

Regression of Instance Boundary by Aggregated CNN and GCN

Yanda Meng¹, Wei Meng¹, Dongxu Gao¹, Yitian Zhao², Xiaoyun Yang³,
Xiaowei Huang⁴, and Yalin Zheng¹(✉)

¹ Department of Eye and Vision Science, Institute of Life Course and Medical Sciences, University of Liverpool, Liverpool, UK
yalin.zheng@liverpool.ac.uk

² Cixi Institute of Biomedical Engineering, Ningbo Institute of Industrial Technology, Chinese Academy of Sciences, Ningbo, China

³ China Science IntelliCloud Technology Co., Ltd, Shanghai, China

⁴ Department of Computer Science, University of Liverpool, Liverpool, UK

Abstract. This paper proposes a straightforward, intuitive deep learning approach for (biomedical) image segmentation tasks. Different from the existing dense pixel classification methods, we develop a novel multi-level aggregation network to directly regress the coordinates of the boundary of instances in an end-to-end manner. The network seamlessly combines standard convolution neural network (CNN) with Attention Refinement Module (ARM) and Graph Convolution Network (GCN). By iteratively and hierarchically fusing the features across different layers of the CNN, our approach gains sufficient semantic information from the input image and pays special attention to the local boundaries with the help of ARM and GCN. In particular, thanks to the proposed aggregation GCN, our network benefits from direct feature learning of the instances' boundary locations and the spatial information propagation across the image. Experiments on several challenging datasets demonstrate that our method achieves comparable results with state-of-the-art approaches but requires less inference time on the segmentation of fetal head in ultrasound images and of optic disc and optic cup in color fundus images.

Keywords: Regression, Semantic Segmentation, CNN, GCN, Attention, Aggregation

1 Introduction

The accurate assessment of anatomic structures in biomedical images plays an important role in the management of many medical conditions or diseases. For instance, fetal head (FH) circumference in ultrasound images is a critical indicator for prenatal diagnosis and can be used to estimate the gestational age and to monitor the growth of the fetus [25]. Similarly, the size of the optic disc (OD) and optic cup (OC) in color fundus images is of great importance for the diagnosis of glaucoma, an irreversible eye disease [35]. Manual annotation of this

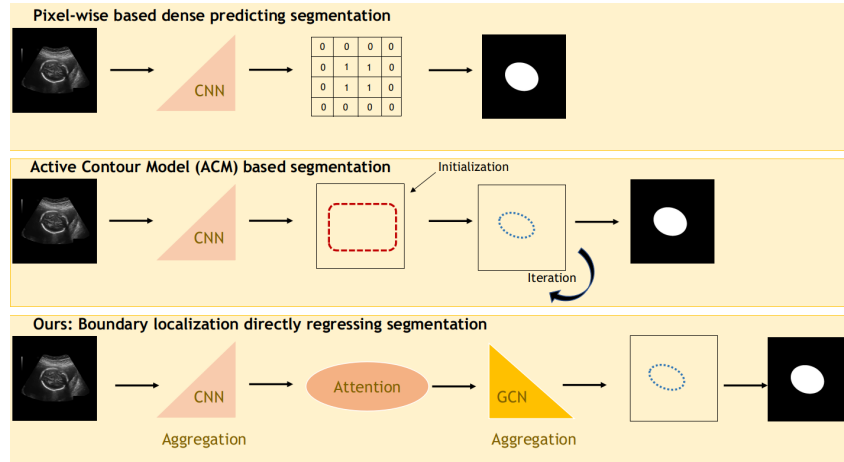


Fig. 1: Three different segmentation paradigms by deep learning. Top row: pixel-wise based methods [14,22,6] that classify each pixel into objects or background. Middle row: active contour based methods [32,9] that need iterative optimization in action to find the final contours. Bottom row: our proposed method that directly regresses the locations of object boundaries by information aggregation through CNN and GCN, enhanced by an attention module.

kind of structures by delineating their boundaries in clinics is unrealistic as it is costly, time consuming, labor intensive, and subject to human experience and errors. Automatic segmentation of biomedical images is believed to be able to help improve the efficiency of workflow in clinical scenarios. Inspired by the way clinicians annotate images, we propose an aggregated network to solve the segmentation tasks through directly regressing the locations of objects' boundaries, and demonstrate the effectiveness of the network in the segmentation of FH in ultrasound and OD & OC in color fundus images, respectively.

The biomedical image semantic segmentation task remains a challenging problem in the field of computer vision. The commonly-used deep learning-based semantic segmentation methods [22,6,41] (top row of Fig. 1) classify each pixel of an image into a category or class. These methods benefit from Convolution Neural Networks (CNN)'s excellent ability to extract high-level semantic features. Being a part of the understanding of scenes or global contexts, these methods need to learn the object location, object boundary, and object category from the high-level semantic information and local location information [31]. However, they suffer from the loss of local location information at the pixel-level [8], because a large receptive field corresponds to a small feature map, and this dilemma has increased the difficulties of dense prediction tasks. In order to solve this problem, approaches in [50,4] either maintain the resolution of the input image with dilated convolution, or capture sufficient receptive fields with pyramid pooling modules. The insights behind these methods indicate that the spatial in-

formation and the receptive field are both important to achieving high accuracy. However, it is hard to meet these two requirements simultaneously with CNN [45]. In particular, it is often challenging to maintain enough spatial information of the input image.

To address the aforementioned challenges, we follow a straightforward and intuitive methodology that human operators take to segment objects and regard segmentation as a regression task. Compared with the preserving abstraction of spatial details [50,4], we use a combination of CNN, ARM, and GCN to directly regress the boundary locations of the instances in the Euclidean space. Our method is different from the recent polygon-based active contour models (ACM) methods [32,9,20] (middle row of Fig. 1), which need to initialize the boundaries and iteratively find the final object boundaries for a new image. On the contrary, we directly supervise the model to learn the precise location of boundaries and produce the boundaries without iteration during inference. Compared with the pixel-wise based methods, our method needs to learn and extract more spatial information to regress the location directly. To address this issue, the local spatial information propagation nature of GCN is exploited. GCN has recently been applied to many low-level tasks, such as scene understanding [29], semantic segmentation [6], and pose estimation [52], because GCN can propagate the information through neighbor nodes (short range) and hence allow the model to learn local spatial correlation structure.

