# Golden Cudgel Network for Real-Time Semantic Segmentation

Guoyu Yang[1], Yuan Wang[2], Daming Shi[1*], Yanzhong Wang[1]

[1]Shenzhen University [2]Zhejiang University of Technology

gyyang2024@mails.szu.edu.cn, sherlockwang@zjut.edu.cn, dshi@szu.edu.cn

wangyanzhong2022@email.szu.edu.cn

## Abstract

*Recent real-time semantic segmentation models, whether single-branch or multi-branch, achieve good performance and speed. However, their speed is limited by multi-path blocks, and some depend on high-performance teacher models for training. To overcome these issues, we propose Golden Cudgel Network (GCNet). Specifically, GCNet uses vertical multi-convolutions and horizontal multi-paths for training, which are reparameterized into a single convolution for inference, optimizing both performance and speed. This design allows GCNet to self-enlarge during training and self-contract during inference, effectively becoming a "teacher model" without needing external ones. Experimental results show that GCNet outperforms existing state-of-the-art models in terms of performance and speed on the Cityscapes, CamVid, and Pascal VOC 2012 datasets. The code is available at* https://github.com/gyyang23/GCNet.

## 1. Introduction

Semantic segmentation is a critical task in computer vision that classifies each pixel in an image into specific categories. It is vital in fields like autonomous driving [2], medical image analysis [24], and environmental monitoring [37]. While deep learning has significantly improved the performance of semantic segmentation models, they still face challenges with real-time image analysis, limiting their use in applications such as autonomous driving. Therefore, ongoing research and optimization of these algorithms are essential for practical deployment.

To meet the demands of real-time applications, numerous semantic segmentation models have emerged in recent years, emphasizing both performance and inference speed. One of the earliest, ERFNet [25], reduced parameters and computational load by redesigning ResNet [17] blocks using 1D convolutional kernels and skip connec-
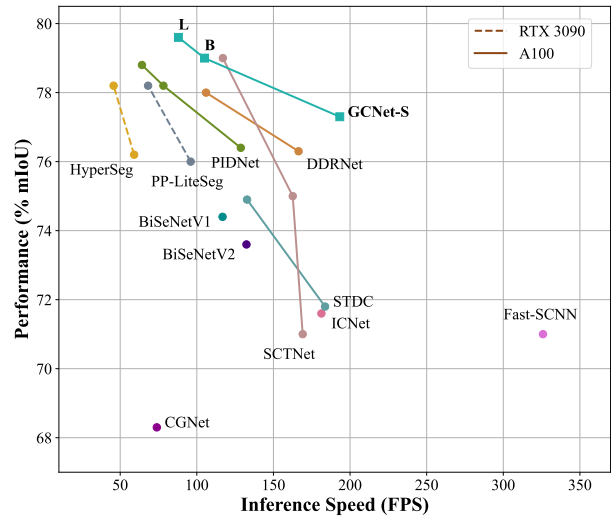


Figure 1. The trade-off between inference speed and performance for real-time semantic segmentation models on the Cityscapes validation set.

tions. However, its single-branch architecture limits its ability to learn spatial information from high-resolution features. In response, BiSeNetV1 [34], BiSeNetV2 [35], and DDRNet [21] introduced two-branch architectures, with one branch focusing on spatial details and the other on deep semantic information. PIDNet [32] expanded this idea with a three-branch design, adding a branch for boundary information. Although multi-branch models have proven effective for enhancing spatial detail capture, SCTNet [33] takes a different approach by proposing a single-branch model that improves semantic representation through knowledge distillation from a high-performance transformer segmentation model [31].

While models mentioned above [21, 32, 33] have achieved impressive speed and performance, their use of Residual Blocks [17] and Conv-Former Blocks [33] can hinder inference speed. As shown in Figure 2, residual connection increase memory access frequency, and the complexity of Conv-Former Block, similar to transformer ar-
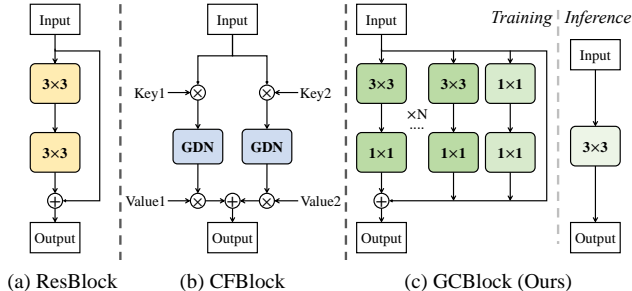
---

*Corresponding author.

Figure 2. A comparison of the proposed GCBlock with multi-path blocks: (a) Residual Block [17], used by model [21, 32, 33]. (b) Conv-Former Block [33], used by model [33]. (c) GCBlock, a block that is scalable in both vertical and horizontal directions.

chitectures, also affects memory access and inference efficiency. In contrast, single-path blocks [28] are better suited for real-time segmentation due to their unidirectional structure, which allows for faster processing and reduced memory overhead. Additionally, although SCTNet uses a single-branch architecture for inference, it relies on a high-performance teacher model during training, which raises costs. The necessity of pre-trained weights from the teacher model on a specific dataset further complicates training and increases resource demands.

Based on the aforementioned description, two key issues currently exist: (1) the decrease in inference speed caused by multi-path blocks and transformer-like convolution structures, and (2) the reliance on high-performance teacher models during training. To address issue (1), we introduce the Golden Cudgel Block (GCBlock), which employs a configuration of vertical multi-convolutions and horizontal multi-paths during training. During inference, these convolutions and paths are reparameterized into a single $3 \times 3$ convolution. This design seeks to simultaneously leverage the training advantages of multi-path blocks, such as avoiding issues of gradient vanishing and explosion, and the inference advantages of single-path blocks. To solve issue (2), we propose the Golden Cudgel Network (GCNet), which stacks multiple GCBlocks. During training, GCNet enlarges itself to enhance its learning capacity, while during inference, it contracts to improve speed. From a broader perspective, GCNet can effectively transform itself into a "teacher model" without the need for an external one.

To evaluate the performance and speed of the proposed GCNet, we conducted experiments on three datasets: Cityscapes [9], CamVid [1], and Pascal VOC 2012 [14]. The results indicate that GCNet achieves a superior balance between segmentation performance and inference speed. As shown in Figure 1, compared to other state-of-the-art real-time semantic segmentation models, GCNet demonstrates enhanced performance and faster speed. The main contributions of this paper are as follows:

- We propose the Golden Cudgel Network family, a model designed to enhance performance through self-enlargement and increase inference speed through self-contraction without any loss in performance.
- We introduce the Golden Cudgel Block, which leverages the training advantages of multi-path blocks and the inference advantages of single-path blocks through reparameterization.
- Experiments conducted on three public datasets demonstrate that the proposed GCNet achieves a superior balance between performance and speed compared to other state-of-the-art real-time semantic segmentation models. In particular, without ImageNet pre-training, GCNet reaches 77.3% mIoU and 193.3 FPS on the Cityscapes.

## 2. Related Work

### 2.1. High-performance Semantic Segmentation

High-performance semantic segmentation models primarily focus on segmentation quality and accuracy, often at the expense of inference speed. FCN [26], one of the earliest deep learning-based models, achieved pixel-level segmentation by replacing fully connected layers in CNNs with convolutional layers. However, its heavy downsampling results in a loss of spatial detail. To mitigate this, the DeepLab series [3–5] integrated dilated convolutions [36] with varying rates to enhance the receptive field while preserving spatial resolution. Recently, transformer-based segmentation models [6, 7, 31] have emerged, leveraging self-attention to capture long-range dependencies and improve performance. Despite their effectiveness, these models remain unoptimized for inference speed, limiting their suitability for real-time applications.

