

A Multi-Branch Hybrid Transformer Network for Corneal Endothelial Cell Segmentation

Yinglin Zhang^{1,6}, Risa Higashita^{1,5*}, Huazhu Fu⁷, Yanwu Xu⁸, Yang Zhang¹,
Haofeng Liu¹, Jian Zhang⁶, and Jiang Liu^{1,2,3,4*}

¹ Department of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen 518055, China

² Cixi Institute of Biomedical Engineering, Chinese Academy of Sciences, China

³ Guangdong Provincial Key Laboratory of Brain-inspired Intelligent Computation, Department of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen 518055, China

⁴ Research Institute of Trustworthy Autonomous Systems, Southern University of Science and Technology, Shenzhen 518055, China

⁵ Tomey Corporation, Nagoya 451-0051, Japan

⁶ Global Big Data Technologies Centre, University of Technology Sydney, NSW, Australia

⁷ Inception Institute of Artificial Intelligence, UAE

⁸ Intelligent Healthcare Unit, Baidu, Beijing 100085, China

k-chen@tomey.co.jp

Abstract. Corneal endothelial cell segmentation plays a vital role in quantifying clinical indicators such as cell density, coefficient of variation, and hexagonality. However, the corneal endothelium’s uneven reflection and the subject’s tremor and movement cause blurred cell edges in the image, which is difficult to segment, and need more details and context information to release this problem. Due to the limited receptive field of local convolution and continuous downsampling, the existing deep learning segmentation methods cannot make full use of global context and miss many details. This paper proposes a Multi-Branch hybrid Transformer Network (MBT-Net) based on the transformer and body-edge branch. Firstly, We use the convolutional block to focus on local texture feature extraction and establish long-range dependencies over space, channel, and layer by the transformer and residual connection. Besides, We use the body-edge branch to promote local consistency and to provide edge position information. On the self-collected dataset TM-EM3000 and public Alisarine dataset, compared with other State-Of-The-Art (SOTA) methods, the proposed method achieves an improvement.

Keywords: Corneal endothelial cell segmentation · Deep learning · Transformer · Multi-branch.

1 Introduction

Corneal endothelial cell abnormalities may be related to many corneal and systemic diseases. Quantifying corneal endothelial cell density, the coefficient of

variation, and hexagonality have essential clinical significance [2]. Cell segmentation is a crucial step to quantify the above parameters. Nevertheless, manual segmentation is time-consuming, laborious, and unstable. Therefore, an accurate and fully automatic corneal endothelial cell segmentation method is essential to improve diagnosis efficiency and accuracy.

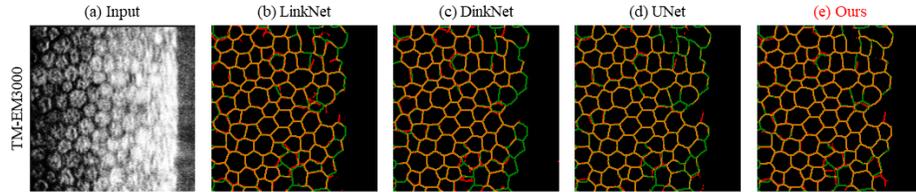


Fig. 1. Segmentation results on TM-EM3000.(a) is the input equalized corneal endothelium cell image.(b), (c), (d), (e) are the segmentation results of LinkNet, DinkNet, UNet, and our method. The red line represents the prediction result, and the green line represents ground truth, orange when the two overlap.

The main challenge of accurate segmentation is the blurred cell edges, which are difficult to segment, as shown in Fig.1, and needs more details and context information to release this problem. UNet [12] captures contextual semantic information through the contracting path and combines high-resolution features in the contracted path with upsampled output to achieve precise localization. UNet++ [18] optimizes it through a series of nested, dense skip connections to reduce the semantic gap between the encoder and decoder’s feature maps. Fabijńska [7] first applied UNet to the task of corneal endothelial cell segmentation. Vigueras-Guillén et al. [15] applied the complete convolution method based on UNet and the sliding window version to the analysis of cell images obtained by SP-1P Topcon corneal endothelial microscope. Fu et al. [8] proposed a multi-conetxt deep network, by combining prior knowledge of regions of interest and clinical parameters. However, due to the limited receptive field of local convolution and continuous downsampling, they cannot make full use of the global context and still miss many details.

