

PortraitNet: Real-time Portrait Segmentation Network for Mobile Device

Leonardo Pesce 223100001
CIE6004: Image Processing and Computer Vision

Abstract—The paper "PortraitNet: Real-time Portrait Segmentation Network for Mobile Device" [1] introduces PortraitNet, a real-time portrait segmentation model designed for mobile devices. It utilizes a lightweight U-shape architecture with two auxiliary losses during training, specifically the boundary loss and consistency constraint loss, to improve accuracy and robustness. PortraitNet outperforms existing methods in accuracy and efficiency, particularly in generating sharper boundaries and handling challenging lighting conditions. Additionally, it can process 224×224 RGB images at 30 FPS on an iPhone 7 in real-time.

I. INTRODUCTION

The paper discusses the growing importance of real-time portrait segmentation, driven by the demand for mobile applications that enable background editing of portrait images. It highlights the challenges posed by the unique characteristics of portrait images, such as the need for precise segmentation in the presence of ambiguous boundaries and varying illumination conditions. Existing methods have prioritized accuracy over efficiency, making them unsuitable for real-time mobile applications.

To address these challenges, the authors propose PortraitNet, a specialized convolutional neural network designed for real-time portrait segmentation on mobile devices. PortraitNet utilizes a specific network architecture with a $32\times$ downsampling rate to enhance efficiency and employs a U-shape architecture for better segmentation. The network is trained using auxiliary losses, including a boundary loss to improve boundary accuracy and a consistency constraint loss to enhance robustness under varying illumination.

The results demonstrate that PortraitNet achieves high accuracy, with 96.62% accuracy on the EG1800 dataset and 93.43% on the SupervisePortrait dataset. Furthermore, it can process images at 30 frames per second on an iPhone 7 with an input size of 224×224 , making it suitable for real-time mobile portrait segmentation.

II. METHOD

A. Model Structure

Let's now have a look at the model implemented in the paper 1. The model is based on the MobileNet-v2 [2], which is an architecture specifically designed for mobile devices (low computational power). The MobileNetV2 design employs an inverted residual structure, which differs from traditional residual models by utilizing slender bottleneck layers for both input and output within the residual block. Instead of

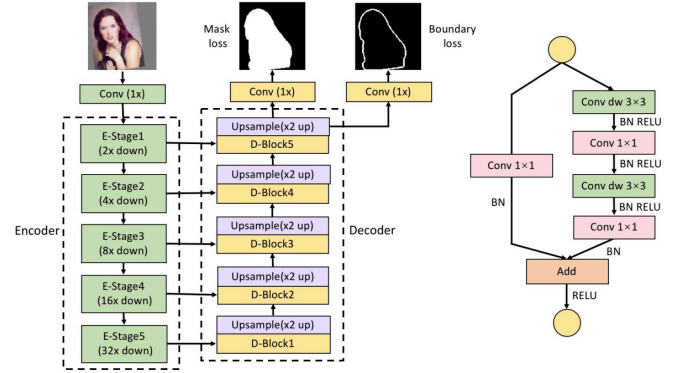


Fig. 1: Structure of the model: on the left the U shape, while on the right the deconvolution block

expanded representations in the input, MobileNetV2 employs lightweight depthwise convolutions to filter features in the intermediate expansion layer.

The structure is U-shape with a $32\times$ downsampling in the encoder phase that combined with the residual architecture extracts global and spatial information about the image.

B. Auxiliary Losses

To make up for the lower complexity of the model (since it needs to run in real-time) two auxiliary losses are introduced.

1) *Boundary Loss*: Given that, portrait segmentation needs sharper edges and that the edge represents less than 10% of the labeled dataset we use the focal loss to balance this fact. Focal loss L_e is combined with cross-entropy loss L_m in the following way:

$$L = L_m + \lambda L_e \quad (1)$$

$$L_m = - \sum_{i=1}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i)) \quad (2)$$

$$L_e = - \sum_{i=1}^n ((1 - p_i)^\gamma y_i \log(p_i) + p_i^\gamma (1 - y_i) \log(1 - p_i)) \quad (3)$$

Where γ is the weight of the boundary loss, and the probability p_i is computed as

$$p_i(z_j) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad (4)$$

With z the original output of potraitNet and K the number of ground-truth classes.

