# ExtremeC3Net: Extreme Lightweight Portrait Segmentation Networks using Advanced C3-modules

Hyojin Park
Seoul National University
wolfrun@snu.ac.kr

Lars Lowe Sjösund
Clova AI, NAVER Corp
lars.sjosund@navercorp.com

YoungJoon Yoo
Clova AI, NAVER Corp
youngjoon.yoo@navercorp.com

Jihwan Bang
Search Solutions, Inc
jihwan.bang@navercorp.com

Nojun Kwak
Seoul National University
nojunk@snu.ac.kr

## Abstract

*Designing a lightweight and robust portrait segmentation algorithm is an important task for a wide range of face applications. However, the problem has been considered as a subset of the object segmentation problem. bviously, portrait segmentation has its unique requirements. First, because the portrait segmentation is performed in the middle of a whole process of many real-world applications, it requires extremely lightweight models. Second, there has not been any public datasets in this domain that contain a sufficient number of images with unbiased statistics. To solve the problems, we introduce a new extremely lightweight portrait segmentation model consisting of a two-branched architecture based on the concentrated-comprehensive convolutions block. Our method reduces the number of parameters from $2.1M$ to $37.7K$ (around 98.2% reduction), while maintaining the accuracy within a 1% margin from the state-of-the-art portrait segmentation method. In our qualitative and quantitative analysis on the EG1800 dataset, we show that our method outperforms various existing lightweight segmentation models. Second, we propose a simple method to create additional portrait segmentation data which can improve accuracy on the EG1800 dataset. Also, we analyze the bias in public datasets by additionally annotating race, gender, and age on our own. The augmented dataset, the additional annotations and code are available in https://github.com/HYOJINPARK/ExtPortraitSeg .*

## 1. Introduction

Designing algorithms working on face data has been considered as an important task in the computer vision field and many sub-areas such as detection, recognition, key-
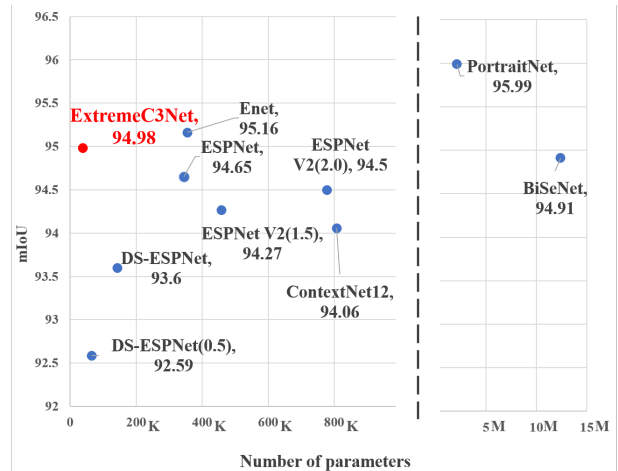


Figure 1. Accuracy (mIoU) vs. complexity (number of parameters) on the EG1800 validation set. Our proposed ExtremeC3Net has high accuracy with small complexity.

point extraction are actively studied. Among them, the portrait segmentation also has been highly used in industrial environments such as background editing, security checks, and face resolution enhancement [21, 32].

Because researchers have considered the portrait segmentation problem as a subset of semantic segmentation, most portrait segmentation algorithms have employed general semantic segmentation algorithms [18, 29, 13] trained on a portrait segmentation dataset. However the portrait segmentation comes with some unique obstacles.

The first thing is the small number of images in the dataset. The EG1800 dataset [21], a popular public portrait segmentation dataset, contains only around 1,300 training images. Also, we observed that this datsaet has large biases with respect to attributes such as race, age, and gender from

our additional annotation. For general semantic segmentation tasks, data augmentation methods like random noise, translation, and color changing are utilized to overcome the small dataset size. However, it is clear that these methods are not sufficient for resolving the small size and imbalance of portrait segmentation datasets.

Second, portrait segmentation is usually used just as one of several steps in real-world applications. Since many of these applications run on mobile devices, the segmentation model needs to be lightweight to ensure real-time speeds. Researchers have developed plenty of lightweight segmentation methods, but most of them are still not lightweight enough. A few recent examples are PortraitNet [32] with $2.1M$ parameters, ESPNetV2 [11] with $0.78M$ parameters, and ESPNet [18] with $0.36M$ parameters. In general, when deploying a portrait segmentation model on a mobile device, the smaller the number of parameters without degradation of accuracy the better. Furthermore, one is often limited in what operations are available when deploying onto embedded systems. Therefore, it can be good to use as few different types of operations as possible.

The contributions of the work can be summarized as follows: (1) We propose a new extremely lightweight segmentation model having less than $0.05M$ parameters as well as having competitive segmentation accuracy with PortraitNet [32], without using deconvolution operations. Our model only uses $1.8\%$ of the number of parameters of our baseline PortraitNet, 37.7K compared to 2.1M, but the accuracy degradation is just about $1\%$ on the EG1800 dataset as can be seen in Figure 1. (2) We introduce a simple and effective data generation method to get enough number of dataset and to alleviate the imbalance of the dataset. For performance enhancement, we generated 10,448 images for training using our generation method. Also, we additionally annotated the public dataset EG1800 according to race, gender, and age. In doing so, we detected strong biases to specific attributes. From the experiments, we found that the proposed data augmentation can enhance the segmentation accuracy and also enrich the balance among different attributes.

## 2. Related Work

**Data Augmentation:** Data augmentation techniques have become an important factor for successful training of various deep networks. Many researches have shown that the augmentation alleviates over-fitting as well as enlarging the number of data samples. In the fields of classification and detection, not only baseline augmentation methods such as cropping and flipping, but also novel patch and region-based method such as CutOut [6] and CutMix [30] have been proposed. In semantic segmentation, basically two augmentation methods are used. One is image calibration methods including rotation, flipping and cropping of the im-

ages. The other is image filtering methods controlling image attributes such as brightness, contrast and color. In addition to the basic augmentation methods, since the labeling cost is high, many segmentation methods additionally use crawled images with the label from relatively simple segmentation algorithms, similar to those in text detection and localization studies [7, 1]. Jin *et al.* [9] crawled images from the web and labeled them with dense conditional random field (CRF). Similarly, [19, 16, 25] applied weakly-supervised method to enlarge the datasets automatically.

