

Waseda University Master Thesis

**Cascaded Fully Convolutional Networks for Object Boundary
Detection**

44161542-3: YUAN JIANG

Master (Engineering)

Professor Jinglu HU, Supervisor

Neurocomputing System
Information Architecture

The Graduate School of Information, Production and Systems

February 2018

ABSTRACT

In computer vision, object boundary is defined as the enclosed set of edge pixels by which an object is tightly surrounded. Given the object boundary, the shapes and locations of the objects in an image can be easily recognized. However, object boundary detection is still a challenging problem due to the multi-scale object problem, the interference caused by the local edges with less semantic information, and the large variance of backgrounds and scenes. In this thesis, we develop a cascaded framework to improve the performance of the present edge detection algorithms in object boundary detection. Within our framework, present networks can be cascaded one-by-one and trained end-to-end to gain more semantic information from the inputs. We test our method in two scenes: neuronal boundary detection in Electron Microscopy (*EM*) images, and object boundary detection in natural images. Massive experiments and analyses show that the proposed cascaded fully convolutional networks can exactly outperform the competitors with the help of the cascaded structure.

Keywords: Fully Convolutional Networks, Object Boundary Detection, Neuronal Segmentation, Cascaded Structures

ACKNOWLEDGMENTS

Contents

Table of contents	i
List of figures	iii
List of tables	v
1 Introduction	1
1.1 Background and Motivation	1
1.2 Application Scenes	3
1.2.1 EM Images	3
1.2.2 Natural Images	4
1.3 Organization of the thesis	5
2 Related Work	7
2.1 Traditional edge detection	7
2.1.1 Hand-crafted feature based methods	8
2.1.2 Deep learning based methods	9
2.2 Object Boundary Detection in Natural Images	10
2.3 Neuronal Boundary Detection in EM Images	11
3 Methodology	15
3.1 Network Structure	15
3.2 Training Phase	17
3.2.1 Formulation of the Multi-recursive-input	18
3.2.2 Loss Function	19
3.3 Testing Phase	21
3.4 Model Interpretation	21
3.4.1 Cascaded Architecture <i>vs.</i> Single-stage Architecture	21
3.4.2 Multi-recursive-input <i>vs.</i> Single-recursive-input	23
3.4.3 End-to-end Training <i>vs.</i> Stepwise Self-tuning	23

4 Experiments	25
4.1 Neuronal Boundary Detection in EM Images	25
4.1.1 Evaluation Metric	26
4.1.2 Mouse Piriform Cortex Dataset	26
4.1.3 ISBI 2012 EM Segmentation Dataset	29
4.2 Object Boundary Detection in Natural Images	32
4.2.1 Metrics	32
4.2.2 PASCAL VOC Contour Dataset	32
4.3 Control Experiments for Model Interpretation	35
4.3.1 Evaluations on Cascaded Architecture vs. Single-stage Architecture	36
4.3.2 Evaluations on Multi-recursive-input vs. Single-recursive-input	37
4.3.3 Evaluations on End-to-end Training vs. Stepwise Self-tuning	38
5 Conclusions	41
Bibliography	43

List of Figures

1.1	Object-level boundary detection is different from the traditional local edge detection, where the former mainly focuses on detecting the high-level semantic boundaries. The first image contains the annotations, which are marked with red lines, for object-level boundary detection.	2
2.1	Illustration of VD2D3D, a recursive deep network with the stepwise training. . . .	12
3.1	Three types of cascade-like networks: (a) cascaded network with single-recursive input; (b) cascaded network with multi-recursive inputs; (c) single-stage network with recursive inputs to fine-tune itself.	16
3.2	Illustration of the proposed cascaded fully convolutional network for object boundary detection.	18
3.3	Multi-recursive-input <i>vs.</i> Single-recursive-input.	22
3.4	Prediction exmaples from different side-outputs and fused outputs in different stages. .	24

4.1 (b)Neuronal boundary detection and (c) segmentation can be easily converted from each other, which we used in the experiments on Mouse Piriform Cortex Dataset[3] and ISBI 2012 EM Segmentation Dataset[5]. (1) By applying the graph-based algorithms such as watershed, we can transfer the boundary prediction into segmentation. (2) By calculating the 2D gradient in the segmentation ground-truth, the boundary annotation can be obtained.	27
4.2 Precision (rand merge)-recall (rand split) curves on Mouse Piriform Cortex Dataset[3]. Our 3-stage network outperforms all the previous works on this dataset.	30
4.3 Examples selected from PASCAL VOC Contour Dataset[1], where ground-truth contours are labeled as red. There are various objects (human, artificialities, animals, plants, etc.) appearing in complex scenes (indoor and outdoor environments, colorful backgrounds, blurs, textural confusions, etc), which significantly increases the difficulty of detecting the object-level boundary.	31
4.4 Precision-recall curves on PASCAL VOC Contour Dataset[1].	35
4.5 Qualitative examples to reveal how cascaded networks improve the high-level boundary detection result. For example, the arrows mean the false positive detections are removed step by step.	37

List of Tables

4.1	Mouse Piriform Cortex Dataset	27
4.2	Rand F-scores on Mouse Piriform Cortex Dataset[3]	29
4.3	The Rand F-scores part from the leaderboard of ISBI 2012 EM Segmentation Challenge[5].	31
4.4	Object boundary detection evaluation comparison on PASCAL VOC Contour Dataset[1]. Our proposed 3-stage cascaded fully convolutional network achieves the new state-of-the-art on this benchmark, with a significant improvement (around 2% over the second)).	34
4.5	Control Experiment 1: Cascaded Architecture <i>vs.</i> Single-stage Architecture	36
4.6	Control Experiment 2: Multi-recursive-input <i>vs.</i> Single-recursive-input	38
4.7	Control Experiment 3: End-to-end Training <i>vs.</i> Stepwise Self-tuning	38

Chapter 1

Introduction

1.1 Background and Motivation

Object boundary detection aims at finding the enclosed boundary in which contains an object. Comparing to traditional edge detection, object boundary detection only considers the boundaries related to objects, despite of the low-level edges which have less semantic information (Fig. 1.1). So the object boundary can be easily applied for instance-level object segmentation, object proposal[1, 2], biomedical engineering[3, 4, 5], etc. While it is important and useful in various vision tasks, there are few literature about the object boundary detection due to three challenges: First, object boundary detection is a multi-scale problem; Second, amount of local edges with less semantic information must be inhibited to gain a high-level boundary map; And the large variance of backgrounds and scenes.

Traditional methods mostly focus on extract single-scale features for edge detection[6, 7, 8, 9]. However, objects has various scales in many applications, resulting in the demand of multi-scale descriptors to detect object boundaries. HED[10] is the first to adopt the supervised multi-scale feature learning to edge detection, where the feature map from each stage of the convolution neural networks is supervised by the ground-truth labels. We follow the multi-scale feature extracting

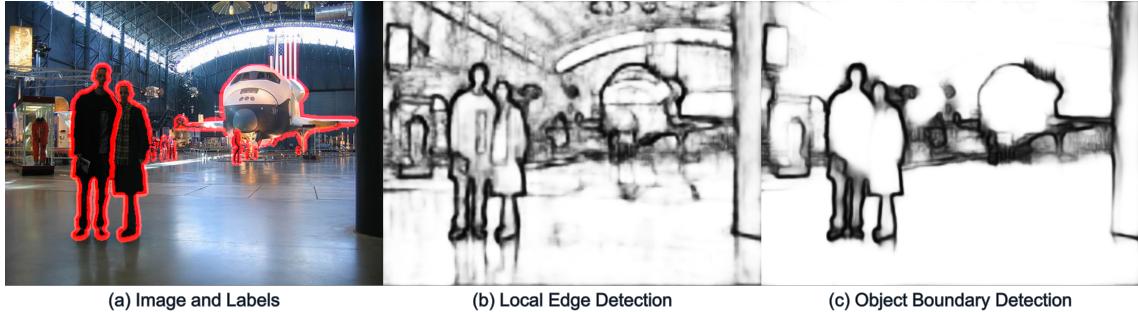


Figure 1.1 Object-level boundary detection is different from the traditional local edge detection, where the former mainly focuses on detecting the high-level semantic boundaries. The first image contains the annotations, which are marked with red lines, for object-level boundary detection.

proposed by HED[10] and extend it to object boundary detection, in which the multi-scale problem is considerably serious.

What's more, the target of traditional edge detection is to identify the discontinuities of image brightness, regardless of the semantic information which indicate whether the edges enclose an object in it. This is the root why some traditional edge detection methods fail to detect the object boundary. We propose the cascaded fully convolutional networks (*FCNs*)[11] with holistic training and testing to overcome the shortages of traditional edge detectors. By recursively connect the FCNs rather than simply stacked the convolution layers, our networks have the ability to inhibit the meaningless edges inside the objects and can be easily interpreted.

Another issue for object boundary detection is the large variance of backgrounds and scenes. Our networks draw inspiration from the related works about neuronal boundary detection in Electron Microscopy (*EM*) images, so we firstly make an effort to apply them on two mainstream datasets for neuronal boundary detection and neuronal segmentation in EM images. To extend the cascaded framework into common cases, we then test our networks in a large natural image dataset, PASCAL Contour[1].

We highlight three points in the thesis. (1) extending the multi-scale holistic network of HED[10] to the multi-scale multi-recursive network for object boundary detection. (2) cascaded

network structure with end-to-end training, which is easy to be integrated with other fully convolutional networks. (3) and systematic experiments in two applications to show the reliability and interpretability of our architecture.

1.2 Application Scenes

1.2.1 EM Images

With the development of Electron Microscopy (*EM*), neuroscientists and cognitive scientists are able to study the connection between structures and functions of neurons from the high-resolution EM images[12]. The first step for the study is often to reconstruct the structure of neuronal circuits. However it would take thousands of hours for an expert to annotate a large number of EM images, which motivates researchers to explore some automatic reconstructing algorithms. To reconstruct the neuronal circuits, we firstly reconstruct the 2D structure of neurons from a serials of EM images produced by a serial section EM. Then the stack of 2D neuron structure will be concatenated into the 3D neuron structure.

As the basic step, 2D neuron reconstruction has drawn a lot of attentions, especially those adopting the powerful deep neural networks[13, 3, 4]. All these approaches convert the reconstruction problem into a 2D boundary detection issue, since the later is easier to be modeled as a binary classification problem and thus there are many literatures and tools can be used to solve it. So long as the neuron boundary is decided, the neuronal circuit can be segmented using simple graph-based methods such as the watershed algorithm[14].

