

LiteSeg: A Novel Lightweight ConvNet for Semantic Segmentation

Taha Emara, Hossam E. Abd El Munim, Hazem M. Abbas

Computer and Systems Engineering Department, Faculty of Engineering, Ain Shams University,
Cairo, Egypt

Email: {taha@emaraic.com, hossameldin.hassan@eng.asu.edu.eg, hazem.abbas@eng.asu.edu.eg}

Abstract—Semantic image segmentation plays a pivotal role in many vision applications including autonomous driving and medical image analysis. Most of the former approaches move towards enhancing the performance in terms of accuracy with a little awareness of computational efficiency. In this paper, we introduce LiteSeg, a lightweight architecture for semantic image segmentation. In this work, we explore a new deeper version of Atrous Spatial Pyramid Pooling module (ASPP) and apply short and long residual connections, and depthwise separable convolution, resulting in a faster and efficient model. LiteSeg architecture is introduced and tested with multiple backbone networks as Darknet19, MobileNet, and ShuffleNet to provide multiple trade-offs between accuracy and computational cost. The proposed model LiteSeg, with MobileNetV2 as a backbone network, achieves an accuracy of 67.81% mean intersection over union at 161 frames per second with 640×360 resolution on the Cityscapes dataset.

Index Terms—semantic image segmentation, atrous spatial pyramid pooling, encoder decoder, and depthwise separable convolution.

I. INTRODUCTION

Semantic image segmentation is defined as the assigning of every pixel in a given image to a specific categorical label. Semantic segmentation [1]–[4], and similar to image classification [5]–[7] and object detection [8], [9], has seen considerable progress due to the employment of deep learning architectures, especially convolutional neural networks (CNN). This progress has resulted in a much better quality of real-world applications, such as autonomous driving, medical diagnosis [10], and aerial image segmentation [11].

Despite the high accuracy achieved by recent proposed architectures [2], [3] for semantic segmentation, they are not computationally efficient especially for the applications that are needed to be run on edge devices, such as autonomous driving cars, robots, or augmented reality kits. Numerous attempts have been investigated in providing lightweight semantic segmentation architectures, such as ERFNet [12], ESPNet [13], Enet [14], CCC [15], and DSNet [16]. Some of these lightweight architectures attempts aimed at obtaining real-time performance with a considerable reduction in network parameters which significantly causes a loss in accuracy measures [13]–[15], [17]. Other methods paid more attention to both accuracy and real-time performance which leads to gain a better real-time performance when compared to complex networks and a better accuracy than the first group [12], [16].

Semantic segmentation. The Fully Convolutional Network

(FCN) [1] is a pivotal approach which paves the way to employ deep learning methods into the semantic segmentation problem. In FCN, a classification model such as, GoogleNet [7] or VGG [18] was used as an encoder to extract features from tested images and then these feature maps were upsampled to pixelwise dense predictions by cascaded layers of unpooling and deconvolution operations. Accuracy of this architecture was improved by using skip architecture, in which semantic information from deep layer and spatial information from the earlier layers were combined to get better results. Despite the breakthrough of FCN architecture, it has suffered from low resolution prediction. Many variant of FCN are proposed to solve this problem. For example, the work in [19] proposed a multi-scale network which employed a different three scale to generate a fine, high resolution predictions. Another solution was proposed by [20], in which a more complex deconvolution network was used to produce high resolution predictions instead of the used one in [1] which used a single bilinear interpolation layer. A different approach [21] employed dilated convolution to increase the receptive field without any increase the in number of parameters and computational cost, followed by bilinear interpolation layers to scale up the feature maps to the input image size. Then a conditional random field (CRF) [22] was used as a post processing to refine the result image. PSPNet [3], Deeplabv3 [4], Deeplabv3+ [2] capture information at multiple scales by either applying pooling operations with different kernel size and they called it pyramid pooling module (PPM) or employing dilated convolution with different rate and that was called Atrous Spatial Pyramid Pooling (ASPP).