**Convolution Factorization:** Convolution factorization dividing the convolution operation into several stages has been used to reduce computational complexity. In Inception [23, 24, 22], several convolutions are performed in parallel, and the results were concatenated. Then, a $1 \times 1$ convolution is used to reduce the number of channels. Xception [5], MobileNet [8] and MobileNetV2 [20] use the depth-wise separable convolution *(ds-Conv)*, which performs spatial and cross-channel operations separately to decrease computation. ResNeXt [27] and ShuffleNet [33] applied a group convolution to reduce complexity. In segmentation, many lightweight segmentation algorithms [18] also adopt the convolution factorization methods to reduce the number of parameters.

**Segmentation:** PortraitFCN+ [21] built a portrait dataset from Flickr and proposed a portrait segmentation model based on FCN [10]. After that, PortraitNet proposed a novel portrait segmentation model with higher accuracy than PortraitFCN+ with real-time execution time. Also, there are many lightweight segmentation models. Enet [13] was the first architecture designed for real-time segmentation. Later, ESPNet [18] improved both speed and performance by introducing an efficient spatial pyramid of dilated convolutions. ERFNet [17] used residual connections and factorized dilated convolutions into two asymmetric dilated convolutions. ContextNet [14] and FastSCNN [15] designed two network branches each for global context and detailed information, respectively. Similarly, BiSeNet [29] proposed a two-paths network for preserving spatial information as well as acquiring a large enough receptive field. There are several works which further reduced model parameters by combining the dilated convolution and the depth-wise separable convolution. ESPNetV2 [11] applied the group point-wise and depth-wise dilated separable convolutions to learn representations and showed better performance at various tasks. C3 [12] resolved the degradation of accuracy from naive combination of dilated convolution and depth-wise separable convolution by proposing a concentrated-comprehensive convolution (C3). In C3, two asymmetric convolutions are added in front of the dilated depth-wise convolution. In this paper, we propose an advanced version of the C3 module which has a much smaller model size to satisfy the aforementioned requirement of the portrait seg-

| | Input | B1 | B2 | B3 | | Input | B1 | B2 | B3 |
|---|---|---|---|---|---|---|---|---|---|
| L1 | $112 \times 112 \times 27$ | 1 | 2 | 3 | L5 | $56 \times 56 \times 56$ | 2 | 4 | 8 |
| L2 | $56 \times 56 \times 48$ | 1 | 3 | 4 | L6 | $56 \times 56 \times 56$ | 2 | 4 | 8 |
| L3 | $56 \times 56 \times 99$ | 1 | 3 | 5 | L7 | $56 \times 56 \times 56$ | 2 | 4 | 8 |
| L4 | $56 \times 56 \times 56$ | 2 | 4 | 8 | L8 | $56 \times 56 \times 56$ | 2 | 4 | 8 |

Table 1. Detailed setting of the dilation ratio. L denotes layer and B denotes block in each C3-module.

mentation.

## 3. Method

In this section, we explain a simple data generation framework to solve the lack of dataset in two situations. Also, we introduce our advanced C3-module with a well-designed combination of dilation ratios and propose an extremely lightweight segmentation model based on it.

### 3.1. Advanced C3-module

Stacking multiple dilated convolutional layers increases the receptive field and has been used to boost performance in many segmentation models. For further reducing model complexity, depth-wise separable convolutions are widely used. However, concentrated-comprehensive convolutions block (C3) [12] pointed out that the over-simplified operation performed by the depth-wise separable dilated convolution causes severe performance degradation due to loss of information contained in the feature maps. When we combine a dilated convolution with depth-wise separable convolution for our basic module, we also observe the same degradation. To mitigate this problem [12] designed the C3-block, which is composed of a concentration stage and a comprehensive convolution stage. The concentration stage compresses information from the neighboring pixels to alleviate the information loss. It does not use the standard square depth-wise convolution but instead uses an asymmetric depth-wise convolution for reducing complexity.

In this paper, we use the C3-module, but unlike in [12], we reduce the number of C3-blocks from 4 to 3, as shown in Figure 2(b). We also re-designed the dilation ratio for each C3-module, based on insights about neural network kernel properties from Zeiler *et al*. [31]. A kernel which is close to the input image extracts common and local features, while one close to the classifier extracts more class-specific and global features. The models C3 [12] and ESPNet [18] use the same dilation ratio combination for every module in their architectures, and do not take this kernel property into account. We design these ratios more carefully, using small ratios in modules close to the input, and larger ratios for later modules. Table 1 shows the detailed dilation ratio settings for each layer.

### 3.2. ExtremeC3Net Architecture

In this part, we introduce the architecture of our ExtremeC3Net for segmentation. We use a two-branched architecture for higher model efficiency, similar to BiseNet [29], ContextNet [14], and Fast-SCNN [15]. The authors of [14] and [15] argue that a two-branch architecture has different properties compared to the common encoder-decoder structure.

The encoder-decoder structure extracts features using an encoder network and recovers the original resolution using a decoder network, which often uses deconvolution operations and concatenation with the encoder feature maps. The amount of computation of the decoder is increased proportionally to the resolution reduced in the encoder.

The two-branch architecture is composed of a deep network and a shallow network running in parallel. The deep network learns complex and global context features and the shallow network preserves spatial detail. Due to this preservation of detail, there is less need to have a sophisticated decoder module, as we can use the simpler bilinear upsampling operation for resolution recovery. This removes the need to use the deconvolution operation, and reduces the number of parameters.

As shown in Figure 2, we design a two-branched architecture consisting of a CoarseNet-branch and a FineNet-branch. The CoarseNet-branch extracts deep feature embeddings and the FineNet-branch has responsibility for preserving spatial detail. Finally, the feature maps from the two branches are combined using element-wise addition.

The CoarseNet-branch is constructed from a series of advanced C3-modules and gives a rough segmentation. It reduces the size of the feature map to a quarter by first applying a convolution with stride 2 and an advanced C3-module. After that, seven C3-modules sequentially produce feature maps without the downsampling operation. Each C3-module has a different combination of dilation ratios, as mentioned above, for making better features. We concatenate the downsampled input image onto the feature maps at two different points in the branch. ESPNet showed that this practice can increase the information flow and accuracy.

The FineNet-branch creates more exact boundary lines by only downsampling the feature maps by a factor 2. This branch is kept shallow, as the large feature size would cause the computational complexity to grow rapidly with network depth. First, it reduces the feature map size by applying a convolution with stride 2. Then it applies a C3-module to capture spatial detail information.

