# COMPLEMENTARY NETWORK WITH ADAPTIVE RECEPTIVE FIELDS FOR MELANOMA SEGMENTATION

*Xiaoqing Guo, Zhen Chen, Yixuan Yuan*

Department of Electrical Engineering, City Univeristy of Hong Kong, Hong Kong SAR, China

## ABSTRACT

Automatic melanoma segmentation in dermoscopic images is essential in computer-aided diagnosis of skin cancer. Existing methods may suffer from the hole and shrink problems with limited segmentation performance. To tackle these issues, we propose a novel complementary network with adaptive receptive filed learning. Instead of regarding the segmentation task independently, we introduce a foreground network to detect melanoma lesions and a background network to mask non-melanoma regions. Moreover, we propose adaptive atrous convolution (AAC) and knowledge aggregation module (KAM) to fill holes and alleviate the shrink problems. AAC explicitly controls the receptive field at multiple scales and KAM convolves shallow feature maps by dilated convolutions with adaptive receptive fields, which are adjusted according to deep feature maps. In addition, a novel mutual loss is proposed to utilize the dependency between the foreground and background networks, thereby enabling the reciprocally influence within these two networks. Consequently, this mutual training strategy enables the semi-supervised learning and improve the boundary-sensitivity. Training with Skin Imaging Collaboration (ISIC) 2018 skin lesion segmentation dataset, our method achieves a dice co-efficient of 86.4% and shows better performance compared with state-of-the-art melanoma segmentation methods[1].

***Index Terms***— Melanoma segmentation, adaptive receptive fields, semi-supervised learning

## 1. INTRODUCTION

Melanoma is the most dangerous form of skin cancer, accounting for a large percentage of skin cancer deaths [7]. Fortunately, if detected early, melanoma survival rate exceeds 95% [7]. Dermoscopy is an imaging technique to visualize deep levels of the skin and it is widely applied to diagnose melanoma. However, manually reviewing dermoscopy images is an error-prone and time-consuming work even for professional dermatologists. In this regard, the development of

[1] https://github.com/Guo-Xiaoqing/Skin-Seg



**Fig. 1**. Illustrations of (a) hole problem, (b) shrink problem. Each group includes the original image, ground truth and prediction of U-Net [5] from left to right.

computational support systems for automated segmentation and analysis of dermoscopy images is highly desirable.

Automated melanoma segmentation remains to be challenging due to the huge variations of melanomas in terms of shape, color and texture. Moreover, some samples may contain artifacts, such as hairs, ruler marks and color calibrations, blurring melanoma lesions. Many algorithms have been proposed to tackle these challenges [3, 6, 8]. Yuan et al. [8] incorporated Lightness channel from CIELAB color space and three channels of HSV space together for the melanoma segmentation. Sarker et al. [6] presented a melanoma segmentation model with negative log likelihood and end point error loss functions to preserve sharp boundaries. Li et al. [3] proposed a dense deconvolutional network with hierarchical supervision to capture local and global contextual information for melanoma segmentation. Although existing methods have achieved significant success, they still suffer from the hole (Fig.1 (a)) and shrink (Fig.1 (b)) in predictions. The relatively low contrast between melanoma and non-melanoma regions confuses the network and causes the appearance of holes. The fuzzy boundaries lead to the shrinking prediction and further decrease the sensitivity of prediction.

