

# 1 U-RISC: an ultra-high-resolution 2 electron microscopy dataset 3 challenging existing deep learning 4 algorithms

5 **Ruohua Shi<sup>1,2†</sup>, Wenyao Wang<sup>1†</sup>, Zhixuan Li<sup>2</sup>, Liuyuan He<sup>2</sup>, Kaiwen Sheng<sup>1</sup>, Lei  
6 Ma<sup>1,2</sup>, Kai Du<sup>\*3</sup>, Tingting Jiang<sup>\*2</sup>, Tiejun Huang<sup>1,2,3</sup>**

**\*For correspondence:**

[kai.du@pku.edu.cn](mailto:kai.du@pku.edu.cn) (KD);  
[ttjiang@pku.edu.cn](mailto:ttjiang@pku.edu.cn) (TTJ)

†These authors contributed equally  
to this work

7 <sup>1</sup>Beijing Academy of Artificial Intelligence Institution, Beijing, China; <sup>2</sup>Department of  
8 Computer Science and Technology, Peking University, Beijing, China; <sup>3</sup>Institute for  
9 Artificial Intelligence, Peking University, Beijing, China

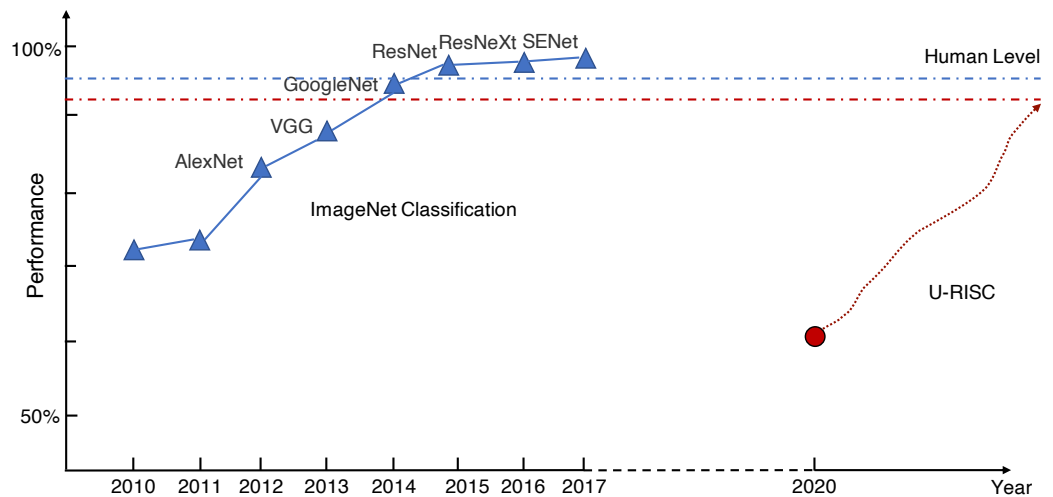
---

11 **Abstract** Connectomics is a developing field aiming at reconstructing the connection of the  
12 neural system at nanometer scale. Computer vision technology, especially deep learning methods  
13 used in image processing, has promoted connectomic data analysis to a new era. However, the  
14 performance of the state-of-the-art methods still falls behind the demand of scientific research.  
15 Inspired by the success of ImageNet, we present the U-RISC, an annotated **U**ltra-high **R**esolution  
16 **I**mage **S**egmentation dataset for **C**ell membrane, which is the largest cell membrane annotated  
17 Electron Microscopy (EM) dataset with a resolution of 2.18nm/pixel. Multiple iterative annotations  
18 ensured the quality of the dataset. Through an open competition, we reveal that the performance  
19 of current deep learning methods still has a considerable gap with human-level, different from ISBI  
20 2012, on which the performance of deep learning is close to human. To explore the causes of this  
21 discrepancy, we analyze the neural networks with a visualization method, attribution analysis. We  
22 find that in U-RISC, it requires a larger area around a pixel to predict whether the pixel belongs to  
23 the cell membrane or not. Finally, we integrate currently available methods to provide a new  
24 benchmark (0.67, 10% higher than the leader of competition, 0.61) for cell membrane  
25 segmentation on U-RISC and propose some suggestions in developing deep learning algorithms.  
26 The U-RISC dataset and the deep learning codes used in this paper will be publicly available.