We propose an aggregated GCN decoder with graph vertices sampling from sparse to dense, which contributes to globally propagate the spatial relationship information across the whole image. This will provide greater representational power and more sufficient information propagation than previous segmentation methods based on Conditional Random Fields or Markov Random Fields [2,30]. Thus, we can directly regress explicit boundary location with the Euclidean space coordinate representation. This strategy addresses the concerns of most recent works [43,44], which share the similar idea but convert the Euclidean space representation into polar representation, and regressing the low-level distance between the center point and boundary points. They found that CNN cannot regress the Euclidean space coordinate representation of the boundary well as some more noise may be added, and the CNN may not maintain enough spatial information [43,44]. Our proposed aggregation GCN can handle this issue well, and our experiment results prove that. Besides, those methods' performance may suffer from the low-quality of center point, so, Xie *et al.* [43] utilized center sample methods to classify and selected high-quality center points to improve the segmentation result. In contrast, our methods can directly regress the boundary location without any further center selection process. As for the proposed CNN aggregation mechanism, some low-level features are unnecessarily over-extracted while object boundaries are simultaneously under-sampled. In order to extract more useful and representative features, we apply the ARM working as a filter between CNN encoder and GCN decoder, which cooperates with the GCN to gain more effective semantic and spatial features, especially the boundary location information from CNN.

In summary, this work makes the following contributions:

- We take a straightforward and intuitive approach to (biomedical) image semantic segmentation and regard it as a direct boundary regression problem in an end-to-end fashion.
- We propose aggregating mechanisms on both CNN and GCN modules, to enable them to reuse and fuse the contextual and spatial information. The additional attention mechanism helps the GCN decoder to gain more useful semantic and spatial information from the CNN encoder.
- We apply a new loss function tailored for object boundary localization that will help to make update step size adaptive to the error values during the training stage.

It is envisaged that the proposed framework may serve as a fundamental and strong baseline in future studies of biomedical semantic segmentation tasks.

2 Related Work

2.1 Pixel-based Methods

Fully Convolution Neural Networks (FCNs) [31] and U-Net architectures [38] are widely used in semantic segmentation tasks [22,6]. These methods are aimed at extracting more spatial information or extending the receptive field that is of pivotal importance in semantic segmentation tasks. However, it is still difficult to capture longer-range correspondence between pixels in an image [48].

Aggregation module In order to gain global contextual dependencies of an image, methods like [50,53,47,41] proposed to fuse multi-scale or multi-level features through aggregating across semantic and spatial feature domains. Zhao *et al.* [50] proposed a pyramid network that utilizes multiple dilated convolution blocks [46] to aggregating global feature maps on different scales. Other approaches such as Deeplab methods [4,5,6] exploited parallel dilated convolution with different rates to extract features at an arbitrary resolution and preserve the spatial information. However, it is still hard to efficiently learn the discriminative feature representation as many low-level features are unnecessarily over-extracted. Therefore, these aggregation methods may result in an excessive use of information flow.

Attention mechanism Alternatively, some other algorithms exploited the benefits of attention mechanism to integrate local discriminative representation and global contextual features. For example, DANet and CSNet [15,34] used the attentions in spatial and channel dimensions respectively to adaptively integrate local features with their global dependencies. Furthermore, Zhao *et al.* proposed the point-wise spatial attention network [51], which connected each position in the feature map with all the others through self-adaptive attention maps to harvest local and long-range contextual information flexibly and dynamically. In this work, an ARM module is also used to supervise our model to learn discriminate features from input images.

2.2 Polygon-based Methods

Instead of assigning each pixel with a class, some recent methods [32,9,20,43,44] started to predict the position of all vertices of the polygon around the boundary of the target objects. The recent work [43,44] used polar coordinates to represent object contours. Both methods achieved comparable results with pixel-based segmentation methods in instance segmentation tasks. Also, the combination of FCNs and Active Contour Models (ACMs) [27] has been exploited. Some methods formulated new loss functions that were inspired by the ACMs principles [7,21] to tackle the task of ventricle segmentation in cardiac MRI. Other approaches used the ACMs as a post-processor of the output of an FCN, for example, Marcos *et al.* [32] proposed a Deep Structured Active Contours model that combined ACMs and pre-trained FCNs to learn the energy surface of the reference map. These ACM-based methods achieved state-of-the-art performance in many segmentation tasks. However, there are still two main limitations. First, the contour curve must be initialized, while the initialized curve is far away from the ground truth, it may be insufficient to optimize or make an inference. Second, due to the iterative inference mechanism of ACMs, they require a relatively longer running time during training and testing.

2.3 GCNs in Segmentation

GCNs have been applied to image segmentation tasks recently, as they can propagate and exchange the local short-range information through the whole image to learn the semantic relations between objects [39,48]. In 2D image semantic segmentation tasks, Li *et al.* proposed a Dual Graph Convolutional Network (DGCNet) [48], which applied two orthogonal graphs frameworks to compute the global relational reasoning of the whole image and the reasoning process can help the whole network to gain rich global contextual information. Another work [39] proposed by Shin *et al.* shared the similar idea, and utilized GCN to learn the global structure of the shape of the object, which reflected the connectivity of neighbouring vertices. Apart from using GCN to learn global contextual information from 2D input, our approach also exploits spatial and local location information. Compared with a recent similar work [33], our method further exploit the relations between low-level and much more high-level vertex information in GCN decoder and perform a ‘skip up sampling’ in terms of Graph convolutions between two layers. This operation helps our model further extract feature correlations among different layers.