### 2.2. Real-time Semantic Segmentation

Real-time semantic segmentation models aim to simultaneously enhance speed and performance. However, optimizing speed through adjustments to depth, width, or structure can compromise performance, necessitating strategies to balance both. Architecturally, these models fall into two categories: single-branch and multi-branch architectures.

**Single-branch Architecture:** As single-branch architecture models, they [22, 25, 30] typically balance performance and speed by optimizing backbone blocks or designing lightweight segmentation heads. However, with the rise of multi-branch models that offer enhanced performance and speed, single-branch models face challenges. Despite this, numerous studies continues to advance single-branch models. For example, STDC [15] addresses the speed slow addition of branches in BiSeNet [34] by introducing the Detail Aggregation Module, which retains spatial information through a single-branch approach. SCTNet [33] utilizes knowledge distillation with SegFormer [31] as the teacher
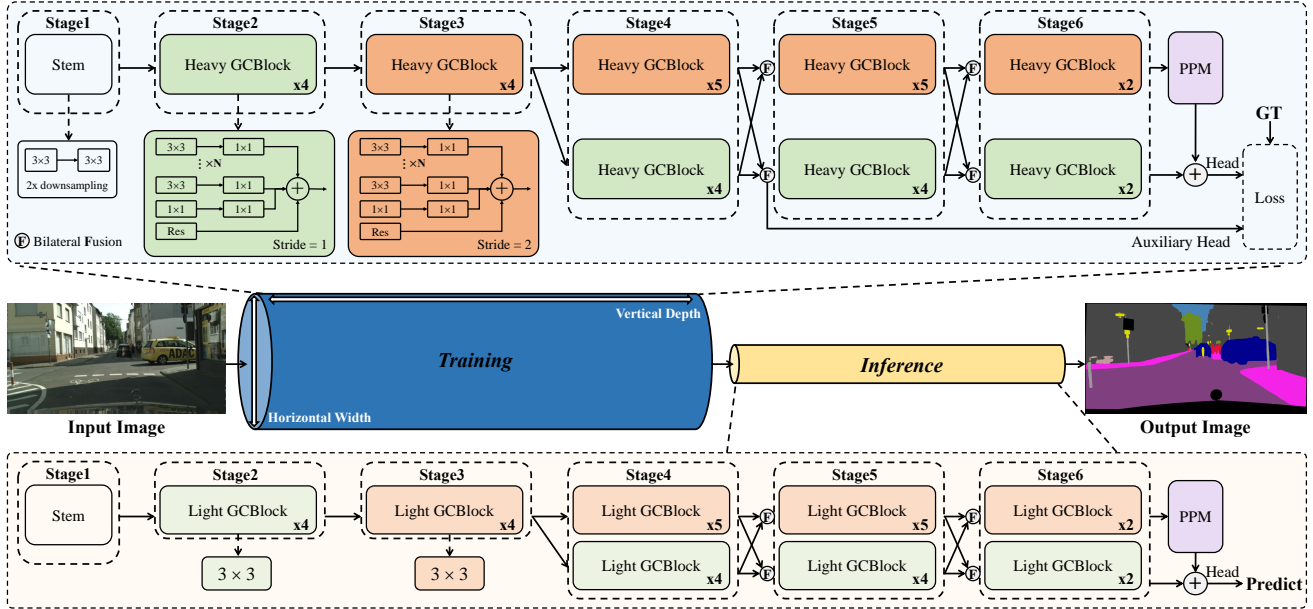
Figure 3. The overall architecture of GCNet. After feature flow into two branches, the upper branch corresponds to the semantic branch, while the lower branch corresponds to the detail branch. The orange box indicates that the first block within the GCBlocks has a stride of 2, while the remaining blocks have a stride of 1. The green box signifies that all GCBlocks maintain a stride of 1. PPM refers to the Deep Aggregation Pyramid Pooling Module [21].

model, employing a Conv-Former Block to bridge the semantic gap between CNN and transformer features, allowing for more effective learning of rich semantic information.

**Multi-branch Architecture:** Single-branch architecture models use skip connections to link feature maps between the backbone and head, preserving spatial details. In contrast, multi-branch models add branches to learn both spatial details and boundary information. BiSeNetV1 [34] and V2 [35] propose a two-branch architecture where one branch focuses on deep semantic information and the other on spatial details, with no weight sharing between them. Conversely, models like Fast-SCNN [23] and DDRNet [21] share some low-level backbone weights. Fast-SCNN extracts shallow features and splits them into two branches for feature retention and global feature extraction, while DDR-Net performs multiple bilateral fusions to efficiently merge information. PIDNet [32] introduces a three-branch architecture to capture spatial detail, deep semantic information, and boundary information.

Whether single-branch models or multi-branch models, although they continually compete in architectural design to optimize performance and speed, their inference speeds are limited by multi-path blocks, as shown in Figure 2. In contrast, our GCNet, based on a two-branch architecture, enhances its learning capability during training by vertically increasing convolutions and horizontally expanding paths. During inference, it improves inference speed by reparameterizing multi-convolution multi-path blocks into a single

convolution, all without sacrificing performance. Although some studies [11–13] have employed similar methods to enhance model performance and speed in image classification tasks, such approaches have been seldom explored in semantic segmentation tasks.

## 3. Methodology

In this section, we first present the overall framework of GCNet. Then, we describe the training structure of the GCBlock and its reparameterization into a single convolution during inference. Finally, we detail the supervisory strategy and loss function used in model training.

### 3.1. Golden Cudgel Network

To address the issue of increased training costs associated with the use of teacher models, we propose the Golden Cudgel Network (GCNet). We establish a two-branch architecture as baseline model and subsequently enhance it by deepening the model with additional convolution and broadening it with multiple paths, thereby augmenting its learning capacity. After training, we reduce the model's convolution and the number of paths to improve inference speed, a process that does not compromise performance. From a macro perspective, the self-enlarged GCNet functions as the "teacher model", while the self-contracted GCNet serves as the "student model".

As shown in Figure 3, GCNet is a two-branch archi-

Table 1. Details of the different versions of GCNet. H represents height, W represents width, and C represents channels. For GCNet-S, C is set to 32, while for GCNet-M and GCNet-L, C is set to 64. In rows $S_4$, $S_5$, and $S_6$, the left side of each column indicates the dimensions of the semantic branch, while the right side indicates the dimensions of the detail branch.

