

An Inception Inspired Deep Network to Analyse Fundus Images

Fatmatulzehra Uslu¹

¹ Bursa Technical University, Electrical-Electronic Engineering Department, Bursa, Turkey
fatmatulzehra.uslu@btu.edu.tr

Abstract

A fundus image usually contains the optic disc, pathologies and other structures in addition to vessels to be segmented. This study proposes a deep network for vessel segmentation, whose architecture is inspired by inception modules. The network contains three sub-networks, each with a different filter size, which are connected in the last layer of the proposed network. According to experiments conducted in the DRIVE and IOSTAR, the performance of our network is found to be better than or comparable to that of the previous methods. We also observe that the sub-networks pay attention to different parts of an input image when producing an output map in the last layer of the proposed network; though, training of the proposed network is not constrained for this purpose.

1. Introduction

Deep learning methods have been shown to produce state of the art performance for many image analysis problems [10]. Their performance can be associated with the use of deeper networks with very large number of parameters, residual connections [7] and regularisation techniques such as dropout [15]. However, the performance of medical image analysis problems is limited by amounts of available labelled data.

The segmentation of retinal vessels in fundus images may be required to extract topology or to measure bifurcation angles or vessel widths [5, 20, 19]. Although vessels appear to be darker than surroundings in a fundus image, they may be confused with pathologies or the border of the optic disc or the field of view (FOV). Also, the central light reflex, a bright stripe running along a vessel centreline, may cause segmentation methods to perceive it as background between two thin vessels [5].

Many studies using deep networks for vessel segmentation have trained a single network to label the location of vessels in a pixel-wise manner [9, 8]. However, when the image background is cluttered or there is large variations on the appearance of vessels, the capacity of a deep network, which is often limited by the amount of labelled data, may need to be used for eliminating uneven illumination, large noise and other structures such as the optic disc and pathologies, in addition to precisely finding the location of vessels.

This study presents a deep network for vessel segmentation in fundus images. The architecture of the network is inspired by inception modules [18]. The proposed network consists of three sub-networks, each with a different filter size. The decision given by each sub-network is combined in the last layer of the proposed network to produce a single decision. We evaluate the performance of the proposed network on two fundus image datasets generated by different imaging modalities: DRIVE and IOSTAR datasets. Based on our experiments, we observe that each sub-network seems to specialise at a different region of an

input image, without giving any supervision for it: pixels inside FOV versus those outside FOV, non-vessel pixels inside FOV and vessel pixels inside FOV.

2. Related Work

2.1. U-net

U-net is one of the best known architectures in medical image segmentation [14]. The architecture has two paths: one for encoding the input image and the other one for decoding the corresponding segmentation map. Two paths were connected with skip connections to improve gradient flow through the network. Skip connections also provide access to features produced at early layers in generation of the segmentation maps due to the concatenation of the features at encoding path and those at decoding path. In the architecture, there are nine convolutional layers. Each convolutional layer contains two 3×3 filters. The filters are followed by pooling along the encoding path to half the grid size. They are preceded by upsampling along the decoding path to double the grid size.

2.2. Inception Modules

Inception modules were proposed by Szegedy *et al.* [18, 17], which were used in “GoogleNet”. The inception modules were designed to be micro-networks [18], which can be located at desired depths of a macro-network. The main characteristic of inception modules is that they contain a range of filter types in parallel; a basic version consists of 1×1 , 3×3 , 5×5 filters. This structure provides rich feature set for the next convolutional layer of a deep network by leading to better performance with a small parameter number, which was proven by the performance of “GoogleNet” [18] in ImageNet Large-Scale Visual Recognition Challenge 2014 (ILSVRC14). In order to reduce the parameter number of the inception modules, one may use 1×1 filters to decrease the channel number, as illustrated in **Fig. 1(a)**, or replace large filters with small ones; a 5×5 filter is factorised to two 3×3 filters in **Fig. 1(b)**.

2.3. Residual Connections

He *et al.* integrated residual connections into deep networks and showed an increase in the performance of the networks due to the residual connections facilitating better gradient flow through the layers of the networks [7]. As shown in **Fig. 1(c)**, the output of a convolutional layer with a residual connection becomes the sum of the output of the convolutional layer ($f(x)$) and its input (x).