3 Method

3.1 Graph Representation

The manually annotated object boundaries are extracted from the binary image and equally sampled into N vertices with the same angle interval $\Delta\theta$ (e.g.

$N = 360$, $\Delta\theta = 1^\circ$). The geometric center of the boundary represents the center vertex. We describe the object contour with vertices and edges as $B = (V, E)$, where V has $N + 1$ vertices in the Euclidean space, $V \in \mathbb{R}^{N \times 2}$, and $E \in \{0, 1\}^{(N+1) \times (N+1)}$ is a sparse adjacency matrix, representing the edge connections between vertices, where $E_{i,j} = 1$ means vertices V_i and V_j are connected by an edge, and $E_{i,j} = 0$ otherwise. Every two continuous vertices on the contour are connected with an edge and are both connected to the center vertices with another two edges to form a triangle. For the OD and OC segmentation, their contours are sampled separately while the geometric centre of the OC is shared as the centre vertex. Thus, there are 360 triangles and 361 vertices for instances in FH images and 720 triangles and 721 vertices for OD and OC images. For more details, please refer to the supplementary material.

We directly use the coordinates in the Euclidean space to represent all the vertices and exploit the semantic and spatial correspondence between the inputs' instance and boundaries. Besides, our boundary representation method is not sensitive to the center point as the boundary does not have too many correlations with the center point.

3.2 Graph Fourier Transform & Convolution

According to [10], the normalized Laplacian matrix is $L = I - D^{-\frac{1}{2}}ED^{-\frac{1}{2}}$, where I is the identity matrix, and D is a diagonal matrix that represents the degree of each vertex in V , such that $D_{i,i} = \sum_{j=1}^N E_{i,j}$. The Laplacian of the graph is a symmetric and positive semi-definite matrix, so L can be diagonalized by the Fourier basis $U \in \mathbb{R}^{N \times N}$, such that $L = U\Lambda U^T$. The columns of U are the orthogonal eigenvectors $U = [u_1, \dots, u_n]$, and $\Lambda = \text{diag}([\lambda_1, \dots, \lambda_n]) \in \mathbb{R}^{N \times N}$ is a diagonal matrix with non-negative eigenvalues. The graph Fourier transform of the vertices representation $x \in \mathbb{R}^{N \times 3}$ is defined as $\hat{x} = U^T x$, and the inverse Fourier transform as $x = U\hat{x}$. The spectral graph convolution of i and j is defined as $i * j = U((U^T i) \odot (U^T j))$ in the Fourier space. Since U is not a sparse matrix, this operation is computationally expensive. To reduce the computation, Defferrard *et al.* [12] proposed that the convolution operation on a graph can be defined in Fourier space by formulating spectral filtering with a kernel g_θ using a recursive Chebyshev polynomial [12]. The filter g_θ is parametrized as a Chebyshev polynomial expansion of order K , such that

$$g_\theta(L) = \sum_{k=1}^K \theta_k T_k(\hat{L}) \quad (1)$$

where $\theta \in \mathbb{R}^K$ is a vector of Chebyshev coefficients, and $\hat{L} = 2L/\lambda_{max} - I_N$ represents the rescaled Laplacian. $T_k \in \mathbb{R}^{N \times N}$ is the Chebyshev polynomial of order K , that can be recursively computed as $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$ with $T_0 = 1$ and $T_1 = x$. Therefore, the spectral convolution can be defined as

$$y_j = \sum_{i=1}^{F_{in}} g_{\theta_{i,j}}(L)x_i \quad (2)$$

where x_i is the i -th feature of input $x \in \mathbb{R}^{N \times F_{in}}$, which has F_{in} features, with $F_{in} = 2$ in this work and $y \in \mathbb{R}^{N \times F_{out}}$ is the output. The entire filter operation is computationally faster and the complexity drops from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ [3].

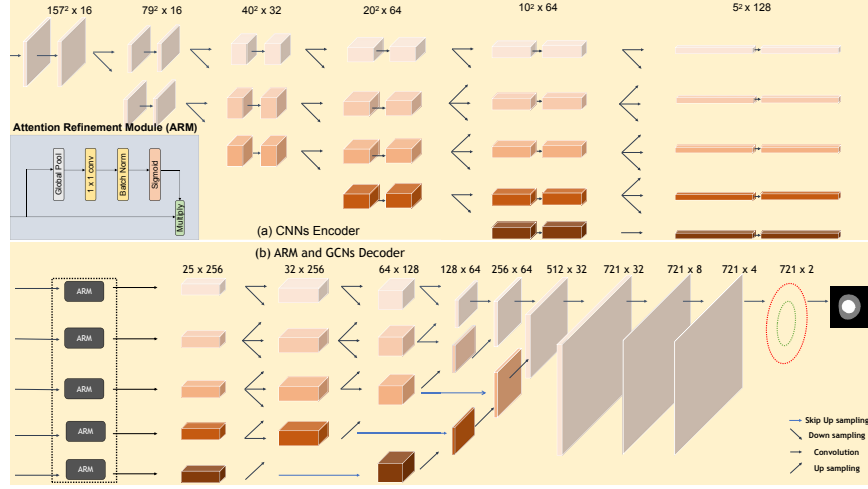


Fig. 2: Overview of our proposed network structure. The size of feature maps of the CNN encoder and vertex maps of the GCN decoder for each stage (columns) are shown. In the CNN encoder, the horizontal black arrow represents CNN convolutional operations that are achieved by a standard CNN Residual Block [24] with kernel size 3×3 , stride 1, followed by a Batch Normalization (BN) layer [26] and Leaky ReLU as the activation function. The down-sampling is conducted by setting stride size as 2, the lower level feature is bi-linearly up-sampled by a factor 2. In the GCN decoder, down-sampling and up-sampling are conducted by graph vertices sampling, which is described in Section 3.3, and the horizontal black arrow represents residual graph convolution (ResGCN) blocks [28] with polynomial order 4. The horizontal blue arrow achieves ‘skip up sampling’ with vertices number four times up sampled in terms of graph vertices sampling method via retained vertices. In this figure, the example is for OD and OC segmentation, and for FH segmentation, the convolution operation will be the same. Still, the feature map and vertex map size will be different because of different input size and number of contours of instances.

