

SLSDeep: Skin Lesion Segmentation Based on Dilated Residual and Pyramid Pooling Networks

Md. Mostafa Kamal Sarker^{1,*}, Hatem A. Rashwan¹, Farhan Akram², Syeda Furraka Banu³, Adel Saleh¹, Vivek Kumar Singh¹, Forhad U H Chowdhury⁴, Saddam Abdulwahab¹, Santiago Romani¹, Petia Radeva⁵, and Domenec Puig¹

¹ Department of Computer Engineering and Mathematics, Universitat Rovira i Virgili, 43007 Tarragona, Spain.

² Imaging Informatics Division, Bioinformatics Institute, 30 Biopolis Street, # 07-01 Matrix, 138671, Singapore.

³ Department of Technology and Engineering Management, Rovira i Virgili University, 43007 Tarragona, Spain.

⁴ Liverpool School of Tropical Medicine, Liverpool L3 5QA, UK.

⁵ Department of Mathematics, University of Barcelona, 08007 Barcelona, Spain.

Abstract. Skin lesion segmentation (SLS) in dermoscopic images is a crucial task for automated diagnosis of melanoma. In this paper, we present a robust deep learning SLS model, so-called SLSDeep, which is represented as an encoder-decoder network. The encoder network is constructed by dilated residual layers, in turn, a pyramid pooling network followed by three convolution layers is used for the decoder. Unlike the traditional methods employing a cross-entropy loss, we investigated a loss function by combining both Negative Log Likelihood (NLL) and End Point Error (EPE) to accurately segment the melanoma regions with sharp boundaries. The robustness of the proposed model was evaluated on two public databases: ISBI 2016 and 2017 for skin lesion analysis towards melanoma detection challenge. The proposed model outperforms the state-of-the-art methods in terms of segmentation accuracy. Moreover, it is capable to segment more than 100 images of size 384×384 per second on a recent GPU.

Keywords: skin lesion segmentation melanoma, deep learning, dilated residual networks, pyramid pooling

1 Introduction

According to the skin Cancer Foundation Statistics, the percentage of both melanoma and non-melanoma skin cancers has been increasing rapidly over the last few years [19]. Dermoscopy, non-invasive dermatology imaging methods, can help the specialists to inspect the pigmented skin lesions and diagnose malignant melanoma at an initial-stage [12]. Even the professional dermatologist can not

* Corresponding Author: mdmostafakamal.sarker@urv.cat

properly classify the melanoma only by relying on their perception and vision. Sometimes, human tiredness and other distractions during visual diagnosis can also yield high number of false positives[4]. Therefore, a Computer-Aided Decision system (CAD) is needed to assist the dermatologists to properly analyze the dermoscopic images and accurately segment the melanomas. Many attempts of melanoma segmentation have been proposed in the literature. An overview of the different melanoma segmentation techniques is presented in [26]. However, this task is still a challenge, since the dermoscopic images has various characteristics including different sizes and shapes, fuzzy boundaries, different colors, and the presence of hair [8]. In last few decades, many approaches have been proposed

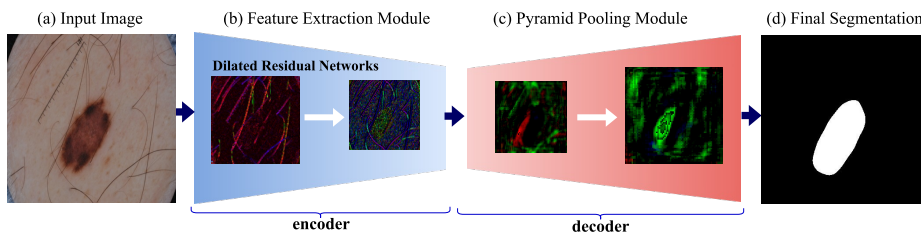


Fig. 1. Architecture of the proposed skin lesion segmentation network. Negative Log Likelihood (NLL) and End Point Error (EPE)

to cope with the aforementioned challenges. Most of these methods are based on thresholding, edge-based and/or region-based active contour models, clustering and supervised learning [5]. However, these methods are unreliable when dermoscopic images are inhomogeneous and/or lesions have fuzzy or blurred boundaries[5]. Furthermore, their performance relies on efficient pre-processing algorithms, such as filtering, illumination correction and hair removal, which badly affect the generalizability of these models.