The transformer has been proved to be an effective method for establishing long-range dependencies. Vaswani et al. [14] proposed a transformer structure system for language translation tasks through a complete attention mechanism to establish the global dependence of input and output among time, space, and levels. Prajit et al. [11] explored the use of the transformer mechanism on visual classification tasks, replacing all spatial convolutional layers in ResNet with stand-alone self-attention layers. However, local self-attention will still lose part of the global information. Wang et al. [16] establish a stand-alone attention layer by using two decomposed axial attention blocks, to reduce the number of parameters and calculations, and allow performing attention in a larger or even global range.

Some previous works obtain better segmentation results by taking full advantage of edge information. Chen et al. [1] proposed the deep contour-aware network, using a multi-task learning framework to study the complementary information of gland objects and contours, which improves the discriminative capability of intermediate features. Chen et al.[5] improved the network output by learning the reference edge map of CNN intermediate features. Ding et al. [6] proposed to use boundary as an additional semantic category to introduce boundary layout constraints and promote intra-class consistency through the boundary feature propagation module based on unidirectional acyclic graphs.

We need to preserve more local details and make full use of the global context. In this paper, we propose a Multi-Branch hybrid Transformer Network(MBT-Net). At first, we apply a hybrid residual transformer feature extraction module to give full play to the advantages of convolution block and transformer block in terms of local details and global semantics. Specifically, we use the convolutional block to focus on local texture feature extraction and establish long-range dependencies over space, channel, and layer by the transformer and residual connection. Besides, we define the corneal endothelial cell’s segmentation task more entirely from the perspective of edge and body. Body-edge branches provide precise edge location information and promote local consistency. The experimental results show that the proposed method is superior to other state-of-the-art methods and has achieved better performance on two corneal endothelial datasets.

2 Method

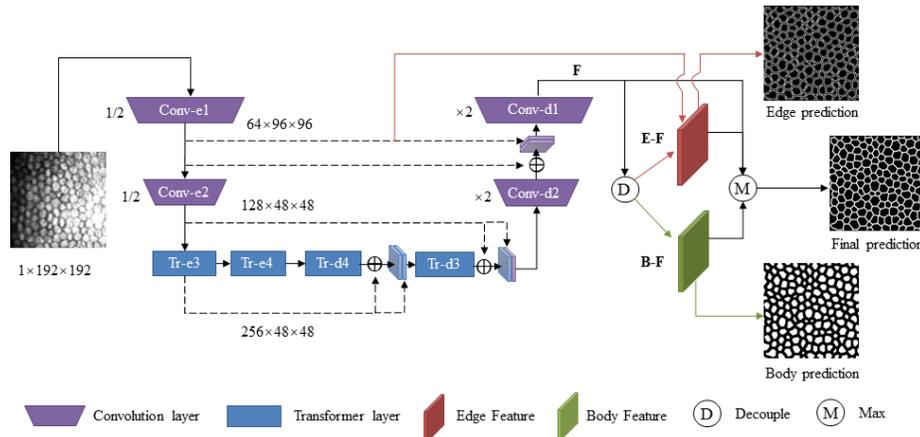


Fig. 2. The pipeline of multi-branch hybrid transformer network. Conv-e1, Conv-e2, Conv-d1 and Conv-d2 represent encoder and decoder layer based on convolution block. Tr-e3, Tr-e4, Tr-d3, and Tr-d4 are based on transformer blocks.

In this paper, we propose the MBT-Net, as shown in Fig. 2. Firstly, the feature F of equalized corneal endothelium cell image is extracted by the hybrid

residual transformer encoder-decoder module. Each convolution layer contains two basic residual blocks with a kernel size = 3×3 . Each transformer layer contains two residual transformer blocks with kernel size = 1×48 . Then, the feature is decoupled into two parts, body and edge. Also, the edge texture information from Conv-e1 is fused into the edge feature. Finally, we take the maximum response of edge feature E-F, body feature B-F, and feature F to obtain the fused feature to predict the final segmentation result. The training process of these three branches is explicitly supervised.

In this pipeline, the convolutional layer focuses on local texture feature extraction, which retains more details. The residual connection and transformer make full use of the feature map’s global context information in a more extensive range of space, different channels, and layers. The edge perspective helps preserve boundary details, and the body perspective promotes local consistency. The low-resolution feature map of d_{x+1} is refined by features from e_x by concatenating and addition operation.