| Stage | Output Dimension | Quantity of Blocks | |
|---|---|---|---|
| | | GCNet-S/M | GCNet-L |
| Input | H × W × 3 | | |
| $S_1$ | H/4 × W/4 × C | | |
| $S_2$ | H/4 × W/4 × C | 4 | 5 |
| $S_3$ | H/8 × W/8 × 2C | 4 | 5 |
| $S_4$ | H/16 × W/16 × 4C, H/8 × W/8 × 2C | 5, 4 | 5, 5 |
| $S_5$ | H/32 × W/32 × 8C, H/8 × W/8 × 2C | 5, 4 | 5, 5 |
| $S_6$ | H/64 × W/64 × 16C, H/8 × W/8 × 4C | 2, 2 | 3, 3 |

tecture model, with its backbone composed of stacked GCBlocks. After extracting shallow features (stage 3), the features flow into two branches: the semantic branch and the detail branch. The semantic branch is designed to learn deep semantic information, while the detail branch focuses on learning spatial detail information. After passing through stage 4 or stage 5, the features from the semantic branch are fused with those from the detail branch after channel compression and upsampling. Conversely, the features from the detail branch undergo channel expansion and downsampling before being integrated with the features from the semantic branch. This process enhances the richness of the features from both branches. We do not utilize various attention modules [16, 18, 29, 32] for feature fusion, as this would complicate the model and decrease inference speed. Instead, we utilize 3 × 3 convolutions for channel compression and expansion, along with bilinear interpolation for upsampling. After stage 6, the features from the semantic branch are processed through a pyramid pooling module [21] to generate richer feature representations, which are then upsampled and fused with the features from the detail branch before being passed to the segmentation head for prediction. The segmentation head of GCNet consists of a 3 × 3 convolution followed by a 1 × 1 convolution. The 3 × 3 convolution is used to integrate the features from both branches while adjusting the channel dimension to $O_c$, and the 1 × 1 convolution is used to align the number of classes. For instance, if the number of classes is 19, the 1 × 1 convolution will adjust $O_c$ to 19. Additionally, we implement an auxiliary segmentation head during training to further enhance GCNet's performance.

As illustrated in Figure 3, the backbone of GCNet consists of six stages. Stage 1 serves as the Stem, comprising two 3 × 3 convolutions with a stride of 2, which is intended for rapid downsampling of the image. The sub-

sequent stages are made up of stacked GCBlocks. Notably, the orange box in the figure indicates that the first block within the GCBlocks has a stride of 2, while the remaining blocks have a stride of 1; the green box denotes that all GCBlocks maintain a stride of 1. To accommodate various application scenarios, we adjusted both the quantity of stacked GCBlocks and the convolution width of the GCBlocks at each stage, resulting in three distinct versions of GCNet: GCNet-S, GCNet-M, and GCNet-L, as detailed in Table 1. The OC of the segmentation heads for GCNet-S, GCNet-M, and GCNet-L are 64, 128, and 256, respectively.

### 3.2. Golden Cudgel Block

**Structure.** To address speed degradation from the complex structures of multi-path and transformer-like convolution blocks, we propose the Golden Cudgel Block (GCBlock). During training, GCBlock expands into a multi-convolution, multi-path structure to leverage the benefits of these blocks. During inference, it simplifies to a single convolution through reparameterization, enhancing efficiency. Specifically, based on the bottleneck structure [17], we achieve vertical reparameterization of convolutions into a single 3 × 3 convolution by removing activation functions between convolutional layers, keeping only one at the output. During our investigation, we found that after training, the reduction in model parameter values hindered the lossless fusion between the first 1 × 1 convolution and the 3 × 3 convolution within the bottleneck. Thus, we eliminated the first 1 × 1 convolution, retaining one 3 × 3 convolution and one 1 × 1 convolution. Additionally, we introduced Path$_{1\times1\_1\times1}$ (consisting of two 1 × 1 convolutions) and more Path$_{3\times3\_1\times1}$ (comprising one 3 × 3 and one 1 × 1 convolution) to enhance learning capability. Although GCBlock contains multiple convolutions and paths, these components integrate into a single 3 × 3 convolution during inference, ensuring both performance and speed. The structure of GCBlock are illustrated in Figure 3, where N indicating the number of Path$_{3\times3\_1\times1}$. We will detail the process of vertically reparameterizing each path into a single 3 × 3 convolution and then summing them horizontally.

**Conv-Bn to Conv.** To enhance the stability of training, Batch Normalization (BN) is typically applied after convolution to perform channel normalization and linear scaling. Before reparameterizing between convolutions, it is necessary to merge the convolution and BN. Let $\mu$, $\sigma$, $\gamma$, and $\beta \in \mathbb{R}^{C_{out}}$, represent the mean, variance, scaling factor, and bias in BN, respectively, and let $W \in \mathbb{R}^{C_{out} \times C_{in} \times k \times k}$ and $B \in \mathbb{R}^{C_{out}}$ denote the weight and bias of a k × k convolution. $C_{in}$ represents the number of input channels, $C_{out}$ and represents the number of output channels. The merging of the convolution and BN can be computed as follows:

$$W' = \frac{\gamma}{\sqrt{\sigma + \varepsilon}}W, B' = \frac{(B - \mu)\gamma}{\sqrt{\sigma + \varepsilon}} + \beta, \qquad (1)$$

4

where $\varepsilon$ is a constant term used in BN to prevent division by zero, and it is typically set to 1e-5.

**Path$_{3\times3\_1\times1}$ to Conv$_{3\times3}$.** For a concatenation consisting of a $3 \times 3$ convolution followed by a $1 \times 1$ convolution, if there are no batch normalization layers or activation functions in between, they can be merged into a single $3 \times 3$ convolution. Let the weight and bias of the $3 \times 3$ convolution be denoted as $W_{3\times3} \in \mathbb{R}^{C_{out1}\times(C_{in}\times3\times3)}$ and $B_{3\times3} \in \mathbb{R}^{C_{out1}}$, respectively, and the weight and bias of the $1 \times 1$ convolution be denoted as $W_{1\times1} \in \mathbb{R}^{C_{out2}\times(C_{out1}\times1\times1)}$ and $B_{1\times1} \in \mathbb{R}^{C_{out2}}$. Given an input $x$, the output $y$ can be computed as follows:

$$
\begin{aligned}
y &= W_{1\times1} * (W_{3\times3} * x + B_{3\times3}) + B_{1\times1} \\
&= W_{1\times1} \cdot I_2(W_{3\times3} \cdot I_1(x) + B_{3\times3}) + B_{1\times1} \\
&= W_{1\times1} \cdot I_2(W_{3\times3} \cdot I_1(x)) + W_{1\times1} \cdot I_2(B_{3\times3}) + B_{1\times1} \\
&= W_{1\times1} \cdot I_2(W_{3\times3} \cdot I_1(x)) + W_{1\times1} \cdot B_{3\times3} + B_{1\times1},
\end{aligned}
\tag{2}
$$

where $*$ denotes the convolution operation, while $\cdot$ represents matrix multiplication. To perform matrix operations, $I(\cdot)$ denotes the im2col operator, which is utilized to transform the input $x$ into a two-dimensional matrix based on the shape of the convolution weight. For instance, $I(\cdot)$ will convert the input $x$ with dimensions $C_{in}\times H\times W$ into a matrix of shape $(C_{in} \times 3 \times 3) \times (H' \times W')$, according to the dimensions of $W_{3\times3}$. Since $B_{3\times3}$ is a constant sequence, $W_{1\times1} \cdot I_2(B_{3\times3})$ is equivalent to $W_{1\times1} \cdot B_{3\times3}$. Furthermore, since the shape of $W_{1\times1}$ is $C_{out2} \times (C_{out1} \times 1 \times 1)$ and the shape of $W_{3\times3} \cdot I_1(x)$ is $C_{out1} \times (H' \times W')$, we can conclude that:

$$
\begin{aligned}
&W_{1\times1} \cdot I_2(W_{3\times3} \cdot I_1(x)) \\
&= W_{1\times1} \cdot W_{3\times3} \cdot I_1(x).
\end{aligned}
\tag{3}
$$

Subsequently, by substituting Equation 3 into Equation 2, the following result is obtained:

$$
\begin{aligned}
y &= W_{1\times1} \cdot I_2(W_{3\times3} \cdot I_1(x)) + W_{1\times1} \cdot B_{3\times3} + B_{1\times1} \\
&= W_{1\times1} \cdot W_{3\times3} \cdot I_1(x) + W_{1\times1} \cdot B_{3\times3} + B_{1\times1} \\
&= (W_{1\times1} \cdot W_{3\times3}) \cdot I_1(x) + (W_{1\times1} \cdot B_{3\times3} + B_{1\times1}) \\
&= W'_{3\times3} \cdot I_1(x) + B'_{3\times3},
\end{aligned}
\tag{4}
$$

where $W'_{3\times3}$ represents the new weight, while $B'_{3\times3}$ represents the new bias. It is evident that a concatenation consisting of a $3 \times 3$ convolution followed by a $1 \times 1$ convolution can be merged into a new $3 \times 3$ convolution without any loss in performance.