A point-wise convolution is applied to the last feature maps from CoarseNet, to get the same number of channels as from FineNet. Bilinear upsampling increases the output from CoarseNet and FineNet with a factor 4 and 2 respectively. The two full resolution feature maps are then aggregated by element-wise summation.
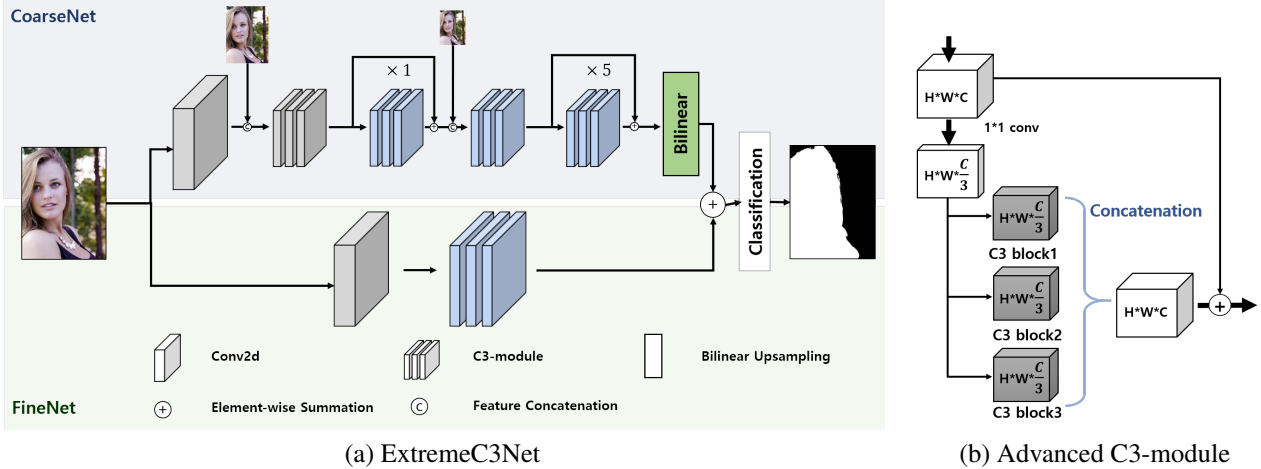
(a) ExtremeC3Net          (b) Advanced C3-module

Figure 2. (a) The model structure of ExtremeC3Net. Gray (green) color represents downsampling (unpsampling) (b) The structure of advanced C3-module consisting of three C3-blocks.



(a) The situation of having human segmentation ground truths      (b) The situation of having only raw images
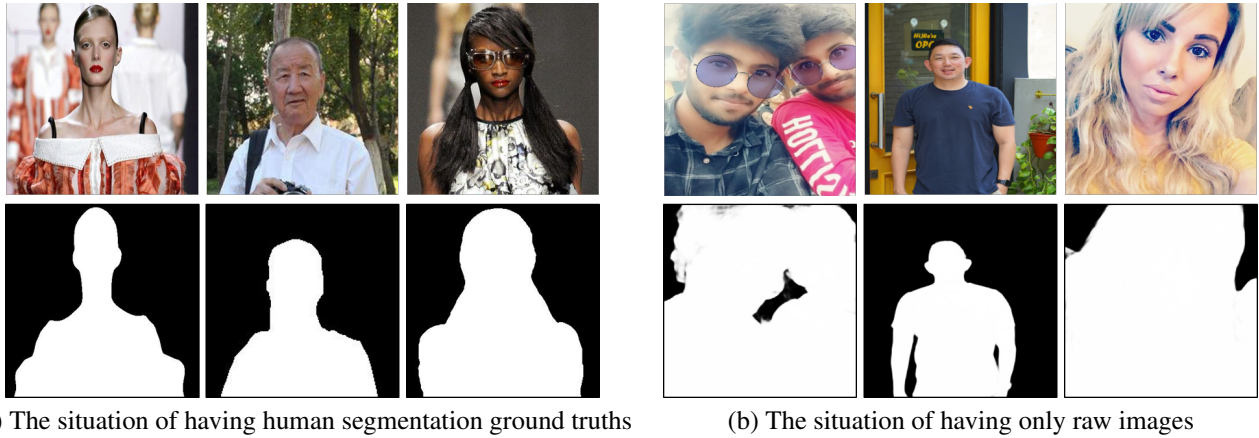
Figure 3. Examples of images and segmentation masks generated by our proposed framework in two situations (a) Data generated from the Baidu human segmentation dataset, by using a face detector. (b) Data generated from raw images by using a segmentation model.

The intersection over union(IoU) score is widely used to evaluate segmentation quality, and can be directly optimized using the Lovász-Softmax loss [3] as a surrogate function. We apply this loss function both as the main segmentation loss and also as an auxiliary loss focused around the boundary area. We define the boundary area as the non-zero part of the difference between the morphological dilation and erosion of the binary ground truth segmentation mask. The final loss is as shown in Equation 1. Lovász denotes a Lovász-Softmax loss and $f$ is a $7 \times 7$ size filter used for the dilation and erosion operations. $\mathcal{P}$ denotes the ground truth, and $\mathcal{B}$ is the boundary area as defined by the morphology operation. The $i(j)$ is an included pixel in the ground truth(boundary area). $y^*$ is a binary ground truth value and $\hat{y}$ is a predicted label from a segmentation model.

$$\mathcal{B} = (f \oplus mask) - (f \ominus mask)$$
$$Loss = \text{Lovász}_{i \in \mathcal{P}}(y_i^*, \hat{y}) + w\text{Lovász}_{j \in \mathcal{B}}(y_j^*, \hat{y}_j) \quad (1)$$