2.1 Residual Transformer Block

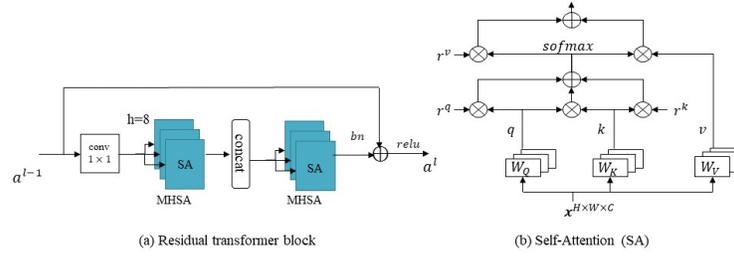


Fig. 3. The residual transformer block contains two 1×1 convolutions, a height-axial and a width-axial Multi-Head Self-Attention block (MHSA). MHSA compute axial self-attention (SA) with eight head. r , W_Q , W_K , and W_V are learnable vectors, where r related to the relative position

The residual transformer block [16] is shown in Fig.3, which contains two 1×1 convolution to control the number of channels to be calculated and a height-axial and a width-axial Multi-Head Self-Attention block (MHSA), which significantly reduces the amount of calculation. This setting allows the transformer layer to consider the global spatial context in feature map size straightly. The axial Self-Attention(SA) module is defined as:

$$y_o = \sum_{p \in N_{1 \times m}(o)} softmax_p(q_o^T k_p + q_o^T r_{p-o}^q + k_p^T r_{p-o}^k)(v_p + r_{p-o}^v) \quad (1)$$

For a given input feature map x , queries $q = W_Q x$, keys $k = W_K x$, values $v = W_V x$ are linear projections of feature map x , where W_Q , W_K , W_V are learnable parameters. $r_{(p-o)}^q$, $r_{(p-o)}^k$, $r_{(p-o)}^v$ measure the compatibility from position p to o in query, key and value. They are also learnable parameters. The $softmax_p$ denotes a softmax function applied to all possible p positions. $N_{1 \times m}(o)$ represents

the local $1 \times m$ square region centered around location o , y_o is the output at position o .

$$a^{l_2} = a^{l_1} + \sum_{i=l_1}^{l_2-1} f(a^i) \quad (2)$$

Besides, all the block used in encoder-decoder is in residual form, which can propagate input signal directly from any low layer to the high layer, optimizing information interaction [9], [10]. Taking any two layers $l_2 > l_1$ into consideration, the forward information propagation process is formulated as Eq.(2).

2.2 The Body, Edge, and Final Branches

The information from the body and edge perspectives is combined to better define corneal endothelial cell segmentation. The body branch provides general shape and overall consistency information to promote local consistency, while the edge branch provides edge localization information to improve the segmentation accuracy of image details.

We decouple the feature F extracted by the hybrid residual transformer encoder-decoder module into $F_{body} = \phi(F)$ and $F_{edge} = F - F_{body}$, where ϕ is implemented by convolution layer. Also, the low level information from encoder Conv-e1 is fused into the edge feature, $F_{edge} = F_{edge} + \psi(F_{e1})$, where ψ is dimension operation. Finally, the above three feature maps are fused into $F_{final} = \varphi(F, F_{edge}, F_{body})$ for final segmentation prediction, where φ represents the maximum response.

The training process of these three branches is explicitly supervised. The three masks used in training are shown in Fig.4. The final prediction mask is the ground truth annotated by experts. The edge prediction mask is extracted from the final prediction mask by the canny operator. The body prediction mask is obtained by inverting the final prediction and then performing Gaussian blurring at the edges.

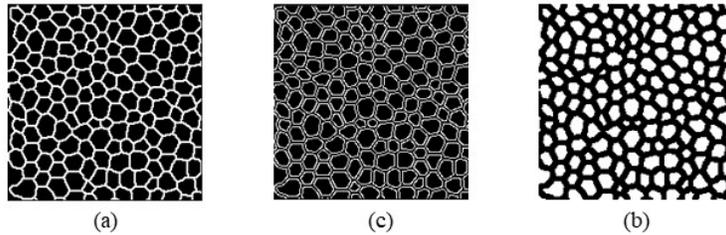


Fig. 4. Three kinds of masks on TM-EM3000.(a) The final prediction mask from the annotation of an expert, (b) The edge prediction mask extracted from (a) through the canny operator, (c) The body prediction mask by relaxing the edge of invert image of (a) with a Gaussian kernel.

