

# GSANet: Semantic Segmentation with Global and Selective Attention

Qingfeng Liu, Mostafa El-Khamy, Dongwoon Bai, Jungwon Lee

## Abstract

This paper proposes a novel deep learning architecture for semantic segmentation. The proposed Global and Selective Attention Network (GSANet) features Atrous Spatial Pyramid Pooling (ASPP) with a novel sparsemax global attention and a novel selective attention that deploys a condensation and diffusion mechanism to aggregate the multi-scale contextual information from the extracted deep features. A selective attention decoder is also proposed to process the GSA-ASPP outputs for optimizing the softmax volume. We are the first to benchmark the performance of semantic segmentation networks with the low-complexity feature extraction network (FXN) MobileNetEdge, that is optimized for low latency on edge devices. We show that GSANet can result in more accurate segmentation with MobileNetEdge, as well as with strong FXNs, such as Xception. GSANet improves the state-of-art semantic segmentation accuracy on both the ADE20k and the Cityscapes datasets.

## I. INTRODUCTION

Semantic segmentation is an important computer vision task that has many industrial applications, like autonomous driving, medical image analysis and mobile phone cameras. There has been significant progress by recent research works that proposed deep neural networks for semantic segmentation. The accuracy of pixel-level semantic segmentation improves with the knowledge of multiscale contextual information. PSPNet [1] performs spatial pyramid pooling (SPP) at several grid scales. Atrous spatial pyramid pooling (ASPP) has been inspired by SPP and utilized in DeepLabV3+ [2]. ASPP collects information from the extracted features at different scales and receptive fields using dilated convolutions with different dilation rates. Another approach, called self-attention [3] gathered momentum, as it has been shown it can capture the important information without being constrained to a local regular grid. Self-attention [4], [5] and its low cost variants [6], [7] have been adopted for semantic segmentation.

We address in this paper two important issues that are crucial to more accurate semantic segmentation, and that have previously received little or no attention. The first issue is how to fuse the captured multiscale contextual information. DeepLabV3+ [2] combines the contextual information using a simple  $1 \times 1$  convolutional filter that cannot model the relative importance of the different scales at the location of interest. The second issue is how to obtain relevant global contextual information. For example, the ASPP deployed in DeepLabV3 deploys global average pooling (GAP) to capture the mean feature as the global contextual information of the input feature, which will be the same at all spatial locations.

This paper proposes a novel architecture called the Global and Selective Attention Network (GSANet) to address these issues. GSANet constitutes of a novel Selective Attention (SA) module for the fusion of the contextual information from the multiple scales, while giving the appropriate relative importance to the different scales at the different spatial locations. The proposed SA-ASPP module aggregates multiscale contextual information from the Feature eXtraction Network (FXN) using atrous convolutions with different dilation rates. The proposed selective attention combines these features by calculating the relative importance of the multiscale features using a condensation and diffusion mechanism. GSANet also features a Global Attention Feature (GAF) module (instead of global average pooling) to provide a more relevant global feature with different global contextual information at the different spatial locations, and reflect the different perceptions of the global information when observed from different locations.

