

Real-Time Semantic Segmentation via Auto Depth, Downsampling Joint Decision and Feature Aggregation

Peng Sun · Jiaxiang Wu · Songyuan Li · Peiwen Lin · Junzhou Huang · Xi Li

Abstract To satisfy the stringent requirements on computational resources in the field of real-time semantic segmentation, most approaches focus on the hand-crafted design of light-weight segmentation networks. Recently, Neural Architecture Search (NAS) has been used to search for the optimal building blocks of networks automatically, but the network depth, downsampling strategy, and feature aggregation way are still set in advance by trial and error. In this paper, we propose a joint search framework, called AutoRTNet, to automate the design of these strategies. Specifically, we propose hyper-cells to jointly decide the network depth and downsampling strategy, and an aggregation cell to achieve automatic multi-scale feature aggregation. Experimental results show that AutoRTNet achieves 73.9% mIoU on the Cityscapes test set and 110.0 FPS on an NVIDIA TitanXP GPU card with 768×1536 input images.

Keywords Real-time semantic segmentation · Neural network search

Peng Sun
Zhejiang University
E-mail: sunpeng1996@zju.edu.cn

Jiaxiang Wu
Tencent AI Lab

Songyuan Li
Zhejiang University

Peiwen Lin
SenseTime Research

Junzhou Huang
University of Texas at Arlington

Xi Li
Corresponding Author
Zhejiang University
E-mail: xilizju@zju.edu.cn

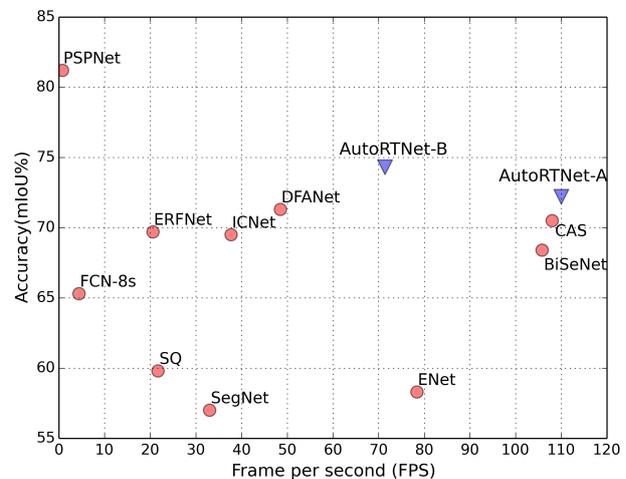


Fig. 1: The inference speed and accuracy for different networks on the Cityscapes test set. Compared with other methods, our AutoRTNet locates in the right-top since it features lower latency with comparable accuracy. **Best viewed in color.**

1 Introduction

Semantic segmentation, a fundamental topic in computer vision, aims at assigning per-pixel semantic labels for images. Recent approaches [54, 6, 7, 55] based on fully convolutional networks [25] have achieved remarkable accuracy on public benchmarks [3, 9, 10]. Such improvements, however, come at the cost of deeper and less efficient networks, which may not be applicable to many real-time systems, e.g. autonomous driving and video surveillance.

To perform fast semantic segmentation with satisfactory accuracy, the design philosophy of real-time segmentation network architectures mainly concentrates on three aspects: 1) building block design [16, 33], which considers the block-level feature representation capacity, computational complex-

ity, and receptive field size; 2) network depth and downsampling strategy [16, 18], which directly affect the accuracy and speed of networks, hence real-time networks favor shallow layers and fast downsampling; and 3) feature aggregation [49, 53], which fuses multi-scale features to compensate the loss of spatial details caused by fast downsampling.

The above hand-crafted networks make huge progress, while they require expertise in architecture design based on laborious trial and error. To relieve this burden, some researchers introduce neural architecture search (NAS) methods [2, 56, 23, 46] into this field, and obtain excellent results [5, 22, 52, 30]. AutoDeepLab [22] and DPC [5] focus on high-quality segmentation instead of real-time applications. To meet the real-time demand, CAS [52] searches a customized architecture by introducing a latency loss function. Though its building block is searched, the network depth, downsampling strategy and feature aggregation way are still set by hand in advance and non-adjustable during the search process. While the aforementioned three aspects are highly correlated and indispensable for a remarkable real-time segmentation network, and these non-adjustable settings increase the difficulties and limitations to find an optimal real-time segmentation architecture (i.e. the best trade-off between performance and speed). These motivate us to explore all of the aspects automatically during the search process.

In this paper, we propose a joint search framework to search for the optimal building blocks, network depth, downsampling strategy, and feature aggregation way simultaneously. Specifically, we propose hyper-cells to jointly decide the network depth and downsampling strategy via a cell-level pruning process in an adaptive manner, and an aggregation cell to fuse features from multiple spatial scales automatically. As for the hyper-cell, we introduce a novel learnable architecture parameter for it, and the network depth and downsampling strategy are fully determined concurrently according to the optimized hyper-cell architecture parameters. As for the aggregation cell, we aggregate multi-level features in the network automatically to effectively fuse the low-level spatial details and high-level semantic context.

We denote the resulting network as **Auto** searched **Real-Time** semantic segmentation network or **AutoRTNet**. We evaluate AutoRTNet on both Cityscapes [9] and CamVid [3] datasets. The experiments demonstrate the superiority of AutoRTNet, as shown in Figure 1, where our AutoRTNet achieves the best accuracy-efficiency trade-off.

The main contributions can be summarized as follows:

- We propose a joint search framework for real-time semantic segmentation that automatically searches for the building blocks, network depth, downsampling strategy, and feature aggregation way simultaneously.
- We propose the hyper-cell to learn the network depth and downsampling strategy in an adaptive manner via

the cell-level pruning process, and the aggregation cell to achieve automatic multi-scale feature aggregation.

- Notably, AutoRTNet has achieved 73.9% mIoU on the Cityscapes test set and 110.0 FPS on an NVIDIA TitanXP GPU card with 768×1536 input images.

2 Related Work

2.1 Semantic Segmentation

High-Quality Segmentation FCN [25] is the pioneer work which has greatly promoted the development of semantic segmentation. Extensions to FCN follow many directions. Encoder-decoder structures [1, 20, 31] combine low-level and high-level features to improve accuracy of semantic segmentation. DRN [50] and DeepLab [6, 7] use dilated convolution to effectively enlarge the receptive field. To capture multi-scale context information, DeepLabV3 [6] and PSPNet [54] propose the pyramid modules. Recently, attention mechanism [41] has been used in segmentation field [11, 51, 55, 17]. These outstanding works are designed for high-quality segmentation, which is inapplicable to real-time applications.

Real-Time Methods Various algorithms have been proposed for real-time semantic segmentation. Some works [45] reduce the computation overheads via restricting the input image size, channel-pruning algorithms [33, 1] are introduced to boost the inference speed, and most real-time methods focus on designing the light-weight and effective network architectures. The design philosophy mainly can be summarized into the following three aspects. And in our work, we fully explore all three aspects simultaneously.

Building block design The building block design [33, 36, 28, 16] requires researchers to give sufficient consideration to the computational complexity, feature representation capacity, and receptive field size, which is essential for real-time semantic segmentation. For example, ENet [33] and DABNet [16] propose light-weight blocks and stack them with different dilation rates to form a whole network. MobileNet and its variants [14, 37] use blocks with depth-wise separable convolution in pursuit of light-weight models.

Network depth and downsampling strategy Different from high-quality segmentation networks using pre-defined backbones (ResNet [13], Xception [8], etc.) as encoders, the network depth and downsampling strategy (i.e. how many blocks or layers in each stage) are determined mostly by hand for real-time segmentation networks (e.g. DABNet [16], DFANet [18], ERFNet [36]), as they directly affect the accuracy and speed of the networks. For pursuing fast inference speed, real-time networks always enjoy shallow layers and perform fast downsampling with factor 16 or 32.

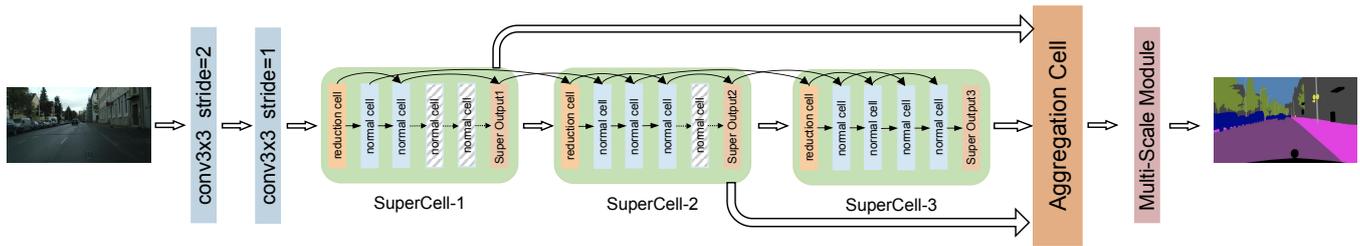


Fig. 2: Illustration of our joint network architecture search framework. The network begins with two convolution layers and contains three hyper-cells which search for the optimal network depth and downsampling strategy via the cell-level pruning process, each hyper-cell contains a reduction cell and n normal cells. The cells marked with the dotted white line are pruned after optimization. The aggregation cell is designed to perform automatic multi-scale feature aggregation effectively, and it seamlessly integrates the outputs of hyper-cells.