**Path$_{1\times1\_1\times1}$ to Conv$_{3\times3}$.** A concatenation of two $1 \times 1$ convolutions can also be merged into a $3 \times 3$ convolution according to Equation 4. The reason is that the first $1 \times 1$ convolution can be viewed as a $3 \times 3$ convolution with a non-zero weight value at the center and zero weight values in the surrounding positions. Notably, to satisfy the conditions in Equation 4, the stride of the first $1 \times 1$ convolution is set to 2 during downsampling, rather than being

applied to the subsequent $1 \times 1$ convolution. Theoretically, the Path$_{1\times1\_1\times1}$ can accommodate an arbitrary number of $1 \times 1$ convolutions; however, experiments have shown that stacking two $1 \times 1$ convolutions yields the best model performance.

**Path$_{residual}$ to Conv$_{3\times3}$.** For residual connection, we first construct a $1 \times 1$ convolution with weight $W_{residual} \in \mathbb{R}^{C_{out}\times(C_{in}\times1\times1)}$ ($C_{out} = C_{in}$) and bias $B_{residual} \in \mathbb{R}^{C_{out}}$ (set to 0) to equivalently replace it. Next, this $1 \times 1$ convolution is transformed into a $3 \times 3$ convolution. Specifically, we traverse $W_{residual}$ using $C_{out}$ and $C_{in}$. When $c_{out}$ equals $c_{in}$ during the traversal, the weight value at that position is set to 1; otherwise, it is set to 0. After completing the traversal, the resulting $1 \times 1$ convolution is equivalent to the residual connection. Since a $1 \times 1$ convolution is a special case of a $3 \times 3$ convolution, it is straightforward to convert the $1 \times 1$ convolution into a $3 \times 3$ convolution. In Path$_{residual}$, we use a BN layer, so it is also necessary to integrate the converted convolution with the BN. Notably, when the stride of the GCBlock is set to 2, the residual connection is not used.

**Multipath to Single Convolution.** After reparameterizing each path into $3 \times 3$ convolutions, parallel sets of $n$ weights $W \in \mathbb{R}^{C_{out}\times(C_{in}\times3\times3)}$ and $n$ biases $B \in \mathbb{R}^{C_{out}}$ are obtained. Let $W_i$ represent the $i$-th weight and $B_i$ the $i$-th bias, where $i$ is less than or equal to $n$. Given an input $x$, the output $y$ from multiple paths is calculated as:

$$
\begin{aligned}
y &= \sum_{i=1}^{n} \left( W_i * x + B_i \right) \\
&= \sum_{i=1}^{n} \left( W_i \cdot I_i(x) + B_i \right) \\
&= \left( \sum_{i=1}^{n} W_i \right) \cdot I(x) + \sum_{i=1}^{n} B_i \\
&= W' * x + B'
\end{aligned}
\tag{5}
$$

where $I_i$ is the $i$-th im2col operator. Since the shape of $W_i$ is consistent, all $I_i$ are the same. This allows for direct summation of these $3 \times 3$ convolutions, enabling the acquisition of a new $3 \times 3$ convolution without performance loss.

### 3.3. Deep Supervision and Loss Function

Previous research [21, 32, 35] has demonstrated that adding an auxiliary segmentation head to appropriate positions in the model during training can improve the model's performance without increasing inference time. Building on prior work, we introduced an auxiliary segmentation head in GC-Net, which is removed during inference. As illustrated in Figure 3, after the bilateral fusion in the stage 4, the features from the detail branch are passed to the auxiliary segmentation head for loss calculation. Notably, the structure of the auxiliary segmentation head is identical to that of the

Table 2. Ablation study on Path$_{1\times1\_1\times1}$ for GCNet-S. "Number" indicates the number of convolutions used in the path, "Memory" refers to the GPU memory utilized during training, and "Time" represents the training duration in hours.

| Number | Memory | Time | mIoU |
|--------|--------|------|------|
| 0 | 20.58 GiB | 4.0 h | 76.1 |
| 1 | 21.87 GiB | 4.5 h | 76.6 |
| 2 | 24.61 GiB | 5.0 h | 76.7 |
| 3 | 27.31 GiB | 5.4 h | 76.4 |

primary segmentation head. Based on the losses from both the segmentation head and the auxiliary segmentation head, the total loss can be expressed as:

$$L = L_{sh} + \alpha L_{ash}, \tag{6}$$

where $L_{sh}$ represents the loss from the segmentation head, $L_{ash}$ denotes the loss from the auxiliary segmentation head, and $\alpha$ is the weight coefficient, which is set to 0.4 in GCNet. To effectively handle imbalanced data and difficult samples, we employed OHEM Cross Entropy [27] as the loss function, consistent with previous work [32].

## 4. Experients

### 4.1. Datasets and Implementation Details

To demonstrate the effectiveness of GCNet, we conducted experiments on three public datasets: Cityscapes [9], CamVid [1], and Pascal VOC 2012 [14]. During training on Cityscapes, we did not utilize ImageNet [10] pretrained weight and instead trained from scratch using four A100 GPUs. For training on CamVid and Pascal VOC 2012, we fine-tuned the model using weight from Cityscapes. During inference, we used a single A100 GPU. More details regarding the datasets, training settings, inference settings, and evaluation metrics can be found in Appendix A.

### 4.2. Ablation Studies

**Number of Convolutions in Path$_{1\times1\_1\times1}$.** To validate the effectiveness of Path$_{1\times1\_1\times1}$, we conducted an ablation study by varying the number of $1 \times 1$ convolutions in this path. Specifically, we established a baseline consisting of a Path$_{3\times3\_1\times1}$ and a Path$_{residual}$, and then incrementally increased the number of $1 \times 1$ convolutions in Path$_{1\times1\_1\times1}$ from 0 to 4. As shown in Table 2, both the memory usage and training time increased with the number of convolutions, calculated based on four A100 GPUs. At 0 convolutions (baseline), the model only achieved 76.1% mIoU. However, as the number of convolutions increased, the mIoU improved, reaching peak performance with two convolutions. Consequently, we set the number of $1 \times 1$ convolutions in this path to 2.
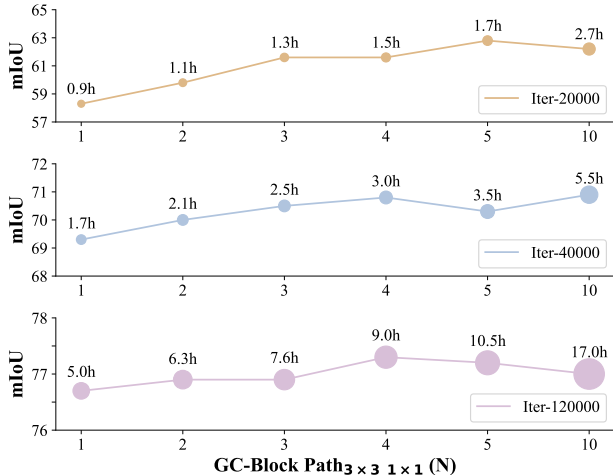


Figure 4. Ablation study on Path$_{3\times3\_1\times1}$ for GCNet-S. "N" indicates the number of Path$_{3\times3\_1\times1}$, while "Iter-20000" signifies that 20000 iterations were completed, and so on.