3. Method

Inspired by the architecture of the inception modules, we present a network with three parallel sub-networks, each with a

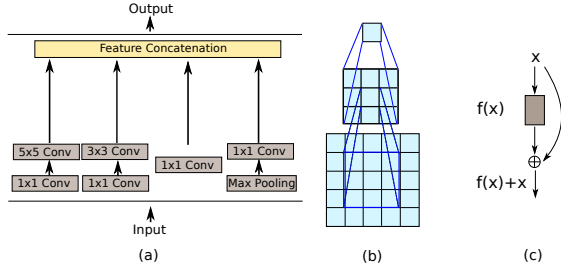


Figure 1. (a) The inception network with pooling [18, 17] (b) Realisation of 5×5 filter with two 3×3 filters (c) A convolutional layer with a residual connection. (Best viewed in color.)

different filter size. The outputs of the sub-networks are combined with a final convolutional layer to allow them both to be jointly trained and to produce a single vessel mask for an input image. We expect the sub-networks to act as three experts giving a joint decision on the generation of vessel masks.

Fig. 2 illustrates the general structure of the proposed architecture. We use U-net as a sub-network. Each sub-network contains modules with one of predetermined filter types, which are 1×1 , 3×3 and 5×5 filters. In order to increase gradient flow through the network, we integrate residual connections into the modules [7]. In contrast to U-net, we use one module, in the place of two convolution layers [14], at each grid size and padded features prior to applying filters to preserve their size. In the final layer of each sub-network, we reduce the number of filters to 1 by using 1×1 convolution, which can be viewed as each sub-network producing its own object mask. Eventually, three features generated by the sub-networks are combined with 1×1 convolution on the top layer of the proposed network to produce a single object mask. We train the network as a regressor, with Euclidean loss between vessel masks synthesised by the network and ground truth images.

Similar to the U-net [14], we double filter sizes along the encoding path of the proposed network and halve them along its decoding path. When downsampling features along the encoding path, we use 2×2 max pooling with a stride of 2 for the sub-network with a 1×1 filter and apply convolution with a stride of 2 pixels for the other two sub-networks. When upsampling features along the decoding path, we use bilinear interpolation after reducing channel number of incoming features by a factor of 2 with 1×1 filters to lighten computation burden, in contrast to U-net [14] where channel reduction was performed after upsampling.

Each convolutional layer of our network is preceded by batch normalisation and followed by RELU activation function, apart from the one in layer 11, which is followed by sigmoid function. Apart from the modules with 1×1 filters, we use 1×1 filters in each module regardless of its filter type to reduce the number of incoming filters, prior to applying filters specific to each sub-network. Similar to inception modules [18, 17], we factorise 5×5 filters to two 3×3 filters to reduce the parameter number of the proposed network, which is 250, 933 in total. **Table 1** shows the number of filters for each layer of sub-networks regardless of filter type.

Table 1. Filter numbers for the layers of the sub-network with 1×1 filters.

Layer No	Filter No	Filter No in Upsampling Layer
1	8	
2	16	
3	32	
4	64	
5	128	64
6	64	32
7	32	16
8	16	8
9	8	
10	1	

4. Material and Experimental Setup

4.1. Material

The segmentation performance of the proposed network was assessed on a well known fundus image dataset, the DRIVE¹ [16] and a recently released fundus image dataset, IOSTAR². The DRIVE was captured by Canon CR5 with FOV of 45° from generally healthy people. The number of images in the dataset is 40; of these, 7 show the signs of diabetic retinopathy. The resolution of images is 768×584 pixels. In order to allow a fair performance comparison, the dataset is divided into two: one for training and the other one for testing. Each set contains 20 images. Manually traced vessel maps and FOV masks are also provided with the dataset.

IOSTAR dataset [22] contains 30 Scanning Laser Ophthalmoscopy (SLO) images, which were captured by an EasyScan camera with a FOV of 45° . The resolution of images is 1024×1024 pixels. FOV masks and binary vessel masks are included in the dataset. We used the first 20 images as the training and validation sets and evaluated the performance of the proposed method on the the last 10 images.