2.3 Loss Function

$$Loss = \lambda_1 L_b(\hat{y}_b, y_b) + \lambda_2 L_e(\hat{y}_e, y_e) + \lambda_3 L_f(\hat{y}_f, y_f) \quad (3)$$

In this paper, we jointly optimize the body, edge, and final losses, as shown in Eq.(3), where $\lambda_1, \lambda_2, \lambda_3$ are hyper parameters to adjust the weight of three different losses. As the final prediction is the output we finally use to compare with ground truth, we give it a higher weight than edge and body branch. In our experiment, we set $\lambda_1 = 0.5, \lambda_2 = 0.5, \lambda_3 = 1.2$. y_b, y_e, y_f represent the ground truth of body, edge and final prediction respectively, and $\hat{y}_b, \hat{y}_e, \hat{y}_f$ are corresponding prediction from model. The binary cross entropy loss is used, as shown in Eq.(4).

$$L = \frac{1}{N} \sum_i [y_i \ln \hat{y}_i + (1 - y_i) \ln(1 - \hat{y}_i)] \quad (4)$$

Where N represents the total number of pixels, y_i denotes target label for pixel i , \hat{y}_i is the predicted probability.

3 Experiments

3.1 Datasets and implementation Details

TM-EM3000 contains 184 images of corneal endothelium cell and its corresponding segmentation ground truth, with size = 266×480 , collected by specular microscope EM3000, Tomey, Japan. To reduce the interference of lesions and artifacts and build a data set with almost the same imaging quality, we select a patch with a size of 192×192 from each image. This dataset is manually annotated and reviewed by three independent experts. We split it into the training set 155 patches, the validation set 10 patches, and the test set 19 patches.

Alizarine Dataset is collected by inverse phase-contrast microscope (CK 40, Olympus) at $200 \times$ magnification [13]. It consists of 30 images of corneal endothelium acquired from 30 porcine eyes and its corresponding segmentation ground truth, with image size = 768×576 , and mean area assessed per cornea = $0.54 \pm 0.07 \text{ mm}^2$. Since each image in this dataset is only partly annotated, we select ten patches of size 192×192 from each image to have 300 patches in total. And then split it into the training set 260 patches, validation set 40 patches. The training set and validation set do not overlap.

Implementation Details. We use the RMSprop optimization strategy during model training. The initial learning rate is $2e-4$, epochs = 100, batch size = 1. The learning rate optimization strategy is ReduceLRonPlateau, and the network input size is 192×192 . All the models are trained and tested with PyTorch on the platform of NVIDIA GeForce TITAN XP.

3.2 Comparison with SOTA methods

We compare performance of the proposed method with LinkNet [3], DinkNet [17], UNet [12], UNet++[18] and TransUNet[4] on **TM-EM3000** and **Alizarine** dataset. We use dice coefficient(DICE), F1 score(F1), sensitivity(SE), and

Table 1. Quantitative evaluation of different methods. The proposed method achieves the best performance.

Model	TM-EM3000				Alisarine			
	DICE	F1	SE	SP	DICE	F1	SE	SP
LinkNet34 [3]	0.711	0.712	0.719	0.941	0.766	0.801	0.805	0.956
DinkNet34 [17]	0.717	0.718	0.724	0.944	0.767	0.805	0.821	0.953
UNet [12]	0.730	0.743	0.763	0.945	0.775	0.811	0.814	0.960
UNet++ [18]	0.728	0.739	0.775	0.938	0.773	0.811	0.850	0.947
TransUNet [4]	0.734	0.742	0.769	0.941	0.783	0.821	0.866	0.948
Proposed	0.747	0.747	0.768	0.946	0.786	0.821	0.877	0.944

specificity(SP) as evaluation indicators, where DICE and F1 are most important because they are related to the overall performance.