To address the hole and shrink problems mentioned before, we propose a complementary network consisting of a foreground segmentation network and a background segmentation network. With the fact that the dependency of two networks is crucial, we propose a mutual loss to optimize the complementary network collaboratively. In this way, our network is sensitive to boundary and can effectively deal with shrink problem. Additionally, we propose AAC to explicitly control the receptive field for incorporating local and context information, and KAM to convolve shallow features by adaptive receptive field kernels learned from deep features. With AAC and KAM, our model can expand the segmented region to fit the ground truth lesion and fill holes.
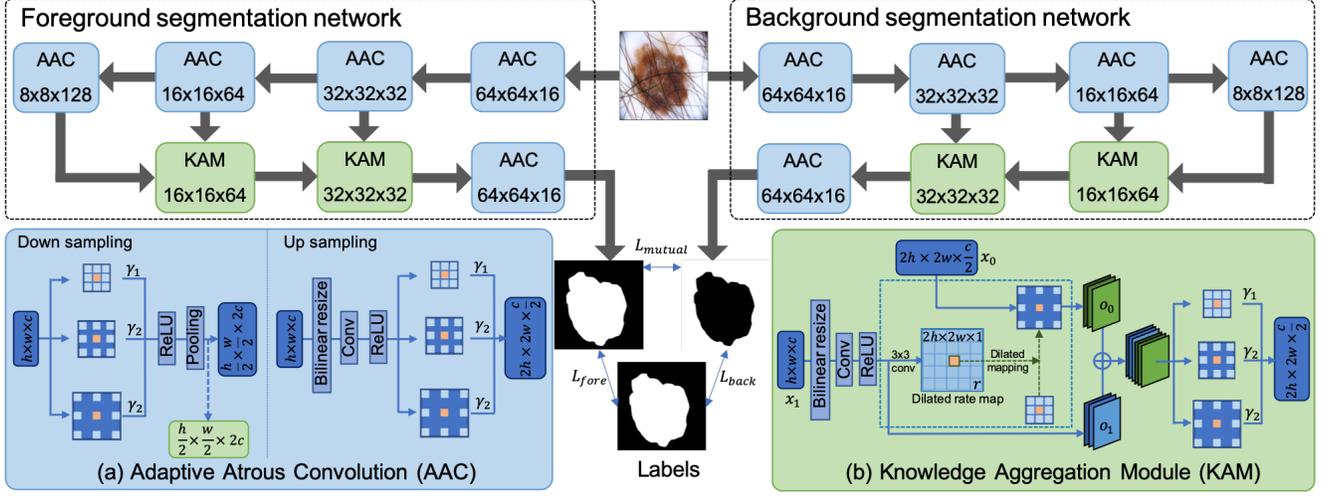
**Fig. 2**. Illustration of the proposed complementary segmentation network, including foreground and background networks.

## 2. METHOD

In this paper, we propose a complementary network for melanoma segmentation in Fig.2. Our model is composed of two networks, foreground and background segmentation networks. For each segmentation network, the input image is firstly passed through the contracting path with four downsampling AACs. The sequences of AACs extract multi-scale features from low-level contexture details to high-level semantic structures. Then, extracted features from different stages are incorporated by KAMs, which can aggregate informative features and suppress noises. After that, an upsampling AAC is introduced to further extend the size of feature maps. Finally, an upsampling layer and a fully connected layer are used to predict the class result per pixel. The foreground segmentation network and background network are trained collaboratively and jointly by minimizing the proposed mutual loss together with their individual foreground and background loss.

### 2.1. Adaptive Atrous Convolution

Atrous convolution [3, 6] allows us to explicitly enlarge the receptive field of filters. However, it is tricky to manually choose the dilated rate, and the fixed-size receptive field is the main limitation of atrous convolution. We propose the AAC module to alleviate information loss caused by sparse kernels and adaptively enlarge the receptive field as in Fig.2 (a). We first utilize various dilated rates to enlarge the receptive field of kernels. Considering different sizes and shapes of objects in images, the importance of feature maps with different receptive fields may not be equal. Therefore, we assign an importance score for each output of dilated convolution, and AAC can be formulated as:

$$g[i,j] = \sum_{k=1}^{K} \sum_{m=1}^{M} \sum_{n=1}^{N} \gamma_k \cdot f[i + r_k \cdot m, j + r_k \cdot n] \cdot h[m,n], \quad (1)$$

where $f[i,j]$ and $g[i,j]$ are input and output feature maps of AAC. $h[m,n]$ denotes the convolution kernel with $M \times N$ size. $r_k$ is dilated rate that equals to $k$, and $K$ is set as 3. $\gamma_k$ denotes the important score for dilated convolution with rate $r_k$. All $\gamma$ values are initialized as $\frac{1}{3}$ and updated every iteration during training phases. Thus, each layer can adjust to its appropriate receptive fields gradually. Finally, the weighted output feature maps are added together for further processing. With receptive fields being adaptively enlarged, the context information can be incorporated into local information, and the problem of holes and shrinks can be relieved.