Real-Time Segmentation. Most of the mentioned approaches are not efficient for real-time applications as they employed large backbone networks such as GoogleNet [7], Xception [18], or ResNet [6], or employed a large CNN architectures for both the encoder or decoder sides. This has lead to having a large number of parameters to be tuned and a large floating point operations (FLOPS), even though they are efficient form accuracy perspective. Many approaches have been proposed to deal with this problem, e.g., ERFNet [12] employed a residual connection and depthwise separable convolution to increase receptive field to achieve high accuracy with a reasonable performance. Alternatively, ESPNet [13] proposed an efficient module called efficient spatial pyramid

arXiv:1912.06683v1 [cs.CV] 13 Dec 2019

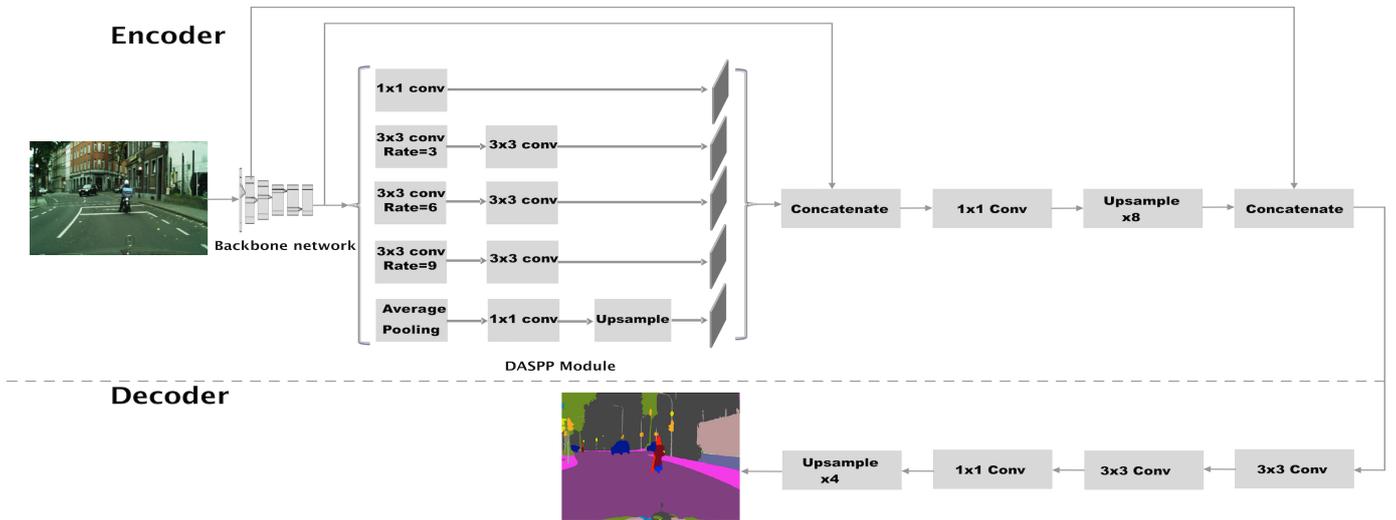


Fig. 1. General LiteSeg diagram including the encoder module with its components backbone network and DASPP module, and the decoder module. Encoder module takes an input image and generates a high dimensional feature vector. The decoder module restores the spatial information from this feature vector.

(ESP), which uses point wise convolution and spatial pyramid of dilated convolution. ESPnet along with Enet provide a lightweight architectures but with a degradation in accuracy. RTSeg [17] provided a decoupled encoder-decoder architecture which allows to plug any encoder (i.e., VGG16 [18], MobileNet [23], ShuffleNet [24], ResNet18 [6]) or decoder (i.e., UNet [10], Dilation [25], SkipNet [1]) architectures independently. They have found out that using SkipNet architecture along with MobileNet and ShuffleNet provided the best trade-off between accuracy and performance.

Motivated by the encoder-decoder architecture, Atours Spatial Pyramid Pooling (ASPP), dilated convolution, and depthwise separable convolution, we design a novel architecture called LiteSeg which is capable of adapting any backbone network. This capability would allow a variety in trade-offs between computational cost and accuracy to fit multiple needs by choosing different backbone networks.

In summary, our main contributions are:

- LiteSeg, a real time competitive architecture is presented and tested with three different backbone networks, Darknet19 [8], MobileNetV2 [23], and ShuffleNet [24], achieving performance 70.75%, 67.81%, and 65.17%, respectively on Cityscapes dataset.
- A new deeper version of ASPP module is adapted to improve the results along with using long and short residual connection.