Feature aggregation The fast downsampling of real-time networks easily results in the loss of spatial details. Thus, multi-scale feature aggregation [49, 53, 18] has been proposed to remedy the loss of spatial details. ICNet [53] proposes a image cascade network with multi-scale inputs. BiSeNet [49] decouples the network into context and spatial paths to make a right balance between the accuracy and speed. DFANet [18] aggregates multi-scale features from different layers to remedy the loss of spatial details.

2.2 Neural Architecture Search

NAS Overview Neural architecture search (NAS) focuses on automating the network architecture design process. Early NAS methods are time-consuming (e.g. thousands of GPU days) and computationally expensive via reinforcement learning [56, 2, 57, 39] or evolutionary algorithms [29, 35]. Recently, the emergence of differentiable NAS methods [23, 46, 4] has greatly relieved the time-consuming problem while achieving excellent performance. DARTs [23] is the pioneer work for gradient-based NAS, they propose an iterative optimization framework which is based on the continuous relaxation of the architecture representation. SNAS [46] constrains the architecture parameters to approximate one-hot, resolving the inconsistency in optimizing between the performance of derived child networks and converged parent networks. FBNet [42], ProxylessNAS [4], MnasNet [39] propose multi-objective optimization with the consideration of real-world latency.

NAS For Segmentation DPC [5] is the first work for dense image prediction using NAS methods and searches for a multi-scale representation module. The similar work to us is AutoDeepLab [22], they propose a hierarchical search space and search for the downsampling path. Although they also search for the downsampling strategy, the mechanism is extremely different from ours. They design the network level continuous relaxation to learn the downsampling path,

while we search for the downsampling strategy via the cell-level pruning progress. Moreover, they cannot search for the adaptive network depth and feature aggregation way and they focus on high-quality segmentation instead of real-time applications. For real-time requirements, CAS [52] searches for an architecture with customized resource constraints and achieves excellent real-time performance. However, our approach also searches for the adaptive network depth, downsampling strategy and feature aggregation way, which is extremely different from CAS [52].

NAS For Object Detection The combination of multi-scale features is also essential for object detection [21, 24]. In the field of NAS, NAS-FPN [12] and Auto-FPN [48] also search for an architecture that merges features of varying dimensions and are successful at searching for the appropriate combination method. Unlike us, NAS-FPN [12] proposes merging cells and uses RNN controller to select candidate feature layers and a binary operation in each merging cell. Their search space only consists of two binary operations, i.e. sum and global pooling for their simplicity. Auto-FPN [48] searches for an efficient feature fusion module, which search space is specially designed for detection and flexible enough to cover many popular designs of detectors. Thus the search space design, motivation and implementation of above both methods are different from ours.

3 Methods

The joint search framework is shown in Figure 2. We propose the hyper-cell to search for the optimal network depth and downsampling strategy as they directly affect the accuracy and speed of networks. For remedying the loss of spatial details caused by fast downsampling, a novel aggregation cell is proposed for automatic multi-scale feature aggregation. The framework contains two pre-defined convolution layers, three hyper-cells and an aggregation cell. The multi-scale module [52] is subsequently used to extract the

global and local context for final prediction. For real-time demands, we take the real-world latency into consideration during the search process. We begin this section with the differentiable architecture search. Afterwards, we elaborate the proposed hyper-cell and aggregation cell in detail.

3.1 Differentiable Architecture Search

Intra-cell Search space The hyper-cell is the building block of the network, and the cell is the basic component unit of the hyper-cell, as shown in Figure 2. There are two types of cells, normal cells and reduction cells [23, 46]. The reduction cells reduce the feature map size by a factor of 2 for downsampling, and the factor is 1 in normal cells.

A cell is a directed acyclic graph (DAG) consisting of an ordered sequence of N nodes, denoted by $\mathcal{N} = \{x^{(1)}, \dots, x^{(N)}\}$. Each node $x^{(i)}$ is a latent representation (i.e. feature map), and each directed edge (i, j) is associated with some candidate operations (e.g. conv, pooling) in operation set $\mathcal{O}^{(i,j)}$, representing all possible transformations from $x^{(i)}$ to $x^{(j)}$. Each cell has two inputs (the outputs of the previous two cells) and one output (the concatenation of all the intermediate nodes in the cell). The structure of cell is shown on the right in Figure 3. Each intermediate node $x^{(j)}$ is computed based on all of its predecessors:

$$x^{(j)} = \sum_{i < j} \tilde{o}^{(i,j)}(x^{(i)}) \quad (1)$$

where $\tilde{o}^{(i,j)} \in \mathcal{O}^{(i,j)}$ is the optimal operation at edge (i, j) .

In order to determine the optimal operation $\tilde{o}^{(i,j)}$ at edge (i, j) , we represent the intra-cell search space with a set of one-hot random variables from a fully factorizable joint distribution $p(M)$ [46]. Specifically, each edge (i, j) is associated with a one-hot random variable $M^{(i,j)}$. We use $M^{(i,j)}$ as a mask to multiply all the candidate operations $\mathcal{O}^{(i,j)}$ at edge (i, j) , and the intermediate node $x^{(j)}$ is given by:

$$x^{(j)} = \sum_{i < j} \sum_{o \in \mathcal{O}} m_o^{(i,j)} \cdot o^{(i,j)}(x^{(i)}) \quad (2)$$

where $m_o^{(i,j)} \in M^{(i,j)}$ and $m_o^{(i,j)}$ is a random variable in $\{0, 1\}$, it is evaluated to 1 if operation $o^{(i,j)}$ is selected.

To make $p(M)$ differentiable, we use Gumbel Softmax technique [15, 27] to relax the discrete sampling distribution to be continuous and differentiable:

$$M^{(i,j)} = f_{\alpha^{(i,j)}}(G^{(i,j)}) = \text{softmax}((\log \alpha^{(i,j)} + G^{(i,j)})/\lambda) \quad (3)$$

where $M^{(i,j)}$ is the softened one-hot random variable for operation selection at edge (i, j) , $\alpha^{(i,j)}$ is the intra-cell architecture parameter at edge (i, j) , $G^{(i,j)} = -\log(-\log(U^{(i,j)}))$

is a vector of Gumbel random variables, $U^{(i,j)}$ is a uniform random variable in the range $(0, 1)$. λ is the temperature of softmax, and as λ approaches 0, $M^{(i,j)}$ approximately becomes one-hot. The technique of using Gumbel Softmax makes the entire intra-cell search differentiable [43, 42, 46] to both network parameter w and architecture parameter α .

For the candidate operation set \mathcal{O} , we collect the operations as follows:

- zero operation
- skip connection
- 3×3 max pooling
- 3×3 conv
- 3×3 conv, repeat 2
- 3×3 separable conv
- 3×3 separable conv, repeat 2
- 3×3 dilated separable conv, dilation=2
- 3×3 dilated separable conv, dilation=4
- 3×3 dilated separable conv, dilation=2, repeat 2

Intra-cell Latency Cost For the operation selection of cells towards real-time network, we take real-world latency into consideration. Specifically, we build a GPU-latency lookup table [4, 39, 42, 52] that records the inference time cost of each candidate operation. The latency of each operation is measured in micro-second on a TitanXP GPU. During the search process, we associate a cost $lat_o^{(i,j)}$ with each candidate operation $o^{(i,j)}$ at edge (i, j) , thus the latency cost of cell p is formulated as:

$$lat_p = \sum_{(i,j)} \sum_{o \in \mathcal{O}} m_o^{(i,j)} \cdot lat_o^{(i,j)} \quad (4)$$

where $m_o^{(i,j)} \in M^{(i,j)}$ and $M^{(i,j)}$ denotes the softened one-hot random variable at edge (i, j) . By using the pre-built lookup table and above sampling process, the latency loss is also differentiable with respect to $m_o^{(i,j)}$.