3.3 Graph Vertices Sampling

To achieve multi-level aggregated graph convolutions on different vertex resolutions, we follow [36] to form a new topology and neighbour relationships of vertices. More specifically, we use the permutation matrix $Q_d \in \{0, 1\}^{m \times n}$ to down-sample m vertices, $m = 360$ or 720 in our work. Q_d is gained by iteratively

decreasing vertices, which uses a quadratic matrix to keep the approximations of the surface error [17]. The down-sampling is a pre-processing, and the discarded vertices are saved with barycentric coordinates. We conduct up-sampling with another transformation matrix $Q_u \in \mathbb{R}^{m \times n}$. The up-sampled vertices V_u can be obtained by a sparse matrix multiplication, i.e., $V_u = Q_u V_d$, where V_d are down-sampled vertices.

3.4 Proposed Aggregation Network

Our novel aggregation graph regression network is motivated by fusing features hierarchically and iteratively [47,53,41], which consists of an image context encoder, an attention refinement module and a vertex location decoder. Both the encoder and decoder contain aggregation mechanisms through up-samplings and down-samplings, which provide improvements in extracting the full spectrum of semantic and spatial information across stages and resolutions. Besides, the attention module plays an essential role to guide the feature learning and refine the output from the CNN encoder, then passes to the GCN decoder through multi-paths. In Section 5.3, our ablation study demonstrates that the proposed aggregation module helps to extract more useful information, and the attention module helps to refine the extracted features from the encoder to guide feature learning better.

Semantic Encoder Fig. 2 (a) shows the detailed structure of our semantic encoder, which maintains high-resolution representations by connecting low-to-high resolution convolutions in parallel, where multi-scale fusions are repeated across different levels (rows). Our encoder is designed to lessen the spatial information loss and extract a wider spectrum of semantic features through different receptive fields. The encoder takes input images of shape $314 \times 314 \times 3$ (Fundus OD & OC images) or $140 \times 140 \times 1$ (Ultrasound FH images), with operations of up-sampling and down-sampling. The aggregation block can extract and reuse more features across various resolutions and scales, which helps to reduce spatial information loss during the encoding process.

Attention Module: We propose an Attention Refinement Module (ARM) to refine the features from the outputs of the encoder. As Fig. 2 (a) & (b) shows, ARM contains five attention blocks, and each block employs global average pooling to capture global context through the different channels, and conducts an attention tensor to lead the emphasis of feature learning through a convolution layer followed by a BN layer and sigmoid as the activation function. For the filter, the kernel size is 1×1 , and the stride is 1. This design can refine the output features of each stage in the Semantic Encoder, which easily integrates the global context information.

Spatial Decoder The decoder takes refined multi-paths outputs from the attention module, then employ ResGCN blocks [28] through different stages and levels, which has been shown that as layers go deeper, ResGCN blocks can prevent vanishing gradient problems. As Fig. 2 (b) shows, our decoder fuses and reuses the features extracted by ResGCN blocks through different stages. Benefits from the graph Vertices sampling, our decoder can regress the location of

the vertices from sparse to dense, which allows the ResGCN blocks to hierarchically extract spatial location information from refined outputs of the attention module. For each ResGCN Block, it consists of 4 graph convolution layers, and each graph convolution layer is followed by a Batch Normalization layer [26] and Leaky ReLU as the activation function. After ResGCN blocks and graph vertices up-samplings, the number of vertices is up-sampled from 25 to 721, and each vertex is represented by a vector of length 32. Different from [33], Our decoder further explored the relations between low and high level resolution of vertices features, which improves the performance and is shown in Section 4. At last, three graph convolution layers are added to generate 2D object contour vertices, which reduces the output channels to 2, as each contour vertex has two dimensions: x and y . With the output from the decoder, we connect every two consecutive vertices on the boundary to form a polygon contour as the final segmentation result.

3.5 Loss Function

L2 and L1 loss have been widely used in regression tasks, such as object detection [19,23] and human pose estimation [42]. However, it is difficult for the L1 loss to find the global minimization in the late training stage without fine-tuning of the learning rate. L2 loss is sensitive to outliers which may result in unstable training in the early training stage.

In this work we solve segmentation as a contour vertices location regression problem. Following Wing-loss [13] and Smooth-L1 loss [18], we adopt a new loss function, Fan-loss (Fig. 3) that can take small update steps when reaching small range errors in the late training stage and can remain stable training during the early training stage. This loss function is defined as:

$$L(x) = \begin{cases} W[e^{|x|/\epsilon} - 1] & \text{if } |x| < W \\ |x| - C & \text{otherwise} \end{cases} \quad (3)$$

Where W is non-negative and decide the range of the non-linear part, ϵ limits the curvature between $(-W, W)$ and $C = W - W[e^{|w|/\epsilon} - 1]$ connects the linear and non-linear parts. After several evaluation experiments, the parameter W is set to 8 and ϵ to 5 for FH segmentation and $W = 6$, $\epsilon = 5$ for OD & OC segmentation. For the OD & OC segmentation tasks, we integrate a weight mask and assign more weights to the vertices that belong to the OC, as OC is usually difficult to segment because of poor image quality or low color contrast.

4 Experiments

4.1 Datasets

We evaluate our approach with two major types of biomedical images on two segmentation tasks respectively: fundus images of retinal for OD & OC segmentation, and ultrasound images of the fetus for FH segmentation.