2) *Consistency Constraint Loss*: In general, in a segmentation process, hard labels (binary labels) are used. Given our complexity requirements, since it has been shown that for small models soft labels help in the training process, a new method to generate soft labels is proposed (usually generating soft labels requires a teacher model that needs extensive training). Given an input A a second input A' is generated simply by a texture enhancement transformation (change in brightness, contrast, sharpness, noise...). At this point, both images are passed through the model obtaining the output B and B' . B' is worse than B because it was generated by A' , a lower quality version of A . In the end, B is used as a higher quality soft label for B' and the KL-divergence is used to compute the loss between the two:

$$L = L'_M + \alpha \times L_c \quad (5)$$

$$L_c = \frac{1}{n} \sum_{i=1}^n q_i \times \log \frac{q_i}{q'_i} \times T^2 \quad (6)$$

$$L'_m = - \sum_{i=1}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i)) - \sum_{i=1}^n (y_i \log(p'_i) + (1 - y_i) \log(1 - p'_i)) \quad (7)$$

Where α is used to balance the two losses, T to smooth the output, p_i and p'_i are:

$$p_i(z_j) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad p'_i(z'_j) = \frac{e^{z'_j}}{\sum_{k=1}^K e^{z'_k}} \quad (8)$$

and q_i and q'_i are:

$$q_i(z_j) = \frac{e^{\frac{z_j}{T}}}{\sum_{k=1}^K e^{\frac{z_k}{T}}} \quad q'_i(z'_j) = \frac{e^{\frac{z'_j}{T}}}{\sum_{k=1}^K e^{\frac{z'_k}{T}}} \quad (9)$$

Adding this loss results in higher accuracy and robustness under different lighting conditions.

C. Pseudo-code

We will now analyze the pseudo-code of training. The model training follows the typical structure. The dataset was split into 3 with train, test, and final evaluation. The first two are used during the epochs of the model training phase to compute new parameters of the network and check if the values are better than before. The last is used to compute the final accuracy of the model.

Algorithm 1 pseudocode for model training

- 1: computes the learning rate
 - 2: executes a training cycle
 - 3: executes a testing phase to update the loss
 - 4: save the model if the loss is the lowest between all the cycles
-

The structure of a training cycle is as follows.

Algorithm 2 pseudocode for the training cycle

- 1: **for** every element in train set **do**
 - 2: i=image, ie=image enhanced, m=mask, e=edge
 - 3: mi, ei = network(i) ▷ mask i, edge i
 - 4: mie, eie = network(ie) ▷ mask ie, edge ie
 - 5: loss += softmax_loss(mi, m)
 - 6: loss += focal_loss(ei, e) * edge_ratio
 - 7: loss += softmax_loss(mie, m)
 - 8: loss += focal_loss(eie, e) * edge_ratio
 - 9: loss += KL_loss(mi, mie)
 - 10: loss += KL_loss(ei, eie)
 - 11: computes the gradient and updates the parameters
 - 12: **end for**
-

The first loss is the one between the mask computed on the image and the real mask, the second is the loss computed between the predicted edges on the image and the real edges. The next two are the same but the prediction was made with the enhanced image. The last two are the Kullback-Leibler Divergence loss between the mask and edge predicted on the real image and the mask and edge predicted on the enhanced image.

The test cycle is simply composed by the sum of all the Intersection Over Union (IoU) computed on the test set.