Neuronal boundary detection in EM images is a special case of object boundary detection. There are several commons between neuronal boundary detection and object boundary detection, for example, both them focus on the high-level boundary, with less consideration about the low-level local edges. In neuronal boundary detection, to be specific, the low-level local edge includes the ones from background and EM noises, and also the ones produced by confounding structures

like nucleus, mitochondria, etc. Hence, traditional edge detectors[10, 6, 7] often fail to detect the boundary in EM images (Fig. 1.1). Deep neural networks are also widely used to detect neuronal boundary and have gained promising progress[13, 3], while lack of Multi-recursive-input vs. Single-recursive-input. Different from all these works, we use a lot of experiments not only to show the performance of the proposed cascaded fully convolutional networks in neuronal boundary detection, but also to explain how the proposed architecture removes the insider local edges and background noises, and strengthens the boundary with low-contrast.

We compare our method with competitors on two mainstream datasets for neuronal segmentation, the Piriform dataset[3] and the ISBI 2012 Challenge dataset[5].

1.2.2 Natural Images

Although our idea is firstly motivated by the problem we met in neuronal boundary detection, we also take lots of attentions in extending the proposed models into more commonly used applications, such as the object boundary detection in natural images. Given the object boundary map of a natural image, we can easily obtain the object proposal, object segmentation[15, 16], object detection[17, 18], and recognize object using shape based classification methods[19, 20, 21].

Object boundary detection in natural images is still a difficult problem due to the fact that objects vary in colors, brightness, scales, shapes, textures and gestures. One can hardly recognize the boundaries of all kinds of objects in natural images with only one of the features. Thus, many traditional edge detection methods[6, 7] do not perform well when directly applying them for high-level object boundary detection.

Performance benchmark is done on the famous PASCAL VOC dataset[17]. PASCAL VOC is one of the most famous object detection and segmentation challenge held once a year. We collect the training and validation data from PASCAL VOC 2007, to PASCAL VOC 2012. Ground-truth boundary maps are obtained from the instance-level segmentation annotations provided by the official. Comparison results show our cascaded network outperforms not only the traditional edge

detectors but also the deep networks without cascade.

1.3 Organization of the thesis

The rest of this thesis is structured as follows: Chapter 2 introduces some related works in traditional edge detection methods and neuronal boundary detection methods, concluding their progresses and shortages. In Chapter 3, we describe the proposed cascaded fully convolutional networks, giving an insight into why we design such a architecture and how it works. To evaluate the performance and verify the interpretation, Chapter 4 tests our models on three datasets, two for neuronal boundary detection / segmentation in EM images, and one for object boundary detection in natural images. The thesis ends in Chapter 5 with conclusions and future works of this research.

Chapter 2

Related Work

Comparing with traditional edge detection, object boundary detection focuses on the high-level semantic information about the objects appeared in the image. It is useful in many applications and has many advantages to the local edges, however, object boundary still lacks of researches in literature. So we will start from traditional edge detection methods to seek the recommendation. On the other hand, neuronal boundary detection in EM images has achieved considerable progresses. Deep neural networks perform well in the neuronal boundary detection, inspiring us to explore a more powerful and explainable one to improve the boundary detection results for not only the EM images, but also the natural images.

2.1 Traditional edge detection

Edge detection is a fundamental task with a long history in computer vision. Contrast to object boundary detection, local edge detection concerns all the pixels with color, brightness, texture, and other properties changing sharply in local parts of an image.

2.1.1 Hand-crafted feature based methods

Canny[6] can be seen as the pioneer of computationally extracting edges from digit images. Canny uses image gradient to get the edges. Firstly, the input RGB image should be converted to a greyscale one, whose 2D gradient map will be next calculated using the Sobel[7] filters. Then the non-maximum suppression is applied in the gradient map to get the thinned edges. After that Canny proposed to adopt an empirically selected double threshold to filtering the edges, where the edges higher than the high threshold will be marked as the "strong edges", the edges lower than the low threshold will be marked as background, and the rest of edges are so called the "weak edges". Canny descriptor use only the brightness gradient of the image, but has almost no ability to utilize the color and texture information of the image, much less the semantic information.

The followers [22, 15, 23] tend to combine brightness, color and texture gradients to make a better use of RGB inputs. All these methods are patch-based learning methods. In training phase, random sampled mini patches are cropped from original input image. Features of these patches are extracted using various cues mentioned above. Then the features of training samples will be fed into a two-class classifier, such as Support Vector Machine and Random Forest, to judge whether the central pixel of the sample is an edge pixel or not. In test phase, a one-step sliding window will be applied for a dense probability map about edge. To detect edges with more semantic meaning, the *gPb*[15] computes gradients at different scales, but resulting in a huge increase of computational load. Lim et al.[8] rather propose a structured feature learning pipeline for an effective and fast detection.

All these methods choose hand-crafted features based on local cues, with few design for object-level boundary detection. Experiments in Chapter 4 will show their shortages in object boundary detection tasks.

2.1.2 Deep learning based methods

Deep Learning, especially the Convolutional Neural Networks (*CNNs*), developed rapidly in recent years. The CNNs showed their power firstly in the large-scale object recognition challenges. AlexNet, proposed by Krizhevsky et al.[24] in 2012, improved the top-5 rate in ILSVRC-2010 and ILSVRC-2012 competition by around 8%, which drives researchers' trying Arbelaez2011 on applying the deep feature extraction for most of the computer vision tasks. VGGNet[25], GoogLeNet[26] and ResNet[27] are the followers of [24]. All these works tend to find the ways to train a deeper or wider network so that the abstract object-level feature can be better learned and extracted. In the next years, CNNs are also widely applied for object detection[28, 18, 29], object segmentation[16], and object contour detection[9, 30, 31].

Similar to [32], Shen et al.[9] cluster the edge structures into several subclasses and convert the two-class classification task to a multi-class one. By this way, the deep networks are able to learn different model parameters for each subclass, and the results are more interpretable. [9] utilizes the deep features to capture high-level object information and outperforms all the traditional hand-crafted feature based methods. However, in testing phase, the features extracted by fully connected layer must be fed into a structured random forest classifier[8] for the final contour detection result. The two-step learning limits the performance of deep networks. What's more, due to the existing fully connected layers, inputting the original images with various scales is forbidden. So the authors[9] have to adopt the patch-based training and testing, leading to a low detection speed.

In 2014, Fully Convolutional Networks (*FCNs*)[11] are proposed to handle with the per-pixel annotation problem, such as the object segmentation in the famous PASCAL VOC Challenge[17]. Developed from VGGNet, FCNs remove the last fully connected layers of it and replace with 1×1 convolution layers, which releases the limitation of the input size. A non-parameter deconvolution layer is followed by the last convolution layer to up-sample the low-resolution feature maps into probability maps with same resolution of the input. Finally, after properly cropping, the output probability maps will have the same size with the input. When training, these probability maps

can be used to calculate the pixel-wise loss with ground-truth labels. FCNs enable the end-to-end training and testing of CNNs for such a pixel-wise prediction task. With FCNs, the precision and speed will not be restricted by the sliding window scheme.

Features from different stages of the convolution neural networks have the multi-scale information[10]. Xie et al.[10] propose to make a full use of these information by extracting features from each stage of the FCNs, then up-sampling them to the same size with the input, which results in multi-scale outputs called side-outputs. All the side-outputs and the fusion of these side-outputs will contribute on the loss, where both the low-level stages with high resolution but little semantic information, and the high-level stages with low resolution but much semantic information can be trained easier. Benefited from the holistic training and multi-scale features, HED[10] achieves great performance in the famous edge detection benchmark, BSDS 500[15], improving the F-measure[15] from 0.756[9] to 0.782. Note that the F-measure of human being's annotation is 0.8.

2.2 Object Boundary Detection in Natural Images

Although HED[10] and other edge detection methods[9, 8, 15] have achieved considerable success in relatively low-level edge detection, they perform bad when applied for the high-level "*object-only*" boundary detection[1]. Fully Convolutional Encoder-Decoder Network (*CEDN*)[1] is the pioneer to distinguish object-level boundary detection from traditional local-level edge detection. CEDN uses VGG-16[25] as the encoder and designs an approximately symmetric but light-weighted decoder. During training, it fixes the model parameters of the encoder and only updates the parameters of the decoder.

CEDN outperforms all the competitors on the large PASCAL VOC Contour Dataset[1] in object boundary detection benchmark. However, it still has several problems: (1) the encoder network does not benefit from the training data at all; (2) to achieve the performance in the leaderboard, it requires all the training data to be processed for 30 times, due to the comparatively deep architec-

ture.

2.3 Neuronal Boundary Detection in EM Images

Neuronal segmentation is the fundamental step to learn the structures of neural cells and analyze the relationship between the structures and the functions[12, 33, 34]. In early days, however, it would take an expert with professional knowledge lots of time to recognize the segment of neuronal circuits[33] in massive EM images, which motivated researchers to find a automatic algorithm to annotate the neuronal segments. People find the key to the neuronal segmentation is to detect the membranes of the cells, while suppressing the local edges within them[35, 36, 37, 38]. So that graph-based algorithms such as watershed[39, 40, 14] and graph cut[41] can be directly used to convert the boundary to segmentation.

To detect neuronal boundary is to extract membranes in an EM image. In this degree, Neuronal Boundary Detection can be seen as a special case of Object Boundary Detection, where the "*object*" in the former is only the neuron circuit.

Deep learning based methods are able to meet the demand of different applications, by learning the parameters of not only the classifiers, but also the feature extractors. With the rapid development of deep networks, adopting them into the hard neuronal segmentation tasks becomes popular among researchers. Comparing with traditional hand-designed features, deep features need less experts' knowledge to design the feature extractors and can make a full use of the data, which is proved to be effective in neuronal boundary detection tasks[13, 5, 4, 42, 3]. As the winner of the famous ISBI 2012 EM Segmentation Challenge[43], Ciresan et al.[13] successfully apply the deep neural networks with a stack of *convolution-subsampling* units for membrane pixel detection. Ronneberger et al.[5] propose U-net, which reuses the features of the layers in the encoder part by merging them with the layers in the decoder part, to detect neuronal boundary and segment EM circuits. With the benefit of fully convolutional network structure and multi-level feature con-

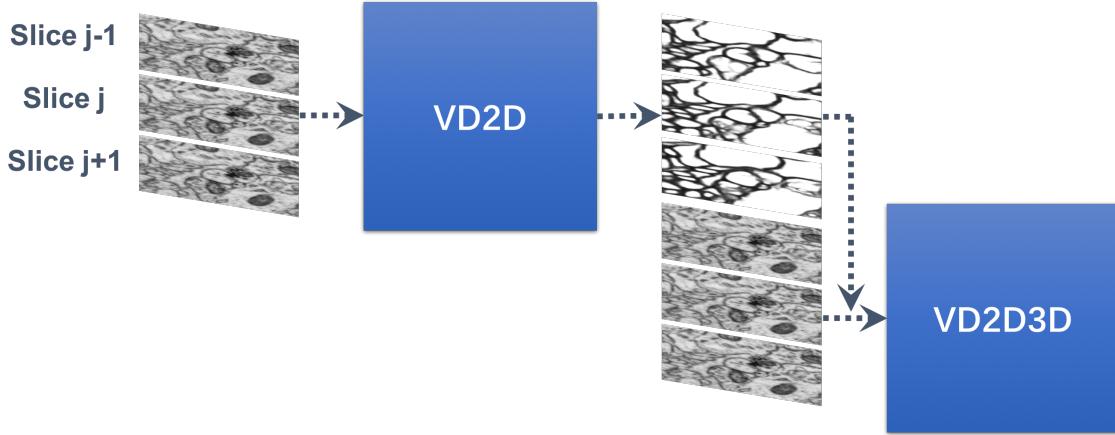


Figure 2.1 Illustration of VD2D3D, a recursive deep network with the stepwise training.

catenating, U-net shows the efficiency in neuronal segmentation by outperforming the competitors in ISBI 2012 by a large margin. With the proposal of ResNet[27] in object recognition and detection, Quan et al.[4] and Fakhry et al.[42] try to implement the deep residual units in their fully convolutional networks, which is also proved to be effective in neuronal boundary detection issues.