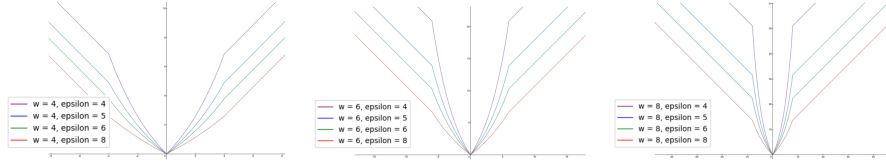


Fig. 3: The loss function plotted with different parameter settings, where w controls the non-linear part and epsilon (ϵ) limits the curvature.

Fudus OD & OC images: 2068 images from five datasets are merged together. 190 fundus images are randomly selected as the retina test dataset, the rest 1878 fundus images are used for the training. Considering the negative influence of non-target areas in fundus retina images, we first localize the disc centers by detector [37] and crop to 314×314 pixels and then transmit into our network. **Refuge** [35] consists of 400 training images and 400 validation images. The pixel-wise OD & OC gray-scale annotations are provided. **Drishti-GS** [40] contains 50 training images and 51 validation images. All images are taken centered on OD & OC with a field-of-view of 30 degrees. The annotations are provided in the form of average boundaries. **ORIGA** [49] contains 650 fundus images. The OD & OC boundaries were manually marked by experienced graders from the Singapore Eye Research Institute. **RIGA** [1] contains 750 fundus images from **MESSIDOR** [11] database. The OD and OC are labeled manually by six ophthalmologists and the mean OD and OC are used as the ground truth. **RIM-ONE** [16] contains 169 fundus images, annotated by five different experts. **Ultrasound FH images:** The HC18-Challenge dataset are used which contains 999 two-dimensional (2D) ultrasound images with size of 800×540 pixels collected from the database of Radboud University Medical Center [25]. We apply zero-padding to each image to 840×840 pixels, and then resize into 140×140 as the input image, then we randomly select 94 images as the test dataset, and the model is trained on the rest 905 images.

4.2 Implementation Details

To augment the dataset, we randomly rotating the input image of training dataset for both segmentation tasks. To be specific, the rotation ranges from -15 to 15 degree. We randomly select 10% of training dataset as the validation dataset. We use stochastic gradient descent with a momentum of 0.9 to optimize the Fan-loss. The number of graph vertices for FH is sampled to 361, 256, 128, 64, 32, 25 crosses five stages with Graph Vertices Sampling introduced in Section 3.3. We trained our model for 300 epochs for all the experiments, with a learning rate of $1e-2$ and decay rate of 0.997 every epoch. The batch size is set as 48. All the training processes are performed on a server with 8 TESLA V100 and 4 TESLA P100, and all the test experiments are conducted on a local workstation with Geforce RTX 2080Ti.

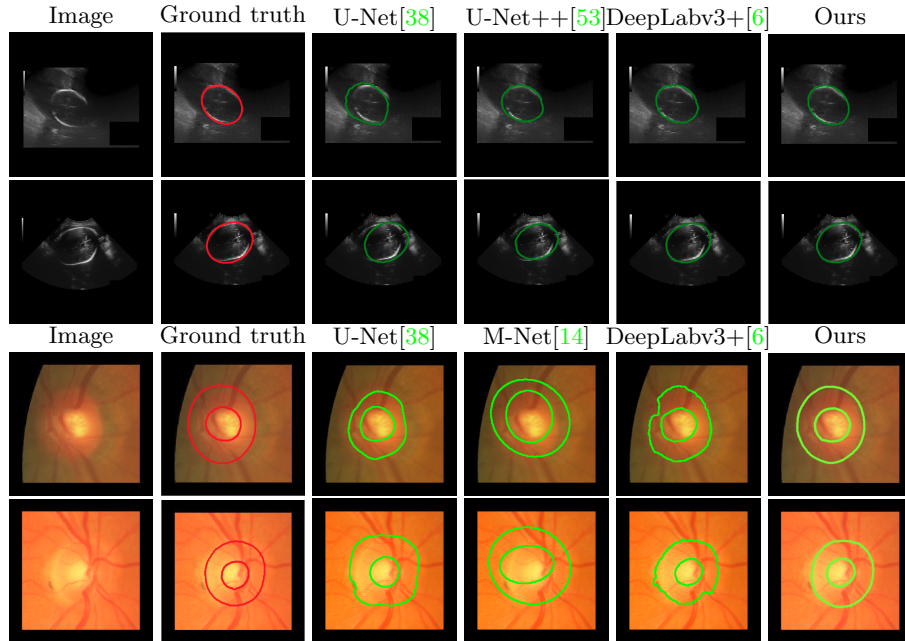


Fig. 4: Qualitative results of segmentation on the testing images of the fundus dataset and HC18-Challenge [25]. Top two rows are the ultrasound FH segmentation results, and the bottom two rows are the fundus OD & OC segmentation results.

5 Results

In this section, we present our experimental results on the OD & OC and FH segmentation task in comparison to other state-of-the-art methods. We compare our model with other state-of-the-art methods, including U-Net [38], PolarMask [43], M-Net [14], U-Net++ [53], DANet [15], DARNet [9], DeepLabv3+ [6], CGRNet [33] through running their open public source code. Dice score and Area Under the Curve (AUC) are used as the segmentation accuracy metrics. The results of an ablation study are shown in order to demonstrate the effectiveness of the proposed aggregation mechanism, attention mechanism and loss function, respectively.