3.2 Adaptive Network Depth and Downsampling

Hyper-Cell Search Space The network depth and downsampling strategy directly affect the accuracy and speed of networks for real-time semantic segmentation. To adjust them in an adaptive manner, we formulate the two design-making processes as a single cell-level pruning process. Specifically, we propose a hyper-cell, as shown in Figure 3, which consists of a reduction cell and n normal cells. We introduce $n+1$ edges to connect each cell with the hyper-cell's output, and associate them with the learnable architecture parameter β . The intra-cell architecture parameters α of n normal cells are shared in the same hyper-cell.

We determine the depth of each hyper-cell by limiting that only one edge can be activated for each hyper-cell, and all cells behind this activated edge can be pruned safely. Each specific edge in hyper-cell s is associated with a one-hot random variable $U^s = (u_1^s, u_2^s, \dots, u_{n+1}^s)$ from a fully

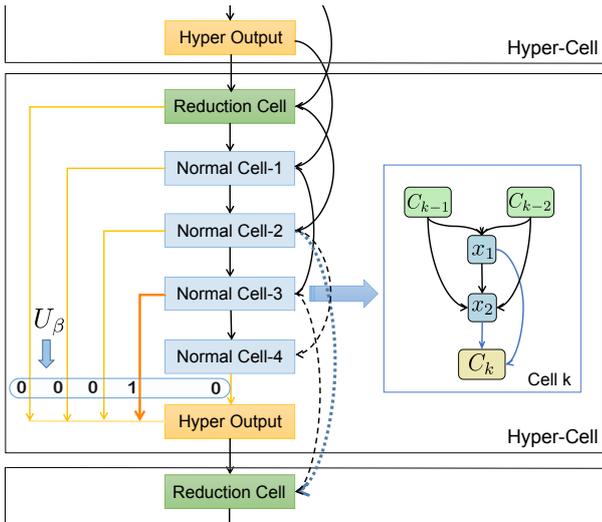


Fig. 3: Illustration of our hyper-cell. The hyper-cell consists of a reduction cell and n normal cells and $n + 1$ edges with architecture parameter which encodes the depth of the hyper-cell. The structure of cell is shown on the right in this figure.

factorizable joint distribution $P(U)$. The U^s works as a mask during the training process, and the output of the hyper-cell s is designed as:

$$HyperOut^{(s)} = \sum_{p=1}^{n+1} u_p^s \cdot (C_p^s) \quad (5)$$

where C_p^s is the output of p -th cell in hyper-cell s , u_p^s represents the random variable in $\{0, 1\}$ of p -th edge of hyper-cell s . Similarly, we adopt the Gumbel Softmax based sampling process to make the training process differentiable.

$$U^s = f_{\beta^s}(G^s) = \text{softmax}((\log \beta^s + G^s)/\lambda) \quad (6)$$

where U^s is the softened one-hot random variable for edge selection of hyper-cell s , β^s is the architecture parameter of hyper-cell s . G^s and λ are similar to the ones in equation (3). The hyper-cell architecture parameter β we introduced can be effectively optimized together with the network parameter w , intra-cell architecture parameter α in the same round of back-propagation. After stacking hyper-cells to form a whole network, the network depth and downsampling strategy can be fully explored concurrently according to the architecture parameter β .

To better explain the cell-level pruning process, we give an example as follows. In the initial phase, let's say we have five cells (one reduction cell and four normal cells) and each cell in hyper-cell keeps its original inputs and outputs. As shown in Figure 3, if the fourth edge is activated currently (i.e. U is $\{0, 0, 0, 1, 0\}$), the Normal Cell-4 will be pruned in this iteration, and the output of this hyper-cell is the output of Normal Cell-3. At the same time, the reduction cell in

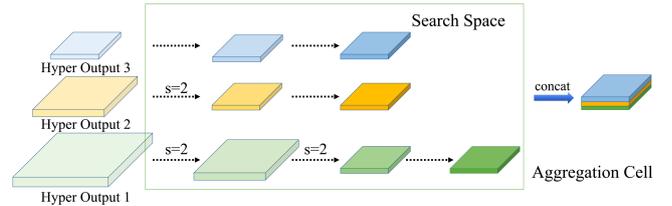


Fig. 4: Overview of the aggregation cell for automatic multi-scale feature aggregation. The aggregation cell contains E edges (dotted arrows), each edge is equipped with some candidate operations. The “s=2” means stride = 2. **Best viewed in color.**

next hyper-cell $s + 1$ takes the outputs of hyper-cell s and Normal Cell-2 in hyper-cell s as its inputs, to stick to the “two-input” principle of the cell. The learning and adjusting like this go through the entire searching phase.

By introducing the architecture parameter β in the proposed hyper-cell, we can dynamically adjust and search for the network depth as well as the downsampling strategy for real-time semantic segmentation.

Network Latency Cost We define the set of cells in all hyper-cells in the initial phase as P , after optimization, the number of the set is reduced and the new set is marked as \bar{P} . For the current architecture (α, β) containing several hyper-cells, the total latency excludes the pruned cells and can be calculated as:

$$Latency(\alpha, \beta) = \sum_{p \in \bar{P}} lat_p \quad (7)$$

where \bar{P} is the set of remaining cells in all hyper-cells of architecture (α, β) . The lat_p is the latency of cell p . Thus the total loss function is formulated as:

$$L((\alpha, \beta), w) = CE((\alpha, \beta), w) + \gamma \log(Latency((\alpha, \beta))) \quad (8)$$

where $CE((\alpha, \beta), w)$ is the cross-entropy loss of architecture (α, β) with network weight w and γ controls the magnitude of latency term (i.e. balance the trade-off between accuracy and speed).

3.3 Network-Level Auto Feature Aggregation

For remedying the loss of spatial details in real-time segmentation networks due to fast downsampling, we propose the aggregation cell to automatically aggregate features by optimal operations from different levels in the network. The aggregation cell seamlessly integrates the outputs of above

hyper-cells, and the outputs of the early hyper-cells compensate for the loss of spatial details.

The structure of the proposed aggregation cell is shown in Figure 4. The aggregation cell takes three hyper-cells’ outputs with different resolutions as its inputs, thus the aggregation cell is designed to combine multi-scale features (i.e. low-level spatial details and high-level semantic context). The aggregation cell is designed as a directed acyclic graph consisting of M nodes and E edges, each node is a latent representation (i.e. feature map) and each directed edge is associated with some candidate operations. As shown in Figure 4, each edge’s stride is set to 1, unless explicitly specified by “s=2” (stride 2), which acts as the downsampling connection. The output of the aggregation cell is designed as the concatenation of the final feature maps from three hyper-cells. We use the same sampling and optimization process as intra-cell search in Section 3.1 to optimize the aggregation cell’s architecture parameter.

Given the candidate operation set, the aggregation cell also efficiently enlarges the receptive field of the network. Moreover, the aggregation cell is designed for effectively improving the segmentation accuracy, thus we introduce it without the latency constraint. For the operation set of the aggregation cell, we collect following 5 kinds of operations:

- 1×1 conv, repeat 2
- 3×3 conv, repeat 2
- 3×3 dilated separable conv, dilation=2, repeat 2
- 3×3 dilated separable conv, dilation=4, repeat 2
- 3×3 dilated separable conv, dilation=8, repeat 2

4 Experiments

To verify the effectiveness and superiority of our joint search framework, we compare our AutoRTNet with other state-of-the-art methods on two challenging benchmark datasets: Cityscapes [9] and CamVid [3]. Moreover, we conduct a series of ablation studies to verify the effectiveness of the proposed hyper-cell and aggregation cell. Finally, we provide an in-depth analysis about the detailed architecture of AutoRTNet.

4.1 Implementation Details

Searching For the searching process, the whole network contains three hyper-cells and the initial cell numbers in these hyper-cells are $\{5, 10, 10\}$, respectively. The intermediate node number of the cell is set to 2. The initial channel number is 8, and the channels are $\times 3$ when downsampling in reduction cells. The search process, which is conducted on the Cityscapes dataset, runs 150 epochs with mini-batch size 16, which takes approximately 16 hours with 16 TitanXP GPU cards. Similar to FBNet [42], we postpone the training

of the hyper-cell architecture parameter β by 50 epochs to warm-up weight w and intra-cell architecture parameter α . The α and β are optimized by Adam, with initial learning rate 0.001, momentum (0.5, 0.999) and weight decay $1e-4$. The w is optimized using SGD with momentum 0.9, weight decay $1e-3$, and cosine learning scheduler that decays learning rate from 0.025 to 0.001. For Gumbel Softmax, we empirically set the initial temperature λ in equation (3) and (6) as 3.0, and gradually decrease to the minimum value of 0.03. We set the node number M and edge number E as 7 in the aggregation cell.

