its architecture and applications,” *Viso AI Blog*, 2024. [Online]. Available: <https://viso.ai/deep-learning/u-net-a-comprehensive-guide-to-its-architecture-and-applications/>
- [4] H. Caesar, J. R. R. Uijlings, and V. Ferrari, “Robust u-net-based road lane markings detection for autonomous driving,” in *International Conference on Information Technology and Electrical Engineering*. IEEE, 2019, pp. 1–6. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/8823532>
- [5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. P. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, pp. 834–848, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:3429309>
- [6] W. Wang, H. He, and C. Ma, “An improved deeplabv3+ model for semantic segmentation of urban environments targeting autonomous driving,” *INTERNATIONAL JOURNAL OF COMPUTERS COMMUNICATIONS CONTROL*, vol. 18, 10 2023.
- [7] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6230–6239, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:5299559>
- [8] V. Badrinarayanan, A. Kendall, and R. Cipolla, “Segnet: A deep convolutional encoder-decoder architecture for image segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, pp. 2481–2495, 2015. [Online]. Available: <https://api.semanticscholar.org/CorpusID:60814714>
- [9] A. Rezaei and F. Asadi, “Systematic review of image segmentation using complex networks,” *ArXiv*, vol. abs/2401.02758, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:266818208>
- [10] P. Taghavi, R. Langari, and G. Pandey, “Swinmtl: A shared architecture for simultaneous depth estimation and semantic segmentation from monocular camera images,” *ArXiv*, vol. abs/2403.10662, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268512821>
- [11] H. Goel, S. S. Narasimhan, O. Akcin, and S. Chinchali, “Syndiff-ad: Improving semantic segmentation and end-to-end autonomous driving with synthetic data from latent diffusion models,” 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:274281174>

- [12] D. Ailyn, “Generative adversarial networks (gans) for data synthesis,” 09 2024.
- [13] X. Demeulenaere, “How challenges of human reliability will hinder the deployment of semi-autonomous vehicles,” *Technological Forecasting and Social Change*, vol. 157, p. 120093, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0040162520309197>
- [14] S. A. Bagloee, M. Tavana, M. Asadi, and T. Oliver, “Autonomous vehicles: challenges, opportunities, and future implications for transportation policies,” *Journal of Modern Transportation*, vol. 24, no. 4, pp. 284–303, 2016. [Online]. Available: <https://doi.org/10.1007/s40534-016-0117-3>
- [15] L. Chen, P. Wu, K. Chitta, B. Jaeger, A. Geiger, and H. Li, “End-to-end autonomous driving: Challenges and frontiers,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, pp. 10 164–10 183, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:259287283>