Recently, deep learning methods applied to image analysis, specially Convolutional Neural Networks (CNNs) have been used to solve the image segmentation problem[15]. These CNN-based methods can automatically learn features from raw pixels to distinguish between background and foreground objects to attain the final segmentation. Most of these approaches generally are based on encoder-decoder networks [15]. These networks learn to map the features of an image to a segmented image. The encoder networks are used for extracting the features from the input images, in turn the decoder ones used to construct the segmented image. The U-net network proposed in [18] has been particularly designed for biomedical image segmentation based on the concept of Fully Convolutional Networks(FCN) [15]. The U-net model reuses the feature maps of the encoder layers to the corresponding decoders and concatenates them to upsampled (via deconvolution) decoder feature maps called “skip-connections”. The U-Net model for SLS outperformed many classical clustering techniques [14].

In addition, the deep residual network (ResNet) model [24] is a 50-layers network designed for segmentation tasks. ResNet blocks are used to boost the overall depth of the networks and allow more accurate segmentation depending on more significant image features. Moreover, Dilated Residual Networks (DRNs) proposed in [23] increase the resolution of the ResNet blocks output by replacing a subset of interior subsampling layers by dilation [22]. DRNs outperform the normal ResNet without adding algorithmic complexity to the model. DRNs are able to represent both tiny and large image features. Furthermore, Zhao et. al. [27] proposed a Pyramid Pooling Network (PPN) that is able to extract additional contextual information based on a multi-scale scheme.

Inspired by the success of the aforementioned deep models for semantic segmentation, we propose a model combining skip-connections, dilated residual and pyramid pooling networks for SLS with different improvements. In our model, the encoder network depends on DRNs layers, in turn the decoder depends on a PPN layer along with their corresponding connecting layers. More features can be extracted from the input dermoscopic images by combining DRNs with PPN, in turn it also enhances the performance of the final network. Finally, our SLS segmentation model uses a new loss function, which combines Negative Log Likelihood (NLL) and End Point Error (EPE) [1]. Mainly, cross-entropy is used for multi-class segmentation models, however it is not as useful as NLL in binary class segmentation. Thus, in such melanoma segmentation, we propose to use NLL as a loss function. In addition, for preserving the melanoma boundaries, EPE is used as a content loss function. Consequently, this paper aims at developing an automated deep SLS model with two main contributions:

- An encoder-decoder network for efficient SLS without any pre- and post-processing algorithms based on dilated residual and pyramid pooling networks to enclose coarse-to-fine features of dermoscopic images.
- A new loss function that is a combination of Negative Log Likelihood and End Point Error for properly detecting the melanoma with sharp edges.

2 Proposed Model

2.1 Network Architecture

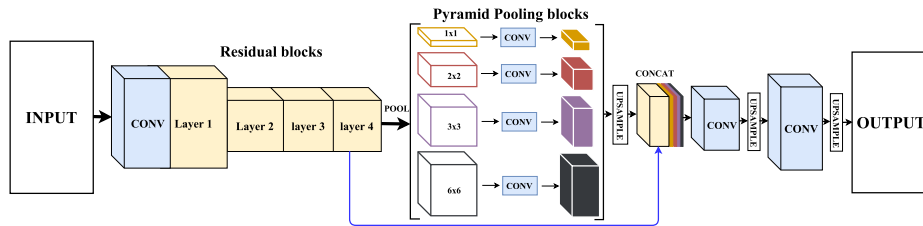


Fig. 2. Architecture of the Encoder-decoder network.