Algorithm 3 pseudocode for the testing cycle

- 1: **for** every element in test set **do**
 - 2: i=image and m=mask
 - 3: mi, ei = network(i) (mask i, edge i)
 - 4: total_IoU += IoU(mi, m)
 - 5: **end for**
 - 6: **return** total_IoU
-

III. RESULTS

The model trained for 1000 epochs with a batch size of 64, a learning rate of 0.001, and a weight decay of 0.0005. The final model accuracy (tested on 130 images) was 94.97%, the accuracy metric was IoU on the real and computed masks. In the following graph, we can see the training loss curve, which follows a typical inverse function decreasing sharply at the beginning and reducing the convergence speed gradually.



Fig. 2: example 1

We now run our model on the dataset provided EG1800 and test its accuracy and results. For every example there are 4 images, the first and the third contain the predicted mask and edge of the person, while the second and the fourth contain the real mask and edge.

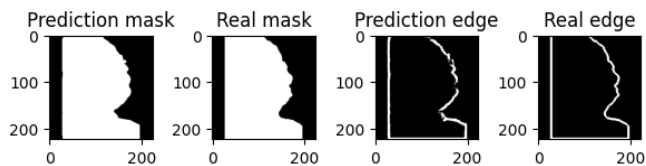


Fig. 3: example 1

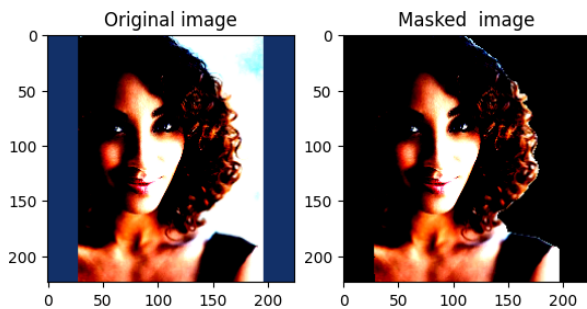


Fig. 4: mask result of example 1

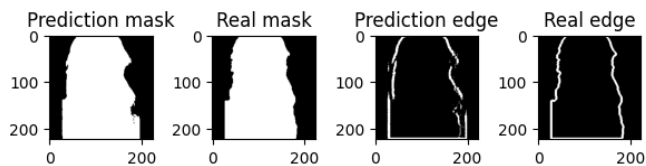


Fig. 5: example 2

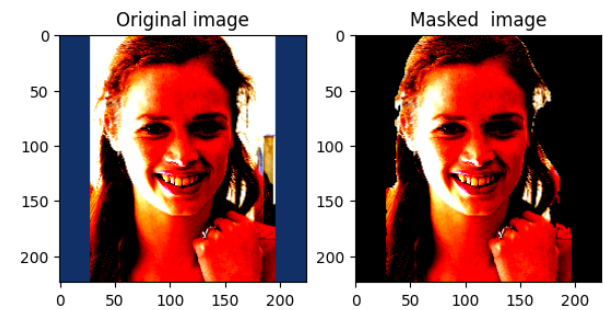


Fig. 6: mask result of example 2

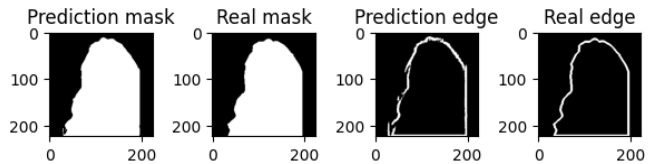


Fig. 7: example 3

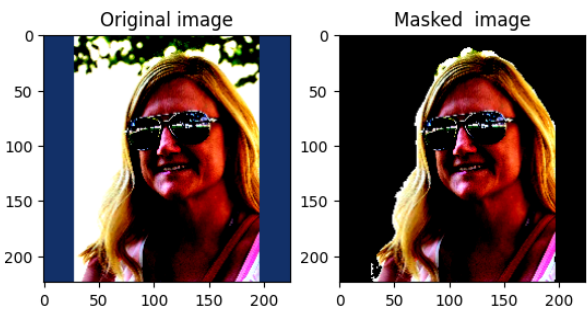


Fig. 8: mask result of example 3

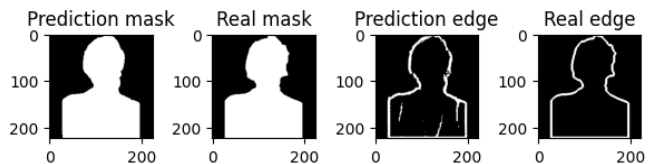


Fig. 9: example 4

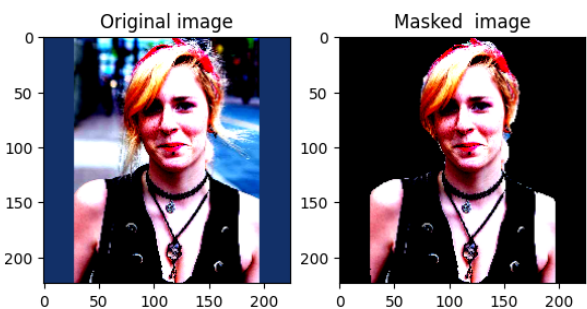


Fig. 10: mask result of example 4

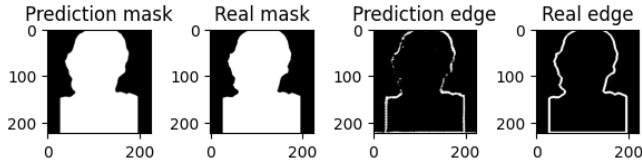


Fig. 11: example 5

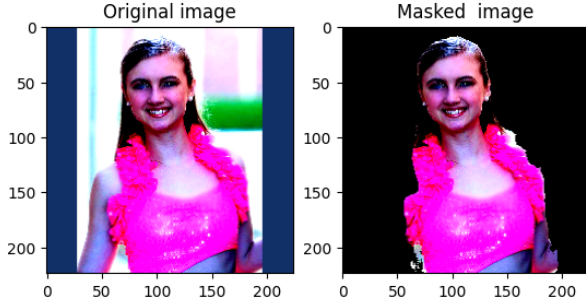


Fig. 12: mask result of example 5

REFERENCES

- [1] S.-H. Zhang, X. Dong, H. Li, R. Li, and Y.-L. Yang, "Portraitnet: Real-time portrait segmentation network for mobile device," *Comput. Graph.*, vol. 80, pp. 104–113, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:131924823>