However, most of these works engraft the state-of-the-art network architectures proposed in other vision works as a "*black-box*" to boost the performance of neuronal boundary detection. It is necessary to seek an interpretable and effective approach, so that it can be theoretically comprehensive and enlighten the relative works in the future.

Most related to our work is Very Deep 2D-3D network (*VD2D3D*) proposed by Lee et al.[3]. A Very Deep 2D (*VD2D*) network is firstly trained using patches randomly cropped from the images. The model of *VD2D* will be then used to produce the predictions of a stack of images produced by the serial section EM. [3] believes the adjacent slices will help to detection the boundary, so *VD2D3D* concatenates the raw image and the prediction of its adjacent slices, produced by *VD2D*, as the new training materials of *VD2D* and recursively fine-tunes last several layers of *VD2D*. We find in experiments that [3] has three main problems: (1) it lacks of multi-scale features to better represent the complex neural membranes; (2) stepwise training limits the learning ability of the recursive deep networks; (3) patch-based testing relies on the sliding window pipeline, which

seriously slow down the testing speed. Inspired by [3], we develop a cascaded fully convolutional network with multi-recursive inputs and train it end-to-end. The proposed network overcomes the shortages of [3] to achieve the state-of-the-art performance on a neuronal segmentation benchmark[3] and a large natural object boundary detection benchmark[5].

Chapter 3

Methodology

In this part, we will introduce the proposed Cascaded Fully Convolutional Networks for object boundary detection. First, the cascaded network structure with multi-recursive inputs is presented. Then we will describe our end-to-end training, where every part of the cascaded network can benefit from the training data. Next, in Section 3.3, we will demonstrate how to apply the learned fully convolutional model for the accurate yet fast detection. In the last section of this chapter, we spend a lot of ink to lead you into the insight of the proposed architecture, explaining how it works.

3.1 Network Structure

As shown in Fig. 3.1, we considered several methods to assemble our cascaded networks. Fig. 3.1(a) shows the structure of single-recursive-input cascaded fully convolutional networks, which feeds the final output of the former sub-network into the next sub-network, recursively. Each of these sub-networks is the fully convolutional network with the same structure but unshared model parameters. Fig. 3.1(b) illustrates the cascaded structure with multi-recursive inputs, our final choice. In this design, the predictions from not only the last layer but also the side-output layers[10] will be the recursive inputs of the next stage, so we name it as multi-recursive-input design. With

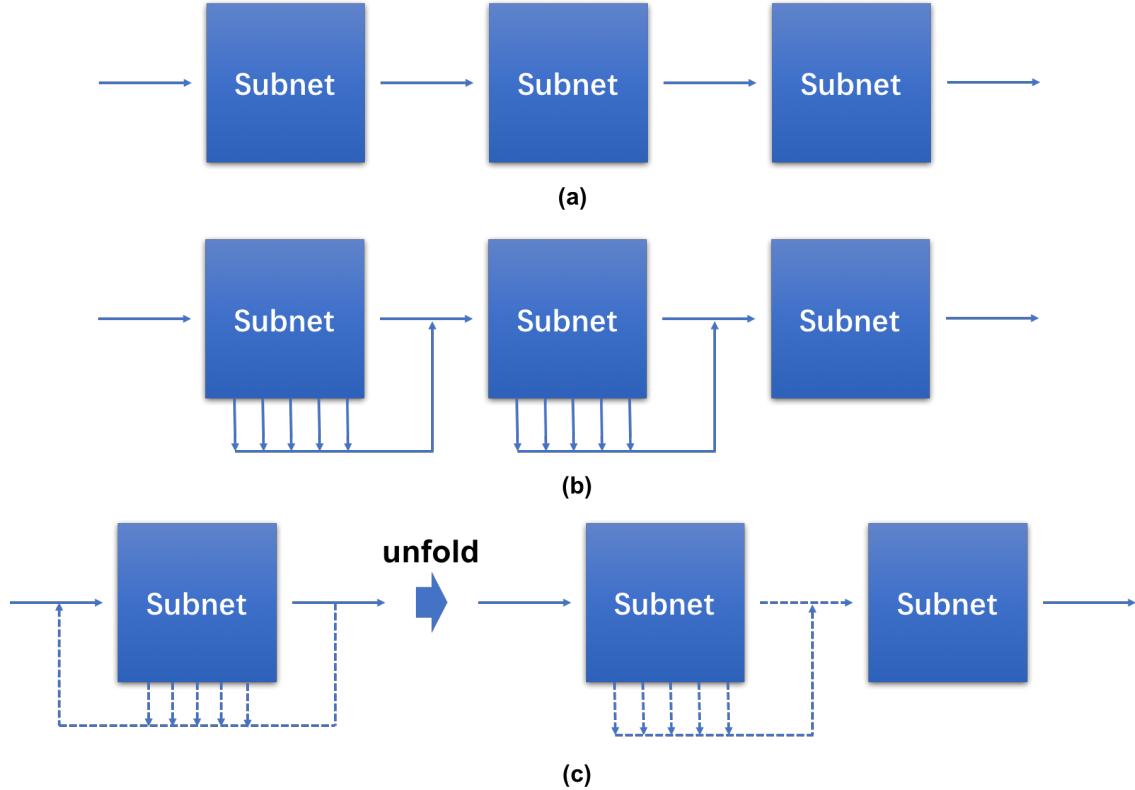


Figure 3.1 Three types of cascade-like networks: (a) cascaded network with single-recursive input; (b) cascaded network with multi-recursive inputs; (c) single-stage network with recursive inputs to fine-tune itself.

contrary to the single-recursive-input one, the proposed architecture makes the full use of features with multiple scales extracted from the former sub-network. Fig. 3.1(c) is another design for constructing such a multi-stage network, where there is only one sub-network, but fine-tuned time after time by recursively feeding the predictions into itself. To unfold the recursive connection of the left side of Fig. 3.1, we can get the equivalent network as shown in the right side. The unfolded version is much similar to the proposed architecture. However, we train all the parameters of sub-networks jointly in the proposed cascaded network, while Fig. 3.1(c) fixes the parameters of former sub-networks and only fine-tune the ones of the last sub-networks. Experiments in Chapter 4 show our design is more effective and efficient than the others.

To be more specific, our network used in object boundary detection can be illustrated as Fig.

3.2. It consists of several stages, three in the default setting, each of which is a fully convolutional network. By default, each sub-network is built based on the holistically-nested network[10], which is converted from the famous VGG-16[25] network. So the subnet has five stages of the convolution-pooling structure, of which the receptive fields are increased with the stepwise dilated strides. For each level, the side-output is obtained by up-sampling the feature from the last convolution layer of the level and activating it using a sigmoid layer. We mark the side-outputs of level 1 to 5 as $S_1^{(p)}, S_2^{(p)}, S_3^{(p)}, S_4^{(p)}, S_5^{(p)}$, respectively, where $p \in 1, 2, 3$ denotes the order of the sub-network.

Benefited from the receptive field sizes increasing, side-outputs are able to capture features in different scales naturally, without input scaling. Multi-scale features are then fused using the 1×1 convolution filters to produce the final predictions for the stage 1, 2, 3, named $F^{(1)}, F^{(2)}, F^{(3)}$, respectively. And the best prediction we used in benchmarks is the fusion of $F^{(1)}, F^{(2)}$ and $F^{(3)}$. As mentioned above, there are abundant multi-scale information stored in the side-outputs, so we concatenate all the side-outputs with fused prediction and the raw image along with the color channel as the input of the next stage, which called the multi-recursive-input schema.

We choose HED[10] as the sub-network in each stage to carry on the multi-scale feature extraction and holistic deep supervision proposed by [10]. Here, the deep supervision refers to use the ground-truth to supervise the prediction of not only the last output, but also the side-outputs[10], which avoids gradient vanishing in some degree and allows shallow layers efficiently learn the parameters from data.

3.2 Training Phase

Object Boundary Detection can be formulated as a pixel-wise two-class classification problem. Given a training image $X = (X(j), j = 1, 2, \dots, |X|)$ with totally $|X|$ pixels over the spatial dimensions, our goal is to classify each $\tilde{Y}(j)$, the label of pixel $X(j)$, as a boundary pixel ($\tilde{Y}(j) = 1$) or not

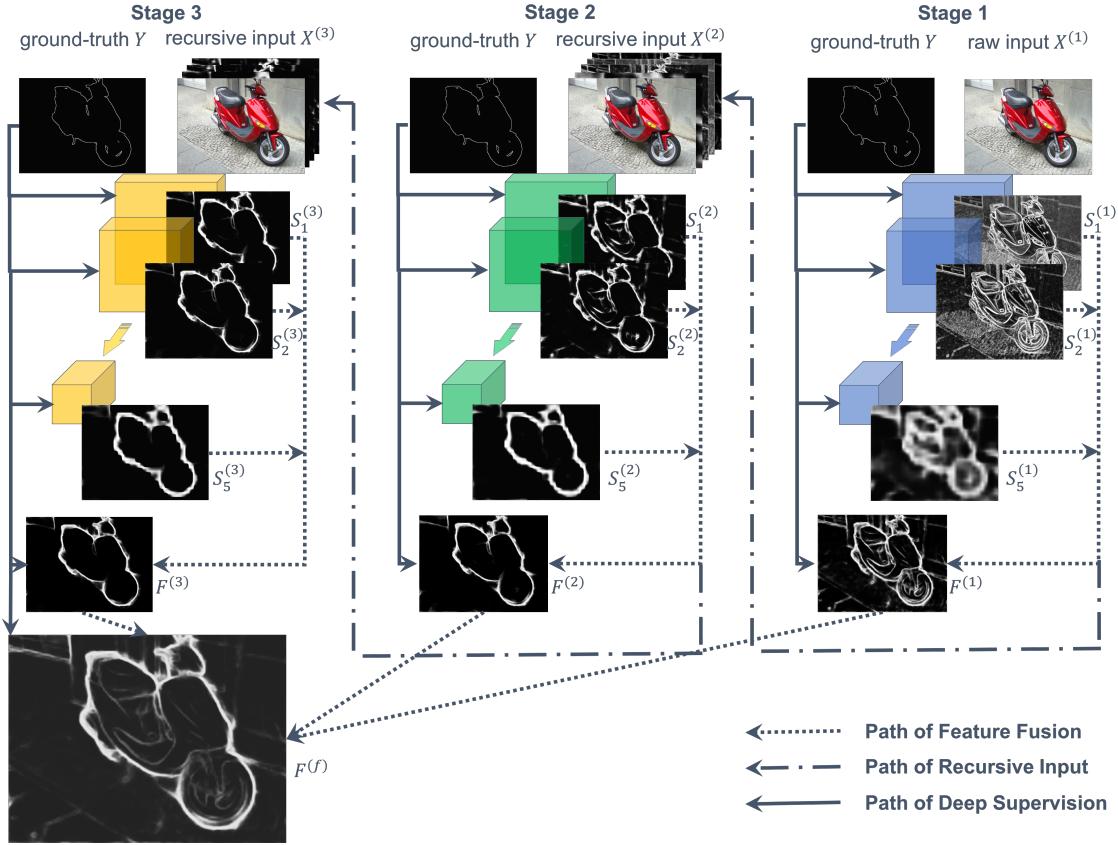


Figure 3.2 Illustration of the proposed cascaded fully convolutional network for object boundary detection.