### 2.2. Knowledge Aggregation Module

Previous methods usually use the "skip-connection" to concatenate multi-level features directly for accurate segmentation [3, 5, 6]. However, the equal weights for different channels of concatenated features are deficient due to the redundant information. Moreover, it is necessary to incorporate context with local information to fill holes since the context information may provide rich clues for local prediction. Therefore, we propose KAM to amalgamate and distill multi-level features as in Fig.2 (b). Assume feature maps from last layer are $x_1 \in \mathbb{R}^{h \times w \times c}$ and features from the contracting path are $x_0 \in \mathbb{R}^{2h \times 2w \times \frac{c}{2}}$, $x_1$ is first convoluted to $o_1 \in \mathbb{R}^{2h \times 2w \times \frac{c}{2}}$. Then, a $3 \times 3$ convolution layer, rate learning layer, is applied to $o_1$ to learn dilated rate map $r \in \mathbb{R}^{2h \times 2w \times 1}$. Each pixel in $r$ indicates the dilated rate of convolution kernel at the corresponding position, and the dilated rate map controls the scaling of receptive fields for each position individually. The weight of the rate learning layer is initialized as $N(0, \sigma^2)$ with $\sigma \ll 1$, and the bias is initialized as ones. This initialization method makes the convolution kernel start from the standard convolutions and gradually adjust to the appropriate dilated rate. Through dilated mapping, the learned dilated rate are applied to each

position for convolving $x_0$, which can be formulated as:

$$o_0[i,j] = \sum_{m=1}^{M} \sum_{n=1}^{N} x_0[i + r[i,j] \cdot m, j + r[i,j] \cdot n] \cdot h[m,n]. \quad (2)$$

When the coordinates $(i + r[i,j] \cdot m, j + r[i,j] \cdot n)$ are not at grids, bilinear interpolation is utilized for approximation as in [2]. Finally, the output $o_0$ of adaptive dilated convolution is concatenated with $o_1$ for further processing. Compared with the standard convolution operator, this dilated mapping enables flexible-size receptive fields according to the semantic information. Moreover, KAM incorporates context information from deep semantic features to texture details, which can alleviate the hole problem and suppress noises.

## 2.3. Joint Loss Function

### 2.3.1. Foreground and Background Loss.

We use focal loss [4] and Jaccard loss [8] to prevent the problem of foreground-background class imbalance. Assume $y_{ij}$ is the ground truth of pixel $j$ in image $i$, its corresponding probability predicted for correct class can be calculated by $p_{ij}^f = \frac{e^{W_{y_{ij}}^\top x_{ij} + b_{y_{ij}}}}{\sum_{k=1}^{2} e^{W_k^\top x_{ij} + b_{ij}}}$, where $W$ and $b$ are weights and bias in the fully connected layer. Then loss function of foreground network can be formulated as :

$$L_{fore} = -\sum_{i=1}^{N} \sum_{j=1}^{n} \frac{(1 - p_{ij}^f)^2 \log(p_{ij}^f)}{N \times n} + \frac{p_{ij}^f y_{ij}}{p_{ij}^f + y_{ij} - p_{ij}^f y_{ij}}, \quad (3)$$

where $N$ denotes mini-batch size, and $n$ denotes the number of pixels in a dermoscopy image. The first term in $L_{fore}$ is focal loss, while the second one is Jaccard loss.

In background segmentation network, $1 - y_{ij}$ is ground truth of background pixel $j$ in image $i$. Similar to foreground loss, background loss function is calculated as:

$$L_{back} = -\sum_{i=1}^{N} \sum_{j=1}^{n} \frac{(1 - p_{ij}^b)^2 \log(p_{ij}^b)}{N \times n} + \frac{p_{ij}^b (1 - y_{ij})}{(1 - y_{ij}) + p_{ij}^b y_{ij}}. \quad (4)$$

This loss function can alleviate class imbalance problem, and avoid additional procedures to re-balance pixels from melanoma region and background.

### 2.3.2. Mutual Loss

To exploit complementary information among the foreground network and background network, we introduce a constraint by minimizing the similarity of predictions from two networks. In particular, we utilize Jensen-Shannon (JS) divergence to measure the similarity of predictions from two networks and propose an exclusion loss to enforce the predicted segmentations from two networks mutually exclusive. The mutual loss thus can be formulated as:

$$L_{mutual} = \frac{1}{N \times n} \sum_{i=1}^{N} \sum_{j=1}^{n} \frac{p_{ij}^f}{2} \log \frac{2p_{ij}^f}{p_{ij}^f + (1 - p_{ij}^b)} +$$
$$\frac{1 - p_{ij}^b}{2} \log \frac{2(1 - p_{ij}^b)}{p_{ij}^f + (1 - p_{ij}^b)} + \frac{2p_{ij}^0 p_{ij}^1}{p_{ij}^0 + p_{ij}^1}, \quad (5)$$

where the first and second terms in $L_{mutual}$ are JS loss, and the third one is the exclusion loss. By minimizing the mutual loss, distributions of $p_{ij}^f$ and $1 - p_{ij}^b$ are constrained to be similar, and the overlap of predictions from two networks tends to be minimized. Mutual loss enhances the complementary information within these two networks. Moreover, the prediction around the boundary is prone to agree with ground truth, and the problem of shrinks will be alleviated.