Fig.1 shows the architecture of the proposed SLSDeep model with DRNs [28] and PPN [10]. The network contains two-fold architecture: encoder and decoder. Regarding the encoder phase, the first layer is a 3×3 convolutional layer followed by 3×3 max pooling with stride 2.0 that generates 64 feature maps. This layer uses Relu as an activation and batch normalization to speed-up the training steps with a random initialization. Following, four pre-trained DRNs blocks are then used to extract 256, 512, 1024 and 2048 feature maps, respectively as shown in Fig.2. The first, third, and fourth DRNs layers are with stride 1.0, in turn the second one is with stride 2.0. Thus, the size of final output of encoder is 1/8 of the input image (e.g. in our model, the input image is in 384×384 and the output feature maps of the encoder is 48×48). For global contextual prior, average pooling is used before feeding to fully connected layers in image classification [21]. However, it is not sufficient to extract necessary information from our skin lesion images. Therefore, we do not use average pooling at the end of the encoder and directly fed the output feature maps to the decoder network.

On the other hand, for the decoder network, we use the concept of PPN for producing multi-scale (coarse-to-fine) feature maps and then all scales are concatenated together to get more robust feature maps. PPN use a hierarchical global prior of variant size feature maps in multi-scales with different spatial filters as shown in Fig.2. In this paper, the used PPN layer extracts feature maps using four pyramid scales with rescaling sizes of 1×1 , 2×2 , 3×3 and 6×6 . A convolutional layer with a 1×1 kernel in every pyramid level is used for generating 1024 feature maps. The low-dimension feature maps are then upsampled based on bilinear interpolation to get the same size of the input feature maps. The input and four feature maps are finally concatenated to produce 6144 feature maps (i.e., 4×1024 feature maps concatenated with the input 2048 feature maps). Sequentially, two 3×3 convolutional layers are followed by two upsampling layers. Finally, a softmax function (i.e. normalized exponential function) is utilized as logistic function for producing the final segmentation map. A ReLU activation with batch normalization is used in the two convolutional layers [11]. Moreover, in order to avoid the overfitting problem, the dropout function with a ratio of 0.5 [20] is used before the second upsampling layer.

The skip connections between all layers of the encoder and decoder were tested during the experiments. However, the best results were provided when only one skip connection was done between the last layer of the encoder and the output of PPN layer of the decoder. The architecture of the encoder and decoder is given in details with the supplementary materials.

2.2 Loss Function

Most of the traditional deep learning methods commonly employ cross-entropy as a loss function for segmentation [18]. Since the melanoma is mostly a small part of a dermoscopic image, the minimization of cross-entropy tends to be biased towards the background. To cope with this challenge, we propose a new loss function by combining objective and content losses: NLL and EPE, respectively. In order to fit a log linear probability model to a set of binary labeled classes, the

NLL is the objective function of the proposed model to minimize. Let $v \in \{0, 1\}$ be a true label for binary classification and $p = Pr(v = 1)$ a probability estimate, the NLL of the binary classifier can be defined as:

$$L_{log}(v, p) = -\log Pr(v|p) = -(v \log(p) + (1 - v) \log(1 - p)). \quad (1)$$

Regarding the content of the loss function, we have also computed an absolute error aiming at maximizing the Peak Signal-to-Noise Ratio (PSNR) by preserving the object boundaries. The used EPE loss follows a classical approach that the generated mask, tu is pixel-wise compared with the corresponding ground-truth, v . The EPE error can be defined by [1]:

$$L_{epe} = \sqrt{(u_0 - u_1)^2 + (v_0 - v_1)^2} \quad (2)$$

where u_0 and u_1 are the first derivatives of u in x and y directions, and v_0 and v_1 are the first derivatives of v .

Hence, our final loss, which combines both NLL and EPE, can be defined as:

$$L_{total} = L_{log} + \alpha L_{epe} \quad (3)$$

where $\alpha < 1$ is a weighted coefficient. In this work, we use $\alpha = 0.5$.