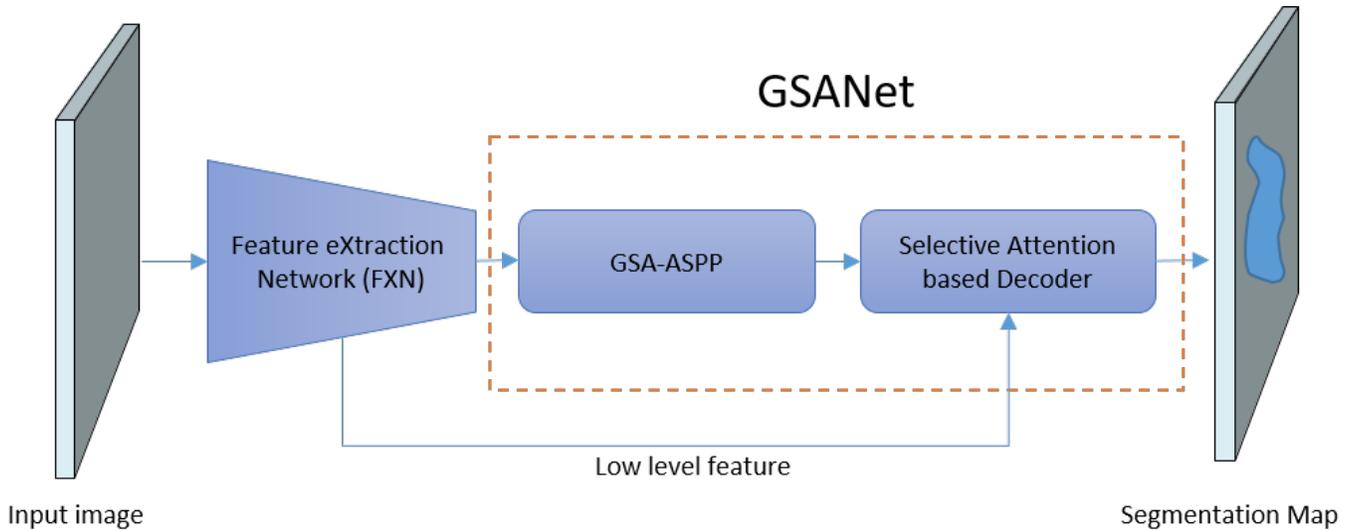


Fig. 1: The system architecture of the proposed GSA Net.

Another contribution is Sparsemax-GAF that deploys sparsemax normalization (instead of the popular softmax normalization) when calculating the GAF so as to suppress the noisy long range contextual information and boost the more prominent global information. The proposed GSA-ASPP calculates the sparsemax-GAF as the global features, and combines this global feature with the multiscale features extracted at different dilation rates using the SA. Please note that the selective attention module and the GAF module are not limited to the context of ASPP. They can serve to enhance any network modules that requires calculation of global contextual information and the fusion of multiscale information. Moreover, we propose the selective attention decoder (SA-Dec), that uses selective attention to combine the GSA-ASPP output together with some low-level features extracted by the FXN for better decoding.

Other contributions of this paper is the investigations of segmentation networks for best performance, as well as for efficiency. With the Xception network as the FXN, we show that the GSA Net provides the state-of-the-art (SOTA) accuracy on public datasets ADE20k and Cityscapes. We are also first to benchmark the performance of the light network MobileNetEdgeTPU [8] (MNEdge), that has been optimized for fast performance on the edge devices, in the semantic segmentation task. We show that with both the MNEdge FXN and the Xception FXN, GSA Net provides more accurate semantic segmentation predictions than the popular DeepLabV3+ framework.

The rest of this paper is organized as follows. Section II describes our proposed GSA Net and its novel components: the SA module in Section II-A, the sparsemax-GAF module in Section II-B, and the SA-Dec in Section II-C. Our experimental setup, implementation details, results, and ablation studies are presented in Section III. Conclusions and discussions are made in Section IV.

## II. GLOBAL AND SELECTIVE ATTENTION NETWORK

The proposed GASNet architecture is demonstrated in Fig. 1. First, deep features are extracted using the FXN. Any deep neural network which showed good accuracy on the classification task, can be used as the FXN. In this paper, we show results when the FXN is either the Xception or the MNEdge network. The extracted deep features are processed by the GSA-ASPP module, shown in Fig. 2, which is composed of two sub-modules, namely the SA-ASPP and the sparsemax-GAF module, illustrated in Fig. 2 and Fig. 3, respectively.

### A. Selective Attention ASPP (SA-ASPP)

The ASPP has multiple atrous filters with different dilation rates, as well as global filter, to aggregate dense information from the FXN features at different scales and receptive fields. This is crucial for pixel-