### 2.3.3. Extension to Semi-supervised Learning

The proposed complementary network can be extended to semi-supervised learning. Under the semi-supervised learning setting, we use foreground loss and background loss for labeled data and compute mutual loss for all training data. Denote the labeled and unlabeled data as $\mathcal{L}$ and $\mathcal{U}$, the total loss function can be represented as

$$L_{Total} = \underset{x \in \mathcal{L}}{L_{fore}} + \underset{x \in \mathcal{L}}{L_{back}} + \underset{x \in \mathcal{D}}{L_{mutual}}, \quad (6)$$

where $\mathcal{D} = \mathcal{L} \cup \mathcal{U}$. Therefore, the complementary network can be not only optimized with pixel-level annotation, but also supervised by dual networks collaboratively and jointly.

## 3. EXPERIMENTS AND RESULTS

We evaluated the proposed method on 2018 ISIC skin lesion segmentation dataset [1], which includes 2594 annotated dermoscopic images. Fourfold cross validation was adopted for the evaluation. The performance of segmentation was evaluated by accuracy (AC), dice coefficient (DI), Jaccard index (JA) and sensitivity (SE). We implemented our model with TensorFlow, and NVIDIA TITAN XP GPU was used for training acceleration. Adam was chosen for optimization with $\beta_1 = 0.5$ and $\beta_2 = 0.999$. Each mini-batch includes 4 samples in training phases. The learning rate was initialized as 0.000001 and dropped by 0.1 every 40 epochs.

We first analyzed the influence of KAM, and showed the learned dilated rate maps at $32 \times 32$ scale as in Fig.3. From the heat map, it is clear that the receptive field is expanded at the hole region, and the hole is disappeared in the final predictions, indicating the proposed KAM is able to solve the hole problem effectively. The receptive field around the boundary is slightly enlarged, which makes the prediction expand to fit the ground truth. Therefore, context information of holes and boundaries is provided to help predict the class of local pixels, which alleviates hole as well as shrink problem and leads to better segmentation performance.

Then we assessed the qualitative performance of our complementary network by comparing it with state-of-the-art methods [5, 6, 8]. We implemented these methods on our dataset and visualized four examples in Fig. 4. It is clear that our complementary network can appropriately fill holes, while holes existed in the results of methods [5, 8]. Moreover, segmentation results obtained by the complementary network
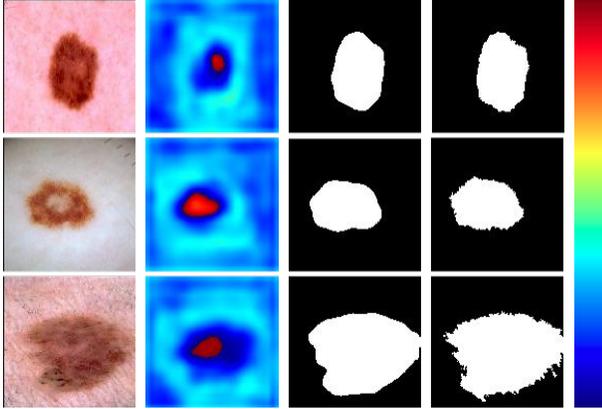
**Fig. 3**. Each row includes the original image, dilated rate map, predictions and ground truth from left to right. Note that red in heat map denotes a larger receptive field.
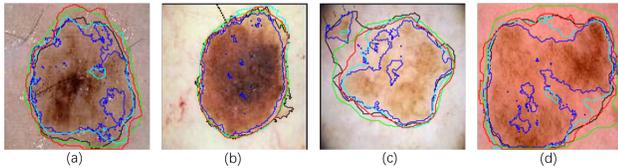


**Fig. 4**. Examples of complementary network results in comparison with other methods. The ground truth is denoted in black. Results of [5], [6], [8] and ours are denoted in blue, cyan, green, and red, respectively.

are better than that of other methods with smallest distance to the ground truth.