5.1 Optic Disc & Cup Segmentation

The retinal dataset we used is merged from five different fundus OD & OC images datasets. In terms of different dataset sources, they may contain different annotation standards for ground truths by different doctors. However, our model still achieve good performance, which shows the robustness and generalizability of our model. Fig. 4 shows some qualitative results. We achieve 0.9697 and

Methods \ Tasks	OC		OD			FH		HD(mm)
	Dice Score	AUC	Dice Score	AUC	vCDR	Dice Score	AUC	
U-Net [38]	0.9016	0.9186	0.9522	0.9648	0.0674	0.9625	0.9688	1.79
M-Net[14]	0.9335	0.9417	0.9230	0.9332	0.0488	-	-	-
U-Net++ [53]	0.9198	0.9285	0.9626	0.9777	0.0469	0.9701	0.9789	1.73
DANet [15]	0.9232	0.9327	0.9654	0.9726	0.0450	0.9719	0.9786	1.69
DARNet [9]	0.9235	0.9339	0.9617	0.9684	0.0455	0.9719	0.9790	1.52
PolarMask [43]	0.9238	0.9366	0.9670	0.9782	0.0419	0.9723	0.9780	1.66
DeepLabv3+ [6]	0.9308	0.9406	0.9669	0.9779	0.0467	0.9779	0.9819	1.58
CGRNet [33]	0.9246	0.9376	0.9688	0.9784	0.0438	0.9738	0.9796	1.58
Our method	0.9255	0.9385	0.9697	0.9791	0.0421	0.9746	0.9801	1.47

Table 1: Segmentation results on retina test dataset for OD & OC and on HC18-Challenge [25] for FH. The performance is reported as Dice score (%), AUC (%), mean absolute error of Hausdorff distance (HD) for FH and mean absolute error of the vertical cup-to-disc ratio (vCDR) for OD & OC. The top three results in each category are highlighted in bold.

0.9255 Dice similarity score on OD & OC segmentation respectively, which are comparable with other pixel-wise based state-of-the-art methods even without any bells and whistles (e.g. multi-scale training, ellipse fitting, longer training epochs, etc.). Tab. 1 provides the results of ours and the other methods. As for the inference speed, our model uses 64.1 milliseconds (ms) per image that is faster than PolarMask [43] (72.1 ms) and DeepLabv3 [6] (323.9 ms). In the supplementary material, we also show some ‘failed’ cases compared with the ground truth. According to the comments from an anonymous expert at the Liverpool Reading Center, our model produces more accurate results than the ground truth. This highlights the potential issue of imperfect ground truth in many deep learning applications.

5.2 Fetal Head Segmentation

Tab. 1 and Fig. 4 shows the quantitative and qualitative results, our model achieves 0.9746 Dice similarity score and 0.9801 % AUC, which outperforms DARNet [9] and DANet [15] by 0.3%. Our model (59.1ms) is faster than PolarMask [43] (65.5 ms) and Deeplabv3+ [6] (290.3ms) for per image inference.

5.3 Ablation Study

We investigate the effect of each component in our proposed model. All the ablation experiments are performed with the same setting as section 4.2 described. The performance in the form of Dice score and AUC are reported in Fig. 5, Tab. 2 and 3. The best performance in each experiment is highlighted in bold. For more qualitative results, please refer to the supplementary material.

Ablation on Parameters of Loss Function We perform Experiments to evaluate the effect of parameter settings of Fan-loss function. When $w = 6$, ϵ

Tasks Loss Function	OC		OD		FH	
	Dice Score	AUC	Dice Score	AUC	Dice Score	AUC
L1	0.9111	0.9259	0.9546	0.9639	0.9505	0.9688
L2	0.9105	0.9210	0.9551	0.9666	0.9440	0.9568
Smooth-L1 [18]	0.9088	0.9114	0.9523	0.9655	0.9394	0.9454
Fan-Loss						
weight mask = 0	0.9184	0.9220	0.9618	0.9739		
weight mask = 3	0.9221	0.9337	0.9649	0.9769		
weight mask = 5	0.9255	0.9385	0.9697	0.9791	0.9746	0.9801
weight mask = 7	0.9175	0.9240	0.9624	0.9720		
weight mask = 9	0.9107	0.9213	0.9600	0.9705		

Table 2: Performance comparisons between different loss function and weight mask parameter settings on the OD & OC segmentation and the FH segmentation respectively. For weight mask = 5, our model achieves best performance on the OD & OC segmentation.

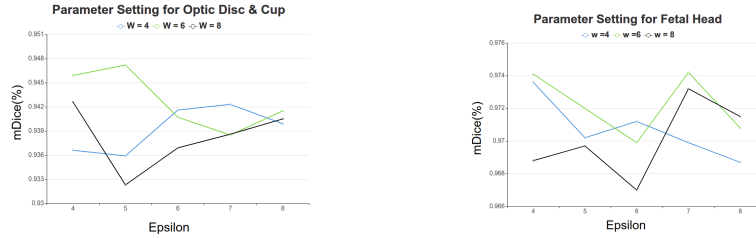


Fig. 5: A comparison of different parameter settings (w and ϵ) for Fan-loss function, measured in terms of the mean Dice score on the fundus dataset for OD & OC. With $w = 6$, $\epsilon = 5$, our model achieves the best performance (0.9255 & 0.9697). On the HC18-Challenge test dataset [25] for FH segmentation, with $w = 6$, $\epsilon = 7$, our model gains the best results 0.9746). It shows that our network is not sensitive to these parameters as no significantly different results are found.

= 5, our model achieve the best performance on OD & OC segmentation test dataset, and $w = 6$, $\epsilon = 7$, for FH segmentation test dataset. For more details, please refer to Fig. 5.

Ablation on Loss Function We conduct experiments to evaluate the effectiveness of the loss function. We compare with L1, L2, Smooth-L1 [18] loss functions, which are commonly used in the regression problem. Tab. 2 shows the quantitative results on OD & OC and FH segmentation tasks respectively. As illustrated, Fan-loss function attains a superior performance over the other three loss functions. In particular, it achieves a mean Dice score that is 1.6% relatively better than that of L1 loss function on OD & OC and 2.7% relatively better than L1 loss function on FH segmentation. Tab. 2 shows comparing with no-weight mask loss function, our proposed weight mask helps to improve OD & OC segmentation results by 0.79% when weight mask = 5 is used.