### 3.3. Data Generation Method

Annotating data often comes with high costs, and the annotation time per instance varies a lot depending on the task type. For example, Papadopoulos et al. [2] estimate the annotation time per instance for PASCAL VOC to be 20.0 seconds for image classification and 239.7 seconds for segmentation, an order of magnitude difference. To mitigate the cost of annotation for portrait segmentation, we consider a couple of plausible situations: 1) having images with ground truth human segmentation 2) having only raw images. We make use of either an elaborate face detector model (case 1) or a segmentation model (case 2) for gener-

| Method | Parameters | FLOPs(all) (G) | FLOPs [11] (G) | mIoU | Paper [32] |
|---|---|---|---|---|---|
| Enet (2016) [13] | 355 K | 0.703 | 0.346 | 95.16 | 96.00 |
| BiSeNet (2018) [29] | 12.4 M | 4.64 | 2.31 | 94.91 | 95.25 |
| PortraitNet (2019)[32] | 2.08 M | 0.666 | 0.325 | 95.99 | 96.62 |
| ESPNet (2018)[18] | 345 K | 0.665 | 0.328 | 94.65 | |
| DS-ESPNet | 143 K | 0.418 | 0.199 | 93.60 | |
| DS-ESPNet(0.5) | 63.9K | 0.296 | 0.139 | 92.59 | |
| ESPNetV2(2.0) (2019)[11] | 778 K | 0.476 | 0.231 | 94.50 | |
| ESPNetV2(1.5) (2019)[11] | 458 K | 0.285 | 0.137 | 94.27 | |
| ContextNet12 (2018)[14] | 838 K | 3.17 | 1.55 | 95.71 | |
| ContextNet12(0.25) | 67.2 K | 0.372 | 0.176 | 93.24 | |
| **ExtremeC3Net(Ours)** | **37.7 K** | **0.286** | **0.128** | **94.23** | |
| **ExtremeC3Net(Ours + generated dataset)** | **37.7 K** | **0.286** | **0.128** | **94.98** | |

Table 2. EG1800 validation results for the proposed ExtremeC3Net and other segmentation models. DS denotes depth-wise separable convolution. FLOPs(all) counted all the operation FLOPs, and FLOPs [11] calculated the number according to ESPNetV2 [11] official code. Performances of $6_{th}$ column were reported in the PortraitNet[32] paper.

ating pseudo ground truths to each situation.

When we have human images and ground truths, we just need a bounding box around the portrait area. We took plenty of images from Baidu dataset [26], which contains 5,382 human full body segmentation images covering various poses, fashions and backgrounds. To get the bounding box and portrait area, we detect the face location of the images using a face detector [28]. Since the face detector tightly bounds the face region, we increase the bounding box size to include parts of the upper body and background before cropping the image and ground truth segmentation.

We also create a second dataset from portrait images scraped from the web, applying a more heavyweight segmentation model to generate pseudo ground truth segmentation masks. This segmentation model consists of a DeepLabv3+ [4] architecture with a SE-ResNeXt-50 [27] backbone. The model is pre-trained on ImageNet and fine-tuned on a proprietary dataset containing around 2,500 fine grained human segmentation images. The model is trained for general human segmentation rather than for the specific purpose of portrait segmentation. Despite this the model works well for the portrait segmentation task, and can be used for acquiring extra training data.

Finally, human annotators just check the quality of each pseudo ground truth image, removing obvious failure cases. This method reduces the annotation effort per instance from several minutes to 1.25 seconds by transforming the segmentation task into a binary classification task. Examples of the generated dataset are shown in Figure 3.

# 4. Experiment

We evaluated the proposed method on the public dataset EG1800 [21], which collected images from Flickr with manually annotated labels. The dataset has a total of $1,800$

images and is divided into $1,500$ train and $300$ validation images. However, we could access only 1,309 images for train and 270 for validation, since some of the URLs are broken. We built additional 10,448 images from the proposed data generation method mentioned in Section 3.3.

We trained our model using ADAM optimizer with initial learning rate to $1e^{-3}$, batch-size to 60, weight decay to $5e^{-4}$. We trained the model with total 600 epochs, and image resolution was set to $224 \times 224$. We used a two-stage training method for training the proposed method. For the first 300 epochs, we only trained the CoarseNet-branch. Then, with intializing the the best parameters of the CoarseNet-branch from the previous step, we trained the overall ExtremeC3Net model for an additional 300 epochs. We evaluated our model followed by various ablations using mean intersection over union (mIoU) and compared with a SOTA portrait segmentation model including other lightweight segmentation models.

Finally, we demonstrated the results of our additional annotation about detailed attributes of faces. The dataset covers many different images, but the analysis below shows that the dataset is biased to the specific races or ages. From Figure 5, we can see that Caucasian occupies majority portion in the dataset. Also, the age group is mainly biased to 'Youth and Middle-aged' group. We give a guideline which attribute is more critical to segmentation accuracy from our analysis and accuracy improvement from our data generation method.