**More Path$_{3\times3\_1\times1}$.** To validate the effectiveness of Path$_{3\times3\_1\times1}$, we conducted an ablation study by varying the number of this path. The baseline comprised a Path$_{3\times3\_1\times1}$, a Path$_{1\times1\_1\times1}$, and a Path$_{residual}$. Although theoretically the number of Path$_{3\times3\_1\times1}$ can be increased indefinitely, extensive research indicates that overly wide or deep models do not necessarily lead to improved performance; our experiments corroborate this finding. As shown in Figure 4, the model's mIoU begins to increase with the addition of Path$_{3\times3\_1\times1}$, peaking at N=4 or N=5. However, at N=10, the model's performance deteriorates, with training time increasing by over 80% compared to N=4. Additionally, we observed that with fewer iterations, increasing the number of Path$_{3\times3\_1\times1}$ can significantly enhance performance (a 4.5% increase at 20000 iterations). Considering the training cost and model performance, we set differ-
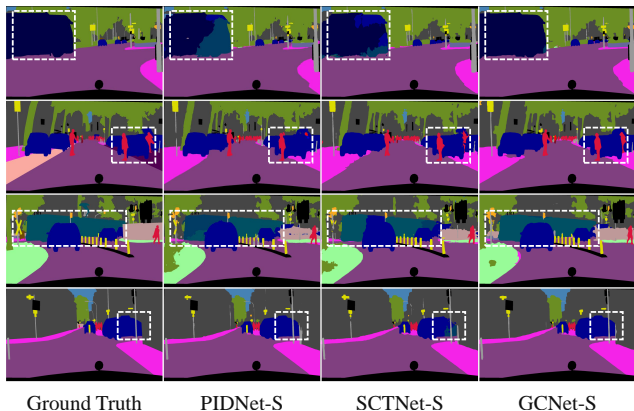


| Ground Truth | PIDNet-S | SCTNet-S | GCNet-S |

Figure 5. Visualization of different segmentation models on the Cityscapes validation set.

Table 3. Comparison on the Cityscapes validation set. "*" indicates that the model was reimplemented and retested by us (with convolution and batch normalization fused). "ImageNet" indicates whether the model utilizes ImageNet pretrained weight.

| Method | Resolution | GPU | FPS | Params | GFLOPs | ImageNet | mIoU (%) |
|---|---|---|---|---|---|---|---|
| ICNet* [38] | $1024 \times 2048$ | A100 | 181.3 | 24.9 M | 74.3 | ✗ | 71.6 |
| Fast-SCNN* [23] | $1024 \times 2048$ | A100 | 325.9 | 1.4 M | 7.3 | ✗ | 71.0 |
| SwiftNetRN-18 [20] | $1024 \times 2048$ | GTX 1080Ti | 39.9 | 11.8 M | 104.0 | ✓ | 75.4 |
| CGNet* [30] | $1024 \times 2048$ | A100 | 73.9 | 0.5 M | 27.5 | ✗ | 68.3 |
| BiSeNetV1* [34] | $1024 \times 2048$ | A100 | 116.8 | 13.3 M | 118.0 | ✓ | 74.4 |
| BiSeNetV2* [35] | $1024 \times 2048$ | A100 | 132.4 | 3.4 M | 98.4 | ✗ | 73.6 |
| STDC1* [15] | $1024 \times 2048$ | A100 | 183.7 | 8.3 M | 67.5 | ✗ | 71.8 |
| STDC2* [15] | $1024 \times 2048$ | A100 | 132.8 | 12.3 M | 93.8 | ✗ | 74.9 |
| HyperSeg-M [19] | $512 \times 1024$ | RTX 3090 | 59.1 | 10.1 M | 7.5 | ✓ | 76.2 |
| HyperSeg-S [19] | $768 \times 1536$ | RTX 3090 | 45.7 | 10.2 M | 17.0 | ✓ | 78.2 |
| PP-LiteSeg-T2 [22] | $768 \times 1536$ | RTX 3090 | 96.0 | - | - | ✓ | 76.0 |
| PP-LiteSeg-B2 [22] | $768 \times 1536$ | RTX 3090 | 68.2 | - | - | ✓ | 78.2 |
| DDRNet-23-Slim* [21] | $1024 \times 2048$ | A100 | 166.4 | 5.7 M | 36.3 | ✗ | 76.3 |
| DDRNet-23* [21] | $1024 \times 2048$ | A100 | 106.0 | 20.3 M | 143.0 | ✗ | 78.0 |
| PIDNet-S* [32] | $1024 \times 2048$ | A100 | 128.7 | 7.7 M | 47.3 | ✗ | 76.4 |
| PIDNet-M* [32] | $1024 \times 2048$ | A100 | 78.2 | 28.7 M | 177.0 | ✗ | 78.2 |
| PIDNet-L* [32] | $1024 \times 2048$ | A100 | 64.2 | 37.3 M | 275.0 | ✗ | 78.8 |
| SCTNet-S-Seg50* [33] | $512 \times 1024$ | A100 | 169.1 | 4.6 M | 7.1 | ✗ | 71.0 |
| SCTNet-S-Seg75* [33] | $768 \times 1536$ | A100 | 168.7 | 4.6 M | 16.0 | ✗ | 74.7 |
| SCTNet-B-Seg50* [33] | $512 \times 1024$ | A100 | 162.6 | 17.4 M | 28.1 | ✗ | 75.0 |
| SCTNet-B-Seg75* [33] | $768 \times 1536$ | A100 | 157.3 | 17.4 M | 63.2 | ✗ | 78.5 |
| SCTNet-B-Seg100* [33] | $1024 \times 2048$ | A100 | 117.0 | 17.4 M | 112.3 | ✗ | 79.0 |
| GCNet-S | $1024 \times 2048$ | A100 | 193.3 | 9.2 M | 45.2 | ✗ | 77.3 |
| GCNet-M | $1024 \times 2048$ | A100 | 105.0 | 34.2 M | 178.0 | ✗ | 79.0 |
| GCNet-L | $1024 \times 2048$ | A100 | 88.0 | 45.2 M | 232.0 | ✗ | 79.6 |

ent values of N for the various versions of GCNet: N=4 for GCNet-S, N=2 for GCNet-M, and N=2 for GCNet-L.