Retraining When the search process is over, the searched network is firstly pretrained on the ImageNet dataset from scratch. We then finetune the network on the specific segmentation dataset (i.e. Cityscapes or CamVid) for 200 epochs with mini-batch size 16. The base learning rate is 0.01 and the poly learning rate policy is adopted with power 0.9, together with momentum is 0.9 and weight decay is 0.0005. Following [44,49], we compute the loss function with the online bootstrapping strategy. Data augmentation contains random horizontal flip, random resizing with scale ranges in [0.5, 2.0], and random cropping into fix size for training.

4.2 Benchmarks and Evaluation Metrics

Cityscapes [9], a public street scene dataset, contains high quality pixel-level annotations of 5000 images with size 1024×2048 and 19,998 images with coarse annotations. 19 semantic classes are used for training and evaluation. CamVid [3] is another public dataset, and it contains 701 images in total. We follow the training/testing set split in [52,3], with 468 training and 233 testing labeled images. These images are densely labeled with 11 semantic class labels. We use three evaluation metrics, including mean of class-wise intersection over union (mIoU), network forward time (Latency), and Frames Per Second (FPS).

4.3 Real-time Semantic Segmentation Results

In this section, we compare the AutoRTNet with other real-time segmentation methods. We run all experiments based on Pytorch 0.4 [34] and measure the latency on an NVIDIA TitanXP GPU card under CUDA 9.0. For fair comparison, we directly quote the reported remeasured or estimated speed results on TitanXP of other algorithms mentioned in [52, 32]. For the AutoRTNet, we report the average inference time through 500 times. In this process, we don’t employ any test augmentation.

Results on Cityscapes. AutoRTNet-A and AutoRTNet-B are searched with latency term weight γ 0.01 and 0.001, respectively. We evaluate them on the Cityscapes test set.

Table 1: Accuracy and speed comparison of our method against other state-of-the-art methods on Cityscapes test set. Methods trained using both fine and coarse data are marked with *. The mark † represents the speed is remeasured by us on Titan XP.

Method	Input Size	mIoU (%)	Latency(ms)	FPS
FCN-8S [25]	512×1024	65.3	227.23	4.4
PSPNet [54]	713×713	81.2	1288.0	0.78
DeepLabV3* [6]	769×769	81.3	769.23	1.3
AutoDeepLab* † [22]	769×769	81.2	303.0	3.3
SegNet [1]	360×640	57.0	30.3	33
ENet [33]	360×640	58.3	12.7	78.4
SQ [40]	1024×2048	59.8	46.0	21.7
ERFNet [36]	512×1024	69.7	48.5	20.6
ICNet [53]	1024×2048	69.5	26.5	37.7
DF1-Seg [47]	768×1536	73.0	29.1	34.4
SwiftNet [32]	1024×2048	75.1	26.2	38.1
ESPNet [28]	512×1024	60.3	8.2	121.7
DFANet [18]	1024×1024	71.3	10.0	100.0
DFANet † [18]	1024×1024	71.3	20.6 †	48.5 †
BiSeNet [49]	768×1536	68.4	9.52	105.8
CAS [52]	768×1536	70.5	9.25	108.0
CAS* [52]	768×1536	72.3	9.25	108.0
AutoRTNet-A	768×1536	72.2	9.09	110.0
AutoRTNet-A*	768×1536	73.9	9.09	110.0
AutoRTNet-B	768×1536	74.3	14.0	71.4
AutoRTNet-B*	768×1536	75.8	14.0	71.4

The validation set is added for training before submitting to online Cityscapes server. Following [52, 49], we scale the resolution of the images from 1024×2048 to 768×1536 as inputs to measure the speed and accuracy. As shown in Table 1, our AutoRTNet achieves the best trade-off between accuracy and speed. AutoRTNet-A yields 72.2% mIoU while maintaining 110.0 FPS on the Cityscapes test set with only fine data and without any test augmentation. When the coarse data is added to the training set, the mIoU achieves 73.9%, which is the state-of-the-art trade-off for real-time semantic segmentation. Compared with BiSeNet [49] and CAS [52] which have comparable speed to us, AutoRTNet-A surpasses them by 3.8% and 1.7% in mIoU on the Cityscapes test set, respectively. Compared with other real-time segmentation methods (e.g. ENet [33], ICNet [53]), our AutoRTNet-A surpasses them in both speed and accuracy by a large margin. Moreover, our AutoRTNet-B achieves 74.3% and 75.8% mIoU (+coarse data) on the Cityscapes test set with 71.4 FPS, which is also the state-of-the-art real-time performance.

Results on CamVid. To validate the transferability of the networks searched by our framework, we directly transfer AutoRTNet-A and AutoRTNet-B, which are obtained on Cityscapes, to the CamVid dataset, as reported in Table 2. With 720×960 input images, AutoRTNet-A achieves 73.5% mIoU on CamVid test set with 140.0 FPS, which

is the state-of-the-art trade-off between accuracy and speed. AutoRTNet-B achieves 74.2% mIoU with 82.5 FPS. We also conduct the architecture search on CamVid ($\gamma = 0.1$) and name the resulting network AutoRTNet-C. Notably, AutoRTNet-C achieves amazing 250.0 FPS while maintaining 68.6% mIoU on the CamVid test set, which surpasses ICNet [53] (67.1% mIoU with 34.5 FPS) and DFANet [18] (64.7% mIoU with 120 FPS) significantly.

Method	mIoU (%)	Latency(ms)	FPS	Parameters (M)
SegNet [1]	55.6	34.01	29.4	29.5
ENet [33]	51.3	16.33	61.2	0.36
ICNet [53]	67.1	28.98	34.5	26.5
BiSeNet [49]	65.6	-	-	5.8
DFANet [18]	64.7	8.33	120	7.8
CAS [52]	71.2	5.92	169	-
AutoRTNet-A	73.5	7.14	140.0	2.5
AutoRTNet-B	74.2	12.1	82.5	3.9
AutoRTNet-C	68.6	4.0	250.0	1.4

Table 2: Results on the CamVid test set with resolution 720×960 .

Parameter Results For many real-time applications on computationally limited mobile platforms, which have restrictive memory constraints, thus model size (number of parameters) is also an important consideration. Table 2 shows the results of our AutoRTNet and other methods on CamVid test set. With only 2.5 million parameters, our AutoRTNet-A achieves impressive accuracy (i.e. 73.5% mIoU) on the CamVid test set, which significantly outperforms existing real-time segmentation networks. The model sizes of AutoRTNet B and C are 3.9M and 1.4M, respectively.

4.4 Ablation Study

The contribution of each component is investigated in following ablation studies on Cityscapes validation set. The latency term weight γ in equation (8) is set to 0.01 and all networks are firstly pretrained on ImageNet in following experiments for fair comparison, if not specially noted.

4.4.1 Comparison with Random Search

Table 3: Comparison with random search on the Cityscapes validation set.

Method	mIoU (%)	Latency (ms)	FPS
AutoRTNet	72.9	9.09	110.0
random search	66.7 ± 2.5	$11.4 \sim 16.2$	$87.5 \sim 61.2$

As discussed in [19,38], NAS is a specialized hyper-parameter optimization problem, while random search is a competitive baseline for the problem. We apply random search to semantic segmentation by randomly sampling ten architectures from our previously-defined search space. The whole search space contains intra-cell operation selection and hyper-cell depth decision, which is significantly challenging for random search to find a satisfactory network. As shown in Table 3, random search achieves average mIoU $66.7\% \pm 2.5\%$ on Cityscapes validation set with ImageNet pretrained, which is substantially lower than our AutoRTNet. The results also demonstrate the effectiveness of our search algorithm.

4.4.2 Hyper-Cell

Robustness Firstly, to verify the robustness of hyper-cell, we set different initial numbers of cells and different random seeds in the initialization phase. The network contains three hyper-cells and the initial cell numbers in hyper-cells are set to $\{a, b, c\}$, after optimization, the numbers of cells remaining in each hyper-cell are $\{\bar{a}, \bar{b}, \bar{c}\}$. As shown in Table 4, the experiments demonstrate that the hyper-cells are insensitive to both initial numbers of cells and random seeds, which verify the robustness of the hyper-cell.

Table 4: The optimization results of hyper-cells with different initial states and different random seeds.