$(\tilde{Y}(j) = 0)$. And we use $\tilde{Y} = (\tilde{Y}(j), j = 1, 2, \dots, |X|)$ to denote the corresponding detected boundary map. Considering the mini-batch with size equal to Z , and the network with P sub-networks and Q levels in each sub-network, we demonstrate our end-to-end training as follows.

3.2.1 Formulation of the Multi-recursive-input

Let $S_p^{(q)}$ be the side-output of the p^{th} level in the q^{th} sub-network, $F^{(q)}$ and $X^{(q)}$ be the fused prediction and training input in the q^{th} sub-network, now we have:

$$X^{(q)} = \begin{cases} X & , q = 1 \\ X \circ S_{q-1}^1 \circ S_{q-1}^2 \circ \cdots \circ S_{q-1}^Q & , q > 1 \end{cases} \quad (3.1)$$

Here \circ represents the concatenating operation along with the color channel.

3.2.2 Loss Function

To demonstrate the loss function of our cascaded network, we firstly present the loss function of side-output $S_p^{(q)}$:

$$\begin{aligned} l_p^{(q)}(\mathbf{W}^{(q)}, \mathbf{w}_p^{(q)}; X^{(q)}, Y) = & -\alpha \sum_{j \in |B|} \log(1 - \text{sigmoid}(s_p^{(q)}(j))) \\ & - (1 - \alpha) \sum_{j \in |\bar{B}|} \log(\text{sigmoid}(s_p^{(q)}(j))) \end{aligned} \quad (3.2)$$

where $\mathbf{W}^{(q)}$, $\mathbf{w}_p^{(q)}$, $|B|$ and $|\bar{B}|$ refer to the model parameters of the q^{th} sub-network, ones of the $S_p^{(q)}$, the sets of boundary and non-boundary ground-truth annotations, respectively. α is the ratio of $|\bar{B}|$ and $|B|$, which cancels out the huge imbalance of positive samples (boundary pixels) and negative samples (non-boundary pixels) in some degree[10].

Let $\beta_p^{(q)}$ be the weight of $S_p^{(q)}$. Given the notation as follows:

$$\begin{aligned} \mathbf{W} &= \{\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}^{(q)}\}, \\ \mathbf{w} &= (\mathbf{w}_p^{(q)}, p \text{ from } 1 \text{ to } P \text{ and } q \text{ from } 1 \text{ to } Q), \end{aligned}$$

the loss of all the side-outputs in the q^{th} sub-network, L_S , is the sum of the losses of all its side-outputs and the fusion output.

$$L_S(\mathbf{W}, \mathbf{w}; X, Y) = \sum_{q=1}^Q \sum_{p=1}^P \beta_p^{(q)} l_p^{(q)}(\mathbf{W}^{(q)}, \mathbf{w}_p^{(q)}; X^{(q)}, Y) \quad (3.3)$$

All the side-outputs in one stage are fused through a 1×1 convolution layer as mentioned in Section 3.1. The fusion of q^{th} sub-network is defined as:

$$F^{(q)} = \sum_{p=1}^P \theta_p^{(q)} S_p^{(q)} \quad (3.4)$$

where $\theta_p^{(q)}$ is the fusion weight of $S_p^{(q)}$.

The fusion loss of q^{th} stage is:

$$\begin{aligned} l_f^{(q)}(\mathbf{W}^{(q)}, \mathbf{w}^{(q)}; X^{(q)}, Y) = & -\alpha \sum_{j \in |B|} \log(1 - \text{sigmoid}(F^{(q)}(j))) \\ & -(1 - \alpha) \sum_{j \in |\bar{B}|} \log(\text{sigmoid}(F^{(q)}(j))) \end{aligned} \quad (3.5)$$

where $\mathbf{w}^{(q)}$ is the set of the parameters for all the side-outputs in stage q .

Similar to Function 3.4, the final prediction is produced by:

$$F^{(f)} = \sum_{q=1}^Q \phi^{(q)} F^{(q)} \quad (3.6)$$

where $\phi^{(q)}$ is the fusion weight of $F^{(q)}$.

The loss of all the fusion outputs is:

$$L_f(\mathbf{W}, \mathbf{w}, \theta, \phi; X, Y) = \sum_{q=1}^Q \beta_f^{(q)} l_f^{(q)}(\mathbf{W}^{(q)}, \mathbf{w}^{(q)}; X^{(q)}, Y) \quad (3.7)$$

where $\theta = (\theta_p^{(q)}, q = 1, 2, \dots, Q, p = 1, 2, \dots, P)$, $\phi = (\phi^{(q)}, q = 1, 2, \dots, Q)$ and $\beta_f^{(q)}$ denotes the loss weight for each sub-network.

Finally, we train all the parameters jointly in an end-to-end framework by minimize the following total loss:

$$L_t(\mathbf{W}, \mathbf{w}, \theta, \phi; X, Y) = L_S(\mathbf{W}, \mathbf{w}; X, Y) + L_f(\mathbf{W}, \mathbf{w}, \theta, \phi; X, Y) \quad (3.8)$$

3.3 Testing Phase

When testing, we input an image X into the first sub-network and obtain the fusion of each sub-network as Function 3.4, step by step. Considering the q^{th} stage as a function $f^{(q)}$, the final prediction \tilde{Y} can be produced by:

$$\tilde{S}_p^{(q)} = f^{(q)}(X^{(q)}, (\mathbf{W}^{(1)})^*, \dots, (\mathbf{W}^{(Q)})^*, (\mathbf{w}_1^{(Q)})^*, \dots, (\mathbf{w}_P^{(Q)})^*) \quad (3.9)$$

$$\tilde{F}^{(q)} = \sum_{p=1}^P (\theta_p^{(q)})^* \tilde{S}_p^{(q)} \quad (3.10)$$

$$\tilde{F}^{(f)} = \sum_{q=1}^Q (\phi^{(q)})^* \tilde{F}^{(q)} \quad (3.11)$$

$$\tilde{Y} = \text{sigmoid}(\tilde{F}^{(f)}) \quad (3.12)$$

3.4 Model Interpretation

3.4.1 Cascaded Architecture vs. Single-stage Architecture

Such a cascaded architecture for object-level boundary detection can be interpreted as follows:

On one hand, our eyes can capture various levels of data due to the information re-organization achieved by the multi-stage vision system in our brain. The proposed cascaded architecture mimics the multi-stage vision system of human-beings, by introducing the recursive links to re-filter the multi-scale features produced by the previous stage. As illustrated in Fig. 3.4, to compare the side-outputs and fusions in row (in the same level), we find the boundary map of each stage is "*clearer*"

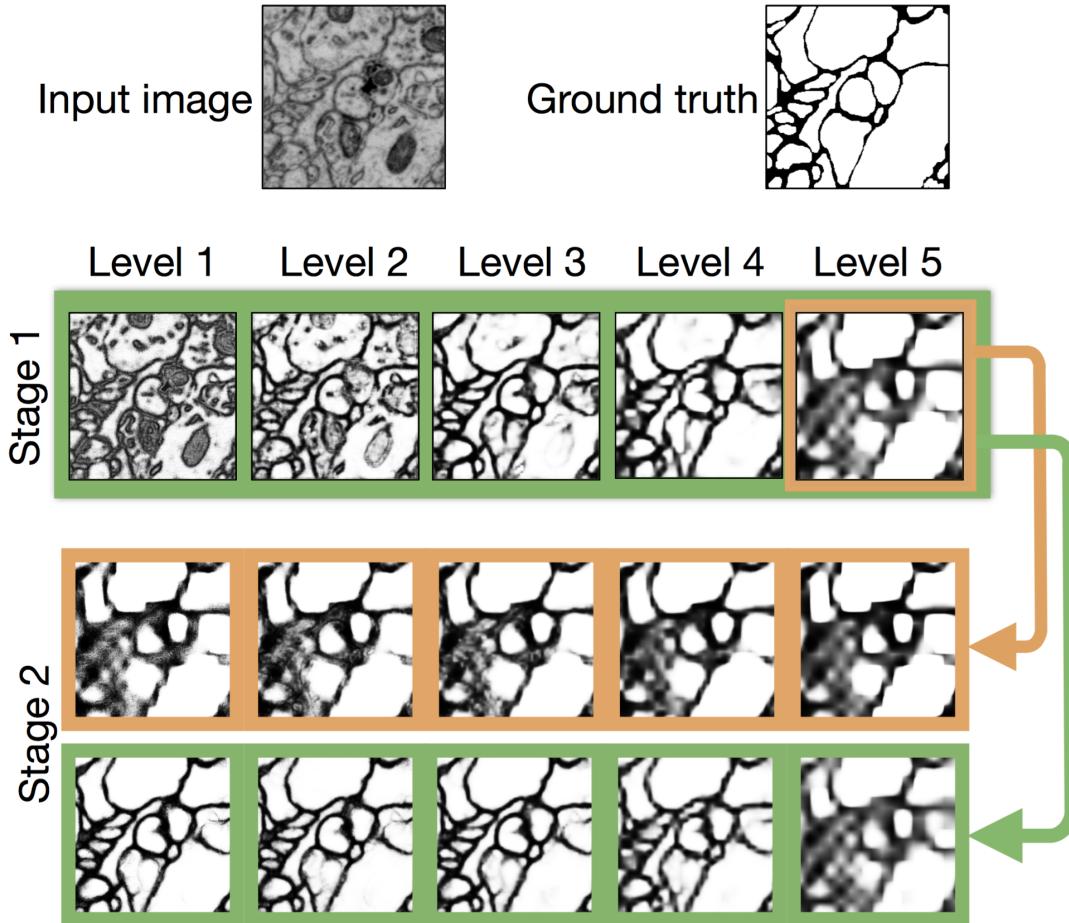


Figure 3.3 Multi-recursive-input vs. Single-recursive-input.

than the former stage. The cascaded network enables the later layers to have a higher-level scope, which breaks the limitation of the receptive field size existing in a single-stage network. Comparing with a considerably deep single-stage network which has the same receptive size, in other words, our cascaded network will not suffer from the gradient vanishing and be much easier to be trained.