Tasks Methods	OC		OD		FH	
	Dice Score	AUC	Dice Score	AUC	Dice Score	AUC
No Aggregation (Encoder + Decoder)	0.9025	0.9065	0.9589	0.9665	0.9567	0.9690
Aggregation	0.9207	0.9303	0.9624	0.9660	0.9700	0.9776
Aggregation + ARM (with CNN decoder)	0.9099	0.9178	0.9529	0.9635	0.9639	0.9758
Aggregation + ARM (Our method)	0.9255	0.9385	0.9697	0.9791	0.9746	0.9801

Table 3: Ablation study on different structure components of the loss function ($w = 6$, $\epsilon = 5$ for FH segmentation and $w = 6$, $\epsilon = 7$ for OD & OC).

Ablation on Angle Interval Experiments are conducted to evaluate the effect of different angle intervals $\Delta\theta$ for vertices sampling. The larger angle interval indicates the smaller number of vertices sampled on the contour. With $\Delta\theta = 1^\circ$, our model achieves best performance on both the FH segmentation and the OD & OC segmentation. The results are shown in supplementary material.

Ablation on Structure Components In this section, we evaluate the effectiveness of our aggregation module, attention module and GCN decoder. First, we compare with no-aggregation structure network, in which we remove all the aggregation parts and attention modules to form a standard encoder-decoder network structure. Then we add aggregated CNN and GCN module to form an aggregation network. To further improve the performance, we design an attention module, and the effect of the attention module is presented in Tab 3. Furthermore, we evaluate the effectiveness of proposed GCN decoder and replace the GCN with CNN, which are the same as we used in the encoder. As illustrated, for the FH segmentation, the proposed aggregation module helps to improve 1.83% on Dice score over the no-aggregation method, the ARM module further improves 0.47%, and GCN decoder further improves 1.11%. For the OD & OC segmentation, the aggregation module improves 1.17 % on average by Dice score, the ARM improves 0.64%, and the GCN decoder improves 1.73%.

6 Conclusion

We propose a straightforward regression method for segmentation tasks by directly regressing the boundary of the instances instead of pixel-wise dense predictions. We have demonstrated its potentials on the segmentation problems of the fetal head and optic disc & cup. In the future work, we will study to extend the proposed model to tackle 3D biomedical image segmentation tasks.

Acknowledgement: Y. Meng thanks the China Science IntelliCloud Technology Co., Ltd for the studentship. D. Gao is supported by EPSRC Grant (EP/R014094/1). We thank NVIDIA for the donation of GPU cards. This work was undertaken on Barkla, part of the High Performance Computing facilities at the University of Liverpool, UK.

References

1. Almazroa, A., Alodhayb, S., Osman, E., Ramadan, E., Hummadi, M., Dlaim, M., Alkatee, M., Raahemifar, K., Lakshminarayanan, V.: Retinal fundus images for glaucoma analysis: the RIGA dataset. In: Medical Imaging 2018: Imaging Informatics for Healthcare, Research, and Applications. vol. 10579, p. 105790B. International Society for Optics and Photonics (2018)
2. Arnab, A., Zheng, S., Jayasumana, S., Romera-Paredes, B., Larsson, M., Kirillov, A., Savchynskyy, B., Rother, C., Kahl, F., Torr, P.H.: Conditional random fields meet deep neural networks for semantic segmentation: Combining probabilistic graphical models with deep learning for structured prediction. *IEEE Signal Processing Magazine* **35**(1), 37–52 (2018)
3. Bruna, J., Zaremba, W., Szlam, A., LeCun, Y.: Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203* (2013)
4. 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. *IEEE Transactions on Pattern Analysis and machine Intelligence* **40**(4), 834–848 (2017)
5. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017)
6. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: ECCV (2018)
7. Chen, X., Williams, B.M., Vallabhaneni, S.R., Czanner, G., Williams, R., Zheng, Y.: Learning active contour models for medical image segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 11632–11640 (2019)
8. Chen, Y., Zhao, D., Lv, L., Zhang, Q.: Multi-task learning for dangerous object detection in autonomous driving. *Information Sciences* **432**, 559–571 (2018)
9. Cheng, D., Liao, R., Fidler, S., Urtasun, R.: Darnet: Deep active ray network for building segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7431–7439 (2019)
10. Chung, F.R., Graham, F.C.: Spectral graph theory. No. 92, American Mathematical Soc. (1997)
11. Decencière, E., Zhang, X., Cazuguel, G., Lay, B., Cochener, B., Trone, C., Gain, P., Ordonez, R., Massin, P., Erginay, A., et al.: Feedback on a publicly distributed image database: the Messidor database. *Image Analysis & Stereology* **33**(3), 231–234 (2014)
12. Defferrard, M., Bresson, X., Vandergheynst, P.: Convolutional neural networks on graphs with fast localized spectral filtering. In: Advances in neural information processing systems. pp. 3844–3852 (2016)
13. Feng, Z.H., Kittler, J., Awais, M., Huber, P., Wu, X.J.: Wing loss for robust facial landmark localisation with convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2235–2245 (2018)
14. Fu, H., Cheng, J., Xu, Y., Wong, D.W.K., Liu, J., Cao, X.: Joint optic disc and cup segmentation based on multi-label deep network and polar transformation. *IEEE transactions on medical imaging* **37**(7), 1597–1605 (2018)
15. Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3146–3154 (2019)