4.2. Experimental Setup

We initially set learning rate to 0.0008 then decreased it with an exponential decay rate of 0.94 during training. We trained the network for 60 epochs for both datasets. We optimised parameters of the proposed network with Adam algorithm, with default values of β_1 and β_2 , by using mini-batches of 64 images. We initialised the parameters of our network with He *et al.*'s technique [6]. The total number of parameters was slightly less than a million.

We randomly cropped 4000 image patches from each image in the training set of the DRIVE. Of these image patches, 400 were used for validation and the rest was allocated for training. The same process was also realised for IOSTAR dataset. The size of image patches for both datasets was set to be 96×96 pixels.

We used color images and applied channel-wise normalisation. The mean contrast for each colour channel was calculated over our training set and subtracted from training, validation and test sets. In order to increase variety on our training sets, we applied color jittering to them, without data augmentation. Final vessel probability maps were generated by combining the output probability maps of the network for input image patches sampled with a stride of 30 pixels at each direction for both DRIVE and IOSTAR datasets.

¹<https://www.isi.uu.nl/Research/Databases/DRIVE/>

²<http://www.retinacheck.org/datasets>

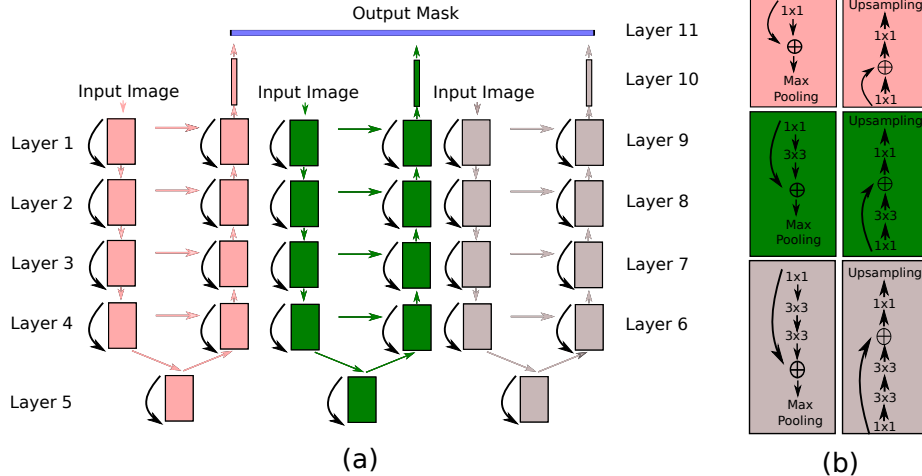


Figure 2. (a) An overview of the proposed architecture. Pink, green and brown boxes respectively show modules with 1×1 , 3×3 and 5×5 filters. Pink, green and brown arrows respectively demonstrate the connection between filters with the same colour code as the arrows and black arrows show residual connections. Blue box illustrates the final convolutional layer of the proposed network. Narrow boxes in pink, green and brown represent filters with 1 channel. The numbers of channels in the other filters are not associated with the widths of the boxes. Input images for each sub-network is the same. (b) Filters inside modules with the same colour code. (Best viewed in color.)

4.3. Evaluation Criteria

In order to binarize final probability maps, we found a threshold with Otsu’s method. In the binary maps, we calculated *accuracy* (*acc*), *sensitivity* (*sens*), *specificity* (*spec*) and *geometric mean* (*g-mean*) [1] by considering pixels inside the FOV masks [5], as follows.

$$acc = \frac{TN + TP}{TP + TN + FN + FP} \quad (1)$$

$$sens = \frac{TP}{TP + FN} \quad (2)$$

$$spec = \frac{TN}{TN + FP} \quad (3)$$

$$g_mean = \sqrt{sens \cdot spec} \quad (4)$$

where *TP*, *TN*, *FN* and *FP* respectively denote true positives, true negatives, false negatives and false positives, which are calculated with the standard approach [5]. A single value for each metric was computed over the complete set of test images. We also calculated the area of the Receiver Operating Characteristic curve (ROC) over final probability maps. This metric gives the discrimination capacity of the network regardless of any threshold.