## 4.1. Evaluation Results on the EG1800 Dataset

We compared the proposed model to PortraitNet[32], which has SOTA accuracy in the portrait segmentation field. Since some sample URLs in the EG1800 dataset are missing, we re-trained the PortraitNet followed the original method in paper and official code from the remaining sam-
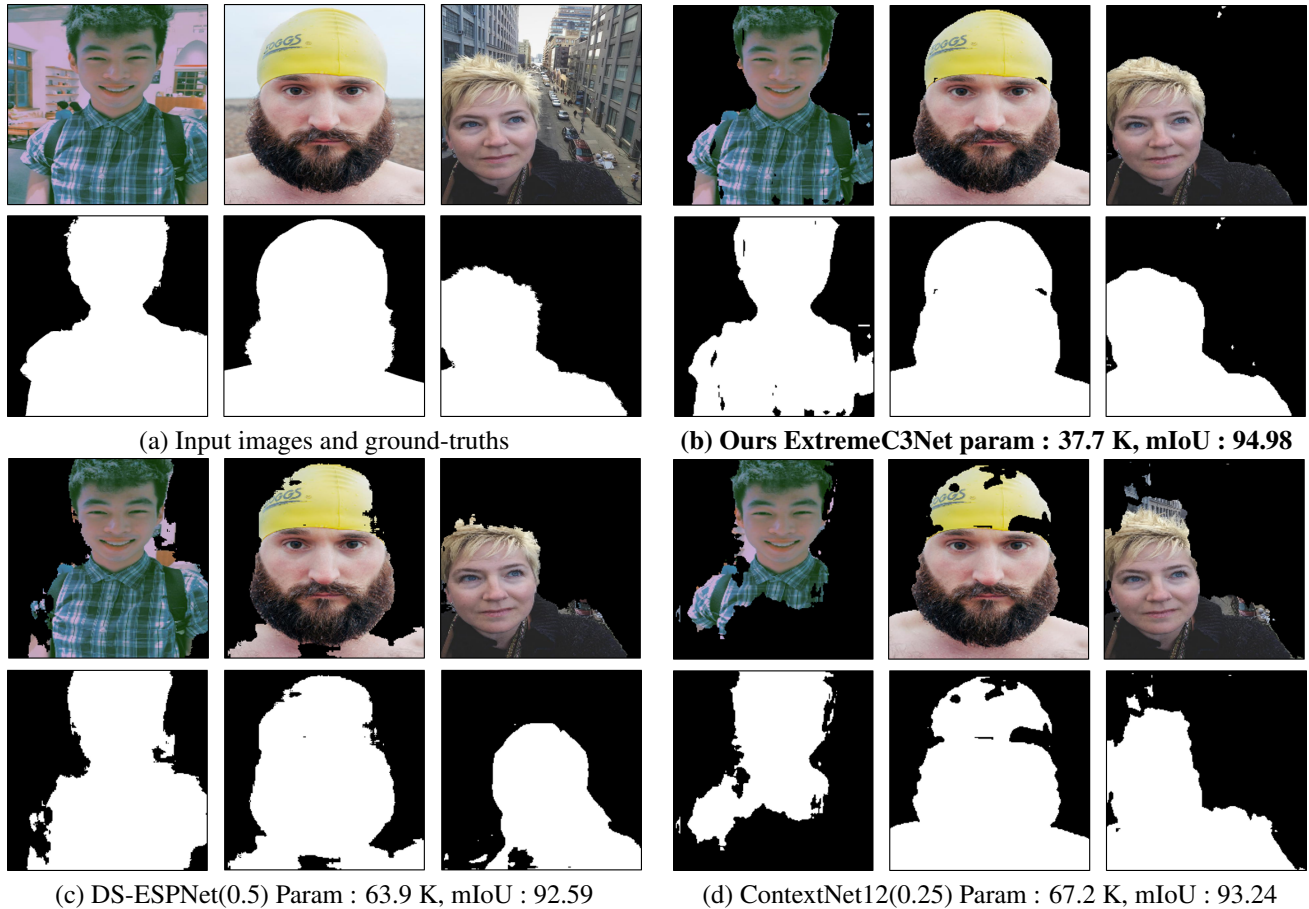
(a) Input images and ground-truths

(b) **Ours ExtremeC3Net param : 37.7 K, mIoU : 94.98**

(c) DS-ESPNet(0.5) Param : 63.9 K, mIoU : 92.59

(d) ContextNet12(0.25) Param : 67.2 K, mIoU : 93.24

Figure 4. Qualitative comparison results on the EG1800 validation dataset.

ples in EG1800 dataset. The PortriatNet compared their work to BiseNet, Enet, and PortraitNet. Therefore, we also re-trained BiSeNet and ENet following the method of PortraitNet for a fair comparison. As shown in Table 2, the accuracy of the re-trained results are slightly decreased due to the reduced size of the training dataset. The aforementioned approaches only counted floating-point operations (FLOPs) of convolution and batch normalization layers, which occupy a large portion of the total FLOPs. However, other operations such as activations and deconvolutions also affect total FLOPs, and these operations are hard to ignore when it comes to the lightweight model case. Therefore, in addition to the FLOPs counting in the official ESPNetV2 code, we also measured the FLOPS including all the operations. The calculation method is described in Supplementary material.

Among the comparison methods, DS-ESPNet has the same structure with ESPNet, with only changing the standard dilated convolutions of the model into depth-wise separable dilated convolutions. For ESPNetV2 (2.0) and ESPNetV2 (1.5), we changed the channel of convolution kernels to reduce the model size. We also reduced the chan-

nel of convolution kernels in DS-ESPNet (0.5) and ContextNet12 (0.25) by half and quarter from the original models to make the model less than 100K parameters and 0.2G FLOPs.

From Table 2, we can see that our proposed method showed comparable or better performance than the other models with less number of parameters and FLOPs. The SOTA PortraitNet showed the highest accuracy in all the experimental results, and has achieved even better performance than the heavier BiSeNet. However, still, PortraitNet requires a large number of parameters, which is a disadvantageous for using it on smaller devices. The proposed ExtremeC3Net has reduced the number of parameters by 98%, and FLOPs is reduced by half compared to PortraitNet, while maintaining accuracy. ESPNet and ESPNet V2 have similar accuracy, but showed a trade-off between the number of parameters and FLOPs. ESPNet V2 has more parameters than ESPNet, but ESPNet needs more FLOPs than ESPNet V2. Enet shows better performance than both models but requires more FLOPs. In comparison, our proposed method has less number of parameters and FLOPs,

| Method | mIoU |
|---|---|
| Baseline with advanced C3-module | 94.09 |
| + Using generated pseudo dataset | 94.27 |
| + Changing cross-entropy loss into Lovász loss | 94.45 |
| + Adding auxiliary loss into boundary area | 94.97 |
| + Applying our methods without pseudo dataset | 94.23 |

Table 3. Ablation study results about our proposed method.

| Method | Position of large ratios | mIoU |
|---|---|---|
| Baseline setting | All 8 layers | 93.16 |
| Advanced setting | L4, L5, L6, L7, L8 | 94.09 |
| Reverse setting | L1, L2, L3 | 92.33 |

Table 4. Ablation study results about adjusting the combination of dilated ratios in the C3-modules.

but better accuracy than ESPNet and ESPNet V2. Our model accuracy is 0.18% less than Enet, but our number of parameters is around 10 times less than Enet. In particular, our ExtremeC3Net has the highest accuracy in an extremely lightweight environment. Both DS-ESPNet(0.5) and ContextNet12(0.25) have a larger number of parameters, but their scores were much lower than our method. Figure 7 shows that the quality of our model is superior to other extremely lightweight models.

We compared the execution speed of the proposed model with SOTA portrait segmentation model PortraitNet on an Intel Core i7-9700 CPU environment with the PyTorch framework. PortraitNet needs 0.119 sec, but ExtremeC3Net takes only 0.062 sec for processing an image. In summary, the proposed ExtremeC3Net showed outstanding performance among the various segmentation model in terms of accuracy and speed.

## 4.2. Ablation Studies

Table 3 shows the accuracy improvement from the various ablations. The baseline model is trained with cross-entropy loss using the EG1800 dataset. Applying Lovász loss, the auxiliary loss, and increasing dataset size all enhance the mIoU. When we used only the Lovász loss and auxiliary loss without the additional 10,448 images, the accuracy was improved slightly. However, when we simultaneously applied our loss methods and the large additional dataset the accuracy was increased from 94.23 to 94.97.