3 Experimental Setup and Evaluation

Database: To test the robustness of the proposed model, it was evaluated on two public benchmark datasets of dermoscopy images for skin lesion analysis: **ISBI 2016** [7] and **ISBI 2017** [9]. The datasets images are captured by different devices at various top clinical centers over the world. In ISBI 2016 dataset, training and testing part contain 900 and 379 annotated images, respectively. The size of the images ranges from 542×718 to 2848×4288 pixels. In turn, ISBI 2017 dataset is divided into training, validation and testing parts with 2000, 150 and 600 images, respectively.

Evaluation Metrics: We used the evaluation metrics of ISBI 2016 and 2017 challenges for evaluating the segmentation performances including Specificity(SPE), Sensitivity(SEN), Jaccard index(JAC), Dice coefficient(DIC) and Accuracy(ACC) detailed in [9] and [7].

Implementation: The proposed model is implemented on an open source deep learning library named PyTorch[16]. For optimization algorithm, we used Adam [13] for adjusting learning rate, which depends on first and second order moments of the gradient. We used a ‘‘poly’’ learning rate policy [6] and selected a base learning rate of 0.001 and 0.01 for encoder and decoder, respectively with a power of 0.9. For data augmentation, we selected random scale between 0.5 and 1.5, random rotation between -10 and 10 degrees. The ‘‘batchsize’’ is set to 16 for training and the epochs to 100. The experiments utilized NVIDIA TITAN X with 12GB memory and its takes approximately 20 hours for train the networks.

Evaluation and results: Since the size of the given images is very large, we resized the input images into 384×384 pixels for training our model. In this

work, we tested different sizes and the 384×384 size yields the best results. In order to separately assess the different contributions of this model, the resulting segmentation for the proposed model with different variations have been computed: (a) The SLSDeep model without the content loss EPE (SLSDeep-EPE), (b) the proposed method with skip connections of all encoder and decoder layers (SLSDeep+ASC) and (c) the final proposed model (SLSDeep) with NLL and EPE loss functions and only one skip connection between the last layer of the encoder and the PPN layer.

Quantitative results on ISBI’2016 and ISBI’2017 datasets are shown in Table 1. Regarding ISBI’2016, we compared the SLSDeep and its variations to the four top methods: ExB, [24], [17] and [25] providing the best results according to [9]. The segmentation results of our model SLSDeep with its variations (SLSDeep-EPE and SLSDeep+ASC) perform much better than the all evaluated methods on the ISBI’2016 with the five aforementioned evaluation metrics. SLSDeep yields the best results among the three variations. In addition, for the DIC score, our model, SLSDeep, improved the results with around 3.5%, while the JAC score was significantly improved with 8%. The SLSDeep yielded results with an overall accuracy of more than 98%.

Table 1. Performance Evaluation on the ISBI Challenges Dataset

Challenges	Methods	ACC	DIC	JAC	SEN	SPE
ISBI 2016	ExB	0.953	0.910	0.843	0.910	0.965
	CUMED[24]	0.949	0.897	0.829	0.911	0.957
	Rahman et. al.[17]	0.952	0.895	0.822	0.880	0.969
	Yuan et. al.[25]	0.955	0.912	0.847	0.918	0.966
	SLSDeep	0.984	0.955	0.913	0.945	0.992
	SLSDeep-EPE	0.973	0.919	0.850	0.890	0.990
	SLSDeep+ASC	0.975	0.930	0.869	0.952	0.979
ISBI 2017	Yuan et. al.[25]	0.934	0.849	0.765	0.825	0.975
	Berseth et. al.[2]	0.932	0.847	0.762	0.820	0.978
	MResNet-Seg[3]	0.934	0.844	0.760	0.802	0.985
	SLSDeep	0.936	0.878	0.782	0.816	0.983
	SLSDeep-EPE	0.913	0.826	0.704	0.729	0.986
	SLSDeep+ASC	0.906	0.850	0.739	0.808	0.905