Random seed setting	Initial phase $\{a, b, c\}$	After optimization $\{\bar{a}, \bar{b}, \bar{c}\}$	mIoU (%)	Frames Per Second (FPS)
2	{5, 10, 10}	{2, 4, 6}	72.9	110.0
2	{5, 10, 15}	{2, 4, 6}	73.0	106.5
2	{5, 15, 15}	{1, 4, 6}	72.5	102.8
2	{10, 10, 10}	{2, 4, 6}	72.8	112.3
1	{5, 10, 10}	{2, 4, 6}	73.0	113.6
3	{5, 10, 10}	{1, 4, 7}	72.8	107.9

Downsampling strategy To demonstrate the superiority of the downsampling strategy searched by hyper-cells, we compare the random downsampling position settings with the searched one. The total cell number is $12(\bar{a}+\bar{b}+\bar{c})$ searched by our framework, we fix the searched cell structures and only random change the downsampling positions (x, y, z) for fair comparison. The (x, y, z) represents the index positions of reduction cells in total 12 cells. After pretraining and retraining, the results in Table 5 demonstrate the superiority of the searched downsampling strategy through hyper-cells. Compared with the random ones, our hyper-cell achieves the best trade-off between accuracy and speed.

Table 5: Comparison to random downsampling strategy.

Downsampling Positions	mIoU(%)	FPS	Downsampling Design Rule
(1, 3, 7)	72.9	106	Hyper Cell
(1, 5, 10)	71.6	90.2	Random
(3, 5, 8)	72.5	75.7	Random
(1, 2, 4)	68.4	125	Random

4.4.3 Aggregation Cell

To demonstrate the effectiveness of the proposed aggregation cell, we conduct a series of experiments with different strategies: a) without multi-scale feature aggregation; b) with random designed aggregation cell using random selected operations from the aggregation cell’s search space; c) with searched aggregation cell (i.e. our AutoRTNet-A). Among them, the result of the random aggregation cell is the average result over ten repeated random experiments and the results are shown in Table 6. Overall, the searched aggregation cell successfully boosts up the mIoU from 69.9% to 72.9% on Cityscapes val set. Particularly, the searched aggregation cell surpasses the random one 1.5% performance gains.

Table 6: Ablation study for effectiveness of the aggregation cell.

Methods	mIoU (%)
a) without aggregation cell	69.9
b) random aggregation cell	71.4
c) AutoRTNet-A	72.9

4.4.4 Hyper-Cell Searching Process

To better analyze how hyper-cell works throughout the whole searching process, we visualize the number of cells of each hyper-cell after the warm-up phase, as depicted in Figure 5. The initial cell numbers are $\{5, 10, 10\}$ and eventually converges to $\{2, 4, 6\}$ in three hyper-cells. The blue lines from top to bottom denote the actual cell numbers according to current architecture parameter β of each hyper-cell, and red curves represent the mathematical expectation values of current cell numbers. We observe that the framework actively explores different cell numbers (i.e. different depths) in each hyper-cell at the early stage of searching, and the expectation values of cell numbers also change gradually. The cell numbers progressively become stable towards the final architecture in the late stage of searching, and the actual cell number lines gradually get close to the expectation curves. Another interesting observation is that hyper-cell #1 finds

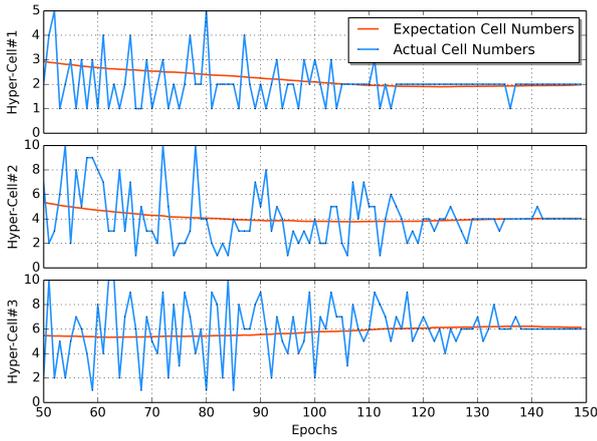


Fig. 5: Illustration of cell numbers in hyper-cells during the search process. Blue lines from top to bottom denote the actual cell number changing in each hyper-cell with the increase of epochs, and red curves represent the mathematical expectation values of the current cell numbers in hyper-cells.

its optimal depth much earlier than the other ones, indicating the search process follows a shallow-to-deep manner as we expected.

4.4.5 Different Latency Settings

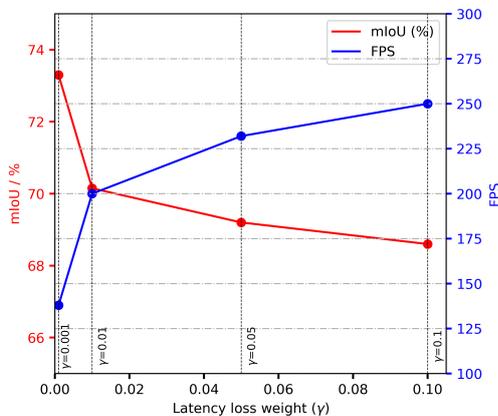


Fig. 6: The results of different latency settings on CamVid dataset.

Our joint search framework searches for the optimal network architectures under different latency settings (i.e. loss weight γ). In Section 4.3, AutoRTNet A and B are searched with $\gamma = 0.01$ and 0.001 on the Cityscapes dataset, respectively, which demonstrates the flexibility of our framework. We also conduct the architecture search on the CamVid dataset with different latency settings, and the results are as depicted in Figure 6. The networks searched with $\gamma = 0.001$,

0.01, 0.05, 0.1 achieve 73.3%, 70.2%, 69.2%, 68.6% mIoU and 138.0, 200.2, 232.1, 250.0 FPS on the CamVid test set, respectively. Notably, our AutoRTNet achieves 250.0 FPS while maintaining 68.6% mIoU, which surpasses ICNet (67.1% mIoU with 34.5 FPS) and DFANet (64.7% mIoU with 120 FPS) significantly.

4.4.6 Insights from Searched AutoRTNet

Finally, we provide an in-depth analysis of the AutoRTNet-A searched by our framework. We use the NAS methods to search the suitable architectures for specific tasks, likewise, we should understand why the searched network works well and it will guide the hand-designed process in turn. We have the following three important observations.

Early Downsampling We notice that the searched downsampling strategy is stable and reasonable. As shown in Table 4, in the first hyper-cell, whether the initialized cell number is 5 or 10, the final number is at most 2 after the optimization. The reason is that the visual information is highly spatially redundant, thus can be compressed into a more efficient representation. Under the latency constraints, the searched downsampling strategy is as we expected and follows the early downsampling [33] priori knowledge.

Suitable receptive field The suitable receptive field size [26] is crucial for semantic segmentation. Too large receptive field may introduce some extra noise or negative interference, and the network lacks of capturing enough context information if it is too small. During the search process, the AutoRTNet continuously adjusts the operation selection to determine the final receptive field, for example, in Figure 7, in the optimized aggregation cell, the operations from the outputs of third hyper-cell always choose the operation with $dil=4$ rather than $dil=2$ or 8 also in the search space. So we should choose right operations for suitable receptive field in hand-designed real-time semantic segmentation networks.

Operation selection The early operations act as good feature extractors, as shown in Figure 7, the selection of operations in early stage always tends to $conv3 \times 3$. The middle and deep layers have the diversity of operation selection. When performing multi-scale feature aggregation in aggregation cell, as shown in Figure 7, we clearly found that the deeper layers enjoy dilated convolution, while the lower layers only prefer the common convolution operations.

4.5 Detailed Time and GPU Information for Fair Comparison

The inference time or FPS is influenced by the GPU device and the input image size of the model. Here we list detailed information of previous approaches in Table 7 for readers

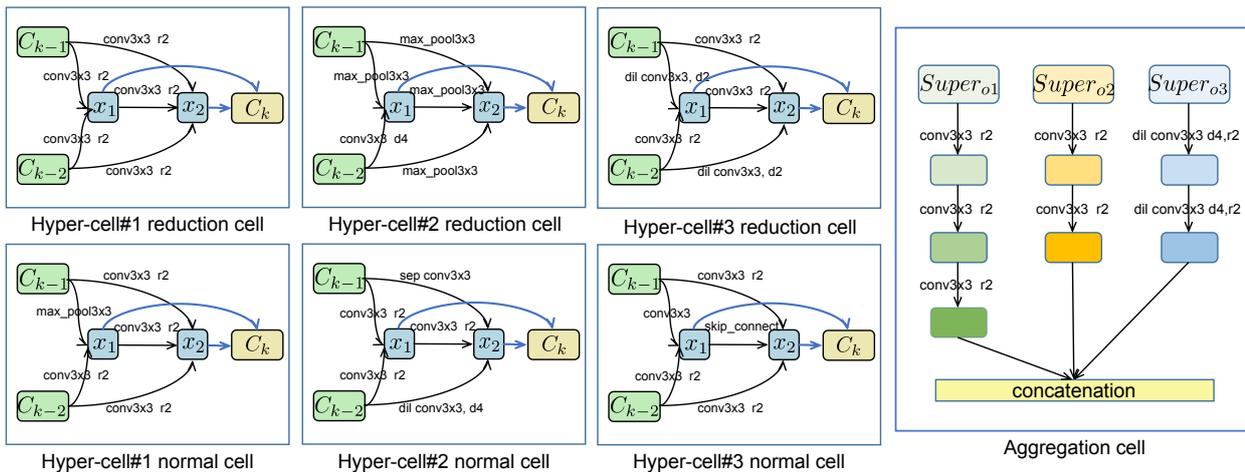


Fig. 7: Illustration of the detailed AutoRTNet-A architecture. The structures of the reduction cells and normal cells in three hyper-cells are shown in the figure respectively. The structure of searched aggregation cell is shown on the right. **Best viewed in color.**

as reference. Our GPU device is Nvidia TitanXP GPU. For fair comparison, we directly quote the reported remeasured or estimated results on TitanXP of other algorithms in CAS [52] and SwiftNet [32] paper. And we remeasure the speed of the methods based on our implementation if the original speed was reported on different GPUs and not mentioned in CAS [52] and SwiftNet [32]. Note that our implementations and speed measurements do not use TensorRT optimizations.