## 4.3. Comparison with State-of-the-art Methods

**Cityscapes.** The experimental results on the Cityscapes dataset are shown in Table 3. To ensure a fair comparison, we made every effort to reproduce each model on the same hardware and tested their inference speeds. The results indicate that our GCNet achieves a better balance between performance and speed. Specifically, GCNet-S achieves 77.3%

mIoU while reaching 193.3 FPS, significantly outperforming other models of similar scale. GCNet-M also performs admirably, surpassing both PIDNet-M and PIDNet-L in terms of both performance and speed, although it is slightly slower than SCTNet-B-Seg100. Furthermore, GCNet-L achieves the highest performance with 79.6% mIoU. As shown in the table, despite GCNet having a higher number of parameters and GFLOPs compared to other models, these factors are not the decisive determinants of inference speed. In fact, inference speed is also influenced by the fre-

Table 4. Comparison on the CamVid test set and the Pascal VOC 2012 validation set. All models were reimplemented by us and retested on the A100. The inference resolution for CamVid was set at $720 \times 960$, while for VOC it was set at $512 \times 2048$. † indicates that ImageNet pretrained weight were used.

| Method | CamVid | | VOC | |
|---|---|---|---|---|
| | FPS | mIoU (%) | FPS | mIoU (%) |
| CGNet [30] | 102.6 | 69.4 | 104.7 | 42.2 |
| BiSeNetV1† [34] | 276.0 | 73.3 | 214.3 | 56.0 |
| BiSeNetV2 [35] | 240.6 | 75.8 | 228.6 | 50.0 |
| STDC1 [15] | 240.8 | 70.2 | 227.0 | 52.9 |
| STDC2 [15] | 149.8 | 71.4 | 137.2 | 56.3 |
| DDRNet-23-Slim [21] | 189.3 | 76.2 | 174.6 | 54.2 |
| DDRNet-23 [21] | 171.1 | 78.2 | 165.4 | 55.8 |
| PIDNet-S [32] | 142.8 | 77.2 | 137.4 | 54.0 |
| PIDNet-M [32] | 128.9 | 78.6 | 123.1 | 57.5 |
| PIDNet-L [32] | 108.4 | 78.7 | 104.0 | 58.2 |
| GCNet-S | 210.6 | 76.6 | 196.9 | 54.9 |
| GCNet-M | 190.0 | 78.5 | 173.4 | 58.2 |
| GCNet-L | 170.4 | 79.1 | 145.1 | 58.8 |

quency of memory accesses within the model. The single-path block design of GCNet reduces the number of memory accesses, thereby enhancing inference speed. We also visualized the segmentation results of PIDNet-S, SCTNet-S-Seg75, and GCNet-S, as illustrated in Figure 5. The figure indicates that GCNet exhibits higher accuracy in pixel classification.

**CamVid.** The experimental results on the CamVid dataset are shown in Table 4. To ensure a fair comparison, we retrained the models listed in the table using the MMSegmentation framework and evaluated their inference speeds. The experimental results indicate that our GCNet-S achieves 76.6% mIoU while also reaching 210.6 FPS. Although PIDNet-S has an mIoU that is 0.6% higher than that of GCNet-S, its FPS is lower by 67.8. Notably, GCNet-L performs exceptionally well at a resolution of $720 \times 960$. Compared to GCNet-S, while its FPS only decreases by 40.2, it achieves 79.1% mIoU, surpassing all models in the PIDNet series in both performance and speed.

**Pascal VOC 2012.** The experimental results on the Pascal VOC 2012 dataset are shown in Table 4. To ensure a fair comparison, we retrained the models listed in the table using the MMSegmentation framework and evaluated their inference speeds. The experimental results indicate that our GCNet-S achieves 54.9% mIoU while also reaching 196.9 FPS. Similar to the results on the CamVid dataset, GCNet-L demonstrates remarkable performance and speed at a relatively low resolution ($512 \times 2048$), even outper-
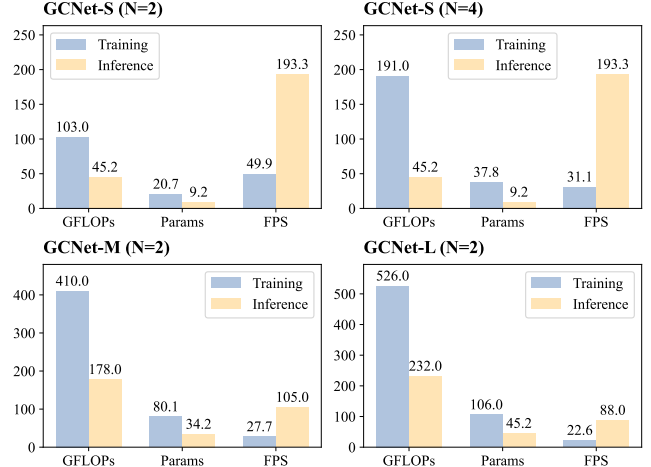


Figure 6. The comparison of GFLOPs, number of parameters, and FPS for GCNet during training versus inference.

forming PIDNet-S. Interestingly, we found that BiSeNetV1 performs better on the VOC dataset, exceeding BiSeNetV2 by 6% in mIoU, with an FPS that is only 14.3 lower. While BiSeNetV1 excels on VOC, the overall performance across the three datasets shows that GCNet outperforms BiSeNetV1.

## 4.4. Comparison of Training and Inference

To visually compare the different configurations of GCNet, we calculated the GFLOPs, number of parameters, and FPS during training versus those during inference, as shown in Figure 6. The figure illustrates that GCNet requires significantly more computational resources during training and operates at a slower speed. However, after self-contracting, both GFLOPs and number of parameters are greatly reduced, while FPS experiences a substantial increase. Notably, the performance of both configurations remains equivalent.

## 5. Conclusion

In this study, we proposed GCNet for real-time semantic segmentation. GCNet achieves stronger performance and faster speed simultaneously through self-enlargement and self-contraction, as confirmed by experiments on three public datasets. Although GCNet operates at high speed, its parameter count and computational complexity are higher than those of other models, which limits its applicability on devices with constrained storage capacity. In the future, we plan to explore more powerful and parameter-efficient reparameterizable structures.

# References

[1] Gabriel J Brostow, Julien Fauqueur, and Roberto Cipolla. Semantic object classes in video: A high-definition ground truth database. *Pattern recognition letters*, 30(2):88–97, 2009. 2, 6, 11

[2] Yingfeng Cai, Lei Dai, Hai Wang, and Zhixiong Li. Multi-target pan-class intrinsic relevance driven model for improving semantic segmentation in autonomous driving. *IEEE Transactions on Image Processing*, 30:9069–9084, 2021. 1

[3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 2

[4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.

[5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. 2

[6] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in neural information processing systems*, 34:17864–17875, 2021. 2

[7] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022. 2

[8] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. 2020. 11

[9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2, 6, 11

[10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 6, 11

[11] Xiaohan Ding, Yuchen Guo, Guiguang Ding, and Jungong Han. Acnet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1911–1920, 2019. 3

[12] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Diverse branch block: Building a convolution as an inception-like unit. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10886–10895, 2021.

[13] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: Making vgg-style convnets great again. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13733–13742, 2021. 3

[14] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010. 2, 6, 11

[15] Mingyuan Fan, Shenqi Lai, Junshi Huang, Xiaoming Wei, Zhenhua Chai, Junfeng Luo, and Xiaolin Wei. Rethinking bisenet for real-time semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9716–9725, 2021. 2, 7, 8, 12

[16] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3146–3154, 2019. 4

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2, 4

[18] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 4

[19] Yuval Nirkin, Lior Wolf, and Tal Hassner. Hyperseg: Patch-wise hypernetwork for real-time semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4061–4070, 2021. 7

[20] Marin Orsic, Ivan Kreso, Petra Bevandic, and Sinisa Segvic. In defense of pre-trained imagenet architectures for real-time semantic segmentation of road-driving images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12607–12616, 2019. 7

[21] Huihui Pan, Yuanduo Hong, Weichao Sun, and Yisong Jia. Deep dual-resolution networks for real-time and accurate semantic segmentation of traffic scenes. *IEEE Transactions on Intelligent Transportation Systems*, 24(3):3448–3460, 2022. 1, 2, 3, 4, 5, 7, 8, 11, 12