16. Fumero, F., Alayón, S., Sanchez, J.L., Sigut, J., Gonzalez-Hernandez, M.: RIM-ONE: An open retinal image database for optic nerve evaluation. In: 2011 24th international symposium on computer-based medical systems (CBMS). pp. 1–6. IEEE (2011)
17. Garland, M., Heckbert, P.S.: Surface simplification using quadric error metrics. In: Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques. pp. 209–216. ACM Press/Addison-Wesley Publishing Co. (1997)
18. Girshick, R.: Fast R-CNN. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1440–1448 (2015)
19. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 580–587 (2014)
20. Gur, S., Shaharabany, T., Wolf, L.: End to end trainable active contours via differentiable rendering. arXiv preprint arXiv:1912.00367 (2019)
21. Gur, S., Wolf, L., Golgher, L., Blinder, P.: Unsupervised microvascular image segmentation using an active contours mimicking neural network. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 10722–10731 (2019)
22. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask R-CNN. In: Proceedings of the IEEE international Conference on Computer Vision. pp. 2961–2969 (2017)
23. He, K., Zhang, X., Ren, S., Sun, J.: Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **37**(9), 1904–1916 (2015)
24. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern recognition. pp. 770–778 (2016)
25. van den Heuvel, T.L., de Bruijn, D., de Korte, C.L., van Ginneken, B.: Automated measurement of fetal head circumference using 2D ultrasound images. *PloS one* **13**(8) (2018)
26. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167 (2015)
27. Kass, M., Witkin, A., Terzopoulos, D.: Snakes: Active contour models. *International journal of computer vision* **1**(4), 321–331 (1988)
28. Li, G., Müller, M., Thabet, A., Ghanem, B.: Can GCNs go as deep as CNNs? arXiv preprint arXiv:1904.03751 (2019)
29. Li, Y., Gupta, A.: Beyond grids: Learning graph representations for visual recognition. In: Advances in Neural Information Processing Systems. pp. 9225–9235 (2018)
30. Liu, Z., Li, X., Luo, P., Loy, C.C., Tang, X.: Semantic image segmentation via deep parsing network. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1377–1385 (2015)
31. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3431–3440 (2015)
32. Marcos, D., Tuia, D., Kellenberger, B., Zhang, L., Bai, M., Liao, R., Urtasun, R.: Learning deep structured active contours end-to-end. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8877–8885 (2018)
33. Meng, Y., Wei, M., Gao, D., Zhao, Y., Yang, X., Huang, X., Zheng, Y.: CNN-GCN aggregation enabled boundary regression for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. p. in press (2020)

34. Mou, L., Zhao, Y., Chen, L., Cheng, J., Gu, Z., Hao, H., Qi, H., Zheng, Y., Frangi, A., Liu, J.: CS-Net: Channel and spatial attention network for curvilinear structure segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 721–730. Springer (2019)
35. Orlando, J.I., Fu, H., Breda, J.B., van Keer, K., Bathula, D.R., Diaz-Pinto, A., Fang, R., Heng, P.A., Kim, J., Lee, J., et al.: REFUGE Challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs. *Medical Image Analysis* **59**, 101570 (2020)
36. Ranjan, A., Bolkart, T., Sanyal, S., Black, M.J.: Generating 3D faces using convolutional mesh autoencoders. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 704–720 (2018)
37. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: Advances in Neural Information Processing Systems. pp. 91–99 (2015)
38. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 234–241. Springer (2015)
39. Shin, S.Y., Lee, S., Yun, I.D., Lee, K.M.: Deep vessel segmentation by learning graphical connectivity. *Medical Image Analysis* **58**, 101556 (2019)
40. Sivaswamy, J., Krishnadas, S., Joshi, G.D., Jain, M., Tabish, A.U.S.: Drishti-GS: Retinal image dataset for optic nerve head (ONH) segmentation. In: 2014 IEEE 11th international symposium on biomedical imaging (ISBI). pp. 53–56. IEEE (2014)
41. Sun, K., Zhao, Y., Jiang, B., Cheng, T., Xiao, B., Liu, D., Mu, Y., Wang, X., Liu, W., Wang, J.: High-resolution representations for labeling pixels and regions. arXiv preprint arXiv:1904.04514 (2019)
42. Toshev, A., Szegedy, C.: Deeppose: Human pose estimation via deep neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1653–1660 (2014)
43. Xie, E., Sun, P., Song, X., Wang, W., Liu, X., Liang, D., Shen, C., Luo, P.: Polar-mask: Single shot instance segmentation with polar representation. arXiv preprint arXiv:1909.13226 (2019)
44. Xu, W., Wang, H., Qi, F., Lu, C.: Explicit shape encoding for real-time instance segmentation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 5168–5177 (2019)
45. Yu, C., Wang, J., Peng, C., Gao, C., Yu, G., Sang, N.: Bisenet: Bilateral segmentation network for real-time semantic segmentation. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 325–341 (2018)
46. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122 (2015)
47. Yu, F., Wang, D., Shelhamer, E., Darrell, T.: Deep layer aggregation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2403–2412 (2018)
48. Zhang, L., Li, X., Arnab, A., Yang, K., Tong, Y., Torr, P.H.: Dual graph convolutional network for semantic segmentation. arXiv preprint arXiv:1909.06121 (2019)
49. Zhang, Z., Yin, F.S., Liu, J., Wong, W.K., Tan, N.M., Lee, B.H., Cheng, J., Wong, T.Y.: ORIGA-light: An online retinal fundus image database for glaucoma analysis and research. In: 2010 Annual International Conference of the IEEE Engineering in Medicine and Biology. pp. 3065–3068. IEEE (2010)

50. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2881–2890 (2017)
51. Zhao, H., Zhang, Y., Liu, S., Shi, J., Change Loy, C., Lin, D., Jia, J.: PSANet: Point-wise spatial attention network for scene parsing. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 267–283 (2018)
52. Zhao, L., Peng, X., Tian, Y., Kapadia, M., Metaxas, D.N.: Semantic graph convolutional networks for 3D human pose regression. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3425–3435 (2019)
53. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: UNet++: A nested U-net architecture for medical image segmentation. In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support, pp. 3–11. Springer (2018)