About the speed gap between the original DFANet [18] and our measurement: The speed of DFANet is reported on TitanX GPU, and also not mentioned in CAS [52] and SwiftNet [32]. So we carefully remeasure the inference time on TitanXP for fair comparison. There still has a speed gap between the original speed and the one we measured, we suspect that this is caused by the inconsistency of the implementation platform. We reimplement the DFANet using official PyTorch [34], and they measure it on their own framework in which the depth-wise separable convolution is more fully optimized.

4.6 Full Quantitative Results on Cityscapes and CamVid Dataset

Here we provide detailed quantitative results of per-class mIoU on the Cityscapes and CamVid datasets. Moreover, we provide the performance of the AutoRTNet on the full-resolution Cityscapes validation set.

4.6.1 Cityscapes Dataset

Compared with other methods, our AutoRTNet-A achieves overall 72.2% mIoU with 110.0 FPS, which is the state-

of-the-art trade-off between accuracy and speed. The per-class accuracy values are shown in Table 8. In comparison with other methods with public per-class accuracy on the Cityscapes test set, our predictions are more accurate in 13 out of 19 classes. AutoRTNet-A achieves slight improvements on the general classes (Road, Sidewalk, Building, Terrain, Car, etc.), while obtaining a significant accuracy improvement on the challenging classes (Truck, Motorbike, Train, Fence, Rider, etc.). AutoRTNet-B achieves 74.3% mIoU on Cityscapes test set with 71.4 FPS.

Cityscapes contains high-resolution 1024×2048 images, which make it a big challenge for real-time semantic segmentation. ICNet [53] focuses on building a practically fast semantic segmentation system with high-resolution image inputs while accomplishing high-quality results. SwiftNet [32] and CAS [52] also perform the experiments on Cityscapes with full-resolution image inputs. In this part, we compare with these methods on the Cityscapes validation set and the results are shown in Table 9. We refer to the speed scaling factors on different GPUs in SwiftNet [32] paper and estimate the speed values of ICNet, SwiftNet, CAS on Titan XP GPU.

AutoRTNet-A Our AutoRTNet-A achieves 75.0% mIoU and delivers 62.7 FPS on full-resolution Cityscapes val set (i.e. 1024×2048). To the best of our knowledge, the real-time performance of AutoRTNet-A outperforms all existing real-time methods. Compared with ICNet, AutoRTNet surpasses it by 5.5% in mIoU with a faster inference speed. Moreover, AutoRTNet outperforms SwiftNet and CAS by 0.6% and 1.0% in mIoU, and has a great advantage in inference speed (i.e. 62.7 FPS vs 38.1 FPS, 62.7 FPS vs 45.2 FPS).

Table 7: The detailed information of our AutoRTNet and other state-of-the-art methods on the Cityscapes test set. Methods trained using both fine and coarse data are marked with *. The mark † represents the speed is remeasured by us.

Method	Input Size	mIoU (%)		Latency (ms) on TitanXP	FPS on TitanXP	Original Results	
		val	test			FPS	GPU
FCN-8S [25]	512x1024	-	65.3	227.23	4.4	-	-
PSPNet [54]	713x713	-	81.2	1288.0	0.78	-	-
DeepLabV3* [6]	769x769	-	81.3	769.23	1.3	-	-
AutoDeepLab* [22]	769x769	-	81.2	303.0	3.3	-	-
SegNet [1]	640x320	-	57.0	30.3	33	-	-
SQ [40]	1024x2048	-	59.8	46.0	21.7	-	Titan X M
ENet [33]	640x320	-	58.3	12.7	78.4	135.4	Titan X
ERFNet [36]	1024x512	-	69.7	48.5	20.6	11.2	TitanX M
ICNet [53]	1024x2048	67.7	69.5	26.5	37.7	30.3	TITAN X(M)
SwiftNet [32]	1024x2048	74.4	75.1	26.2	38.1	34.0	GTX 1080Ti
DF1-Seg [47]	768x1536	74.1	73.0	29.1	34.4	30.7	GTX 1080Ti
ESPNet [28]	1024x512	-	60.3	8.2	121.7	112	TitanX
BiSeNet [49]	768x1536	69.0	68.4	9.52	105.8	105.8	TitanXP
DFANet [18]	1024x1024	-	71.3	20.6†	48.5†	100.0	TitanX
CAS [52]	768x1536	71.6	70.5	9.25	108.0	108.0	TitanXP
CAS* [52]	768x1536	72.5	72.3	9.25	108.0	108.0	TitanXP
AutoRTNet-A	768x1536	72.9	72.2	9.09	110.0	110.0	TitanXP
AutoRTNet-A*	768x1536	74.5	73.9	9.09	110.0	110.0	TitanXP
AutoRTNet-B	768x1536	74.7	74.3	14.0	71.4	71.4	TitanXP
AutoRTNet-B*	768x1536	76.0	75.8	14.0	71.4	71.4	TitanXP

Table 8: Detailed performance comparison of our AutoRTNet-A with other state-of-the-art methods on the Cityscapes test set.

Method	Road	Sidewalk	Building	Wall	Fence	Pole	Traffic light	Traffic sign	Vegetation	Terrain	Sky	Person	Rider	Car	Truck	Bus	Train	Motorcycle	Bicycle	Mean IOU(%)	FPS
SegNet	96.4	73.2	84.0	28.4	29.0	35.7	39.8	45.1	87.0	63.8	91.8	62.8	42.8	89.3	38.1	43.1	44.1	35.8	51.9	57.0	33
ENet	96.3	74.2	75.0	32.2	33.2	43.4	34.1	44.0	88.6	61.4	90.6	65.5	38.4	90.6	36.9	50.5	48.1	38.8	55.4	58.3	78.4
ICNet	97.1	79.2	89.7	43.2	48.9	61.5	60.4	63.4	91.5	68.3	93.5	74.6	56.1	92.6	51.3	72.7	51.3	53.6	70.5	69.5	37.7
ESPNet	97.0	77.5	76.2	35.0	36.1	45.0	35.6	46.3	90.8	63.2	92.6	67.0	40.9	92.3	38.1	52.5	50.1	41.8	57.2	60.3	121.7
ERFNet	97.9	82.1	90.7	45.2	50.4	59.0	62.6	68.3	91.9	69.4	94.2	78.5	59.8	93.4	52.3	60.8	53.7	49.9	64.2	69.7	20.6
BiSeNet	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	68.4	105.8
DFANet	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	71.3	83.1
CAS	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	70.5	108.0
AutoRTNet-A	98.5	84.9	91.4	45.9	53.0	52.2	60.3	67.3	91.5	70.6	93.9	78.2	62.5	95.3	63.7	74.5	63.9	56.8	67.0	72.2	110.0
AutoRTNet-B	98.4	86.6	91.2	52.2	54.9	58.5	63.7	68.4	91.5	71.5	94.9	79.3	61.4	95.4	65.9	78.0	69.8	59.6	69.8	74.3	71.4

AutoRTNet-B Our AutoRTNet-B achieves 76.8% mIoU with 45.6 FPS on full-resolution Cityscapes val set, which is the state-of-the-art real-time performance. Compared with SwiftNet and CAS which have a little bit slower speed than us, our AutoRTNet-B surpasses them by 2.4% and 2.8% in mIoU, respectively.