Furthermore, SLSDeep on the ISBI’2017 provided segmentation results with improvements of 3% and 2% in terms of DIC and JAC scores, respectively. Again, SLSDeep perform better the three top methods of the ISBI’2017 benchmark, [25], [2] and [3], with ACC, DIC and JAC scores. However, [25] yielded the best SEN score with an improvement of 0.9% better than our model. The SLSDeep-EPE and SLSDeep+ASC provided reasonable results, however their results were worse than the three tested methods. SLSDeep-EPE yields the highest SPE,

which is 0.1% and 0.3% more than MResNet-Seg [3] and SLSDeep, respectively. Using the EPE function with the final SLSDeep model significantly improved the DIC and JAC scores of 3% and 5%, respectively, on ISBI'2016 and of 5% and 8%, respectively, with ISBI'2017. In addition, SLSDeep with only one skip connections yields better results than SLSDeep+ASC on both ISBI datasets.

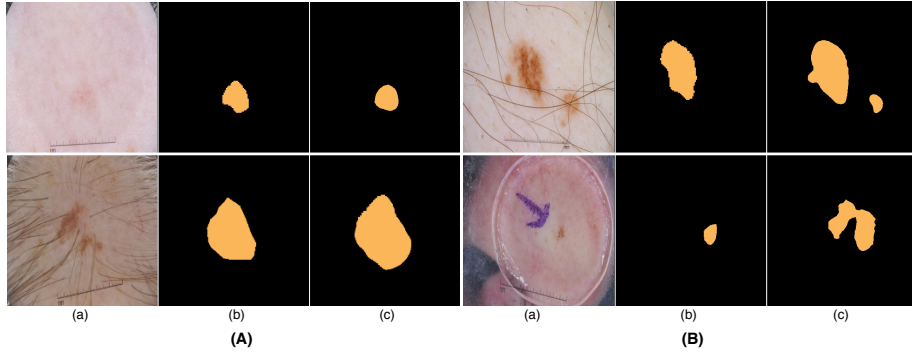


Fig. 3. Segmentation results: (a) *input image*, (b) *ground truth* and (c) *segmentation result*; (A) correct segmentation by our model; (B) incorrect segmentation by our model.

Qualitative results of four examples of the ISBI'2017 dataset are shown in Fig.3. For the first and second examples (on the top-and down-left side), the lesions were properly detected, although the color of the lesion area is very similar to the rest of the skin. In addition, the lesion area was accurately segmented with sharp edges. Regarding to the third example (on the top-right side), SLSDeep properly segmented the lesion area; however a small false region with similar features was also detected. In turn, the last example is very difficult, since the lesion shown in the input image is a very small region. However, the SLSDeep model can segment it, but with bigger size of false negative regions.

4 Conclusions

This paper proposed a novel deep learning skin lesion segmentation model based on training an encoder-decoder network. The encoder network used the dilated ResNet layers with downsampling to extract the features of the input image, in turn convolutional layers with pyramid pooling and upsampling are used to reconstruct the segmented image. This approach outperforms, in terms of skin lesion segmentation, the literature evaluated on two ISBI'2016 and ISBI'2017 datasets. The experiments show that SLSDeep is robust segmentation technique using different evaluation metrics: accuracy, Dice coefficient, Jaccard index and specificity. In addition, the qualitative results show promising skin lesion seg-

mentation. For future work, the proposed model will be explored on different color spaces and applied to other medical applications to prove its versatility.

