

DEEP CONVOLUTIONAL ENCODER-DECODERS WITH AGGREGATED MULTI-RESOLUTION SKIP CONNECTIONS FOR SKIN LESION SEGMENTATION

Ahmed H. Shahin^{*} Karim Amer[†] Mustafa A. Elattar^{*}

^{*} Medical Imaging and Image Processing Group, Center for Informatics Sciences, Nile University, Egypt

[†] Ubiquitous and Visual Computing Group, Center for Informatics Sciences, Nile University, Egypt

ABSTRACT

The prevalence of skin melanoma is rapidly increasing as well as the recorded death cases of its patients. Automatic image segmentation tools play an important role in providing standardized computer-assisted analysis for skin melanoma patients. Current state-of-the-art segmentation methods are based on fully convolutional neural networks, which utilize an encoder-decoder approach. However, these methods produce coarse segmentation masks due to the loss of location information during the encoding layers. Inspired by Pyramid Scene Parsing Network (PSP-Net), we propose an encoder-decoder model that utilizes pyramid pooling modules in the deep skip connections which aggregate the global context and compensate for the lost spatial information. We trained and validated our approach using ISIC 2018: Skin Lesion Analysis Towards Melanoma Detection grand challenge dataset. Our approach showed a validation accuracy with a Jaccard index of 0.837, which outperforms U-Net. We believe that with this reported reliable accuracy, this method can be introduced for clinical practice.

Index Terms— skin lesion, segmentation, melanoma, convolutional neural networks, pyramid pooling modules

1. INTRODUCTION

Skin cancer is by far the most common of all cancers. Skin melanoma incidents represent about 1% of skin cancers but causes a large majority of skin cancer deaths. In 2018, it is estimated that more than 90,000 new incidents of melanoma will be diagnosed in the United States alone, resulting in 10,000 deaths [1]. Consequently, the early diagnosis of melanoma greatly increases the prevalence of recovery. The five-year survival rate for early stage melanoma exceeds 95% [2]. These statistics promote the importance of early diagnosis and identification of skin melanoma lesions. One of the essential steps in computerized analysis of dermoscopic images is automatic skin lesion segmentation.

A number of classical image processing techniques like thresholding, edge based, or region-based methods have been

traditionally used for the sake of skin lesion segmentation [3]. However, the existence of different sources of artifacts in images was a serious limitation to the classical techniques. For feature-based techniques, low-level and handcrafted features were not able to accurately segment skin lesions [4].

In the recent years, alongside with the advancements in the computational power of the graphical processing units (GPUs), Convolutional Neural Networks (CNNs) [5] have emerged as one of the most powerful tools in image processing. CNN models have shown a promising performance in multiple domains including medical image analysis [6]. Long et al. [7] introduced a fully convolutional network (FCN) for segmentation task where the input image is encoded using sequence of convolutions and pooling operators producing deep feature maps, and then decoded by a single deconvolutional upsampling layer to produce output feature maps with the original size. Despite the significant contribution over the classical machine learning and image processing methods, the FCN design had a negative effect on the output probability maps as the location information is distorted due to the sequential pooling operations producing coarse segmentation masks.

To alleviate the deficiencies of FCN, Badrinarayanan et al. [8] designed a trainable decoder with multiple deconvolutional layers operating on gradually upsampled feature maps. Ronneberger et al. [6] introduced skip connections to preserve important location information by concatenating the features from the contracting (encoding) path with the corresponding features in the expanding (decoding) path. PSP-Net [9] proposed a pyramid pooling module to aggregate the global context by using parallel pooling layers with differently sized kernels. Consequently, that provided additional contextual information preserved along the sequential convolution and pooling layers.

In skin lesion segmentation context, Yuan et al. [10] utilized the conventional encoder-decoder architecture with a novel Jaccard-index-based loss function to handle the class imbalance in the dermoscopic images. Their FCN model has shown an improvement in the segmentation accuracy on ISIC 2016 dataset for skin lesion analysis [11]. However, their study suffered from several limitations such as failing to achieve reasonable accuracy on some images that have