---

## 28 Introduction

29 Accurate descriptions of neurons and their connections are fundamental to modern neuroscience.  
30 By depicting neurons with the help of Golgi-staining method (*Golgi, 1885*), Cajal could propose  
31 the classic “Neuron Doctrine” more than a century ago (*y Cajal, 1888*), which opened a new era for  
32 modern neuroscience. Nowadays, the development of electron microscopy (EM) has enabled us to  
33 further explore the structural details of the neural system at nanometer (nm) scales (*Kornfeld and  
34 Denk, 2018; Shawn, 2016*) opening up a new field called “Connectomics” that aims to reconstruct  
35 every single connection in the neural system. One milestone of connectomics is the *C.elegans*  
36 project (*White et al., 1986*) which maps all 302 neurons and 7,000 connections in a worm. Recently,  
37 a small piece of human cortex was imaged with a high-speed scanning EM, which maps ~50,000  
38 neurons and ~110,000,000 synaptic connections (*Shapson-Coe et al., 2021*). Connectomic data



**Figure 1.** The history of ImageNet.

39 increases exponentially with a higher resolution of EM and a larger neural tissue volume, even  
40 reaching petabyte (PB) scale (*Shapson-Coe et al., 2021*). Just as it took almost 15 years to complete  
41 the connectome of *C.elegans*, structural reconstruction for higher-level creatures is becoming more  
42 and more daunting with the explosion of connectomic data. Among many bottlenecks, accurate  
43 annotation from large amounts of EM images is the first one that has to be solved.

44 Manual annotation of all the connectomic data is infeasible because of the high annotation cost.  
45 To reduce the burden of manual annotation for humans, one would hope to enable a machine to  
46 annotate the connectomic data with near-human performance automatically. Hopes are higher  
47 today because of the rapid development of deep learning methods. However, even with deep  
48 learning, it still requires tremendous efforts to achieve human-level performance on this challenging  
49 task. There were a few successful experiences to learn from the computer science community to  
50 make the deep learning method fully comparable to humans in connectomics. The success of deep  
51 learning methods highly depends on the amount of training data and the quality of annotation.  
52 Take the task of image classification as an example; ImageNet (*Russakovsky et al., 2015*) has set  
53 up a research paradigm of applying deep learning methods for vision tasks. In 2009, by releasing  
54 a large-scale accurately annotated dataset, ImageNet provided a benchmark (72%) for image  
55 classification. From 2010 to 2017, a challenge called "The ImageNet Large Scale Visual Recognition  
56 Challenge (ILSVRC)" was organized every year. This challenge significantly boosted the development  
57 of deep learning algorithms. Many champions of this challenge have become the milestones of  
58 deep learning methods, such as AlexNet (*Krizhevsky et al., 2012*), VGG (*Simonyan and Zisserman,*  
59 *2014*), GoogleNet (*Szegedy et al., 2015*), ResNet (*He et al., 2016*), etc. As shown in *Figure 1*, deep  
60 learning performance on image classification finally exceeded human level (95%) after eight years of  
61 development. To summarize, there is a roadmap for the success of ImageNet, which includes three  
62 key steps: the first step is to establish a large-scale dataset with high-quality annotation, which is  
63 very important for deep learning. Based on the dataset, the second step is organizing a challenge  
64 that can evaluate algorithms at a large scale and allow researchers to estimate the progress of their  
65 algorithms, taking advantage of the expensive annotation effort. The third step is the design of new  
66 algorithms based on the previous two steps. Each of the three stages is indispensable.