- [2] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation," *CoRR*, vol. abs/1801.04381, 2018. [Online]. Available: <http://arxiv.org/abs/1801.04381>

A. Encountered problems

The method proposed required a long training (at least a long training for my PC without a GPU). The problem was solved by optimizing as much as possible the code and using Google Colab to use the free GPU resources available. Also instead of training for 2000 epochs we only trained for 1000 epochs (which given the curve 2 shouldn't change too much the results).

B. Consideration and possible improvements

The model in itself, given the complexity of it, works very well. Errors typically lie in hair strands as visible in 10 or in lower parts of the body as in 12 (where the arm is cut off). The first one is a problem that is quite challenging to solve since a solution could be increasing the model size to increase the learning capabilities (but this is in contrast with the objective of the paper), while the second one is due to the images being mainly of the upper part of the body, thus a solution could be extending the dataset.

The model accuracy differs slightly from the one illustrated in the paper probably for two reasons, first, they trained for double the amount of epochs, and second, they removed the KL loss on the edge and the softmax loss on the enhanced image to better optimize for the mask prediction.

IV. CONCLUSIONS

The paper discusses the creation of a light and fast neural network architecture based on MobileNet-v2 to segment people inside the picture (with a target on mobile devices). Of fundamental contribution in this paper besides the structure of the network are also the boundary loss and the consistency constraint loss. We analyzed the structure and mathematical understanding of the paper and finally proposed an implementation of it to validate the results proposed by the author.