Table 4 shows the importance of the combination of dilation ratios of the convolutional filters in the extremely lightweight model depending on the position of the filter. Our ExtremeC3Net consists of eight advanced C3-modules. The baseline setting uses the same dilation ratios $d$ in all 3C-modules, $d = [2, 4, 8]$. The advanced setting uses small dilation ratios at the shallow layers, which are close to the input image, and uses large ratios at the deeper layers, which

are far from the input image, as shown in Table 1. The reverse setting applied the dilated ratio of the filters by the opposite direction, which means smaller dilated ratios located in deeper layers. Our advanced setting performed better than baseline setting, and reverse setting showed much lower accuracy than the others. From the results, we can see that making small receptive fields for the local features in the shallow layers is critical for extremely lightweight portrait segmentation.

| | | # Img | Exp1 | Exp2 | Ours |
|---|---|---|---|---|---|
| Race | Caucasian | 214 | 93.67 | 94.51 | **94.81** |
| | Asian | 46 | 94.31 | **94.78** | 93.79 |
| | Black | 11 | 94.74 | 95.05 | **95.84** |
| Gender | Man | 112 | 92.85 | 94.21 | **94.43** |
| | Woman | 159 | 94.54 | 94.81 | **94.91** |
| Age | Child | 42 | 94.47 | 93.86 | **95.10** |
| | Youth | 219 | 93.87 | **94.79** | 94.63 |
| | Senior | 10 | 90.07 | 92.97 | **93.98** |
| Accuracy | | 271 | 94.09 | 94.59 | 94.98 |

Table 5. Validation results of each attribute about different data augmentation methods and our proposed method. The reason why total number of image is not 270 but 271 that there are two people in one validation image. A class of "Child" includes infant or child and a class of "Youth" denotes youth or middle-aged person. "Exp" means experiment.

## 4.3. Analysis of the EG1800 dataset

The public dataset EG1800 [21] contains 1,800 images for portrait stylizing from segmentation results. However, many urls are broken and the total amount of available images is 1,579. Six human annotators labeled the attributes of each portrait image with respect to race, gender, and age, and for images containing more than one person the annotators counted all of them. We divided the dataset into 18 classes as shown in Table 6, and illustrated the dataset bias in Figure 5. Both training and validation sets have a severe bias towards Caucasian and youth and middle-aged people, and we could not find any image containing a senior black woman. Detailed results can be found in the Supplemental materials.

We conducted a comparative experiment to understand how data augmentation methods can resolve the data bias problems. In Table 4.2, Experiment 1 and 2 use the same model structure (ExtremeC3Net trained with EG1800 dataset. However, they were trained with different data augmentation method. For the Experiment 1, we used naive data augmentation methods such as random resizing, random crop, and random horizontal flip. For the Experiment 2, we applied a more sophisticated data augmentation method proposed by PortraitNet, which consisting of the deformation and the texture augmentation methods [32].
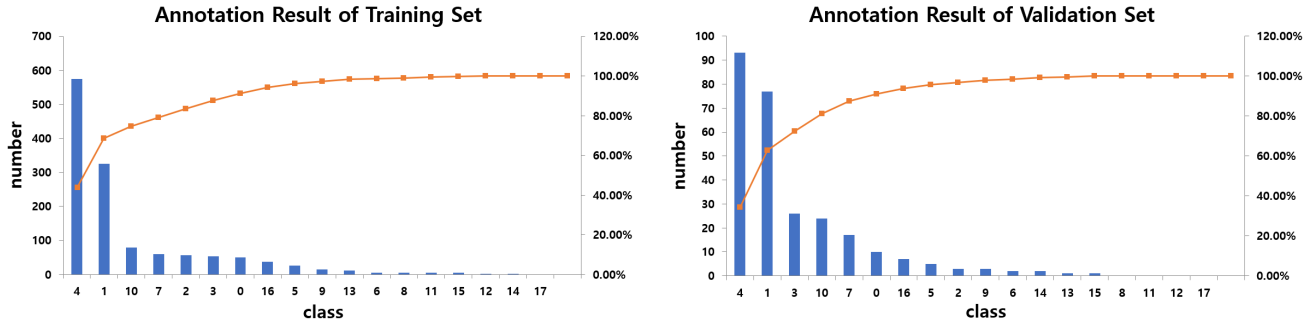
Figure 5. Dataset histogram about the number of each group image. The detailed number is described in supplementary material

| Class | Race | Gender | Age | Class | Race | Gender | Age | Class | Race | Gender | Age |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Caucasian | Man | Child | 6 | Asian | Man | Child | 12 | Black | Man | Child |
| 1 | Caucasian | Man | Youth | 7 | Asian | Man | Youth | 13 | Black | Man | Youth |
| 2 | Caucasian | Man | Senior | 8 | Asian | Man | Senior | 14 | Black | Man | Senior |
| 3 | Caucasian | Woman | Child | 9 | Asian | Woman | Child | 15 | Black | Woman | Child |
| 4 | Caucasian | Woman | Youth | 10 | Asian | Woman | Youth | 16 | Black | Woman | Youth |
| 5 | Caucasian | Woman | Senior | 11 | Asian | Woman | Senior | 17 | Black | Woman | Senior |

Table 6. Classes of attributes. "Child" includes infant or child and "Youth" denotes youth or middle-aged person.



Figure 6. Example images of age group. Row1 : seniors, Row 2: child, youth and middle-aged person from left to right

For the gender and age attributes, the number of images in the dataset has an impact on accuracy. The number of female images is 280 more than the number of male images in the training set, and it seems to make the model overfitted to the female portrait images. The numbers of child and seniors images were lower than youth images, but the accuracy for child images was still high compared to that of youth images. On the contrary, the accuracy of seniors is remarkably lower than the others. From the results, we conjecture that the bias between the races is not that important compared to those of other attributes, but the bias from the gender and age makes meaningful accuracy degradation. Also, the number of images for the child group did not bring a significant imbalanced impact on accuracy, compared to the senior group case.

The sophisticated data augmentation method used in Experiment 2 was effective for improving the accuracy, but it could not solve the imbalance problem completely; the accuracy of seniors was still lower than the other attributes. The possible reason of the phenomenon would be that the distinct features of the seniors in Figure 6, such as wrinkles and ages spots, would be difficult to be covered by the data augmentation method. From the results, the augmented data from the proposed data generation framework was shown to be effective for all the attributes and improved the balanced accuracy of each attribute. The accuracy disparity in an the age groups reduced from $1.82$ (Experiment 2) to $1.12$ (ours).