On the other hand, human can see both the local details and the object-level abstracts from an image, while not confused by the complex multi-scale information. It is mainly because our brain has a multi-stage vision system to filter these information and select the useful information from the right scale[33]. the joint training and feature fusion of the proposed method let data drives the choice from the results of all the stages and levels.

3.4.2 Multi-recursive-input vs. Single-recursive-input

One of the biggest differences between our work with others is the multi-recursive-input. The qualitative comparison can be found in Fig. 3.3, where the single-recursive-input only reuses the large scale features in the next stage and thus loses the ability to accurately localize the boundary.

3.4.3 End-to-end Training vs. Stepwise Self-tuning

To train our network end-to-end, the former stages and the next stages can have an effect on each other. Instead of fixing the model parameters of the previous stages, we let data teach how the network works. With enough training data and the proposed easy-to-train network architecture, our network is able to adjust all the parameters jointly and benefit from it to achieve a better performance than trained stepwise. Refer to Chapter 4 for quantization results and details.

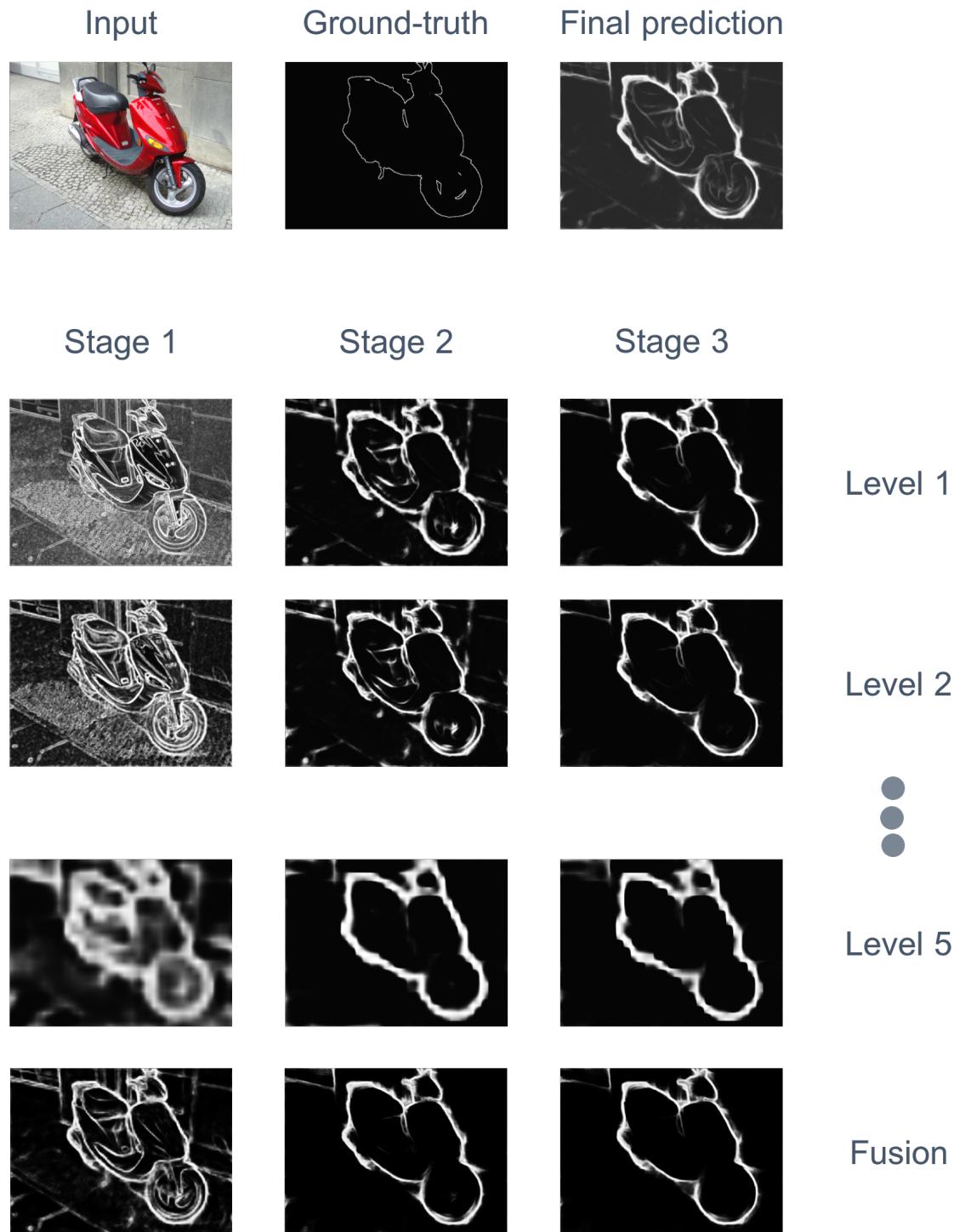


Figure 3.4 Prediction examples from different side-outputs and fused outputs in different stages.

Chapter 4

Experiments

As addressed in Chapter 1, we firstly test our cascaded fully convolutional network in the task of neuronal boundary detection on two mainstream datasets, Mouse Piriform Cortex Dataset[3] and ISBI 2012 EM Segmentation Dataset[5], then extend to adapt with the object boundary detection in natural images on the large object boundary dataset, PASCAL VOC Contour Dataset[1]. Besides, we will give the comparative experiments to explain how the performance is achieved by the proposed network structure.

4.1 Neuronal Boundary Detection in EM Images

Neuronal boundary detection in EM images is usually formed as the segmentation annotation problem (Fig. 4.1). In training phase, we compute the 2D gradient in the segmentation ground-truth to obtain the boundary annotation. And then in testing phase, we apply the graph-based algorithms such as watershed algorithm to transfer the boundary prediction into segmentation, since the evaluation metric is based on the segmentation results. In this experiment, we firstly detect the membrane boundary and then apply graph-based watershed[14] to it for a segmentation result.

4.1.1 Evaluation Metric

We follow the metric proposed in [3], using the Rand F-score to report the performance of the segmentation. As the harmonic average of the Rand merge score and Rand split score, marked as V_{merge}^R and V_{split}^R respectively, the Rand F-score:

$$V_{Fscore}^{Rand} = \frac{2V_{merge}^{Rand}V_{split}^{Rand}}{V_{merge}^{Rand} + V_{split}^{Rand}} \quad (4.1)$$

(4.2)

and

$$V_{merge}^{Rand} = \frac{\sum_{ij} p(i, j)^2}{\sum_i (\sum_j p(i, j))^2} \quad (4.3)$$

$$V_{split}^{Rand} = \frac{\sum_{ij} p(i, j)^2}{\sum_j (\sum_i p(i, j))^2} \quad (4.4)$$

where $p(i, j)$ indicates the number of pixel pairs located in the i^{th} segment of the detected segmentation and the j^{th} segment of the ground-truth segmentation.

Rand split score is high for those results with few split errors, while Rand merge score is high for those results with few merge errors. Generally there will be a trade-off between the Rand split score and the Rand merge score for an image or a dataset. The Rand F-score here is to define the trade-off score, just like the role of the F-score or F-measure in many object detection benchmarks to average the precision metric and recall metric.

The Precision (rand merge)-recall (rand split) curves (PR-curves) can be generated by varying the threshold for boundary binarization[44].

4.1.2 Mouse Piriform Cortex Dataset

Mouse Piriform Cortex Dataset[3] contains 4 stacks of grayscale EM images and their corresponding segmentation annotations(Table. 4.1), where the stack 1 for validation and the stack 2, 3, 4 for training.

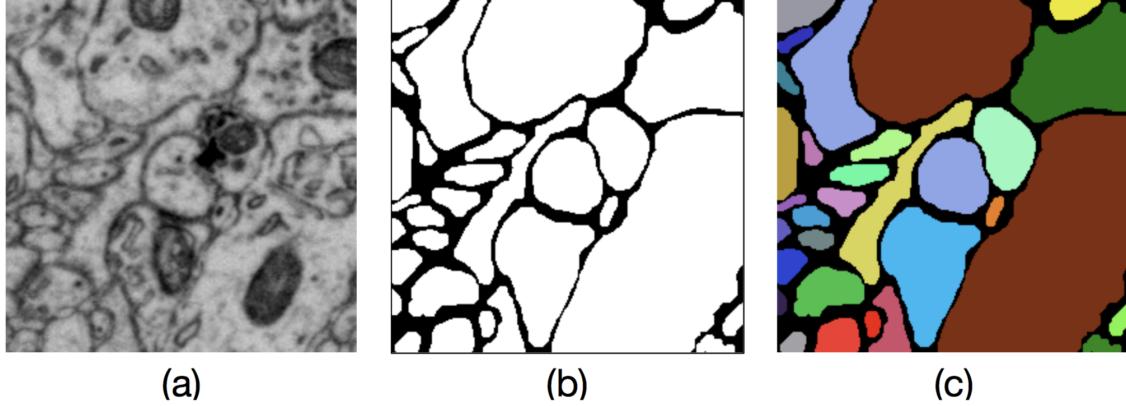


Figure 4.1 (b)Neuronal boundary detection and (c) segmentation can be easily converted from each other, which we used in the experiments on Mouse Piriform Cortex Dataset[3] and ISBI 2012 EM Segmentation Dataset[5]. (1) By applying the graph-based algorithms such as watershed, we can transfer the boundary prediction into segmentation. (2) By calculating the 2D gradient in the segmentation ground-truth, the boundary annotation can be obtained.

Table 4.1 Mouse Piriform Cortex Dataset

	stack 1	stack 2	stack 3	stack 4
Slices	168	170	169	121
Scales	255*255	512*512	512*512	256*256
Usage	Validation	Training	Training	Training

The annotation of EM images is laborious and time consuming[33], leading to the lack of annotated training samples in EM boundary detection. However, the deep networks often require thousands of training data to drive the optimization algorithms. So we augment the raw data to generate sufficient data for training, by rotating the raw image into four directions ($0, \frac{\pi}{2}, \pi, \frac{3\pi}{2}$), then flipping them with three cases (horizontal flipping, vertical flipping, no flipping), and rescaling them to three scales (0.8, 1.0, 1.2), resulting in the training data 36 times the raw images.