67 Following the success of ImageNet, significant progress of EM automatic segmentation was  
68 achieved by the 2012 IEEE International Symposium on Biomedical Imaging (ISBI 2012), which was  
69 the first challenge on EM automatic segmentation with releasing a publicly available dataset (*Arganda-*  
70 *Carreras et al., 2015*). The state-of-the-art (SOTA) method exhibited unprecedented accuracy in

71 EM cellular segmentation on the dataset of ISBI 2012. In particular, the deep learning method  
72 "U-Net" (*Ronneberger et al., 2015*), which was first proposed during the challenge, becomes the  
73 backbone of many SOTA methods in the field. However, today many deep learning methods have  
74 become "exceedingly accurate" and likely saturated at the ISBI 2012 (*Arganda-Carreras et al., 2015*).  
75 In addition, ISBI 2012 images are  $512 \times 512$  pixels with a resolution of  $4 \text{ nm/pixel} \times 4 \text{ nm/pixel}$ ,  
76 while there are many EM images with higher resolution in connectomics because enough high  
77 resolution is essential to unravel the neural structures unambiguously. For instance, 2nm has been  
78 suggested as the historical "gold standard" to identify synapses (*DeBello et al., 2014*), in particular  
79 to identify gap junctions (*Leitch, 1992*) which are common in neural tissues (*Anderson et al., 2009*).  
80 It is not clear if previous classic deep learning methods developed on EM images with relatively  
81 lower resolution can still work well on datasets with higher resolutions.

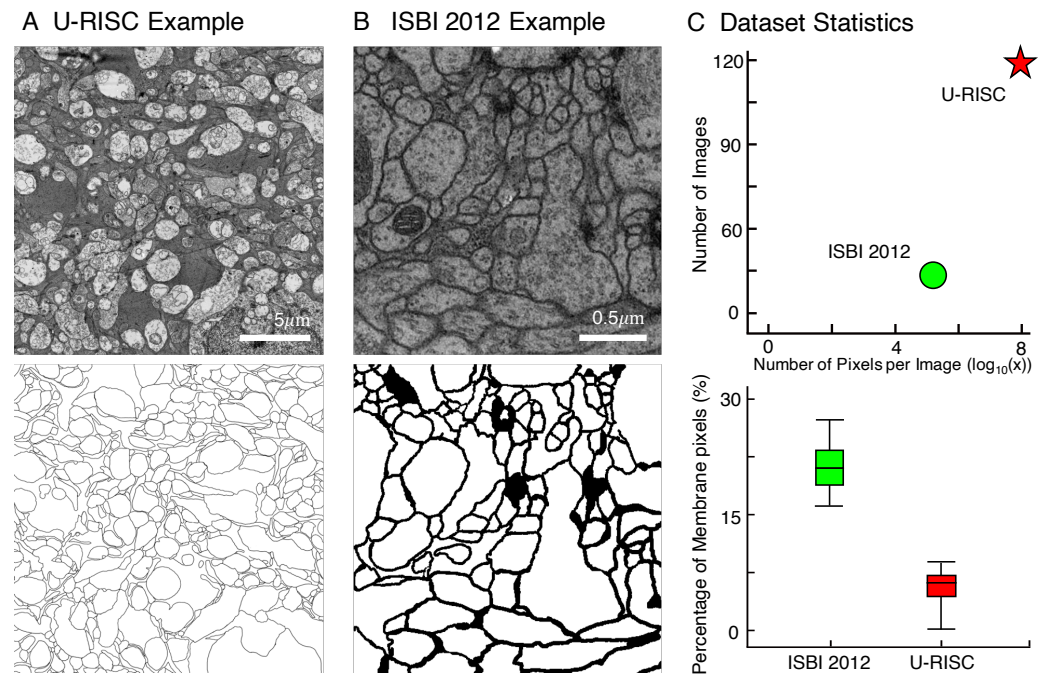
82 Here, to promote the deep learning algorithms in EM datasets, we initiate a new roadmap: We  
83 first annotated the retinal connectomic data, RC1, from rabbit (*Anderson et al., 2011*) and presented  
84 a brand new annotated EM dataset named U-RISC (Ultra-high Resolution Image Segmentation  
85 dataset for Cell membrane). Compared with ISBI 2012, U-RISC has a higher resolution of 2.18  
86 nm/pixel and a larger size of  $9958 \times 9959$  pixels. The precision of annotation was ensured by  
87 multi-steps of iterative verification, costing over 10,000 labour hours in total. Next, based on  
88 U-RISC, a competition of cellular membrane prediction was also organized. Surprisingly, from  
89 448 domestic participants/teams, the top performance of deep learning methods on U-RISC (~  
90 0.6, F1-score) was far below the human-level accuracy ( $> 0.9$ ), in contrast with the near-human  
91 performance of deep learning methods in ISBI 2012. We then made fair comparisons between  
92 ISBI 2012 and U-RISC with the same segmentation methods, including U-Net. The comparison  
93 results confirmed that U-RISC indeed provides new challenges to existing deep learning methods.  
94 U-Net, for example, dropped from 0.97 in ISBI 2012 to 0.57 in U-RISC. To further explore how  
95 these methods work on segmentation tasks, we introduced a gradient-based attribution method,  
96 integrated gradient (*Sundararajan et al., 2017*), to analyze ISBI 2012 and U-RISC. The result showed  
97 that when deciding on whether a pixel belonged to the cell membrane or not, deep learning  
98 methods represented by U-Net would refer to a larger attribution region on U-RISC (about four  
99 times on average) than that on ISBI 2012. It suggests that the deep learning methods might  
100 require more background information to decide the segmentation of the U-RISC dataset. Finally, we  
101 integrated current available advanced methods, combining U-Net and transfer learning recently  
102 introduced (*Conrad and Narayan, 2021*), and provided a benchmark (0.6659), about 10% higher  
103 than the leader board (0.6070), for U-RISC.