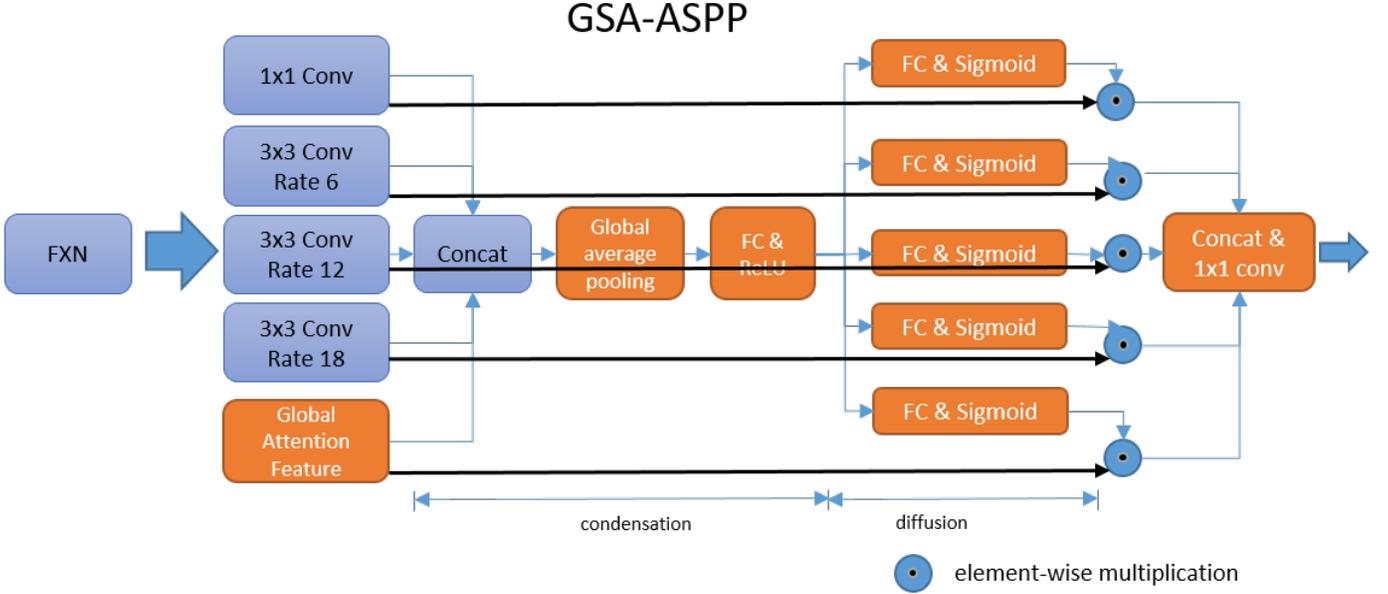


Fig. 2: The architecture of GSA-ASPP. It contains two modules, namely the SA-ASPP and sparsemax-GAF. FXN is the Feature eXtraction Network, such as MobileNetEdgeTPU etc.

level semantic segmentation, where it is need to understand the local and global context surrounding the pixel of interest, in order to classify it. The SA-ASPP is proposed to also use the context information, to give different spatial weights to these aggregated multiscale features, before combining them. First, in the condensation stage, all the  $W \times H \times C$  dense features obtained from the atrous and global filters are concatenated along the feature dimension to form a  $W \times H \times nC$  feature (for the case of  $n-1$  atrous filters and 1 global filter), and then condensed from all scales across the spatial and channel dimensions with global average pooling (GAP) followed by two fully connected layers with a ReLU nonlinearity to form the  $1 \times 1 \times \gamma$  condensate feature. In the diffusion stage,  $n$  channel attentions are obtained from the condensate feature using  $n$  separate fully connected layers, to diffuse the attention back to each branch. Due to the condensation and diffusion mechanism, it is easy to see that the channel attention for each branch does not only use the global contextual information from its own branch, but also borrows information from multiple branches. In Section III-C, we have show that both the condensation and diffusion stages are necessary to obtain better performance. For combining the different multiscale features, the element-wise dot product is applied per channel with the attention values and all the weighted branches are concatenated again to produce the new concatenated feature for further processing.

### B. Global Attention Feature (GAF)

The perception and understanding of the global context is crucial for accurate dense prediction. The ASPP [2] utilizes global average pooling (GAP) to extract the global contextual information that will be uniformly used at different spatial locations. One drawback of GAP is that it treats all pixels similarly. However, pixels at different spatial locations and belonging to different categories should perceive different global contextual information. Our proposed GAF is a global feature and addresses this issue by using self-attention. The global context for each pixel is derived from a linear combination of the features at all the other pixels using a distance metric determined by the degree of similarity between the pixels. If this feature is normalized by the ubiquitous softmax, we will refer to this as softmax-GAF. To give more weights to important information and discard the noisy information, we utilize the sparsemax projection [9] instead of softmax normalization to form the sparsemax-GAF. Sparsemax projection performs Euclidean projection