To train such a considerably large network from scratch is not necessary nor effective, so we follow the framework proposed by a lot of works[11, 10, 44] to initialize model parameters of the first stage with the VGG-16 models pre-trained in ImageNet[45]. To ease the multi-stage training, we train a single-stage model first. Then use it to initialize the first stage of a two-stage network, leaving the second stage is initialized randomly, by which our multi-stage networks can reuse the relatively low-level filters learned by the single-stage network. Note that the proposed cascaded network is still trained end-to-end because we do not fix any parameters in the final training step but train all the parameters jointly instead.

For experiment settings, we use the standard SGD algorithm to optimize the network, with the momentum of 0.9, the weight decay of 2×10^{-4} , and a stepwise learning rate from 1e-8 to 1e-9(after 20000 iterations). By default, there are 5 convolution levels in each subnet. The training for the takes 20000 iterations for the single-stage network, and 10000 iterations for each of the sequential sub-networks.

Quantitative results are shown in Table 4.2 and Fig. 4.2. With stacking the sub-networks, our method can maintain both a high recall and a high precision. Compared with the current state-of-the-art method VD2D3D[3], a deep network with the parameters of several layers fine-tuned by recursively input the prediction into itself, our 2-stage and 3-stage cascaded network can outperform it by around 1% and 1.5% respectively. We argue that such a improvement is meaningful in neuronal boundary detection, since the precise boundary annotation is vital for neural circuits reconstruction and has a low tolerance of errors. The architecture of VD2D3D is much like the Fig. 3.1(c), while the VD2D3D not only feeds the boundary map of the raw input but also generates the predictions of neighbor slices and concatenates all them to enhance the learning (Fig. 2.1). Our networks, taking advantage of the multi-recursive-input and end-to-end training, can outperform the VD2D3D without any additional 3D context information.

Table 4.2 Rand F-scores on Mouse Piriform Cortex Dataset[3]

	Rand Split Score	Rand Merge Score	Rand F-score
N4[13]	0.9010	0.9619	0.9304
VD2D[3]	0.9174	0.9771	0.9463
VD2D3D[3]	0.9555	0.9891	0.9720
Ours(1-stage)	0.9802	0.9576	0.9688
Ours(2-stage)	0.9880	0.9759	0.9819
Ours(3-stage)	0.9815	0.9917	0.9866

4.1.3 ISBI 2012 EM Segmentation Dataset

The other dataset for neuronal boundary detection in EM images is the ISBI 2012 EM Segmentation Dataset. As the dataset for the popular ISBI 2012 EM Segmentation Challenge, ISBI 2012 EM Segmentation Dataset does not offer the ground-truth segmentation for the testing stack, which is reasonable but increases the difficulty of results analyzing and discussion. Besides, the raw training data with annotations are relatively insufficient, where there is only one stack for training with the dimension of $30 \times 512 \times 512$. To train a very deep network using such a few images is theoretically difficult, even though we use the standard data augmentation to enrich the training materials.

Network settings are the same as ones we used in the experiment for Mouse Piriform Cortex Dataset[3]. Standard data augmentation (Section 4.1.2) is considered here, too.

The Rand F-score data from the leaderboard of ISBI 2012 EM Segmentation Challenge are shown in Table 4.3 . Not all the scores are presented since there are many entries in the leaderboard which have not been reported in literature. It worths noting that most of the leading methods

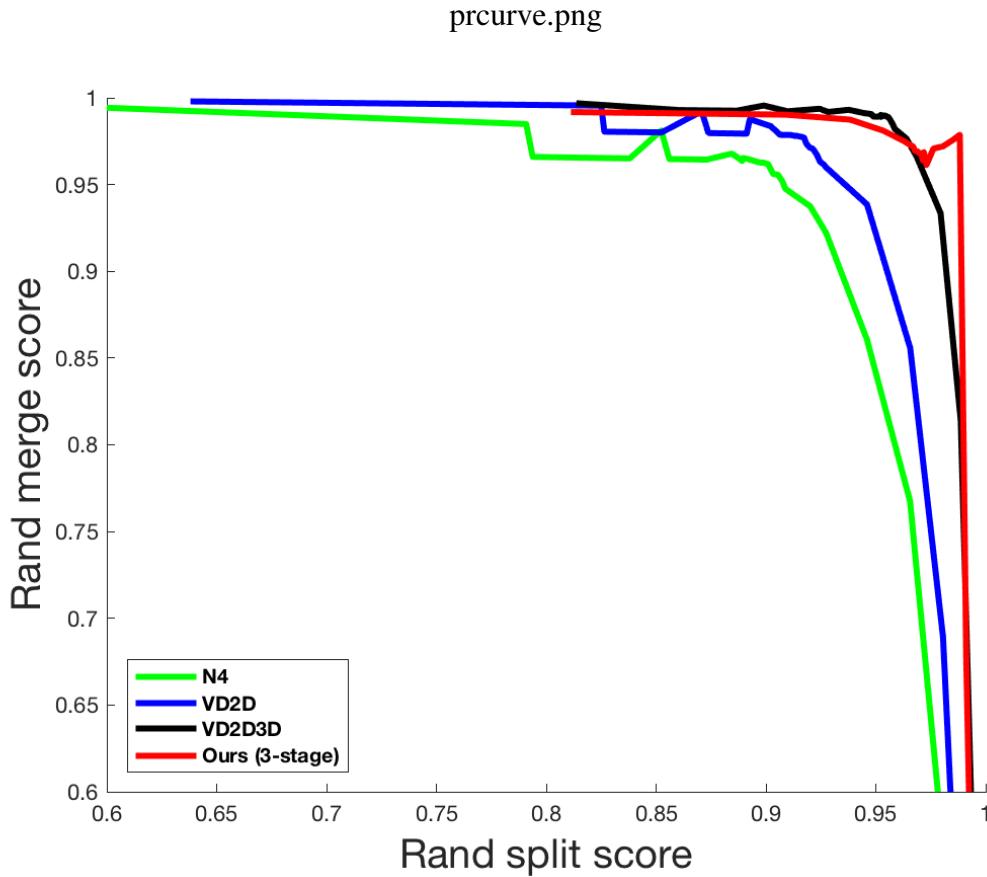


Figure 4.2 Precision (rand merge)-recall (rand split) curves on Mouse Piriform Cortex Dataset[3]. Our 3-stage network outperforms all the previous works on this dataset.

apply post-processing algorithms or assemble several models to improve the ranking. However, the proposed cascaded networks without any post-processing procedures can achieve competitive result to those with post-processing. The post-processing used in each approach is also presented in Table 4.3, where [46] is a post-processing algorithm itself and can be adopted to ours. With a proper post-processing, our network can win the state-of-the-art with more powerful ResNet[27] as its base network. In the future, we will embed the ResNet[27] in our cascaded framework to further improve the performance.

Table 4.3 The Rand F-scores part from the leaderboard of ISBI 2012 EM Segmentation Challenge[5].

	Rand F-score
U-net[5]	0.9727
CUMedVision[47]	0.9768
IAL IC[46]	0.9773
FusionNet[4]	0.9780
Ours(3-stage)	0.9780
PolyMtl[48]	0.9806
Ours(3-stage) + IAL IC	0.9836

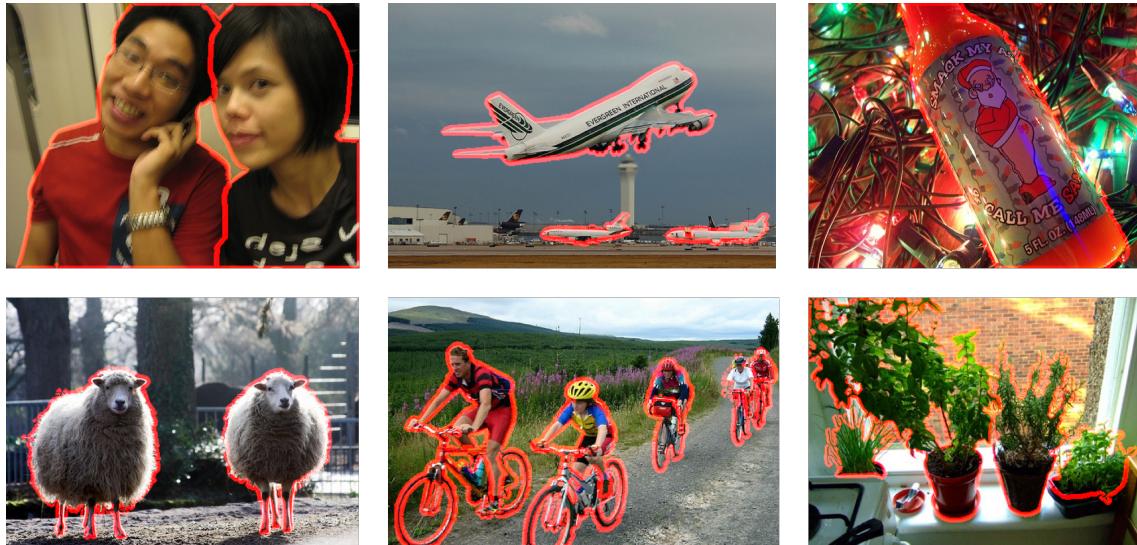


Figure 4.3 Examples selected from PASCAL VOC Contour Dataset[1], where ground-truth contours are labeled as red. There are various objects (human, artificialities, animals, plants, etc.) appearing in complex scenes (indoor and outdoor environments, colorful backgrounds, blurs, textural confusions, etc), which significantly increases the difficulty of detecting the object-level boundary.

4.2 Object Boundary Detection in Natural Images

Our work begins at Neuronal Boundary Detection of EM images but dose not end in it. By applying the proposed multi-stage network to common object boundary detection in natural images, we evaluate the generalization capacity of our network.

4.2.1 Metrics

Precision and Recall are the well-accepted metrics for two-class annotation tasks[44, 49, 1]. Given the notations as follows: (1) True Positive (TP): the times of annotating a positive pixel as positive; (2) False Positive (FP): the times of annotating a negative pixel as positive; (3) True negative (TN): the times of annotating a negative pixel as negative; (4) False positive (FN): the times of annotating a positive pixel as negative; we have Precision (P), Recall (R) and their harmonic average F-measure

$$P = \frac{TP}{TP+FP} \quad (4.5)$$

$$R = \frac{TP}{TP+FN} \quad (4.6)$$

$$F\text{-measure} = \frac{2PR}{P+R} \quad (4.7)$$

$$(4.8)$$

The Precision-Recall Curve (PR-curves) can be also generated by varying the threshold for boundary binarization[44].

4.2.2 PASCAL VOC Contour Dataset

PASCAL VOC Segmentation Datasets are a series of well-accepted object segmentation datasets, including the well-annotated natural photos released from 2007 to 2012. [1] apply DenseCRF[50] to further refine the boundary annotation from the segmentation ground-truth and release a large

instance-level object boundary detection dataset, PASCAL VOC Contour Dataset, with 10582 training images, 1449 testing images and their corresponding boundary annotations. Different from EM images, these natural photos are taken from various scenes and contain kinds of objects (Fig. 4.3). For example, in each photo, the image scale, the number and category of the object are different. There are also many samples with blurred object boundary and confused background on the dataset. All these significantly increase the difficulty of boundary detection on the PASCAL VOC Contour Dataset, while call out a deep network with sufficient capacity to extract the highly abstract object-level features from the data.