References

1. Baker, S., Scharstein, D., Lewis, J., Roth, S., Black, M.J., Szeliski, R.: A database and evaluation methodology for optical flow. *IJCV* 92(1), 1–31 (2011)
2. Berseth, M.: Isic 2017-skin lesion analysis towards melanoma detection. arXiv preprint arXiv:1703.00523 (2017)
3. Bi, L., Kim, J., Ahn, E., Feng, D.: Automatic skin lesion analysis using large-scale dermoscopy images and deep residual networks. preprint arXiv:1703.04197 (2017)
4. Celebi, M.E., Iyatomi, H., Stoecker, W.V., Moss, R.H., Rabinovitz, H.S., Argenziano, G., Soyer, H.P.: Automatic detection of blue-white veil and related structures in dermoscopy images. *CMIG* 32(8), 670–677 (2008)
5. Celebi, M.E., Wen, Q., Iyatomi, H., Shimizu, K., Zhou, H., Schaefer, G.: A state-of-the-art survey on lesion border detection in dermoscopy images. *Dermoscopy Image Analysis* pp. 97–129 (2015)
6. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. arXiv preprint arXiv:1606.00915 (2016)
7. Codella, N.C., Gutman, D., Celebi, M.E., Helba, B., Marchetti, M.A., Dusza, S.W., Kalloo, A., Liopyris, K., Mishra, N., Kittler, H., 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). arXiv preprint arXiv:1710.05006 (2017)
8. Day, G.R., Barbour, R.H.: Automated melanoma diagnosis: where are we at? *Skin Research and Technology* 6(1), 1–5 (2000)
9. Gutman, D., Codella, N.C., Celebi, E., Helba, B., Marchetti, M., Mishra, N., Halpern, A.: Skin lesion analysis toward melanoma detection: A challenge at the international symposium on biomedical imaging (isbi) 2016, hosted by the international skin imaging collaboration (isic). arXiv preprint arXiv:1605.01397 (2016)
10. He, K., Zhang, X., Ren, S., Sun, J.: Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Trans. PAMI* 37(9), 1904–1916 (2015)
11. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: *ICML*. pp. 448–456 (2015)
12. Kardynal, A., Olszewska, M.: Modern non-invasive diagnostic techniques in the detection of early cutaneous melanoma. *Journal of dermatological case reports* 8(1), 1 (2014)
13. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
14. Lin, B.S., Michael, K., Kalra, S., Tizhoosh, H.: Skin lesion segmentation: U-nets versus clustering. arXiv preprint arXiv:1710.01248 (2017)
15. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of CVPR*. pp. 3431–3440 (2015)
16. Paszke, A., Gross, S., Chintala, S., Chanan, G.: Pytorch (2017)
17. Rahman, M., Alpaslan, N., Bhattacharya, P.: Developing a retrieval based diagnostic aid for automated melanoma recognition of dermoscopic images. In: *Applied Imagery Pattern Recognition Workshop (AIPR)*. pp. 1–7. IEEE (2016)

18. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241. Springer (2015)
19. Siegel, R.L., Miller, K.D., Jemal, A.: Cancer statistics, 2017. CA: A Cancer Journal for Clinicians 67(1), 7–30 (2017), <http://dx.doi.org/10.3322/caac.21387>
20. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research 15(1), 1929–1958 (2014)
21. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., et al.: Going deeper with convolutions. CVPR (2015)
22. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122 (2015)
23. Yu, F., Koltun, V., Funkhouser, T.: Dilated residual networks. In: Computer Vision and Pattern Recognition. vol. 1 (2017)
24. Yu, L., Chen, H., Dou, Q., Qin, J., Heng, P.A.: Automated melanoma recognition in dermoscopy images via very deep residual networks. TMI 36(4), 994–1004 (2017)
25. Yuan, Y., Chao, M., Lo, Y.C.: Automatic skin lesion segmentation with fully convolutional-deconvolutional networks. arXiv preprint arXiv:1703.05165 (2017)
26. Zhang, X.: Melanoma segmentation based on deep learning. Computer Assisted Surgery 22(sup1), 267–277 (2017)
27. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: IEEE Conf. CVPR. pp. 2881–2890 (2017)
28. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralla, A.: Scene parsing through ade20k dataset. In: Proceedings of the IEEE Conference CVPR (2017)