We further analyzed the performance of our complementary network by comparing it with methods [5, 6, 8] under the setting of fully supervised and semi-supervised learning, respectively. Specifically, $3^{rd}$ to $6^{th}$ rows show fully supervised learning results with 1945 labeled images, while $8^{th}$ to $11^{th}$ rows report results with 649 labeled and 1296 unlabeled images for our method and results with 649 labeled images for other methods in Table 1. Under fully supervised setting, the proposed method shows superior performance with an improvement of 2.7%, 0.2%, 1.6% in AC, 7.0%, 0.5%, 3.5% in DI and 9.3%, 1.0%, 4.7% in JA compared with the existing deep learning based methods [5, 6, 8], respectively. This result validated the proposed complementary network possesses superior ability to alleviate hole and shrink problems for skin lesion segmentation. Trained with the same labeled data, our semi-supervised method exhibited a significant increment in evaluation criteria compared with [5, 6, 8]. The increment is more distinct than that of fully supervised methods, indicating our method can leverage unlabeled images effectively.

## 4. CONCLUSION

In this paper, we propose a novel complementary network with adaptive atrous convolution (AAC), knowledge aggrega-

**Table 1**. Comparison results for menaloma segmentation.

| **Fully supervised** | | | | |
|---|---|---|---|---|
| Methods | AC (%) | DI (%) | JA (%) | SE (%) |
| Unet [5] | 92.3±0.3 | 79.4±0.4 | 68.3±0.6 | 83.6±0.6 |
| Sarker et al.[6] | 94.8±0.1 | 85.9±0.3 | 76.6±0.6 | **87.9±0.5** |
| Yuan et al.[8] | 93.4±0.2 | 82.9±0.8 | 72.9±0.9 | 85.6±0.6 |
| Our Method | **95.0±0.6** | **86.4±1.3** | **77.6±1.9** | 86.9±1.0 |
| **Semi-supervised** | | | | |
| Unet [5] | 91.0±0.2 | 75.9±0.3 | 63.8±0.3 | 82.3±0.2 |
| Sarker et al.[6] | 93.5±0.1 | 83.0±0.4 | 74.2±0.5 | **87.7±0.5** |
| Yuan et al.[8] | 91.4±0.3 | 77.8±0.9 | 66.4±1.1 | 84.7±0.7 |
| Our Method | **94.4±0.2** | **85.0±0.6** | **75.9±0.9** | 85.0±1.7 |

tion module (KAM) and mutual loss, for melanoma segmentation. Our network is able to alleviate the hole and shrink problems existing in current methods. The proposed complementary network and mutual loss can be further extended to semi-supervised learning, which is significant for medical image segmentation due to the limited annotated data. AAC and KAM can be flexibly transferred to other image segmentation tasks to adaptively enlarge receptive fields and boost the segmentation performance.

## 5. REFERENCES

[1] Codella, N.C.F., Gutman, D., Celebi, M.E., Helba, B., Marchetti, M.A., Dusza, S.W., et al.: Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In: ISBI 2018. pp. 168–172 (2018)

[2] Jaderberg, M., Simonyan, K., Zisserman, A., Kavukcuoglu, K.: Spatial transformer networks. In: NeurIPS. pp. 2017–2025 (2015)

[3] Li, H., He, X., Zhou, F., Yu, Z., Ni, D., Chen, S., Wang, T., Lei, B.: Dense deconvolutional network for skin lesion segmentation. IEEE J. Biomed. Health Inform **23**(2), 527–537 (2019)

[4] Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV. pp. 2980–2988 (2017)

[5] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241. Springer (2015)

[6] Sarker, M.M.K., Rashwan, H.A., Akram, F., Banu, S.F., Saleh, A., Singh, V.K., et al.: Slsdeep: Skin lesion segmentation based on dilated residual and pyramid pooling networks. In: MICCAI. pp. 21–29. Springer (2018)

[7] Siegel, R.L., Miller, K.D., Jemal, A.: Cancer statistics, 2019. CA Cancer J Clin **69**(1), 7–34 (2019)

[8] Yuan, Y., Lo, Y.C.: Improving dermoscopic image segmentation with enhanced convolutional-deconvolutional networks. IEEE J. Biomed. Health Inform **23**(2), 519–526 (2017)