4.6.2 CamVid Dataset

As shown in Table 10, with 720×960 input images, the searched AutoRTNet-A achieves 73.5% mIoU with 140.0 FPS, which is the state-of-the-art trade-off between accuracy and speed on the CamVid test set. In comparison with other

methods, the predictions of our AutoRTNet-A are more accurate in 7 out of the 11 classes. More importantly, the inference speed of AutoRTNet-A achieves 140 FPS, which is very impressive compared with other methods. (e.g. SegNet 29.4 FPS, ENet 61.2 FPS, ICNet 34.5 FPS). The per-class accuracy of AutoRTNet B and C are also shown in Table 10.

4.6.3 Visual Segmentation Results

We provide some visual prediction results on both Cityscapes and CamVid datasets here. As shown in Figure 8 and Figure 9, the columns correspond to input image, ground truth, the prediction of ICNet, and the prediction of our AutoRTNet-

Table 9: Accuracy and speed comparison on the Cityscapes validation set with image resolution 1024×2048 .

Method	Input Size	mIoU (%)	FPS on Titan XP	Original Results	
				FPS	GPU
SQ [40]	1024x2048	59.8	-	-	Titan X(M)
ICNet [53]	1024x2048	69.5	55.6	30.3	Titan X(M)
SwiftNet [32]	1024x2048	74.4	38.1	34.0	GTX 1080Ti
CAS [52]	1024x2048	74.0	45.2	34.2	GTX 1070
AutoRTNet-A	1024x2048	75.0	62.7	62.7	Titan XP
AutoRTNet-B	1024x2048	76.8	45.6	45.6	Titan XP

Table 10: Detailed performance comparison of our AutoRTNet-A with other state-of-the-art methods on the CamVid test set.

Method												Mean IOU(%)	FPS
	Building	Tree	Sky	Car	Sign	Road	Pedestrian	Fence	Pole	Sidewalk	Bicyclist		
SegNet	88.8	87.3	92.4	82.1	20.5	97.2	57.1	49.3	27.5	84.4	30.7	55.6	29.4
ENet	74.7	77.8	95.1	82.4	51.0	95.1	67.2	51.7	35.4	86.7	34.1	51.3	61.2
ICNet	-	-	-	-	-	-	-	-	-	-	-	67.1	34.5
BiseNet-Xception39	82.2	74.4	91.9	80.8	42.8	93.3	53.8	49.7	25.4	77.3	50.0	65.6	-
BiseNet-Res18	83.0	75.8	92.0	83.7	46.5	94.6	58.8	53.6	31.9	81.4	54.0	68.7	-
DFANet	-	-	-	-	-	-	-	-	-	-	-	64.7	120.0
CAS	-	-	-	-	-	-	-	-	-	-	-	71.2	169.0
AutoRTNet-A	88.1	78.4	91.7	93.0	44.0	96.2	66.8	60.8	33.0	88.6	67.4	73.5	140.0
AutoRTNet-B	89.1	79.3	91.7	92.5	41.0	96.8	64.7	66.5	32.2	90.0	72.1	74.2	82.5
AutoRTNet-C	87.5	77.9	91.7	90.6	30.4	95.6	62.2	47.1	25.2	87.5	58.9	68.6	250.0

A. Compared with ICNet, AutoRTNet-A produces more accurate and detailed results with faster inference speed. For example, AutoRTNet-A captures small objects in more details (e.g. traffic light in Figure 8, poles in Figure 9) and generates “smoother” results on object boundaries (e.g. rider, fence in Figure 8, car in Figure 9).

5 Conclusion

In this paper, we propose a novel joint search framework which covers all three main aspects of the design philosophy for real-time semantic segmentation networks. The framework searches for building blocks, network depth, downsampling strategy, and feature aggregation way simultaneously. The hyper-cell is proposed for searching for the network depth and downsampling strategy in an adaptive manner, and the aggregation cell is introduced for automatic multi-scale feature aggregation. Extensive experiments on both Cityscapes and CamVid datasets demonstrate the superiority and effectiveness of our approach.

References

- Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017.
- Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. *arXiv preprint arXiv:1611.02167*, 2016.
- Gabriel J Brostow, Jamie Shotton, Julien Fauqueur, and Roberto Cipolla. Segmentation and recognition using structure from motion point clouds. In *European conference on computer vision*, pages 44–57. Springer, 2008.
- Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. *arXiv:1812.00332*, 2018.
- Liang-Chieh Chen, Maxwell Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jon Shlens. Searching for efficient multi-scale architectures for dense image prediction. In *Advances in Neural Information Processing Systems*, pages 8699–8710, 2018.
- Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings*

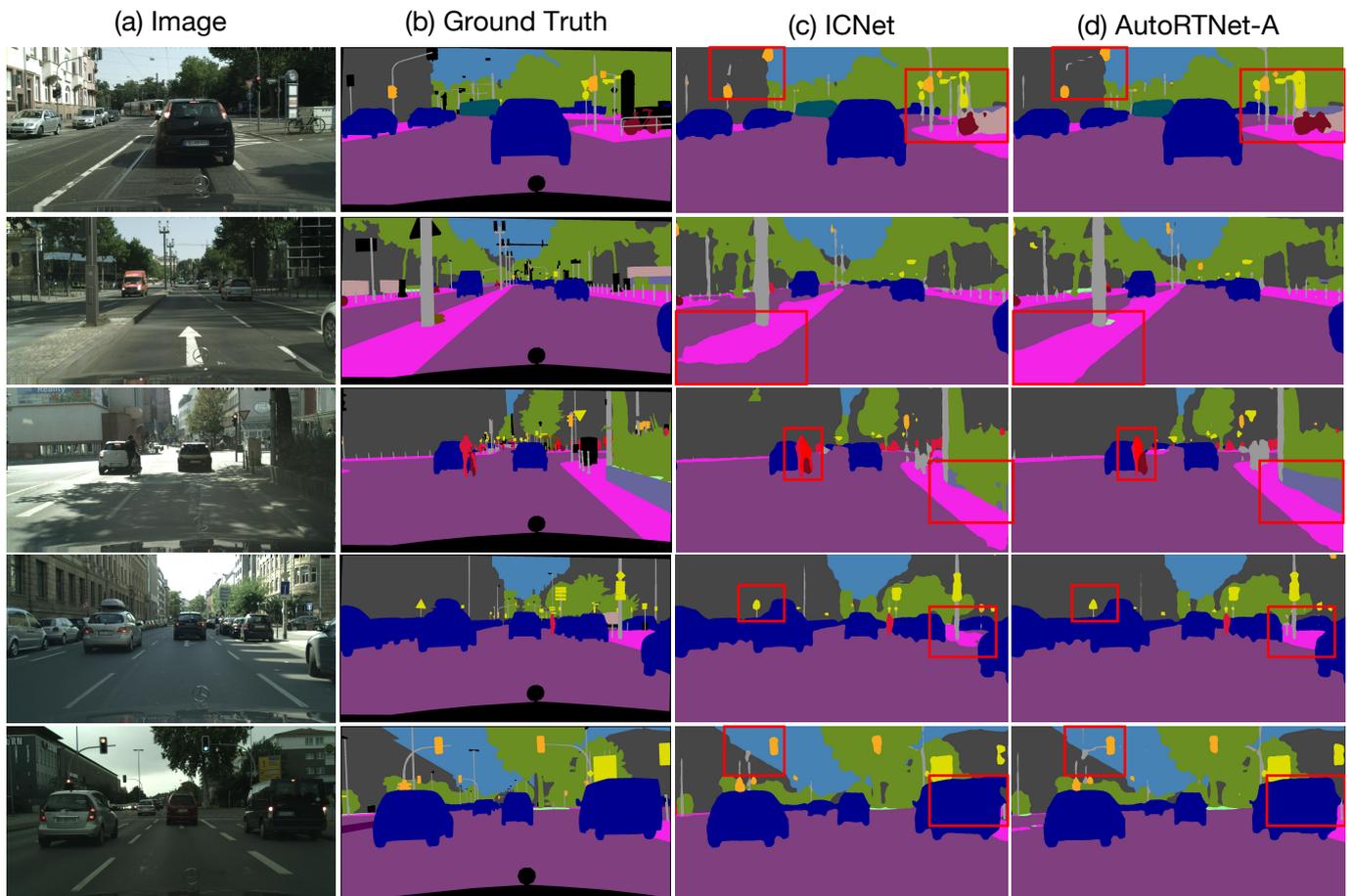


Fig. 8: Visual segmentation results on Cityscapes validation set. (a) Image. (b) Ground Truth. (c) ICNet. (d) AutoRTNet-A.