As shown in Table 1, The proposed method has obtained the best overall performance on both TM-EM3000 and Alisarine data sets. On TM-EM3000, the DICE accuracy and F1 score of our approach are 0.747 and 0.747. On the Alisarine data set, the Dice accuracy and F1 score of our method are 0.786 and 0.821. UNet++[18] is modified from UNet, through a series of nested, dense skip connections to capture more semantic information. However, in general, there is no noticeable improvement observed in this experiment. TransUNet[4] optimized the UNet by using the transformer layer to capture the global context in the encoder part, but our method has achieved better performance. It is mainly due to the following advantages. 1) Long-range dependencies are established through the transformer in both the encoder and decoder. 2) Performing transformer layer on the whole feature map, further reducing the loss of semantic information. 3) The body-edge branch encourages the network to learn more general features and provide edge localization information.

As shown in Fig.5, on the left side of the TM-EM3000 image, the cell boundary is clear. The segmentation performance of different methods is not much different. Nevertheless, on the right side with uneven illumination and the blur cell boundary, the proposed method achieves better segmentation results, closer to the ground truth, and in line with the real situation.

There is no extensive range of fuzzy area in the Alisarine image, and all methods obtained satisfied segmentation results. However, the baselines have varying degrees of loss in details, which lead to discontinuous cell edge segmentation as the white arrow indicated, and the segmentation results do not match the ground truth well as the yellow arrow indicated. The proposed method obtains better segmentation accuracy.

3.3 Ablation Study

The ablation experiment is conducted to explore the transformer’s replacement design, as shown in Table 2. In this process, we gradually replace the encoder-decoder structure’s convolution layer with the transformer from inside to out-

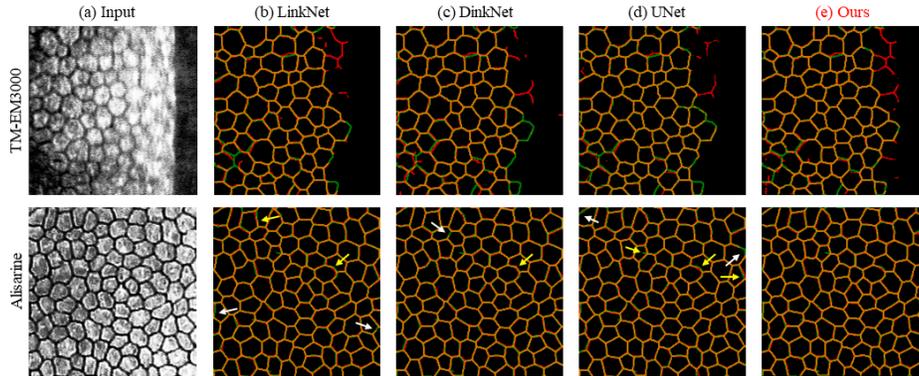


Fig. 5. Qualitative results on TM-EM3000 and the Alisarine Dataset. (a) is the input equalized corneal endothelium cell image. The red line represents the prediction result, and the green line represents ground truth, orange when the two overlap. The white arrow indicates the missing segmentation location, and the yellow arrow indicates the location where the segmentation result does not fit well with the ground truth

side. In the beginning, the model captures more semantic information, and the performance is improved. Then, with the transformer replacing the shallow convolutional layer further, the model starts to lose local information, resulting in the decline of performance. Model 2-2-TR achieves the best balance between local details and global context.

Table 2. Ablation study on the replacement design of transformer. 0-0-TR means no transformer is used. 1-1-TR means e4, d4 is transformer layer. 2-2-TR means e3, e4, and d3, d4 is transformer layer. 3-3-TR means e2, e3, e4, and d2, d3, d4 is transformer layer. 4-4-TR means complete transformer structure.

Model	TM-EM3000				Alisarine			
	DICE	F1	SE	SP	DICE	F1	SE	SP
0-0-TR	0.731	0.737	0.774	0.937	0.776	0.813	0.852	0.948
1-1-TR	0.737	0.746	0.778	0.941	0.778	0.816	0.857	0.948
2-2-TR	0.747	0.747	0.768	0.946	0.786	0.821	0.877	0.944
3-3-TR	0.702	0.714	0.742	0.935	0.777	0.812	0.874	0.940
4-4-TR	0.687	0.707	0.717	0.940	0.769	0.802	0.869	0.936

We also study the influence of transformer and body-edge branch on performance on TM-EM3000 dataset, as shown in Table 3. When neither transformer nor body-edge branches are used, the DICE accuracy and F1 score on TM-EM3000 and Alisarine are 0.720 and 0.733, respectively. After adding the body-edge branch, the performance is improved to 0.731 and 0.737. When the transformer is used, the DICE accuracy and F1 score are 0.736 and 0.741. Using

Table 3. Ablation study on transformer and body-edge branch on TM-EM3000. TR means transformer, and B-E means body-edge branch. When transformer is used, it means 2-2-TR.