## Sparsemax-GAF

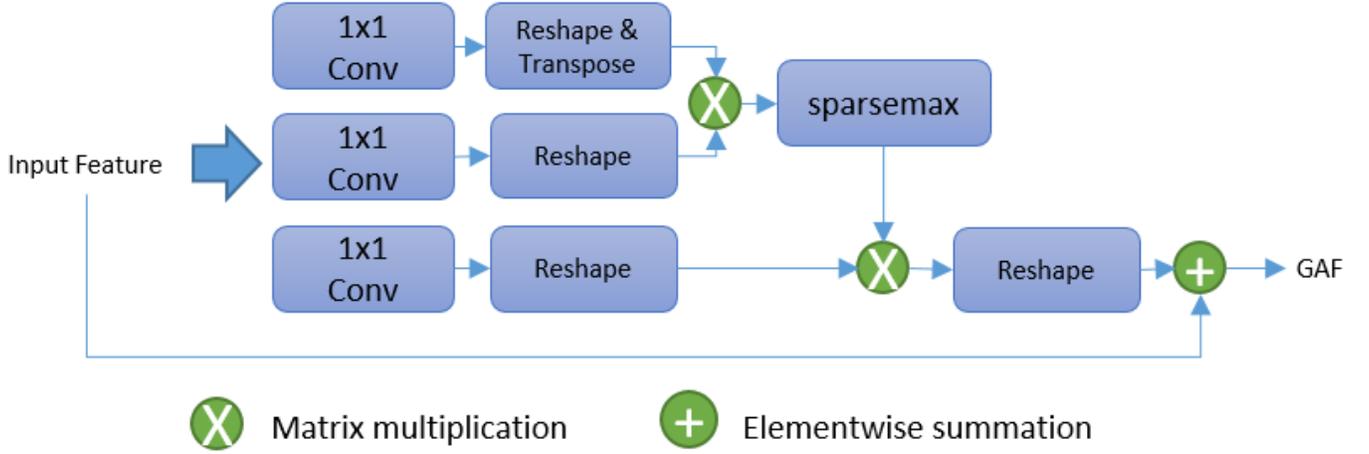


Fig. 3: The architecture of the network proposed to generate the Sparsemax Global Attention Feature (sparsemax-GAF).

of the input attention vector  $\mathbf{z} = [z_1, \dots, z_K]$  onto a probability simplex, as described mathematically below, with the notation that  $z_{(1)} \geq z_{(2)}, \dots, \geq z_{(K)}$  after sorting  $\mathbf{z}$  and

$$\text{sparsemax}(z_i) = \max(0, z_i - \tau(\mathbf{z})), \quad (1)$$

where the threshold  $\tau(\mathbf{z})$  is found by  $\tau(\mathbf{z}) = \frac{(\sum_{j \leq f(\mathbf{z})} z_{(j)})^{-1} - 1}{f(\mathbf{z})}$ ,

$$f(\mathbf{z}) = \max_{k \in \{1, 2, \dots, K\}} \left( 1 + kz_{(k)} > \sum_{j \leq k} z_{(j)} \right).$$

Due to the projection and thresholding, Sparsemax produces sparse probabilities that lead to a selective and more compact attention focus. The sparsemax-GAF strengthens the similarities and amplifies the attentions for similar pixels, while forcing zero attentions on dissimilar pixels and discarding noisy features.

The overall pipeline of the proposed GAF module is demonstrated in Fig. 3. First, the pairwise correlations are calculated for each pixel pair in the input feature map. Then, the raw attention map is calculated, where each value in the map corresponds to the correlation between two pixels. For sparsemax-GAF, the sparsemax projection is applied on each row of the map to obtain the normalized attention map. Third, the normalized attention map is projected onto the original input feature map using matrix multiplication to obtain the attended feature. The GAF is calculated as the sum of the shortcut input feature and the attended feature.

### C. Selective Attention Decoder (SA-Dec)

The proposed SA-Dec builds on top of the DeepLabV3+ decoder. The SA module is inserted to selectively combine, with attention, the different features that are input to the decoder. As shown in Fig. 4, the SA-Dec calculates the selective attention for two features, the low-level feature from the FXN and the multiscale aggregated feature that is output from the GSA-ASPP. The SA module follows the same condensation and diffusion method describes in Sec. II-A to modify the two features, before being processed by the DeepLabV3+ decoder operations constituting concatenation, filtering with convolutional layers, and upsampling.