- of the European conference on computer vision (ECCV), pages 801–818, 2018.
8. François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.
 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.
 10. Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
 11. 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 Conference on Computer Vision and Pattern Recognition*, pages 3146–3154, 2019.
 12. Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V. Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
 13. 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.
 14. Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco An-

- dreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
15. Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
 16. Gen Li and Joongkyu Kim. Dabnet: Depth-wise asymmetric bottleneck for real-time semantic segmentation. In *British Machine Vision Conference*, 2019.
 17. Hanchao Li, Pengfei Xiong, Jie An, and Lingxue Wang. Pyramid attention network for semantic segmentation. *arXiv preprint arXiv:1805.10180*, 2018.
 18. Hanchao Li, Pengfei Xiong, Haoqiang Fan, and Jian Sun. Dfanet: Deep feature aggregation for real-time semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9522–9531, 2019.
 19. Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search. *CoRR*, abs/1902.07638, 2019.
 20. Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1925–1934, 2017.
 21. Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.
 22. Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation.

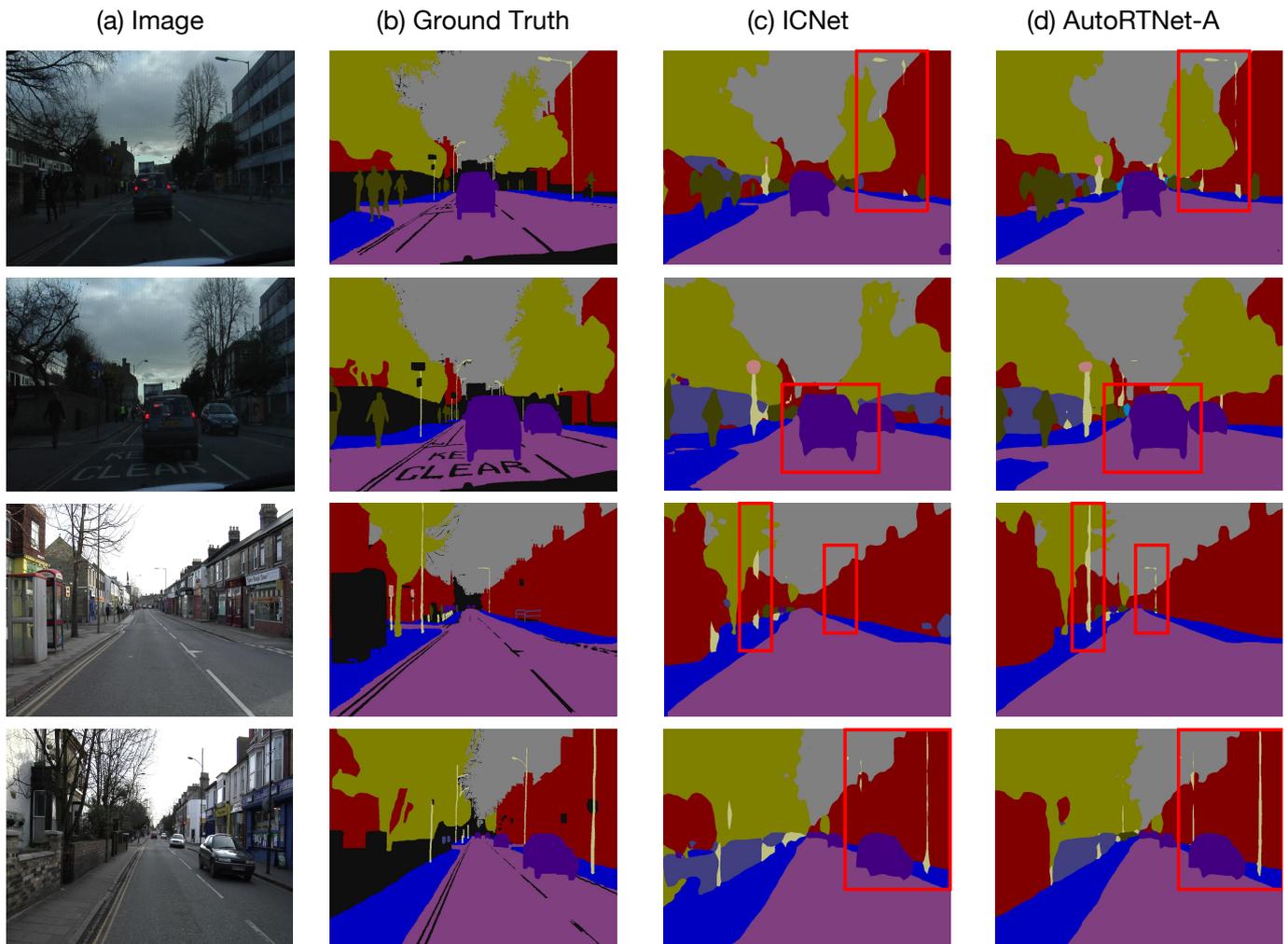


Fig. 9: Visual segmentation results on CamVid test set. (a) Image. (b) Ground Truth. (c) ICNet. (d) AutoRTNet-A.

- In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 82–92, 2019.
23. Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In *ICLR*, 2019.
 24. Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016.
 25. Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
 26. Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In *Advances in neural information processing systems*, pages 4898–4906, 2016.
 27. Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
 28. Sachin Mehta, Mohammad Rastegari, Anat Caspi, Linda Shapiro, and Hannaneh Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 552–568, 2018.
 29. Risto Miiikkulainen, Jason Liang, Elliot Meyerson, Aditya Rawal, Daniel Fink, Olivier Francon, Bala Raju, Hormoz Shahrzad, Arshak Navruzyan, Nigel Duffy, et al. Evolving deep neural networks. In *Artificial Intelligence in the Age of Neural Networks and Brain Computing*, pages 293–312. Elsevier, 2019.
 30. Vladimir Nekrasov, Hao Chen, Chunhua Shen, and Ian Reid. Fast neural architecture search of compact semantic segmentation models via auxiliary cells. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9126–9135, 2019.
 31. Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han. Learning deconvolution network for semantic segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 1520–1528, 2015.
 32. 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 Conference on Computer Vision and Pattern Recognition*, pages 12607–12616, 2019.
 33. Adam Paszke, Abhishek Chaurasia, Sangpil Kim, and Eugenio Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:1606.02147*, 2016.
 34. Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison,

- Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
35. Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4780–4789, 2019.
 36. 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.
 37. Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018.
 38. Christian Sciuto, Kaicheng Yu, Martin Jaggi, Claudiu Musat, and Mathieu Salzmann. Evaluating the search phase of neural architecture search. 2019.
 39. Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2820–2828, 2019.
 40. Michael Trembl, José Arjona-Medina, Thomas Unterthiner, Rupesh Durgesh, Felix Friedmann, Peter Schuberth, Andreas Mayr, Martin Heusel, Markus Hofmarcher, Michael Widrich, et al. Speeding up semantic segmentation for autonomous driving. In *MLITS, NIPS Workshop*, volume 2, page 7, 2016.
 41. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
 42. Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In *CVPR*, pages 10734–10742, 2019.
 43. Bichen Wu, Yanghan Wang, Peizhao Zhang, Yuandong Tian, Peter Vajda, and Kurt Keutzer. Mixed precision quantization of convnets via differentiable neural architecture search. *arXiv preprint arXiv:1812.00090*, 2018.
 44. Zifeng Wu, Chunhua Shen, and Anton Van Den Hengel. High-performance semantic segmentation using very deep fully convolutional networks. 2016.
 45. Zifeng Wu, Chunhua Shen, and Anton van den Hengel. Real-time semantic image segmentation via spatial sparsity. *arXiv preprint arXiv:1712.00213*, 2017.
 46. Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. SNAS: stochastic neural architecture search. In *ICLR*, 2019.
 47. Zheng Pan Jiashi Feng Xin Li, Yiming Zhou. Partial order pruning: for best speed/accuracy trade-off in neural architecture search. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
 48. Hang Xu, Lewei Yao, Wei Zhang, Xiaodan Liang, and Zhenguo Li. Auto-fpn: Automatic network architecture adaptation for object detection beyond classification. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
 49. 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.
 50. Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 472–480, 2017.
 51. Hang Zhang, Han Zhang, Chenguang Wang, and Junyuan Xie. Co-occurrent features in semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 548–557, 2019.
 52. Yiheng Zhang, Zhaofan Qiu, Jingen Liu, Ting Yao, Dong Liu, and Tao Mei. Customizable architecture search for semantic segmentation. In *CVPR*, pages 11641–11650, 2019.
 53. 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.
 54. Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017.
 55. Hengshuang Zhao, Yi Zhang, Shu Liu, Jianping Shi, Chen Change Loy, Dahua Lin, and Jiaya Jia. PSANet: Point-wise spatial attention network for scene parsing. In *ECCV*, 2018.
 56. Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.
 57. Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018.