104 Overall, our contribution in this work lies mainly in the following three parts: (1) we provided the  
105 community a brand new publicly available annotated mammalian EM dataset with known highest  
106 resolution (~ 2.18 nm/pixel) and largest image size ( $9958 \text{ pixel} \times 9959 \text{ pixel}$ ); (2) we organized  
107 a competition and made a comprehensive analysis to reveal the challenges of U-RISC for deep  
108 learning methods; (3) we improved the benchmark with 10% to the F1-score of 0.6659. In discussion,  
109 we proposed further suggestions for improving segmentation methods from perspectives of model  
110 design, loss function design, data processing, etc. We hope our dataset and analysis can help  
111 researchers gain insights into designing more robust methods, which can finally accelerate the  
112 speed of untangling brain connectivity.

## 113 Results

### 114 The largest ultra-high-resolution EM cell membrane segmentation dataset

115 Along with this paper, we proposed a new EM dataset with cell membrane annotated: Ultra-high  
116 Resolution Images Segmentation dataset for Cell membrane (**U-RISC**). To our best knowledge, U-  
117 RISC has the highest resolution among the publicly available annotated EM datasets (see *Figure 2(A)*  
118 as an example). It was annotated upon the rabbit retinal connectomic dataset RC1 (*Anderson et al.,*  
119 *2011*) with a 2.18 nm/pixel resolution at both  $x$  and  $y$  axes. The size of individual image and the



**Figure 2.** Comparison between U-RISC and ISBI 2012. (A and B) An example of U-RISC and ISBI 2012 data includes the raw EM image (top) and the corresponding annotation result (bottom). Black pixels in annotation results represent cellular membranes. (C) (Top) Both the number and size of images in U-RISC surpass those in ISBI 2012. (Bottom) The proportion of annotated pixels,  $5.10\% \pm 2\%$  in ISBI 2012 and  $21.65\% \pm 2\%$  in U-RISC, making latter a more imbalanced dataset.

120 total number of images in U-RISC both exceed the published datasets by far, taking ISBI 2012 as  
121 an example (**Figure 2(C)**) (120 pairs of  $9958 \text{ pixel} \times 9959 \text{ pixel}$  images in U-RISC and 30 pairs of  $512$   
122  $\text{pixel} \times 512 \text{ pixel}$  images in ISBI 2012). One characteristic of U-RISC is that cell membranes only  
123 cover a small area of the images, making it an imbalanced dataset for deep learning (an average of  
124  $5.10\% \pm 2\%$  in U-RISC compared with  $21.65\% \pm 2\%$  in ISBI 2012).