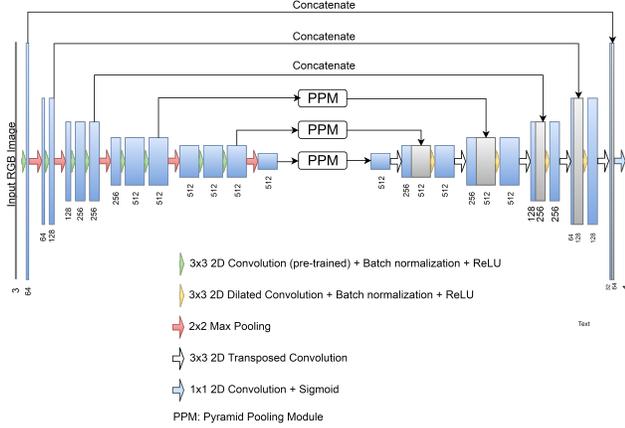


Fig. 1: An illustration of the proposed architecture. Each blue box represents a multi-channel feature map. Number of channels is shown below each box. White boxes represent copied feature maps. The arrows denote the different operations. PPM: Pyramid Pooling Module.

low contrast between the lesion and the skin. Xue et al. [12] explored the use of adversarial learning in the lesion segmentation task.

Inspired by PSP-Net, we propose an adjusted encoder-decoder architecture to overcome the coarse nature of the output segmentation map. Our main idea is to embed the pyramid pooling modules (PPMs) to the skip connections in the deeper convolutional layers. We show that such architecture aggregates the location information from the low convolutional layers, to alleviate the problem of spatial information loss. Our experiments on "ISIC 2018 Skin Lesion Analysis Towards Melanoma Detection" dataset [11] corroborate our hypothesis that this approach refines the output predicted masks when compared to the current state-of-the-art lesion segmentation methods.

2. METHODS

Network Architecture Our proposed network architecture consists of an encoding path and a decoding path as shown in Fig. 1. For the encoding path, we use VGG11 [13] network without the fully connected layers, with an additional batch normalization layer after every convolution. In the decoding path, the features map is upsampled using transposed convolutional layers which double the spatial resolution and reduces the number of channels by half. Additionally, each upsampled map is concatenated with the corresponding feature map from the encoding path through skip connections, then followed by a 3x3 dilated convolution, batch normalization, and ReLU. In all convolutional layers, we use padded convolutions in order to prevent the loss of borders.

As shown in Fig. 1, the two upper skip connections just concatenate the features from the encoding path with those in the decoding path to preserve the spatial information from

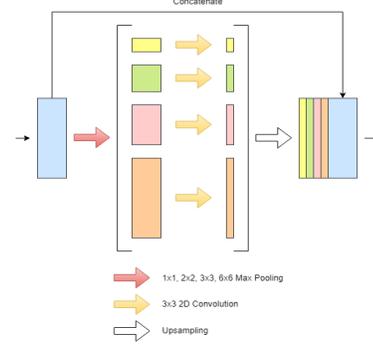


Fig. 2: An illustration of the Pyramid Pooling Module (PPM).

the encoding path. However, the two lower skip connections project the encoded activation maps to a Pyramid Pooling Module (PPM) [9] before concatenation. PPM covers different parts of the feature map through different pooling operations (see Fig. 2). An extra PPM block is used as a bottle neck between the encoding and the decoding path.

There is an output layer after the encoded path which performs a pixel-wise classification. The output layer is a 1x1 convolutional layer with a sigmoid activation function. The entire image is projected to a map of the same size, where each element represents the probability of its correspondence to the foreground (the lesion).

The output probability map is subjected to three successive post processing steps starting with thresholding the probability map at 0.5, then selecting the largest connected component and finally filling the holes of this component producing the final mask.

Training Our model was implemented using PyTorch deep learning framework. Adam optimizer was used for neural network optimization. Dynamic learning rate was used and initialized by 5×10^{-5} then multiplied by 10-1 every 30 training epochs. The number of epochs was initially set to 200 epochs. However, in our experiments the network converged to the best validation accuracy in less than 100 epochs thanks to the pretrained weights. Additionally, the training batch size was set to 16 images. We augmented the training images using series of geometric transformations (horizontal flipping, vertical flipping, rotation, and zooming) to reduce model overfitting. We used Generalized Dice Loss (GDL) [14] as a loss function, which is a modified formula of dice score coefficient (DSC). Unlike DSC, GDL is differentiable and can be used as a loss function in case of imbalanced dataset, as an alternative for the widely used Cross-entropy loss. GDL takes the form:

$$GDL = 1 - \frac{2 \sum_n r_n p_n}{\sum_n r_n + \sum_n p_n} \quad (1)$$

Where: r is the reference ground truth segmentation with pixel values r_n , and p is the predicted probabilities from the CNN with pixel values p_n .