The rest of the paper is organized as follows. Section 2 describes the proposed architecture, LiteSeg, in details. In Section 3, both the accuracy and the computational efficiency of the proposed model is evaluated and the paper is finally concluded in Section 4.

II. METHODS

Here, we will describe our architecture LiteSeg and its new deeper version of ASPP module called Deeper atrous Spatial

Pyramid Pooling (DASPP) module (Figure 1). In addition, the atrous convolution, depthwise separable convolution and long and short residual connection are briefly introduced. Then, Deeplabv3+ [2] which is used as the decoder module will be reviewed.

A. Atrous Convolution

In convolutional architecture, decreasing the receptive field size will result in a spatial information loss that can be attributed to the strided convolution and pooling layers. To overcome this problem, the dilated convolution was used in [21], [25] to increase the receptive field without any reduction in the feature map resolution and an increase in trainable parameters. This allows network to learn global context features across the entire image for refining full-resolution predictions.

B. Depthwise Separable Convolution

Standard convolution is computationally an expensive operation due to the large number of parameters to be tuned and thus the needed FLOPS. To tackle this problem, depthwise separable convolution is a suggested solution to replace standard convolution without compromising the accuracy.

The main idea of depthwise separable convolution is to split both the input and the kernel into channels -they share the same number of channels-, and each input channel will be convolved with the corresponding kernel channel. Then, the pointwise convolution is performed using a 1×1 kernel to project the output of the depthwise convolution into new channel space. Employing depthwise separable convolution [26] was empirically proven to reduce the computational cost with similar or better performance.

C. Long and short residual connection

He et al. [6] proposed a residual learning framework to allow training of very deep networks. Unlike the traditional feedforward neural network, ResNets introduce an identity

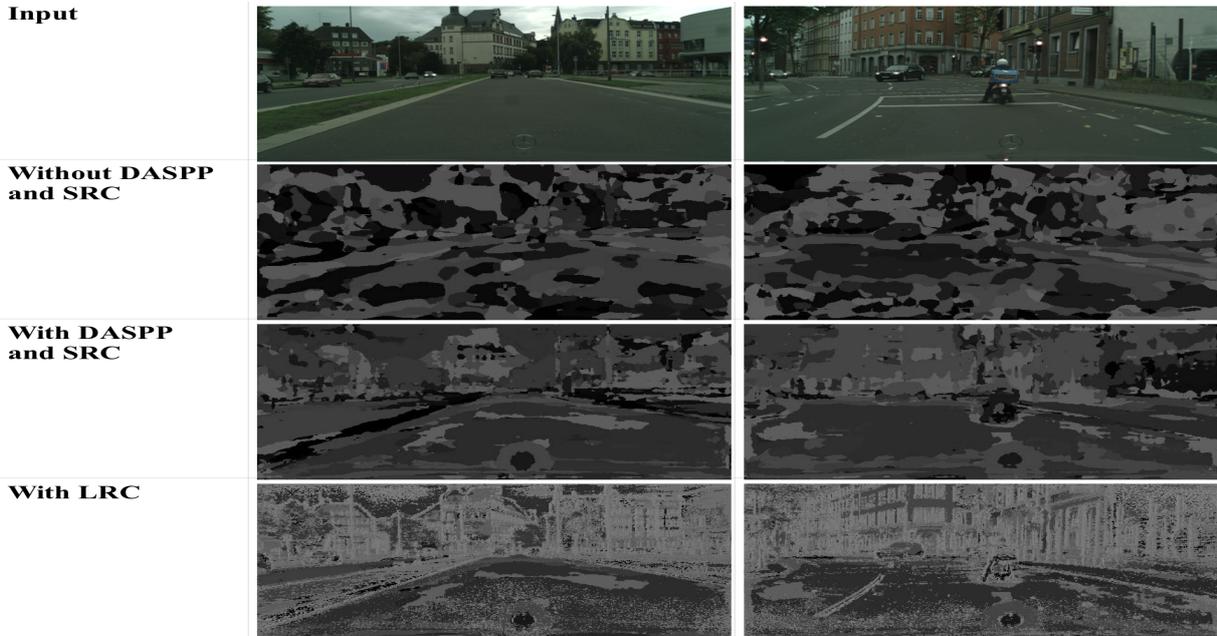


Fig. 2. Visualization of the output after encoder module, to show the effectiveness of short residual connection (SRC), long residual connection (LRC), and DASPP module on our model performance.

shortcut connection. Let X is the input feature map, $F(X)$ is the residual, and $H(X)$ is the output of residual block, the residual learning takes the form $H(X) = F(X) + X$ such that, if there is no residual it will work as identity mapping, that means it can eliminate the effect of the DASPP module if it turns out to be an unnecessary. The resulting learned residual assures that the proposed network would not perform worse than without it.