With such a large dataset, the standard data augmentation with a ratio of 36 is not necessary nor efficient. We follow the light-weight data augmentation mentioned in [1], by randomly cropping four 224×224 patches from the raw image and flipping them once, resulting in $8\times$ training samples. Here, the cropping is mainly for memory efficiency[1] and the mini-batch based training. Original FCNs[11] is not available for the training with multiple training samples in one mini-batch, because these samples may have different scales and can result in error. By fixing the cropped patches into one-size, we are able to enlarge the batch-size, 8 in practice, for a more effective and faster training.

We fine-tune the proposed network on PASCAL VOC Contour Dataset using the same strategy in Section 4.1.2. Firstly, we train a single-stage network with the same structure as HED[10]. Then we use the fine-tuned HED network to initialize the training for our 2-stage and 3-stage cascaded fully convolutional networks. F-measure comparisons are shown in 4.4.

For optimization, we use the standard Stochastic Gradient Descent (*SGD*). Other hyper parameters such as the level of convolution in each subnet, learning rate, momentum and weight decay are same as what we used in the experiments of neuronal boundary detection. As demonstrated in [1], to process all the training samples, named one epoch, takes around 10000 iterations. We train the single-stage, 2-stage and 3-stage networks for 2 epochs, 4 epochs, and 5 epochs respectively. Compared with the CEDN proposed in [1], the training for our cascaded networks costs only half of

Table 4.4 Object boundary detection evaluation comparison on PASCAL VOC Contour Dataset[1]. Our proposed 3-stage cascaded fully convolutional network achieves the new state-of-the-art on this benchmark, with a significant improvement (around 2% over the second)).

	F-measure
SCG[51]	0.36
MCG[51]	0.37
SE[8]	0.37
HED[10]	0.52
Ours(2-stage)	0.53
CEDN[1]	0.57
Ours(3-stage)	0.59

the time for training CEDN, which proves the effectiveness of the proposed recursively end-to-end training.

As we expected, traditional methods for local edge detection such as SCG, its multi-scale version MCG[51], and Structured Random Forest for Edge Detection (*SE*)[8] present poor results on this benchmark. These methods use hand-crafted features which often combine multiple cues like color, brightness, spectrum, to handle different cases. HED[10] has the ability to naturally obtain the hierarchical features due to the adoption of the holistically-nested network, which is also used in our sub-networks. We fine-tune it on PASCAL VOC Contour Dataset and receive a F-measure of 0.52, much better than the traditional methods but lower than our 2-stage network. What's more, our 3-stage network outperforms CEDN[1] by 2% and achieves the state-of-the-art score on this benchmark. PR-curves in Fig. 4.4 illustrate the refinement benefits from the cascaded network structure.

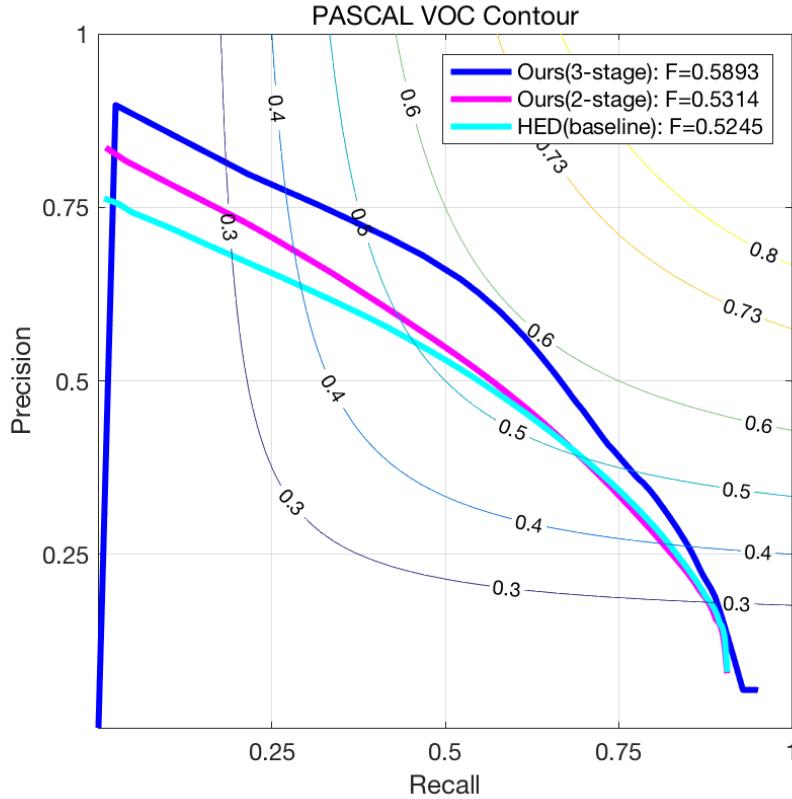


Figure 4.4 Precision-recall curves on PASCAL VOC Contour Dataset[1].

4.3 Control Experiments for Model Interpretation

3

To verify the explanation and interpretation in Section 3.4, we design a series of control experiments and give the neuronal boundary detection evaluation on Mouse Piriform Cortex Dataset[3]. We do not test on the ISBI 2012 Segmentation Dataset[5] since there is no provided annotation for testing images on it.

There are three control experiments, for which we design eight kind of networks with different network structures:

- **Baseline HED[10]** pre-trained on BSDS 500 Dataset[15]
- **Network A** 1-stage (HED[10]) network with 5 levels in one stage, end-to-end training and

Table 4.5 Control Experiment 1: Cascaded Architecture vs. Single-stage Architecture

	Rand F-score
Baseline	0.9680
Network A	0.9688
Network B	0.9739
Network C	0.9819
Network D	0.9866

no recursive input;

- **Network B** 2-stage cascaded network with 4 levels in each stage, end-to-end training and multi-recursive-input;
- **Network C** 2-stage cascaded network with 5 levels in each stage, end-to-end training and multi-recursive-input;
- **Network D** 3-stage cascaded network with 5 levels in each stage, end-to-end training and multi-recursive-input;
- **Network E** 2-stage cascaded network with 5 levels in each stage, end-to-end training and single-recursive-input (only the level 4);
- **Network F** 2-stage cascaded network with 5 levels in each stage, end-to-end training and single-recursive-input (only the level 5);
- **Network G** 2-stage cascaded network with 5 levels in each stage, stepwise training and multi-recursive-input;

4.3.1 Evaluations on Cascaded Architecture vs. Single-stage Architecture

To prove the proposed cascaded architecture effective in boundary detection, we compare our 2-stage (Network C) and 3-stage (Network E) cascaded networks with a 1-stage fine-tuned HED[10] network (Network A) and the baseline HED[10] pre-trained on a local edge boundary detection dataset, BSDS 500[15]. Results show that even a 2-stage cascaded network with only 4 convo-

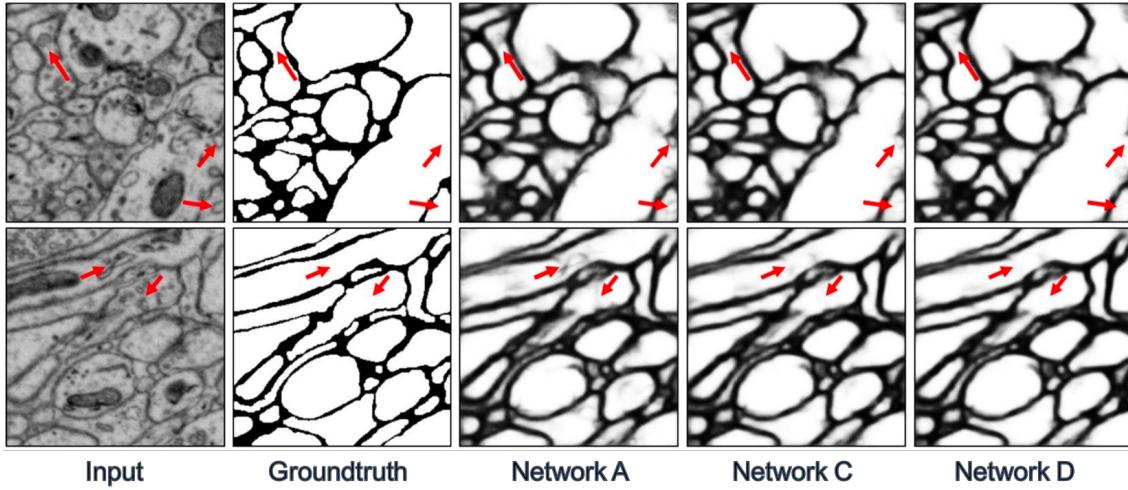


Figure 4.5 Qualitative examples to reveal how cascaded networks improve the high-level boundary detection result. For example, the arrows mean the false positive detections are removed step by step.

lution levels outperforms the baseline and the single-stage fine-tuned network with 5 convolution levels. If we stack more sub-networks while controlling the rest hyper-parameters, we can find the performance boosting from 0.9688 (Network A), to 0.9819 (Network C) and 0.9866 (Network D) step by step (Table 4.5), which indicates that the proposed cascaded structures effective in high-level object boundary detection. The improvement achieved by stepwise removing the confounding local edges inside the objects (circuits), and refining the boundary with low contrast to the non-boundary area, illustrated by Fig. 4.5.

4.3.2 Evaluations on Multi-recursive-input vs. Single-recursive-input

As we mentioned in Section interpretation, Multi-recursive-input plays a important role in the high-level object boundary detection, by providing the multi-scale features in a natural way. In order to evaluate it, we design several 2-stage single-recursive-input networks (Network E, F) and compare them with our 2-stage multi-recursive-input network (Network C). Cause we do not know which level is the best for inputing to the next stage, we input the level 4 (in Network E) and level 5 (in Network F) to the next sub-network respectively. According to Table 4.6, we surprisingly

Table 4.6 Control Experiment 2: Multi-recursive-input vs. Single-recursive-input

	Rand F-score
Network F	0.9410
Network E	0.9656
Baseline	0.9680
Network C	0.9819

Table 4.7 Control Experiment 3: End-to-end Training vs. Stepwise Self-tuning

	Rand F-score
Baseline	0.9680
Network G	0.9762
Network C	0.9819

find the cascaded networks with single-recursive-input perform bad even when comparing with the pre-trained 1-stage network (Baseline). The single-recursive-input from the previous stage has only the features in one scale, either large or small. If we only feed the large-scale features with low-resolution into the next stage, there will be a huge loss in image details. If we only feed the small-scale features into the next stage, the cascaded structure will be deprecated. It is always a hard problem for people to choose a best level as the recursive input, so we highlight the multi-recursive-input again in the cascaded structures.