3. EXPERIMENTS AND RESULTS

Our data was extracted from the ISIC 2018: Skin Lesion Analysis Towards Melanoma Detection grand challenge datasets [11]. The training set consists of 2594 RGB dermoscopic images with spatial resolutions ranging from 576×768 to 6748×4499 . ISIC 2018 Validation and test data ground truth have not yet been released by the time of this paper submission. Thus, we divided the training data into 80% training (2076 images) and 20% validation set (518 images). The validation set was used to evaluate the performance of our approach. In order to guarantee the robustness of the model against different parts of the available data, five-fold cross-validation was used. We resized all the images to 192×256 . However, the reported results were calculated after resampling the output masks to the original sizes. Fig. 3 shows the cross-validation results of our model. We compare our method with U-Net as a baseline against one random validation fold (518 images) of the dataset. Tables 1 & 2 summarize the results of our method compared to U-Net, other approaches on ISIC 2017 dataset (2000 training images and 600 test images) and the results on the official test set (1000 images). We do not include results on ISIC 2016 dataset (900 training images and 375 test images) due to its relatively smaller number of test images. Using a Geforce 1080Ti Nvidia GPU, our method is able to segment around 10 images per second.

4. DISCUSSION

Automatic lesion segmentation is still a challenging task due to the lack of distinctive lesion boundaries to adjacent skin, as well as the existence of various artifacts. In this work, we introduced an encoder-decoder like architecture equipped with PPMs for automatic skin lesion segmentation.

Our model has shown superior, and stable results on the five-fold cross-validation experiments. Compared to the previous approaches on smaller datasets, our method introduced

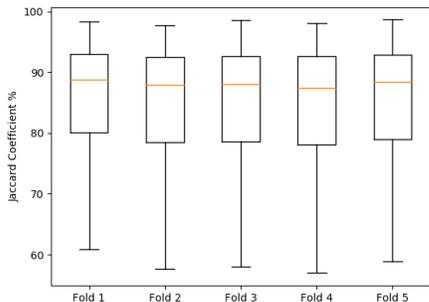


Fig. 3: Five-fold cross-validation results of the proposed model on ISIC 2018 dataset.

Table 1: Results of our method against other methods evaluated on the local validation set. All values are in percentages. JA: Jaccard; DC: dice; SN: sensitivity; and SP: specificity.

Method	Data	JA	DC	SN	SP
Xue et al. [12]	ISIC 2017	78.5	86.7	-	-
Yuan et al. [10]		76.5	84.9	82.5	97.5
U-Net [6]	ISIC 2018	82.6	89.5	91.1	96.8
Proposed Method		83.7	90.3	90.2	97.4

Table 2: Results of our method against other methods evaluated on the official ISIC 2018 test set. The official evaluation metric was Jaccard thresholded at 65%. All values are in percentages.

Method	Thresholded Jaccard
Bissoto et al. [15]	72.8
Hardie et al. [4]	66.3
Proposed Method	73.8

a significant improvement in the segmentation accuracy. ISIC 2018 dataset is the latest publicly available dataset, that extends the previous ISIC 2017 and ISIC 2016 datasets with significantly more training and testing examples. Our model outperformed U-Net on our local validation set, we believe that this improvement is due to the PPMs role in preserving the spatial information along the encoding-decoding process by extracting additional contextual information in the deeper skip connections. Additionally, our method had a higher accuracy when compared to other published methods on the official test set of ISIC 2018 skin lesion analysis towards melanoma detection challenge. Bissoto et al. adopted the conventional U-Net architecture for their submission, while Hardie et al. used Bayesian classifier with handcrafted features and SVM regression for segmentation threshold selection. The difference between our scores on the local validation set and on the official test set is due to the 65% thresholded Jaccard metric used in the test set evaluation. Our model does not rely on heavy augmentations or ensembling that would increase the prediction time. Fig. 4(a,b,c and d) shows success samples of the proposed method. There are entries with higher scores on the official leaderboard of ISIC challenge. However, there are no official presentations for their used methods.

Despite the superior performance in most of the cases, the model needs improvement to handle some failure cases. In Fig. 4(e,f), we can notice the difficulty of differentiating between the lesion and skin areas even for experienced raters, due to the very low contrast between the two classes. One way to further improve the segmentation performance is to use other post-processing techniques such as conditional random field and to combine its parameters in the network training process.