TR	B-E	Dice	F1	SE	SP
✗	✗	0.720	0.733	0.746	0.945
✗	✓	0.731	0.737	0.774	0.937
✓	✗	0.736	0.741	0.786	0.936
✓	✓	0.747	0.747	0.768	0.946

both the body-edge branch and transformer, we improve the performance by 2.7% and 1.4% in total to 0.747 and 0.747.

4 Conclusion

This paper proposes a multi-branch hybrid transformer network for corneal endothelial cell segmentation, which combines the convolution and transformer block’s advantage and uses the body-edge branch to promote local consistency and provide edge localization information. Our method achieves superior performance to various state-of-the-art methods, especially in the fuzzy region. The ablation study shows that both the well-designed transformer replacement and body-edge branches contribute to improved performance.

References

- 0011, C.H., Qi, X., Yu, L., Heng, P.A.: Dcan: Deep contour-aware networks for accurate gland segmentation. CVPR pp. 2487–2496 (2016)
- Al-Fahdawi, S., Qahwaji, R., Al-Waisy, S.A., Ipson, S.S., Ferdousi, M., Malik, A.R., Brahma, A.: A fully automated cell segmentation and morphometric parameter system for quantifying corneal endothelial cell morphology. Computer Methods and Programs in Biomedicine pp. 11–23 (2018)
- Chaurasia, A., Culurciello, E.a.: Linknet: Exploiting encoder representations for efficient semantic segmentation. VCIP pp. 1–4 (2017)
- Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, L.A., Zhou, Y.: Transunet: Transformers make strong encoders for medical image segmentation (2021)
- Chen, L.C., Barron, T.J., Papandreou, G., 0002, M.K., Yuille, L.A.: Semantic image segmentation with task-specific edge detection using cnns and a discriminatively trained domain transform. computer vision and pattern recognition (2016)
- Ding, H., Jiang, X., Liu, Q.A., Magnenat-Thalmann, N., Wang, G.: Boundary-aware feature propagation for scene segmentation. ICCV pp. 6819–6829 (2019)
- Fabijanska, A.: Segmentation of corneal endothelium images using a u-net-based convolutional neural network. Artificial Intelligence in Medicine pp. 1–13 (2018)

8. Fu, H., Xu, Y., Lin, S., Wong, D.W.K., Mani, B., Mahesh, M., Aung, T., Liu, J.: Multi-context deep network for angle-closure glaucoma screening in anterior segment oct. In: International Conference on Medical image computing and computer-assisted intervention. pp. 356–363. Springer (2018)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016)
10. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: ECCV. pp. 630–645. Springer (2016)
11. Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., Shlens, J.: Stand-alone self-attention in vision models. NIPS (2019)
12. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. MICCAI (2015)
13. Ruggeri, A., Scarpa, F., Luca, D.M., Meltendorf, C., Schroeter, J.: A system for the automatic estimation of morphometric parameters of corneal endothelium in alizarine red-stained images. *British Journal of Ophthalmology* pp. 643.0–647 (2010)
14. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, N.A., Kaiser, L., Polosukhin, I.: Attention is all you need. NIPS pp. 5998–6008 (2017)
15. Viguera-Guillén, J.P., Sari, B., Goes, S.F., Lemij, H.G., van Rooij, J., Vermeer, K.A., van Vliet, L.J.: Fully convolutional architecture vs sliding-window cnn for corneal endothelium cell segmentation. *BMC Biomedical Engineering* **1**(1), 1–16 (2019)
16. Wang, H., Zhu, Y., Green, B., Adam, H., Yuille, A., Chen, L.C.: Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In: ECCV. pp. 108–126. Springer (2020)
17. Zhou, L., Zhang, C., Wu, M.: D-linknet: Linknet with pretrained encoder and dilated convolution for high resolution satellite imagery road extraction. *CVPR Workshops* pp. 182–186 (2018)
18. Zhou, Z., Siddiquee, M.R.M., Tajbakhsh, N., Liang, J.: Unet++: A nested u-net architecture for medical image segmentation. *DLMIA/ML-CDS@MICCAI* pp. 3–11 (2018)