5. Results

5.1. Segmentation Performance

Table 2 compares the segmentation performance of the proposed method with that of previous methods. According to the table, our method generates consistent performance for both DRIVE and IOSTAR datasets; though, the datasets were captured with different imaging modalities and so show different characteristics such as larger variation in hue in IOSTAR dataset. The proposed method outperforms previous methods, including other deep leaning based studies, for both datasets with a score of 0.98 on AUC and produces comparable or better performance on the other metrics for DRIVE dataset. For IOSTAR dataset, the proposed network outperforms other state of the art methods with a significant margin on sensitivity, with

a score of 0.81, and G-mean, with a score of 0.89.

Indeed, the high performance of our network seems to be not due to being with large parameter count or extensive pre-processing stage but as a result of its design. The network of Liskowski and Krawiec [9] contains almost 50 times of parameter count that our network includes. Therefore, it requires a far larger size of training dataset, which was obtained with exhaustive data augmentation [9]. Zhang *et al.* [22] and Na *et al.* [11] enhanced the appearance of vessels by applying pre-processing techniques such as homogeneity correction prior to the use of their methods. Moreover, Na *et al.* reported that their method could not reach to the same performance when a homogeneity correction proposed by them [11] was not used, where sensitivity and specificity respectively reduced from 0.76 to 0.75 and from 0.98 to 0.92. It should be noted that our method did not rely on such a mechanism to improve the homogeneity in images; however, it still shows higher performance.

Fig. 3 and Fig. 4 demonstrate segmentation maps with the maximum and minimum G-mean scores, generated by the proposed method for both DRIVE and IOSTAR. As seen in the figures, the proposed method achieves to segment almost the complete vasculature, including many small vessels.

5.2. Activation Map Produced at the Last Layer of Each Sub-Network

This section will investigate roles of sub-networks on the decisions of the proposed network. Because each sub-network possess a different filter type, we expect each sub-network to behave as an expert on specific structures in input images.

Fig. 5 demonstrates activation maps at the final layers of sub-networks and output probability maps generated after combining these activation maps at the final layer of the proposed network, for eight input image patches randomly sampled from DRIVE dataset; each image patch is examined in its corresponding column in the figure. As seen in the figure, activation maps produced by filters after RELU, from top to bottom, respectively seem to focus on non-vessel pixels inside FOV, pix-

Table 2. Vessel segmentation performance comparison on DRIVE and IOSTAR.

Dataset	Year	Method	AUC	Accuracy	Sensitivity	Specificity	G-mean
DRIVE	2019	The proposed method	0.98	0.95	0.81	0.98	0.89
	2017	Orlando <i>et al.</i> [13]			0.79	0.97	0.87
	2016	Liskowski and Krawiec [9]	0.97	0.95	0.75	0.98	0.86
	2016	Oliveira <i>et al.</i> [12]	0.95	0.95	0.86	0.96	0.91
	2015	Li <i>et al.</i> [8]	0.97	0.95	0.76	0.98	0.86
	2015	Wang <i>et al.</i> [21]	0.95	0.98	0.82	0.97	0.89
	2014	Cheng <i>et al.</i> [3]	0.96	0.95	0.72	0.98	0.84
	2013	Fraz <i>et al.</i> [4]		0.94	0.73	0.97	0.84
IOSTAR	2019	The proposed method	0.98	0.96	0.81	0.98	0.89
	2016	Zhang <i>et al.</i> [22]	0.96	0.95	0.75	0.97	0.85
	2017	Na <i>et al.</i> [11]	0.96	0.96	0.76	0.98	0.86

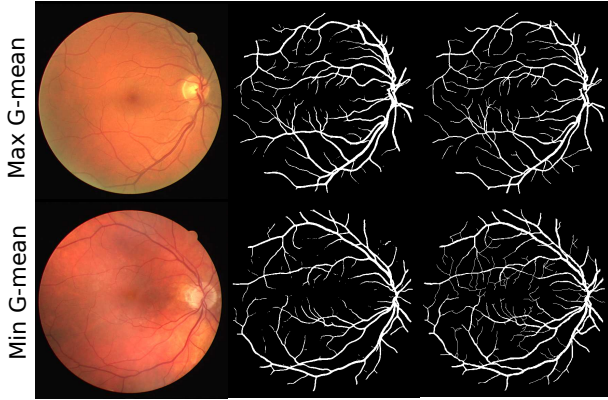


Figure 3. The segmentation performance of the proposed method on DRIVE dataset. Color fundus images are accompanied with vessel maps produced by the proposed method and ground truth segmentation masks respectively. Images from top to bottom respectively belong to 18_test.tif and 07_test.tif. (Best viewed in colour.)