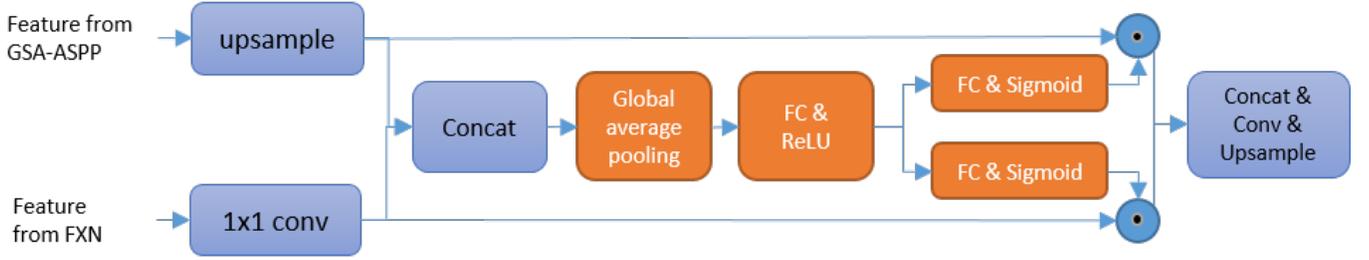


Fig. 4: The network architecture of the proposed Selective Attention Decoder (SA-Dec).

### III. EXPERIMENTS

In this section, extensive experiments are conducted to show the effectiveness of the proposed GSANet semantic segmentation architecture with different feature extraction networks.

#### A. Datasets

We benchmarked the performance of GSANet on two widely used public datasets, namely the Cityscapes and the ADE20k dataset.

**Cityscapes** [10] is specifically created for scene parsing. There are 5k high quality finely annotated images and 20k coarsely annotated images, which are taken on the street road with vehicles. The size of all the images in the dataset is  $2048 \times 1024$ . The finely annotated images are divided into 2975, 500, and 1525 splits for training, validation, and testing, respectively. The dataset contains 30 classes annotations in total, while only 19 classes are used for evaluation.

**ADE20K** [11] dataset is a large-scale dataset used in ImageNet Scene Parsing Challenge 2016. For the challenging version, there are 150 classes and the dataset is divided into 20k, 2k, and 3k images for training, validation, and testing, respectively. Most of the images in the dataset are taken from real life scenes and full of diversity, including objects and scenes of various scales, shapes, and colors etc.. Different from Cityscapes, both scenes and stuff are annotated in this dataset, posing more challenges to participated methods.

#### B. Implementation Details

We demonstrate the effectiveness of the GSANet semantic segmentation architecture with a low-complexity FXN, the MobileNetEdgeTPU (MNEdge), and show it can achieve competitive performance. We also demonstrate that GSANet beats current SOTA performance with the stronger Xception FXN.

The FXNs namely MNEdge and Xception, are pretrained on the ImageNet [12]. We use Stochastic Gradient Descent (SGD) to optimize our network, in which we set the initial learning rate to 0.01 for Cityscapes and 0.007 for ADE20K. During training, the learning rate is decayed according to the “poly” leaning rate policy, where the learning rate is multiplied by  $1 - (\frac{\text{iter}}{\text{max\_iter}})^{\text{power}}$  with  $\text{power} = 0.9$ . For Cityscapes, we randomly crop out half-resolution patches  $512 \times 1024$  from the original images as the inputs. While for ADE20K, we set the crop size to  $480 \times 640$ . For all datasets, we apply random scaling in the range of  $[0.5, 2.0]$ , random horizontal flip as additional data augmentation methods. For the ablation studies below, half resolution image of the Cityscapes dataset, namely  $512 \times 1024$  are taken as the input for models using the MNEdge FXN, and the FLOPs are calculated correspondingly. At the inference stage, we use single scale inference for models using the MNEdge FXN. For Xception based model, we adopted the left-right flipping and multiscale  $[0.75, 1.0, 1.25, 1.5, 1.75, 2.0]$  strategies for inference, following the methodology adopted by the state-of-the-art methods (*cf.* DeepLabV3+) for fair comparisons.

Method	mIoU (%)	# Params (M)	FLOPs (B)
MNEdge + ASPP	70.45	3.15	26.77
MNEdge + SA-ASPP (condensation only)	71.08	3.20	26.74
MNEdge + SA-ASPP (diffusion only)	71.10	3.18	26.76
MNEdge + SA-ASPP	71.21	3.20	26.77
MNEdge + softmax-GAF-ASPP	70.48	3.20	27.14
MNEdge + sparsemax-GAF-ASPP	71.19	3.20	27.18
MNEdge + GSA-ASPP	72.10	3.20	27.18
MNEdge + GSANet	<b>75.07</b>	3.43	37.64

TABLE I: Ablation study of GSANet on Cityscapes validation dataset using MobileNetEdgeTPU (MNEdge). The inference input resolution is half the resolution of the Cityscapes dataset, namely  $512 \times 1024$ .