Fusion and reusing of low-level features -which include color blobs or edges- from bottom layers and high-level features from top layers have been proven to be helpful for high resolution segmentation [27]. This fusion can be done between feature maps from close layers by short residual connection (SRC) and far layers by long residual connection(LRC). These connections act as memory units [28] in the network as they allow preserving the information from the bottom layers to the top layers.

There are two approaches to carry out residual connections, one by element-wise addition [6] and the other by concatenating the feature maps [28]. Here, we employed the concatenation approach as an element-wise addition that require the residual output and the input have the same dimension width, height, and depth instead of the conventional concatenation which requires the same dimension of width and height only. The mismatch in width and height can be maintained by upsampling and optionally a 1×1 convolution can be used to reduce the depth of the features for computational efficiency. It was found out that long skip connection helps to make clearer semantic boundaries and short skip connections with DASPP help in fine tuning the semantic segments and thus providing richer geometrical information (Figure 2).

D. Proposed Encoder

The proposed encoder contains a backbone network architecture which acts as an image classification architecture for feature extraction. These architectures were chosen to meet our performance criteria, so we tested the architecture with different three lightweight models MobileNet, ShuffleNet, and Darknet19. Not only is the type of the backbone network controls the performance, but also the output stride [4] which is defined as the ratio between the input image size and the last feature map of the encoder. Let height H , width W , and depth C be the input image dimension and the outputs of the backbone network are h , w , and c , so the output stride is defined as $os = H \times W / h \times w$. Decreasing the output stride leads to having high resolution feature map and also better results [4] as more spatial information throughout the network is preserved but it comes with computational cost. The output stride of backbone networks is controlled by removing max pool layers and modifying the stride for the last convolution layers. Deeplabv3+ [2] with output stride equal to 16 is the best trade-off between accuracy and computational efficiency. Moreover, they found out that the accuracy can be greatly improved using output stride equal to 8 but with huge computational cost and the computational efficiency can be improved by increasing the output stride to 32 with a compromise in accuracy. Therefore, the proposed backbone network is configured with output stride of 32 for MobileNetV2 [23] and ShuffleNet [24] and output stride of 16 for Darknet19 to achieve different trade-offs between accuracy and speed. DeepLabv3 [4] employs Atrous Spatial Pyramid Pooling (ASPP) module with different dilation rates to capture multi-scale information, following the presented approach in

ParseNet [29]. Here, a new deeper version of ASPP module is proposed (called Deeper Atrous Spatial Pyramid Pooling (DASPP)), by adding standard 3×3 convolution after 3×3 atrous convolutions to refine the features and also fusing the input and the output of the DASPP module via short residual connection. Also, the number of convolution filters of ASPP is reduced from 255 to 96 to gain computational performance.

E. Deeplabv3+ as a Decoder

Deeplabv3+ [2] presented a simplified decoder that is composed of standard 3×3 convolution and upsampling layers. Here, we added another 3×3 convolution layer and reduced the number of filters in all 3×3 convolution from 256 to 96 for computational performance gain. Additionally, the output of the encoder is augmented with low level features from earlier layers of the backbone network via long residual connection. These low level features might have large number of feature maps, and in order to resolve this problem, a 1×1 convolution is utilized to reduce the number of channels of low level feature. Otherwise, with some light backbone networks, there will be no need to apply the 1×1 convolution on low level features because of the low number of channels (e.g., 24 in case of using MobileNet).

III. EXPERIMENTAL RESULTS AND VALIDATION

In our evaluation of the proposed method, the effectiveness of LiteSeg with different backbone networks is empirically tested and the results are compared with the lightweight state-of-the-art architectures on Cityscapes [30] dataset. The performance of the proposed model is measured in terms of mean intersection over union (mIOU), giga floating point operations (GFLOPs), and the number of parameters (Params) in millions.