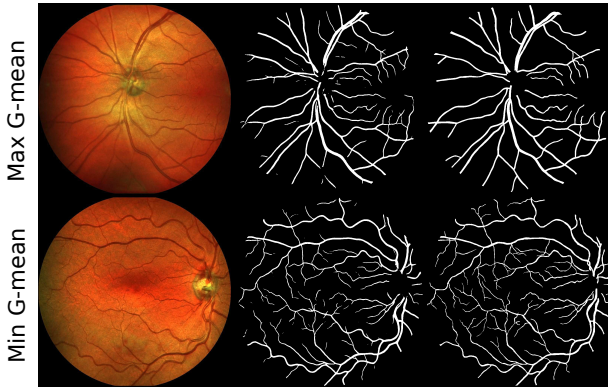


Figure 4. The segmentation performance of the proposed method on IOSTAR dataset. Color fundus images are accompanied with vessel maps produced by the proposed method and ground truth segmentation masks respectively. Images from top to bottom respectively belong to 44.OSN.jpg and 34.ODC.jpg. (Best viewed in colour.)

els outside FOV and vessel pixels inside FOV. Moreover, this specialisation of sub-networks seems to provide better description of image patches regarding the location of vessels and other

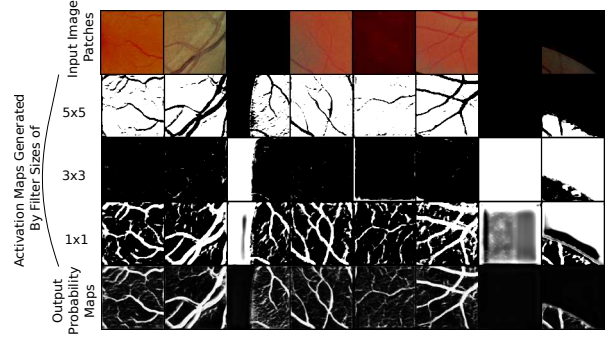


Figure 5. Activation maps generated at the last layers of sub-networks with filter size of 1×1 , 3×3 and 5×5 , demonstrated from the second row to the fourth one. The first and last rows respectively show input image patches and output probability maps generated by the proposed network. Each sub-network seems to become an expert on one task: the pixel classification regarding the position of image patches inside FOV or outside FOV (the sub-network with 3×3 filters), the classification of background pixels inside FOV (the sub-network with 5×5 filters) and the classification of vessel pixels inside FOV (the sub-network with 1×1 filters). Finally, the proposed network combines the activations from the sub-networks in an effective way, without exclusively outputting any of them. Images patches are randomly sampled from DRIVE dataset. (Best viewed in color.)

structures in images.

A detailed examination of the figure reveals that vessels with low contrast, such as arteries with the central light reflection in the second input image patch and tiny vessels in the first and fifth image patches, are well recognised. The proposed network manages to assign larger probabilities to vessels despite large variation in their contrasts and thickness.

6. Conclusion

This paper presented a deep network for medical image segmentation, where the segmentation of a structure may be hindered due to the presence of pathologies, other organs and imaging related problems. The proposed network is a composite of three sub-networks, each with the same architecture but a different filter size. According to the results of experiments carried out in two fundus image datasets, DRIVE and IOSTAR, which are captured by different techniques (CCD camera and SLO), our network outperformed previous studies with a significant

margin on sensitivity and G-mean for IOSTAR dataset, without using any extensive preprocessing to improve the appearance of vessels, in contrast to Zhang *et al.*'s and Na *et al.*'s methods. Also, our network showed better or comparable performance on DRIVE when compared with that of previous methods.