## 5. Conclusion

In this paper, we proposed ExtremeC3Net which is an extremely lightweight two-branched model consisting of CoarseNet and FineNet for solving the portrait segmentation task. The coarseNet produces the coarse segmentation map and the FineNet assists spatial details to catch the boundary information of the object. We also proposed the advanced C3-module for each network which elaborately

adjusts the dilation ratio of the convolutional filters according to their positions. From the experiments on a public portrait segmentation dataset, our model obtained outstanding performance compared to the existing lightweight segmentation models. Also, we proposed a simple data generation framework covering the two situations: 1) having human segmentation ground truths 2) having only raw images. We also analysed the configuration of the public dataset given portrait attributes and studied the effects of the bias among each attribute regarding the accuracy. The additionally labeled samples we generated were shown to be helpful for improving the segmentation accuracy for all the attributes.

# References

[1] Y. Baek, B. Lee, D. Han, S. Yun, and H. Lee. Character region awareness for text detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9365–9374, 2019.

[2] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei. Whats the point: Semantic segmentation with point supervision. In *European conference on computer vision*, pages 549–565. Springer, 2016.

[3] M. Berman, A. Rannen Triki, and M. B. Blaschko. The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4413–4421, 2018.

[4] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018.

[5] F. Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint*, pages 1610–02357, 2017.

[6] T. DeVries and G. W. Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.

[7] A. Gupta, A. Vedaldi, and A. Zisserman. Synthetic data for text localisation in natural images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2315–2324, 2016.

[8] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[9] B. Jin, M. V. Ortiz Segovia, and S. Susstrunk. Webly supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3626–3635, 2017.

[10] 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*, pages 3431–3440, 2015.

[11] S. Mehta, M. Rastegari, L. Shapiro, and H. Hajishirzi. Espnetv2: A light-weight, power efficient, and general purpose convolutional neural network. *arXiv preprint arXiv:1811.11431*, 2018.

[12] 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.

[13] 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.

[14] R. P. Poudel, U. Bonde, S. Liwicki, and C. Zach. Contextnet: Exploring context and detail for semantic segmentation in real-time. *arXiv preprint arXiv:1805.04554*, 2018.

[15] R. P. Poudel, S. Liwicki, and R. Cipolla. Fast-scnn: fast semantic segmentation network. *arXiv preprint arXiv:1902.04502*, 2019.

[16] X. Qi, Z. Liu, J. Shi, H. Zhao, and J. Jia. Augmented feedback in semantic segmentation under image level supervision. In *European Conference on Computer Vision*, pages 90–105. Springer, 2016.

[17] E. Romera, J. M. Alvarez, L. M. Bergasa, and R. Arroyo. Erfnet: Efficient residual factorized convnet for real-time semantic segmentation. *IEEE Transactions on Intelligent Transportation Systems*, 19(1):263–272, 2018.

[18] A. C. L. S. Sachin Mehta, Mohammad Rastegari and H. Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *ECCV*, 2018.

[19] F. Saleh, M. S. Aliakbarian, M. Salzmann, L. Petersson, S. Gould, and J. M. Alvarez. Built-in foreground/background prior for weakly-supervised semantic segmentation. In *European Conference on Computer Vision*, pages 413–432. Springer, 2016.

[20] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation. *arXiv preprint arXiv:1801.04381*, 2018.

[21] X. Shen, A. Hertzmann, J. Jia, S. Paris, B. Price, E. Shechtman, and I. Sachs. Automatic portrait segmentation for image stylization. In *Computer Graphics Forum*, volume 35, pages 93–102. Wiley Online Library, 2016.

[22] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, volume 4, page 12, 2017.

[23] 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 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015.

[24] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.

[25] Y. Wei, X. Liang, Y. Chen, Z. Jie, Y. Xiao, Y. Zhao, and S. Yan. Learning to segment with image-level annotations. *Pattern Recognition*, 59:234–244, 2016.

[26] Z. Wu, Y. Huang, Y. Yu, L. Wang, and T. Tan. Early hierarchical contexts learned by convolutional networks for image segmentation. In *2014 22nd International Conference on Pattern Recognition*, pages 1538–1543. IEEE, 2014.

[27] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 5987–5995. IEEE, 2017.

[28] Y. Yoo, D. Han, and S. Yun. Extd: Extremely tiny face detector via iterative filter reuse. *arXiv preprint arXiv:1906.06579*, 2019.

[29] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 325–341, 2018.

[30] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. *arXiv preprint arXiv:1905.04899*, 2019.

[31] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.

[32] S.-H. Zhang, X. Dong, H. Li, R. Li, and Y.-L. Yang. Portraitnet: Real-time portrait segmentation network for mobile device. *Computers & Graphics*, 80:104–113, 2019.

[33] X. Zhang, X. Zhou, M. Lin, and J. Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. *arXiv preprint arXiv:1707.01083*, 2017.

# Appendix

In this supplementary material, we provide additional results and methods that we could not include due to space limitation. We indicate detailed number of dataset histogram from our additional annotations. We also explain the method for measuring floating-point operations (FLOPs) and show other segmentation examples for qualitative comparison of our method with two extreme lightweight segmentation models based on ContextNet[14] and ESPNet[18].

## Detailed results of additional annotations

Table 7 is detailed numbers of dataset histogram figure in Experiment section. Most of both training and validation sets consist of Caucasian and Youth or Middle-aged person, as shown in Table 7. Even, some of the group images do not exist in the dataset, and the bias of configuration in the dataset is more severe in validation set such as senior Black woman, Black child, and senior Asian man. The total frequency is greater than the total number of data sets (training: 1,309 and validation: 270) because sometimes there is more than one person in the image.

## Additional qualitative comparison results

We show other example results with covering various attributes that we could not illustrate in the Experiment section due to space limitation. Our model has less complexity than other extreme lightweight model but shows better performance as shown in Figure 7.

## FLOPs Calculation

Table 8 shows how we calculated FLOPs for each operation. The following notations are used.