A. Dataset and Computing Environment

The Cityscapes dataset is a large-scale dataset for semantic understanding of urban scenes. It contains 5000 images with fine annotations divided into 2975 images for training, 500 images for validation, and 1525 images for testing. It also contains about 20000 images with coarse annotations that can be used as extra data for fine-tuning the models.

The experiments were carried out on a computer with Intel Core i7-8700 @ 3.2GHZ, 16GB memory, and NVIDIA GTX1080Ti GPU card. This computer runs Ubuntu 18.04 and PyTorch [31] version 0.4.1 with CUDA 9.0 and cudnn 7.0.5.

B. Training Protocol

Stochastic Gradient Descent with Nesterov [32] was used with a momentum value of 0.9 and an initial learning rate of 10^{-7} for ShuffleNet and MobileNetV2 backbone networks and 10^{-8} for Darknet19 backbone network, and a weight decay 4×10^{-5} . We applied multiple learning rate policies where the learning rate changes after every five epochs such that the learning rate of the current epoch is calculated by $initial_learning_rate \times (1 - epoch/max_epochs)^{power}$ with power 0.9.

C. Encoder Options

Baseline Model. First our experiments are conducted with a baseline architecture which employs ASPP and decoder modules from the Deeplabv3+ [2]. This baseline was tested with three different backbone networks, MobileNetV2 [23], ShuffleNet [24], and Darknet19 [8] with output stride of 32 during both training and testing phases. As shown in the first row of Table I, employing Darknet19 as the backbone network for LiteSeg produces an appreciable improvement in the accuracy when compared to MobileNetV2 and ShuffleNet as it is a more efficient classification model [8]. This can be attributed to the fact that the generated features for the decoders make the architecture more efficient as a classifier.

Employing DASPP Module. As shown in the second row of Table I, employing DASPP module along with decreasing output stride from 32 to 16, considerably increases the accuracy of the network by 2.37% when using Darknet19 as a backbone network. It also shows that employing DASPP module along with keeping output stride at 32 for MobileNetV2 and ShuffleNet increases the accuracy of the network by 0.1% and 0.9%, respectively. For DASPP module, we employed dilation rates (3,6,9) for the three 3×3 convolutions in the first layer, and used standard 3×3 convolution in the second layer of convolutions.

Pre-Training on The Coarse Dataset. Due to the lack of finely annotated data for semantic segmentation models, several works [33], [34] found that object-level and image-level labels can improve the result of semantic segmentation models. LiteSeg is trained on the coarse data for 20 epochs and then the trained model is used for training the fine data. The third row of Table I shows that using a trained network on coarsely annotated data improves the accuracy of the network by 0.7%, 1.6%, and 1.3% when using Darknet19, MobileNetV2, and ShuffleNet, respectively.

Multi-Scale Input. Learning network with multi-scale images forces the network to well predict across multiple sizes of input images [8]. Following this strategy, we augmented the dataset with multi-scale input images as our network is a fully convolutional network which makes it accept different dimensions of images. This makes the proposed models efficient for predicting various sizes of input images, as stated in the fourth row of Table I.

Employing Depthwise Separable Convolution. Not only does the use of depthwise separable convolution instead of using standard convolution in our network reduce the FLOPs as stated in Table II, but it also improves the accuracy of the network by 0.5%, 0.7%, and 0.7% when using Darknet19, MobileNetV2, and ShuffleNet, respectively when evaluating images of size 1024×2048 , as stated in the fifth rows of Table I.

D. Computational Performance Evaluation

The computational efficiency of the proposed models is assessed here. Both the inference time, which reflects the real-time performance, and number of parameters, which reflects

TABLE I

EVALUATION RESULTS IN MIOU ON THE CITYSCAPES VALIDATION SET USING *LiteSeg* WITH AN INPUT IMAGE SIZE 512×1024 USING DIFFERENT BACKBONE NETWORKS. BASELINE NETWORK IS MINIMAL VERSION OF DEEPLABV3+. **FT**: USING COARSE DATASET. **MS**: MULTI-SCALE TRAINING STRATEGY. **DW**: EMPLOYING DEPTHWISE SEPARABLE CONVOLUTION. RESULTS WITH '*' WERE EVALUATED ON IMAGES WITH SIZES 512×1024 AND 1024×2048 AND LISTED AS 512×1024 ACCURACY/ 1024×2048 ACCURACY.