[22] Juncai Peng, Yi Liu, Shiyu Tang, Yuying Hao, Lutao Chu, Guowei Chen, Zewu Wu, Zeyu Chen, Zhiliang Yu, Yuning Du, et al. Pp-liteseg: A superior real-time semantic segmentation model. *arXiv preprint arXiv:2204.02681*, 2022. 2, 7

[23] Rudra PK Poudel, Stephan Liwicki, and Roberto Cipolla. Fast-scnn: Fast semantic segmentation network. *arXiv preprint arXiv:1902.04502*, 2019. 3, 7, 12

[24] Wenbo Qi, HC Wu, and SC Chan. Mdf-net: A multi-scale dynamic fusion network for breast tumor segmentation of ultrasound images. *IEEE Transactions on Image Processing*, 2023. 1

[25] Eduardo Romera, José M Alvarez, Luis M Bergasa, and Roberto Arroyo. Erfnet: Efficient residual factorized convnet for real-time semantic segmentation. *IEEE Transactions on Intelligent Transportation Systems*, 19(1):263–272, 2017. 1, 2

[26] Evan Shelhamer, Jonathan Long, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(4):640–651, 2017. 2

[27] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 761–769, 2016. 6

[28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2

[29] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11534–11542, 2020. 4

[30] Tianyi Wu, Sheng Tang, Rui Zhang, Juan Cao, and Yongdong Zhang. Cgnet: A light-weight context guided network for semantic segmentation. *IEEE Transactions on Image Processing*, 30:1169–1179, 2020. 2, 7, 8, 12

[31] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34: 12077–12090, 2021. 1, 2

[32] Jiacong Xu, Zixiang Xiong, and Shankar P Bhattacharyya. Pidnet: A real-time semantic segmentation network inspired by pid controllers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19529–19539, 2023. 1, 2, 3, 4, 5, 6, 7, 8, 11, 12

[33] Zhengze Xu, Dongyue Wu, Changqian Yu, Xiangxiang Chu, Nong Sang, and Changxin Gao. Sctnet: Single-branch cnn with transformer semantic information for real-time segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6378–6386, 2024. 1, 2, 7, 12

[34] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 325–341, 2018. 1, 2, 3, 7, 8, 12

[35] Changqian Yu, Changxin Gao, Jingbo Wang, Gang Yu, Chunhua Shen, and Nong Sang. Bisenet v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International Journal of Computer Vision*, 129: 3051–3068, 2021. 1, 3, 5, 7, 8, 11, 12

[36] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016. 2

[37] Jinghua Zhang, Chen Li, Sergey Kosov, Marcin Grzegorzek, Kimiaki Shirahama, Tao Jiang, Changhao Sun, Zihan Li, and Hong Li. Lcu-net: A novel low-cost u-net for environmental microorganism image segmentation. *Pattern Recognition*, 115:107885, 2021. 1

[38] Hengshuang Zhao, Xiaojuan Qi, Xiaoyong Shen, Jianping Shi, and Jiaya Jia. Icnet for real-time semantic segmentation on high-resolution images. In *Proceedings of the European conference on computer vision (ECCV)*, pages 405–420, 2018. 7, 12

## A. More Experimental Details

### A.1. Datasets

**Cityscapes.** Cityscapes [9] is extensively utilized in urban scene understanding and autonomous driving research. It consists of 19 classes and a total of 5000 images, with 2975 designated for training, 500 for validation, and 1525 for testing. The test set is unlabeled, and model predictions must be uploaded to a specific website for evaluation. Each image has been meticulously annotated and possesses a resolution of $1024 \times 2048$ pixels.

**CamVid.** CamVid [1] is the first video dataset to include semantic labels. It comprises 701 images, with 367 designated for training, 101 for validation, and 233 for testing. Each image possesses a resolution of $720 \times 960$ pixels and features 32 class labels, though only 11 are typically used for training and evaluation. Notably, many current research [21, 32, 35] utilize the training and validation sets for training and the test set for evaluation, and we adopt this same strategy.

**Pascal VOC 2012.** Pascal VOC 2012 [14] is primarily utilized for image classification, object detection, and image segmentation tasks. It includes 20 classes along with 1 background class. For the semantic segmentation task, the dataset consists of 2913 images, with 1464 allocated to the training set and 1449 to the validation set. Unlike Cityscapes and CamVid, the resolution of these images varies.

### A.2. Implementation Details

**Computing Platform.** The hardware of the computing platform comprises an AMD EPYC 7742 CPU and four NVIDIA A100 GPUs. The software stack includes Ubuntu 20.04.6, CUDA 11.3, PyTorch 1.12.1, TorchVision 0.13.1, MMEngine 0.10.2, and MMSegmentation 1.2.2 [8]. During the training phase, both GPUs are utilized, while for the inference phase, only a single GPU is used, with the batch size set to 1.

**Training.** Stochastic Gradient Descent (SGD) with a momentum of 0.9 and a weight decay of 0.0005 was employed as the optimizer. Additionally, a polynomial learning rate decay strategy with a power of 0.9 was utilized. For data augmentation, we applied random scaling within the range of 0.5 to 2.0, random cropping, and random flipping with a probability of 0.5. The training iterations, initial learning rate, random crop size, and batch size for the Cityscapes, CamVid, and Pascal VOC 2012 datasets were configured as follows: [120000, 0.01, $1024 \times 1024$, 12], [7800, 0.001, $720 \times 960$, 12], and [24400, 0.001, $512 \times 512$, 16], respectively. Notably, our model was not pretrained on the ImageNet [10] dataset, and the Cityscapes model weights were used during training on the CamVid and Pascal VOC 2012 datasets.
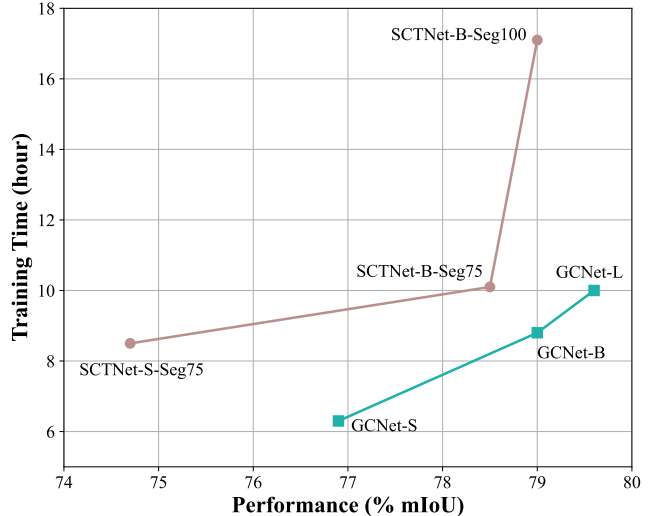


Figure 7. Comparison of training time between GCNet and SCT-Net on the Cityscapes dataset. Four A100s were used for training and the training time was recorded.