We show that each sub-network attuned to specific tasks such as the identification of pixels outside FOV and the classification of vessel pixels from others inside FOV. One may find similarities between the output of our sub-networks and attention maps, which is designed to pay attention to such regions of an image that may facilitate the detection or segmentation of an object of interest [2]. Despite our network consisting of three sub-networks, its total parameter number is far smaller than that of other networks producing similar performance[9]. This allows the proposed network to be easily applied for image segmentation tasks with the limited amount of labelled data.

7. References

- [1] Josephine Akosa. Predictive accuracy: A misleading performance measure for highly imbalanced data. In *Proceedings of the SAS Global Forum*, 2017.
- [2] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3640–3649, 2016.
- [3] Erkang Cheng, Liang Du, Yi Wu, Ying J Zhu, Vasileios Megalooikonomou, and Haibin Ling. Discriminative vessel segmentation in retinal images by fusing context-aware hybrid features. *Machine Vision and Applications*, 25(7):1779–1792, 2014.
- [4] Muhammad Moazam Fraz, A Basit, and SA Barman. Application of morphological bit planes in retinal blood vessel extraction. *Journal of Digital Imaging*, 26(2):274–286, 2013.
- [5] Muhammad Moazam Fraz, Paolo Remagnino, Andreas Hoppe, Bunyarit Uyyanonvara, Alicja R Rudnicka, Christopher G Owen, and Sarah A Barman. Blood vessel segmentation methodologies in retinal images—a survey. *Computer Methods and Programs in Biomedicine*, 108(1):407–433, 2012.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [8] Qiaoliang Li, Bowie Feng, LinPei Xie, Ping Liang, Huisheng Zhang, and Tianfu Wang. A cross-modality learning approach for vessel segmentation in retinal images. *IEEE Transactions on Medical Imaging*, 35(1):109–118, 2015.
- [9] P. Liskowski and K. Krawiec. Segmenting retinal blood vessels with deep neural networks. *IEEE Transactions on Medical Imaging*, 35(11):2369–2380, 2016.
- [10] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen AWM van der Laak, Bram Van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.
- [11] Tong Na, Yitian Zhao, Yifan Zhao, and Yue Liu. Superpixel-based line operator for retinal blood vessel segmentation. In *Annual Conference on Medical Image Understanding and Analysis*, pages 15–26. Springer, 2017.
- [12] Wendeson S Oliveira, Joyce Vitor Teixeira, Tsang Ing Ren, George DC Cavalcanti, and Jan Sijbers. Unsupervised retinal vessel segmentation using combined filters. *PLOS ONE*, 11(2):e0149943, 2016.
- [13] José Ignacio Orlando, Elena Prokofyeva, and Matthew B Blaschko. A discriminatively trained fully connected conditional random field model for blood vessel segmentation in fundus images. *IEEE Transactions on Biomedical Engineering*, 64(1):16–27, 2017.
- [14] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015.
- [15] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [16] Joes Staal, Michael D Abràmoff, Meindert Niemeijer, Max A Viergever, and Bram Van Ginneken. Ridge-based vessel segmentation in color images of the retina. *IEEE Transactions on Medical Imaging*, 23(4):501–509, 2004.
- [17] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, volume 4, page 12, 2017.
- [18] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015.
- [19] Fatmatülzehra Uslu and Anil Anthony Bharath. A multi-task network to detect junctions in retinal vasculature. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 92–100. Springer, 2018.
- [20] Fatmatülzehra Uslu and Anil Anthony Bharath. A recursive bayesian approach to describe retinal vasculature geometry. *Pattern Recognition*, 87:157–169, 2019.
- [21] Shuangling Wang, Yilong Yin, Guibao Cao, Benzhenq Wei, Yuanjie Zheng, and Gongping Yang. Hierarchical retinal blood vessel segmentation based on feature and ensemble learning. *Neurocomputing*, 149:708–717, 2015.
- [22] Jiong Zhang, Behdad Dashtbozorg, Erik Bekkers, Josien PW Pluim, Remco Duits, and Bart M ter Haar Romeny. Robust retinal vessel segmentation via locally adaptive derivative frames in orientation scores. *IEEE Transactions on Medical Imaging*, 35(12):2631–2644, 2016.