Baseline	DASPP	FT	MS	DW	LiteSeg-Darknet	LiteSeg-MobileNet	LiteSeg-ShuffleNet
✓					65.84%	64.70%	60.41%
	✓				68.21%	64.80%	61.3%
	✓	✓			68.94%/71.5%*	66.4%/67.8%*	62.65%/62.2%*
	✓	✓	✓		69.14%/72.3%*	66.49%/69.3%*	63.2%/65.4%*
	✓	✓	✓	✓	69.43%/72.8%*	66.48%/70.0%*	62.45%/66.1%*

TABLE II

EFFECT OF EMPLOYING DEPTHWISE SEPARABLE CONVOLUTION TO REDUCE THE NUMBER OF FLOATING POINT OPERATIONS, INSTEAD OF STANDARD CONVOLUTION. THE UNIT OF ALL LISTED NUMBER IS GIGA FLOATING POINT OPERATIONS (GFLOPs). THEY ARE MEASURED ON IMAGE SIZE 1024×512 .

Convolution type	LiteSeg-Darknet	LiteSeg-MobileNet	LiteSeg-ShuffleNet
Standard Convolution	123.26	18.86	9.36
Depthwise Separable convolution	103.09	4.9	2.75

the memory footprint, are measured. A set 200 images for the burn-in process and 200 images for evaluation are used in the process. Table III compares the proposed models to current state-of-the-art real-time segmentation networks using the same computing environment.

TABLE III

INFERENCE TIME ANALYSIS ON IMAGES WITH RESOLUTION 360×640 AND FULL RESOLUTION 1024×2048 USING OUR MACHINE. DSNET RESULT WAS TAKEN FROM THEIR PAPER, THEY USED NVIDIA GTX 1080TI ON THEIR EXPERIMENTS.

Network	FPS (360x640)	FPS (1024x2048)	Params(in millions)
ErfNet [12]	105	15	2.07
DSNet [16]	100.5	-	0.91
LiteSeg-Darknet (ours)	98	15	20.55
ESPNET [13]	144	25	0.364
LiteSeg-MobileNet (ours)	161	22	4.38
LiteSeg-ShuffleNet (ours)	133	31	3.51

These results clearly show the ability of LiteSeg to generate different lightweight models to manipulate the accuracy and computational efficiency by using different backbone network. For example, using 640×360 input resolution, LiteSeg with MobileNetV2 [23] as a backbone network achieved a speed of 161 FPS which exceeds the speed of ESPNet [13] by 17 FPS on the same machine, while providing an improved accuracy by 7.51%.

E. Cityscapes Benchmark Results

The models with the best result on the validation set are selected and compared with the results of the proposed model when experimented in the test set. The results are then uploaded to the official benchmark of Cityscapes dataset. As shown in Table IV, we compare our result on the test set with other state-of-the-art real-time models for semantic image segmentation. Although the LiteSeg-DarkNet19 has a high GFLOPs compared with ERFNet, it has improved

the accuracy of ERFNet and DSNet by 2.75% and 1.45%, respectively, just with a sacrifice of 7 FPS for ERFNet and 2.5 FPS for DSNet (Table III).

TABLE IV

PERFORMANCE OF OUR PROPOSED LITESEG AND SIMILAR ARCHITECTURES ON CITYSCAPES TEST SET. FOR RESULTS WITH '**', GFLOPs IS COMPUTED ON IMAGE RESOLUTION 640×360 .

Model	GFLOPs	Class mIOU	Category mIOU
SegNet* [35]	286.03	56.1%	79.1%
ESPNet [13]	9.67	60.3%	82.2%
ENet [14]	8.52	58.3%	80.4%
ERFNet [12]	53.48	68.0%	86.5%
SkipNet-ShuffleNet [17]	4.63	58.3%	80.2%
SkipNet-MobileNetNet [17]	13.8	61.5%	82.0%
CCC2 [15]	6.29	61.96%	nan
DSNet [16]	nan	69.3%	86.0%
LightSeg-MobileNet (ours)	4.9	67.81%	86.79%
LightSeg-ShuffleNet (ours)	2.75	65.17%	85.39%
LightSeg-DarkNet19 (ours)	103.09	70.75%	88.29%

TABLE V

CATEGORY RESULTS OF OUR LITESEG MODELS ON CITYSCAPES TEST SET. ALL NUMBER REPRESENT THE MIOU.