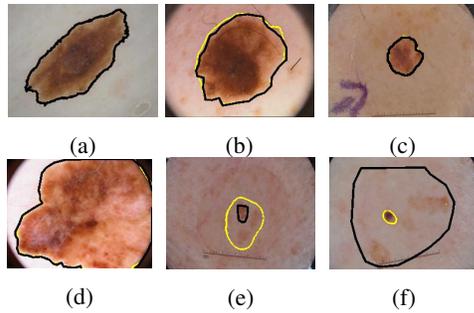


Fig. 4: a, b, c and d: Success samples of our method. e and f: Failure samples of our method. Ground truth label is shown in black contour and the output predicted mask is in yellow.

5. CONCLUSIONS

In this work, we presented a fully automatic algorithm based on CNNs for skin lesion segmentation from dermoscopic images. We proposed the combination of encoder-decoder architectures with the pyramid pooling modules. Our method does not require any preprocessing or gray-scale conversion. According to the reported results, the method has shown its robustness against other methods and various image artifacts. We believe that this model can generalize well by showing reliable performance on other medical segmentation problems.

Acknowledgement In memory of Mohamed Samy Ahmed Kassem (September 1993 – May 2018).

6. REFERENCES

- [1] Rebecca L. Siegel, Kimberly D. Miller, and Ahmedin Jemal, “Cancer statistics, 2018,” *CA: A Cancer Journal for Clinicians*, vol. 68, no. 1, pp. 7–30, 1 2018.
- [2] “Survival rates for melanoma skin cancer, by stage,” <https://www.cancer.org/content/cancer/en/cancer/melanoma-skin-cancer/detection-diagnosis-staging/survival-rates-for-melanoma-skin-cancer-by-stage/>, October 10, 2018.
- [3] Margarida Silveira, Jacinto C. Nascimento, Jorge S. Marques, Andr R. S. Marcal, Teresa Mendonca, Syogo Yamauchi, Junji Maeda, and Jorge Rozeira, “Comparison of Segmentation Methods for Melanoma Diagnosis in Dermoscopy Images,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 3, no. 1, pp. 35–45, 2 2009.
- [4] Russell C Hardie, Redha Ali, Manawaduge Supun De Silva, and Temesguen Messay Kebede, “Skin Lesion Segmentation and Classification for ISIC 2018 Using Traditional Classifiers with Hand-Crafted Features,” *arXiv preprint arXiv: 1807.07001*, 7 2018.
- [5] Y. Lecun, Lon Bottou, Yoshua Bengio, and Patrick Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [6] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-Net: Convolutional Networks for Biomedical Image Segmentation,” in *MICCAI*, 5 2015, pp. 234–241.
- [7] Jonathan Long, Evan Shelhamer, and Trevor Darrell, “Fully convolutional networks for semantic segmentation,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 6 2015, vol. 39, pp. 3431–3440, IEEE.
- [8] Vijay Badrinarayanan, Alex Kendall, and Roberto 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, 12 2017.
- [9] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia, “Pyramid scene parsing network,” *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, vol. 2017-Janua, pp. 6230–6239, 2017.
- [10] Yading Yuan and Yeh-Chi Lo, “Improving Dermoscopic Image Segmentation with Enhanced Convolutional-Deconvolutional Networks,” *IEEE Journal of Biomedical and Health Informatics*, pp. 1–1, 2017.
- [11] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, and Allan Halpern, “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 *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*. 4 2018, pp. 168–172, IEEE.
- [12] Yuan Xue, Tao Xu, and Xiaolei Huang, “Adversarial learning with multi-scale loss for skin lesion segmentation,” in *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*. 4 2018, number Isbi, pp. 859–863, IEEE.
- [13] Karen Simonyan and Andrew Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” *arXiv preprint arXiv: 1409.1556*, 9 2014.
- [14] Carole H. Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M. Jorge Cardoso, “Generalised Dice Overlap as a Deep Learning Loss Function for Highly Unbalanced Segmentations,” in *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, vol. 10553 LNCS, pp. 240–248. 2017.
- [15] Alceu Bissoto, Fbio Perez, Vincius Ribeiro, Michel Fornaciali, Sandra Avila, and Eduardo Valle, “Deep-Learning Ensembles for Skin-Lesion Segmentation, Analysis, Classification: RECOD Titans at ISIC Challenge 2018,” *arXiv preprint arXiv: 1808.08480*, 8 2018.