Method	mIoU (%)
EncNet [13]	44.65
CCNet [6]	45.22
APNB [7]	45.24
OCNet [5]	45.45
Xception65-DeepLab V3+ [2]	45.65
Xception65-GSANet (ours)	<b>47.20</b>

TABLE II: Comparison to state-of-the-art on the validation set of ADE20K.

### C. Abalation Studies and Comparisons with SOTA Methods

Table I shows the effectiveness of GSANet with a low-complexity FXN, MNEdge. The output stride, namely the ratio of the input image size over the size of the FXN output feature, in these experiments is 16. SA-ASPP (condensation only) means we don't use the diffusion part using multiple FC layers, but use a single FC layer instead to obtain channel attentions. SA-ASPP (diffusion only) means no concatenation and global average pooling are used, and each branch has its own channel attention. Compared to the DeepLabV3+ ASPP, the results show SA-ASPP provides 1.5% point gain in the mean intersection over union (mIoU) accuracy. Table I also shows that sparsemax-GAF with the GSA-ASPP has 0.7% point

Method	mIoU (%)
CCNet [6]	81.3
DANet [4]	81.50
ACFNet [14]	81.46
Xception65-DeepLabV3+ [2]	80.42
Xception65-GSANet (ours)	<b>82.04</b>

TABLE III: Comparison to state-of-the-art on the validation set of Cityscapes.

improvement in mIoU over the GAP in SA-ASPP. We can also observe that the sparsemax-GAF gives better mIoU (0.7% point gain) than softmax-GAF.

Table II and Table III benchmarks GSANet with the strong FXN, Xception-65, on ADE20k and Cityscapes, respectively. Xception65-GSANet provides around 1.6% gain over Xception65-DeepLab V3+ on both the ADE20k and the Cityscapes datasets.

#### IV. CONCLUSIONS AND DISCUSSIONS

This paper proposes a novel semantic segmentation architecture, the Global and Selective Attention Network (GSANet). It enhances all components of the SOTA DeepLabV3+ architecture using attention modules. The multiscale aggregation in the ASPP takes into account the importance of the contextual information using the proposed selective attention which deploys condensation and diffusion modules. More relevant global information is extracted using the proposed sparsemax global attention feature. The decoder also deploys attention with condensation and diffusion to dope its different input features with extrinsic information from the other features. With both the low-complexity MobileNetEdge FXN and the strong Xception FXN, we show that GSANet gives better performance than DeepLabV3+ and achieves the state of art accuracy.

#### REFERENCES

- [1] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia, "Pyramid scene parsing network," in *Proc. CVPR*, 2017, pp. 6230–6239.
- [2] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proc. ECCV*, 2018, pp. 833–851.
- [3] 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*, 2017, pp. 5998–6008.
- [4] 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*, 2019, pp. 3146–3154.
- [5] Yuhui Yuan and Jingdong Wang, "Ocnet: Object context network for scene parsing," *CoRR*, vol. abs/1809.00916, 2018.
- [6] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu, "Ccnet: Criss-cross attention for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 603–612.
- [7] Zhen Zhu, Mengde Xu, Song Bai, Tengpeng Huang, and Xiang Bai, "Asymmetric non-local neural networks for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 593–602.
- [8] Google, "Mobilenetedgegpu," <https://github.com/tensorflow/models/tree/master/research/slim/nets>, 2019.
- [9] Andre Martins and Ramon Astudillo, "From softmax to sparsemax: A sparse model of attention and multi-label classification," in *International Conference on Machine Learning*, 2016, pp. 1614–1623.
- [10] 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 *Proc. CVPR*, 2016, pp. 3213–3223.
- [11] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba, "Scene parsing through ADE20K dataset," in *Proc. CVPR*, 2017, pp. 5122–5130.
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li, "Imagenet: A large-scale hierarchical image database," in *Proc. CVPR*, 2009, pp. 248–255.
- [13] Hang Zhang, Kristin J. Dana, Jianping Shi, Zhongyue Zhang, Xiaogang Wang, Amrbrish Tyagi, and Amit Agrawal, "Context encoding for semantic segmentation," in *Proc. CVPR*, 2018, pp. 7151–7160.
- [14] Fan Zhang, Yanqin Chen, Zhihang Li, Zhibin Hong, Jingtuo Liu, Feifei Ma, Junyu Han, and Errui Ding, "Acfnet: Attentional class feature network for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6798–6807.