Model	Flat	Nature	Object	Sky	Construction	Human	Vehicle
LightSeg-MobileNet	97.90%	91.70%	62.75%	94.62%	90.44%	79.26%	90.88%
LightSeg-ShuffleNet	97.88%	91.22%	57.43%	93.99%	89.69%	77.27%	90.28%
LightSeg-DarkNet19	98.44%	98.44%	65.94%	94.99%	91.67%	81.73%	92.95%

In Table V, the mIOU of the main categories of Cityscapes test set are listed and one can easily observe that the most common categories in the dataset have the highest mIOU score. The results of LiteSeg are displayed in Figure 3 for qualitative analysis against ESPNet [13] and ERFNet [12].

IV. CONCLUSION

In this paper, we proposed LiteSeg, a novel lightweight architecture for semantic image segmentation. The ability of LiteSeg to adapt multiple backbone networks allows for providing multiple trade-offs to fit embedded devices and deep learning workstations. We introduced a new module named DASPP to improve semantic boundaries of captured features from backbone network. The proposed network, LiteSeg, was evaluated with ShuffleNet as a backbone network on the Cityscapes test dataset showing that it is able to achieve 65.17% mIoU at 31 FPS for full image resolution 1024×2048 on a single Nvidia GTX 1080TI GPU.