$F$ : A input feature map
$O$ : A output feature map
$K$ : A convolution kernel
$K_h$ : A height of convolution kernel
$K_w$ : A width of convolution kernel
$H_i$ : A height of input feature map
$W_i$ : A width of input feature map
$C_i$ : A input channel dimension of feature map or kernel
$C_o$ : A output channel dimension of feature map or kernel
$g$ : A group size for channel dimension
$H_o$ : A height of output feature map
$W_o$ : A width of output feature map
$g(\cdot)$ : A non-linear activation function

| Training set | | | Validation set | | |
|---|---|---|---|---|---|
| Class | # Frequency | Cumulative value (%) | Class | # Frequency | Cumulative value (%) |
| 4 | 575 | 3.81% | 4 | 93 | 34.32% |
| 1 | 326 | 28.66% | 1 | 77 | 62.73% |
| 10 | 79 | 32.93% | 3 | 26 | 72.32% |
| 7 | 59 | 36.97% | 10 | 24 | 81.18% |
| 2 | 56 | 80.79% | 7 | 17 | 87.45% |
| 3 | 53 | 82.77% | 0 | 10 | 91.14% |
| 0 | 50 | 83.16% | 16 | 7 | 93.73% |
| 16 | 38 | 87.65% | 5 | 5 | 95.57% |
| 5 | 26 | 88.03% | 2 | 3 | 96.68% |
| 9 | 15 | 89.18% | 9 | 3 | 97.79% |
| 13 | 12 | 95.20% | 6 | 2 | 98.52% |
| 6 | 5 | 95.58% | 14 | 2 | 99.26% |
| 8 | 5 | 95.73% | 13 | 1 | 99.63% |
| 11 | 5 | 96.65% | 15 | 1 | 100.00% |
| 15 | 5 | 96.72% | 8 | 0 | 100.00% |
| 12 | 2 | 97.10% | 11 | 0 | 100.00% |
| 14 | 1 | 100.00% | 12 | 0 | 100.00% |
| 17 | 0 | 100.00% | 17 | 0 | 100.00% |
| Total | 1312 | 100.00% | Total | 271 | 100.00% |

Table 7. Detailed number of dataset histogram of each group image

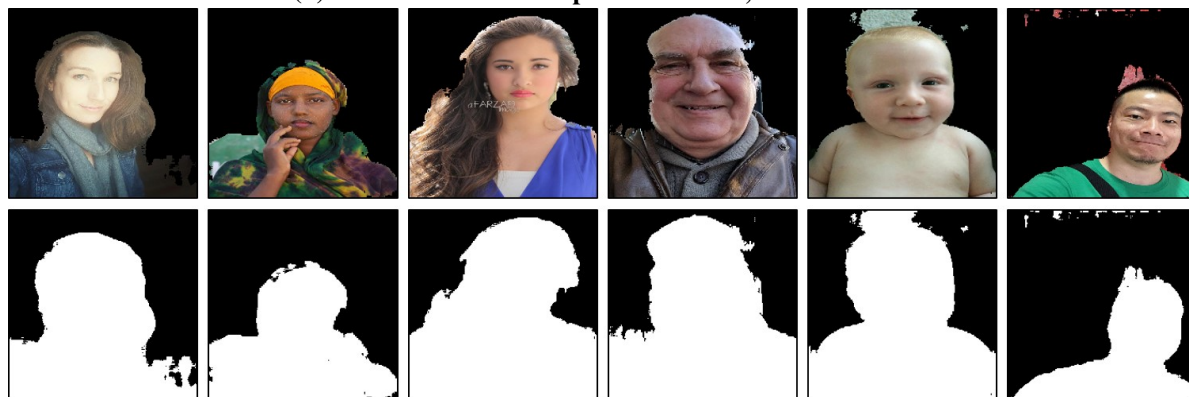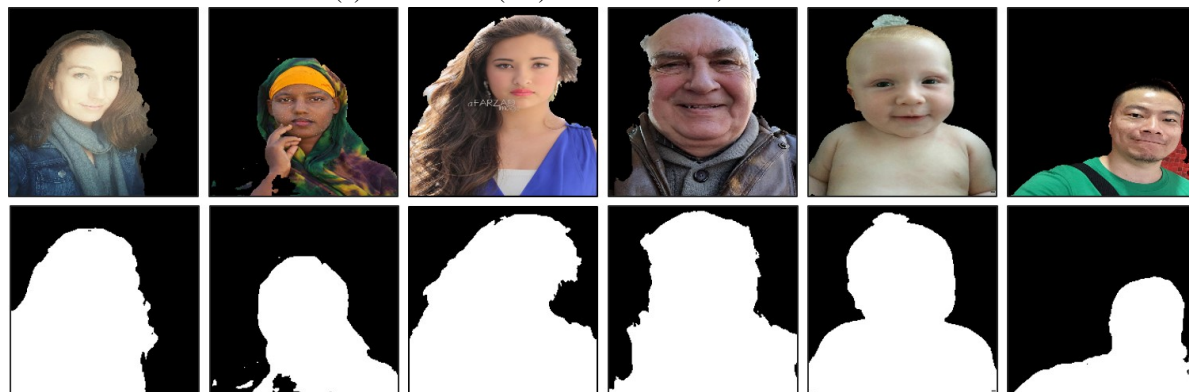| Layer | Operation | Flop |
|---|---|---|
| Convolution | $O = F * K$ | $2 \cdot H_o W_o \cdot K_h K_w \cdot C_i C_o / g$ |
| Deconvolution | $O = F * K$ | $2 \cdot H_i W_i \cdot K_h K_w \cdot C_i C_o / g$ |
| Average Pooling | $O = Avg(F)$ | $H_i \cdot W_i \cdot C_i$ |
| Bilinear upsampling | $f(x,y) = \sum_{i=0}^{1} \sum_{j=0}^{1} a_{ij} x^i y^j$ | $3 \cdot H_i \cdot W_i \cdot C_i$ |
| Batch normalization | $(F - mean)/std$ | $2 \cdot H_i \cdot W_i \cdot C_i$ |
| ReLU or PReLU | $O = g(F)$ | $H_i \cdot W_i \cdot C_i$ |

Table 8. The detail method for calculating FLOPs

(a) Input images and ground-truths

(b) **Ours ExtremeC3Net param : 37.9 K, mIoU : 94.98**

(c) DS-ESPNet(0.5) Param : 63.9 K, mIoU : 92.59

(d) ContextNet12(0.25) Param : 67.2 K, mIoU : 93.24

Figure 7. Additional qualitative comparison results on the EG1800 [21] validation dataset.