**Inference.** The inference image sizes for the Cityscapes, CamVid, and Pascal VOC 2012 datasets were configured at $1024 \times 2048$, $720 \times 960$, and $512 \times 2048$, respectively. Since the image resolution in the Pascal VOC dataset is not fixed, images are adjusted to an approximate resolution of $512 \times 2048$ to maintain the aspect ratio. The inference speed of the same model varies across different CPUs and software environments, even when using the same GPU. To ensure a fair comparison, we made every effort to reimplement the other models.
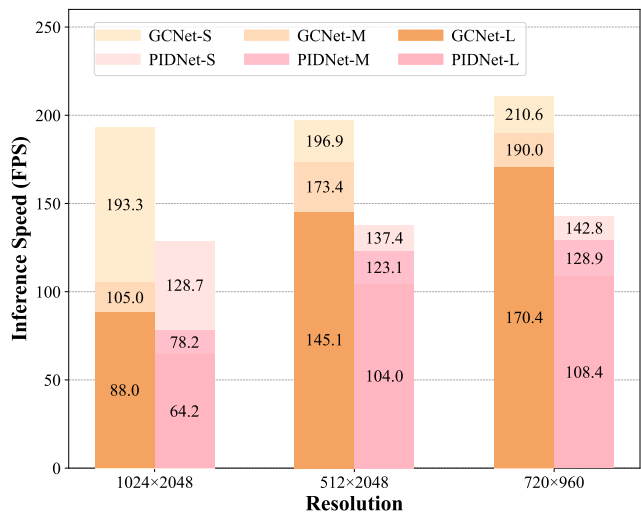


Figure 8. Inference speed of different segmentation models for varying resolutions. The GPU used is A100.

11

Table 5. Comparison of inference speed on varying GPUs using the Cityscape validation set. "*" indicates that the model was reimplemented and retested by us (with convolution and batch normalization fused). "ImageNet" indicates whether the model utilizes ImageNet pretrained weight.

| Method | Resolution | RTX 4080 | RTX 3090 | V100 | A100 | ImageNet | mIoU |
|---|---|---|---|---|---|---|---|
| ICNet* [38] | $1024 \times 2048$ | 92.3 | 108.7 | 76.7 | 181.3 | ✗ | 71.6 |
| Fast-SCNN* [23] | $1024 \times 2048$ | 125.3 | 211.9 | 173.6 | 325.9 | ✗ | 71.0 |
| CGNet* [30] | $1024 \times 2048$ | 38.3 | 53.1 | 48.5 | 73.9 | ✗ | 68.3 |
| BiSeNetV1* [34] | $1024 \times 2048$ | 62.7 | 64.3 | 56.5 | 116.8 | ✓ | 74.4 |
| BiSeNetV2* [35] | $1024 \times 2048$ | 68.7 | 69.6 | 64.4 | 132.4 | ✗ | 73.6 |
| STDC1* [15] | $1024 \times 2048$ | 101.1 | 103.3 | 89.4 | 183.7 | ✗ | 71.8 |
| STDC2* [15] | $1024 \times 2048$ | 75.8 | 73.8 | 64.7 | 132.8 | ✗ | 74.9 |
| DDRNet-23-Slim* [21] | $1024 \times 2048$ | 94.2 | 116.5 | 97.6 | 166.4 | ✗ | 76.3 |
| DDRNet-23* [21] | $1024 \times 2048$ | 54.2 | 53.4 | 47.1 | 106.0 | ✗ | 78.0 |
| PIDNet-S* [32] | $1024 \times 2048$ | 60.1 | 84.2 | 76.3 | 128.7 | ✗ | 76.4 |
| PIDNet-M* [32] | $1024 \times 2048$ | 41.7 | 41.3 | 35.4 | 78.2 | ✗ | 78.2 |
| PIDNet-L* [32] | $1024 \times 2048$ | 30.0 | 31.1 | 27.9 | 64.2 | ✗ | 78.8 |
| SCTNet-S-Seg50* [33] | $512 \times 1024$ | 78.9 | 124.2 | 108.3 | 169.1 | ✗ | 71.0 |
| SCTNet-S-Seg75* [33] | $768 \times 1536$ | 78.7 | 124.0 | 99.7 | 168.7 | ✗ | 74.7 |
| SCTNet-B-Seg50* [33] | $512 \times 1024$ | 77.3 | 119.0 | 97.6 | 162.6 | ✗ | 75.0 |
| SCTNet-B-Seg75* [33] | $768 \times 1536$ | 73.9 | 101.9 | 83.1 | 157.3 | ✗ | 78.5 |
| SCTNet-B-Seg100* [33] | $1024 \times 2048$ | 63.6 | 63.9 | 53.3 | 117.0 | ✗ | 79.0 |
| GCNet-S | $1024 \times 2048$ | 110.1 | 130.9 | 114.1 | 193.3 | ✗ | 77.3 |
| GCNet-M | $1024 \times 2048$ | 47.4 | 50.1 | 44.2 | 105.0 | ✗ | 79.0 |
| GCNet-L | $1024 \times 2048$ | 38.3 | 40.7 | 36.2 | 88.0 | ✗ | 79.6 |

## A.3. Metrics

We adopt mIoU (mean Intersection over Union), number of parameters, GFLOPs (Giga Floating Point Operations), and FPS (Frames Per Second) as metrics. mIoU is commonly used in image segmentation tasks, where it measures the average overlap between predicted results and ground truth annotations, serving as an indicator of model performance. Number of parameters refers to the total number of trainable parameters in the model, providing a measure of its size. GFLOPs quantifies the computational load of the model, reflecting its computational complexity. FPS represents the number of image frames the model processes per second, serving as a metric for inference speed. In general, number of parameters and GFLOPs are not the decisive factors in-

fluencing FPS. This paper focuses on optimizing mIoU and FPS, rather than number of parameters and GFLOPs.

## B. More Experiments

### B.1. Training Time of GCNet and SCTNet

The paper mentions that SCTNet requires a high-performance segmentation model for knowledge distillation training, which is quite time-consuming. To verify that GC-Net demands less training time compared to SCTNet, we recorded the training times of various versions of both models, as shown in Figure 7. Since SCTNet-S/B-Seg50 and SCTNet-S/B-Seg75 share the same training configuration, differing only in inference settings, we only recorded the

training time for SCTNet-S-Seg75 and SCTNet-B-Seg75. The figure reveals that GCNet not only requires less training time but also achieves higher performance than SCTNet. Specifically, SCTNet-B-Seg100 takes 17.1 hours to reach 79.0% mIoU, while GCNet-B achieves this in only 8.8 hours.

## B.2. Inference Speed With Varying Resolution

To provide an intuitive understanding of the inference speed of GCNet at varying resolutions, we visualized its FPS, as shown in Figure 8. The figure reveals that as resolution decreases, the FPS of both GCNet and PIDNet increases significantly, particularly for the M and L versions. Surprisingly, GCNet-M and GCNet-L achieve outstanding FPS at lower resolutions ($512 \times 2048$ and $720 \times 960$), with GCNet-L even surpassing PIDNet-S. This may be attributed to the benefits of reduced memory access enabled by the single-path block, as memory access represents a significant computational cost in computer systems.

## B.3. Inference Speed With Varying GPUs

To demonstrate GCNet's versatility across varying GPUs, we conducted speed tests on the RTX 4080, RTX 3090, V100, and A100, as shown in Table 5. The results reveal that GCNet performs well on both consumer-grade GPUs (RTX 4080 and RTX 3090) and professional-grade GPUs (V100 and A100). Interestingly, smaller models show faster inference speeds on the RTX 3090 compared to the RTX 4080, while larger models perform similarly on both. As a single-branch architecture model, SCTNet is relatively insensitive to changes in lower resolutions, with comparable inference speeds for Seg50 and Seg75. In contrast, multi-branch models show substantial speed improvements with lower resolutions. As illustrated in Figure 8, inference speeds for multi-branch models, especially GCNet, increase significantly as the resolution decreases. We attribute this to the high-resolution branch in multi-branch architectures, which requires the maintenance of larger feature maps and thus more computations. In future work, we plan to further investigate GCNet on lower resolutions using the Cityscapes dataset.