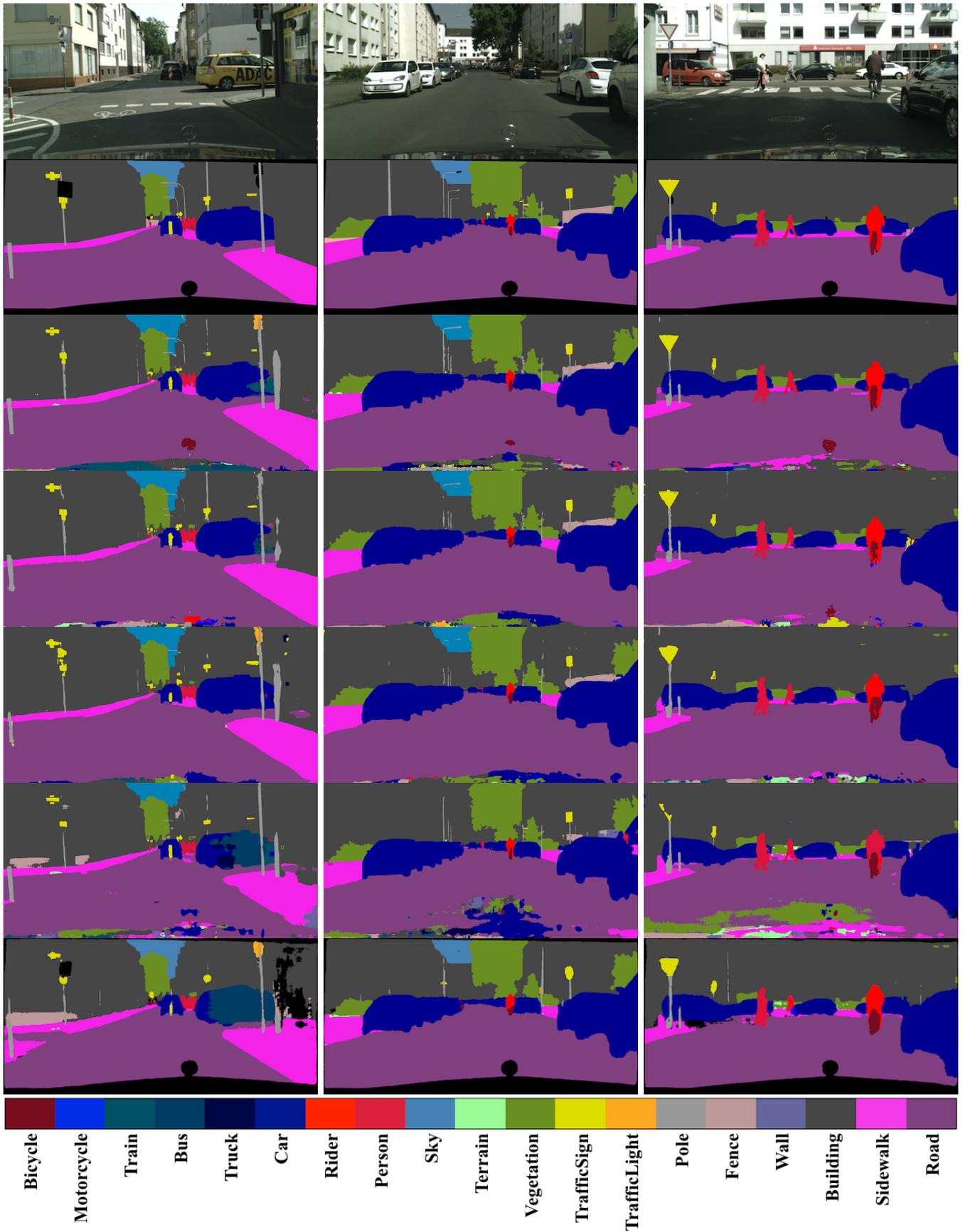


Fig. 3. Visualization results of multiple models on Cityscapes validation set [30]. From top to down 1-Input RGB images; 2-Ground truths; 3-LiteSeg-Darknet predictions; 4-LiteSeg-MobileNet predictions; 5-LiteSeg-ShuffleNet predictions; 6-ERFNet predictions; 7-ESPNet predictions; 8- Color map for Cityscapes classes.

REFERENCES

- [1] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440, 2015.
- [2] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 801–818, 2018.
- [3] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017.
- [4] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017.
- [5] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1492–1500, 2017.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [7] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- [8] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7263–7271, 2017.
- [9] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
- [10] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241, Springer, 2015.
- [11] D. Marmanis, J. D. Wegner, S. Galliani, K. Schindler, M. Datcu, and U. Stilla, "Semantic segmentation of aerial images with an ensemble of cnns," *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 3, p. 473, 2016.
- [12] E. Romera, J. M. Alvarez, L. M. Bergasa, and R. Arroyo, "Efficient convnet for real-time semantic segmentation," in *2017 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1789–1794, IEEE, 2017.
- [13] S. Mehta, M. Rastegari, A. Caspi, L. Shapiro, and H. Hajishirzi, "Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 552–568, 2018.
- [14] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, "Enet: A deep neural network architecture for real-time semantic segmentation," *arXiv preprint arXiv:1606.02147*, 2016.
- [15] H. Park, Y. Yoo, G. Seo, D. Han, S. Yun, and N. Kwak, "Concentrated-comprehensive convolutions for lightweight semantic segmentation," *arXiv preprint arXiv:1812.04920*, 2018.
- [16] W. Wang and Z. Pan, "Dsnnet for real-time driving scene semantic segmentation," *arXiv preprint arXiv:1812.07049*, 2018.
- [17] M. Siam, M. Gamal, M. Abdel-Razek, S. Yogamani, and M. Jagersand, "Rtseg: Real-time semantic segmentation comparative study," in *2018 25th IEEE International Conference on Image Processing (ICIP)*, pp. 1603–1607, IEEE, 2018.
- [18] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [19] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *Proceedings of the IEEE international conference on computer vision*, pp. 2650–2658, 2015.
- [20] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proceedings of the IEEE international conference on computer vision*, pp. 1520–1528, 2015.
- [21] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.
- [22] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected crfs with gaussian edge potentials," in *Advances in neural information processing systems*, pp. 109–117, 2011.
- [23] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018.
- [24] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6848–6856, 2018.
- [25] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [26] V. Nekrasov, C. Shen, and I. D. Reid, "Light-weight refinenet for real-time semantic segmentation," in *BMVC*, 2018.
- [27] Z. Zhang, X. Zhang, C. Peng, X. Xue, and J. Sun, "Exfuse: Enhancing feature fusion for semantic segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 269–284, 2018.
- [28] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- [29] W. Liu, A. Rabinovich, and A. C. Berg, "Parsenet: Looking wider to see better," *arXiv preprint arXiv:1506.04579*, 2015.
- [30] M. Cordts, S. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [31] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. De Vito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.
- [32] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International conference on machine learning*, pp. 1139–1147, 2013.
- [33] J. Xu, A. G. Schwing, and R. Urtasun, "Learning to segment under various forms of weak supervision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3781–3790, 2015.
- [34] P. O. Pinheiro and R. Collobert, "From image-level to pixel-level labeling with convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1713–1721, 2015.
- [35] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.