4.3.3 Evaluations on End-to-end Training vs. Stepwise Self-tuning

In the end, we design a control experiment with only the training strategy varied to see whether our cascaded networks benefit from the end-to-end training or not. As shown in Table 4.7, the cascaded network trained end-to-end (Network C) outperforms the same network but trained stepwise (Network G). The end-to-end training was hardly used to learn such a deep network, due to the consideration of gradient vanishing. However, the holistically-nested sub-networks and multi-

recursive-input release the limitation, which enables the end-to-end training in our deep cascaded networks.

Chapter 5

Conclusions

Object-level boundary is a useful but unexplored topic in computer vision. With a well-detected object boundary, the segment of the object will be easily extracted, which can be applied for object proposal, shape based object recognition and kinds of object detection problems. Based on the previous works in neuronal boundary detection, we develop a series of cascaded fully convolutional networks to detection high-level boundary in EM images. Not end up with the improvement achieved in EM images, we further extend the cascaded architecture to the object boundary detection in natural images, and design several control experiments to interpret how it works. Our model is proved to be both effective and interpretable in various boundary detection scenes, thus easy to be aggregated with other frameworks.

In the future, we will try some other structures for the alternative sub-networks, such as the recent published Richer Convolutional Features for Edge Detection[52], which is now the state-of-the-art on many edge detection benchmarks. Applications of our object-level boundary detection can be also considered, such as the boundary-based object proposal[1].

Bibliography

- [1] J. Yang, B. Price, S. Cohen, and M. Y. H. Lee, “Object contour detection with a fully convolutional encoder-decoder network,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 193–202.
- [2] C. L. Zitnick and P. Dollár, “Edge boxes: Locating object proposals from edges,” in *Proc. of European Conference on Computer Vision*, 2014, pp. 391–405.
- [3] K. Lee, A. Zlateski, A. Vishwanathan, and H. S. Seung, “Recursive training of 2d-3d convolutional networks for neuronal boundary detection,” in *Proc. of The International Conference on Neural Information Processing Systems*, 2015, pp. 3573–3581.
- [4] T. M. Quan, D. G. C. Hildebrand, and W. Jeong, “Fusionnet: A deep fully residual convolutional neural network for image segmentation in connectomics,” *CoRR*, vol. abs/1612.05360, 2016.
- [5] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *Proc. of The International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2015, pp. 234–241.
- [6] J. Canny, “A computational approach to edge detection,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 8, pp. 679–698, 1986.
- [7] I. Sobel, “An isotropic 3 3 image gradient operator,” 02 2014.

- [8] P. Dollar and C. L. Zitnick, “Structured forests for fast edge detection,” in *Proc. of The IEEE International Conference on Computer Vision*, 2013, pp. 1841–1848.
- [9] W. Shen, X. Wang, Y. Wang, X. Bai, and Z. Zhang, “Deepcontour: A deep convolutional feature learned by positive-sharing loss for contour detection,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3982–3991.
- [10] S. Xie and Z. Tu, “Holistically-nested edge detection,” in *Proc. of The IEEE International Conference on Computer Vision*, 2015, pp. 1395–1403.
- [11] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [12] J. W. Lichtman and W. Denk, “The big and the small: Challenges of imaging the brains circuits,” *Science*, vol. 334, pp. 618–623, 2011.
- [13] D. Ciresan, A. Giusti, L. M. Gambardella, and J. Schmidhuber, “Deep neural networks segment neuronal membranes in electron microscopy images,” in *Proc. of The International Conference on Neural Information Processing Systems*, 2012, pp. 2852–2860.
- [14] A. Zlateski and H. S. Seung, “Fusionnet: A deep fully residual convolutional neural network for image segmentation in connectomics,” *CoRR*, vol. abs/1505.00249, 2015.
- [15] P. Arbeláez, M. Maire, C. Fowlkes, and J. Malik, “Contour detection and hierarchical image segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, pp. 898–916, 2011.
- [16] P. O. Pinheiro, R. Collobert, and P. Dollar, “Learning to segment object candidates,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2015.

- [17] M. Everingham, L. J. V. Gool, C. K. I. Williams, J. M. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International Journal of Computer Vision*, vol. 88, pp. 303–338, 2010.
- [18] R. Girshick, “Fast r-cnn,” in *Proc. of The IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.
- [19] X. Wang, B. Feng, X. Bai, W. Liu, and L. J. Latecki, “Bag of contour fragments for robust shape classification,” *Pattern Recognition*, vol. 47, pp. 2116–2125, 2014.
- [20] W. Shen, X. Wang, C. Yao, and X. Bai, “Shape recognition by combining contour and skeleton into a mid-level representation,” in *Proc. of Chinese Conference on Pattern Recognition*, 2014, pp. 391–400.
- [21] W. Shen, Y. Jiang, W. Gao, D. Zeng, and Z. Zhang, “Shape recognition by bag of skeleton-associated contour parts,” *Pattern Recognition Letters*, vol. 83, pp. 321–329, 2016.
- [22] J. Malik, S. Belongie, and J. S. T. Leung, “Contour and texture analysis for image segmentation,” *International Journal of Computer Vision*, vol. 43, pp. 7–27, 2001.
- [23] P. Krähenbühl and V. Koltun, “Learning to propose objects,” in *Proc. of The International Conference on Neural Information Processing Systems*, 2015, pp. 1574–1582.
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Proc. of The International Conference on Neural Information Processing Systems*, 2012, pp. 1106–1114.
- [25] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014.

- [26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [27] X. Z. K. He and J. S. S. Ren, “Deep residual learning for image recognition,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [28] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
- [29] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Proc. of The International Conference on Neural Information Processing Systems*, 2015, pp. 91–99.
- [30] J. Kivinen, C. Williams, and N. Heess, “Visual boundary prediction: A deep neural prediction network and quality dissection,” in *Proc. of The International Conference on Artificial Intelligence and Statistics*, 2014, pp. 512–521.
- [31] W. Shen, X. Wang, Y. Wang, X. Bai, and Z. Zhang, “N⁴-fields: Neural network nearest neighbor fields for image transforms,” in *Proc. of Asian Conference on Computer Vision*, 2014, pp. 536–551.
- [32] J. J. Lim, C. L. Zitnick, and P. Dollar, “Sketch tokens: A learned mid-level representation for contour and object detection,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 3158–3165.
- [33] M. Helmstaedter, “Cellular-resolution connectomics: Challenges of dense neural circuit reconstruction,” *Nature methods*, vol. 10, pp. 501–507, 2013.
- [34] O. Sporns, G. Tononi, and R. Kötter, “The human connectome: A structural description of the human brain,” *PLoS Comput Biol*, vol. 1, p. e42, 2005.

- [35] D. Laptev, A. Vezhnevets, S. Dwivedi, and J. M. Buhmann, “Anisotropic sstem image segmentation using dense correspondence across sections,” in *Proc. of The International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2012, pp. 323–330.
- [36] V. Kaynig, T. J. Fuchs, and J. M. Buhmann, “Geometrical consistent 3d tracing of neuronal processes in sstem data,” in *Proc. of The International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2010, pp. 209–216.
- [37] R. Kumar, A. Vazquez-Reina, and H. Pfister, “Radon-like features and their application to connectomics,” in *Proc. of The IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 2010, pp. 186–193.
- [38] M. Seyedhosseini, R. Kumar, E. Jurrus, R. Giuly, M. Ellisman, H. Pfister, and T. Tasdizen, “Detection of neuron membranes in electron microscopy images using multi-scale context and radon-like features.” in *Proc. of The International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2011, pp. 670–677.
- [39] L. Najman and M. Schmitt, “Geodesic saliency of watershed contours and hierarchical segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, pp. 1163–1173, 1996.
- [40] M. G. Uzunbas, C. Chen, and D. N. Metaxas, “Optree: A learning-based adaptive watershed algorithm for neuron segmentation,” in *Proc. of The International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2014, pp. 97–105.
- [41] Y. Boykov and M. P. Jolly, “Interactive graph cuts for optimal boundary and region segmentation of objects in n-d images,” in *Proc. of The IEEE International Conference on Computer Vision*, 2001, pp. 105–112.

- [42] A. Fakhry, T. Zeng, and S. Ji, “Residual deconvolutional networks for brain electron microscopy image segmentation,” *IEEE Transactions on Medical Imaging*, vol. 36, pp. 447–456, 2017.
- [43] I. Arganda-Carreras, S. C. Turaga, D. R. Berger, D. Ciresan, A. Giusti, L. M. Gambardella, J. Schmidhuber, D. Laptev, S. Dwivedi, J. M. Buhmann, T. Liu, M. Seyedhosseini, T. Tasdizen, L. Kamentsky, R. Burget, V. Uher, X. Tan, C. Sun, T. D. Pham, E. Bas, M. G. Uzunbas, A. Cardona, J. Schindelin, and H. S. Seung, “Crowdsourcing the creation of image segmentation algorithms for connectomics,” *Frontiers in Neuroanatomy*, vol. 9, p. 142, 2015.
- [44] W. Shen, K. Zhao, Y. Jiang, Y. Wang, Z. Zhang, and X. Bai, “Object skeleton extraction in natural images by fusing scale-associated deep side outputs,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 222–230.
- [45] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li, “Imagenet: a large-scale hierarchical image database,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [46] T. Beier, B. Andres, U. Köthe, and F. A. Hamprecht, “An efficient fusion move algorithm for the minimum cost lifted multicut problem,” in *Proc. of European Conference on Computer Vision*, 2016, pp. 715–730.
- [47] H. Chen, X. Qi, J. Cheng, and P. Heng, “Deep contextual networks for neuronal structure segmentation,” in *Proc. of The AAAI Conference on Artificial Intelligence*, 2016.
- [48] M. Drozdzal, G. Chartrand, E. Vorontsov, L. Di-Jorio, A. Tang, A. Romero, Y. Bengio, C. Pal, and S. Kadouri, “Learning normalized inputs for iterative estimation in medical image segmentation,” *CoRR*, vol. abs/1702.05174, 2017.

- [49] W. Shen, K. Zhao, Y. Jiang, Y. Wang, X. Bai, and A. Yuille, “Deepskeleton: Learning multi-task scale-associated deep side outputs for object skeleton extraction in natural images,” *IEEE Transactions on Image Processing*, vol. 26, pp. 5298–5311, 2017.
- [50] P. Krähenbühl, A. Giusti, L. M. Gambardella, and J. Schmidhuber, “Efficient inference in fully connected crfs with gaussian edge potentials,” in *Proc. of The International Conference on Neural Information Processing Systems*, 2011, pp. 109–117.
- [51] P. Arbeláez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik, “Multiscale combinatorial grouping,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [52] Y. Liu, M. Cheng, X. Hu, K. Wang, and X. Bai, “Richer convolutional features for edge detection,” in *Proc. of The IEEE Conference on Computer Vision and Pattern Recognition*, 2017.