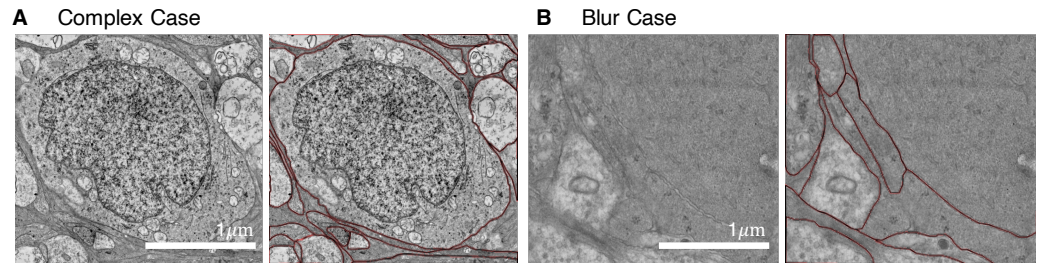
125 We employed an iterative manual annotation procedure to ensure the quality of annotation. Be-  
126 cause of the difficulty of distinguishing cell membrane from organelle membrane, special attention  
127 was paid to exclude organelle membrane from annotation (**Figure 3(A)**). In practical connectomic  
128 research, the image quality can be affected by many reasons like insufficient staining, thick section,  
129 etc. Considering this, we retained several images with low quality in U-RISC to make the dataset  
130 closer to the actual situation. Annotation on these images costs more time and caution (**Figure 3(B)**).  
131 Labeling errors could be detected and then corrected in each round of iteration (**Figure 4**). For  
132 scientific research reasons, the human labeling process is very valuable for uncovering the human  
133 learning process. Therefore, the intermediate annotated results were also reserved for public  
134 release (<https://brain.baai.ac.cn/biodb-rabbit-details.html>).

### 135 Ultra-high resolution EM images segmentation competition

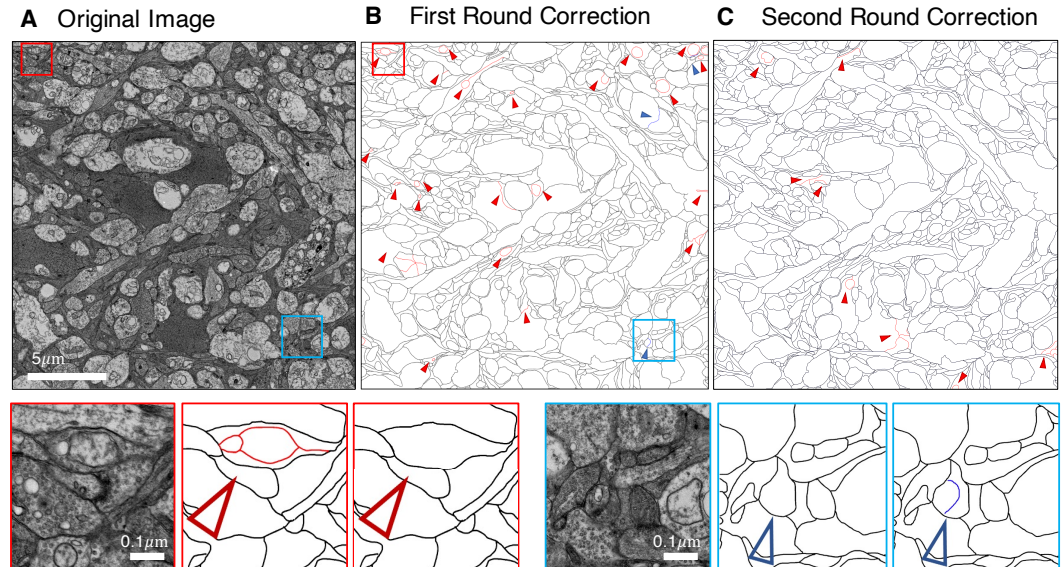
136 To investigate the performance of the deep learning methods on U-RISC and propose a benchmark,  
137 a competition about cellular membrane segmentation was organized by BAAI (Beijing Academy  
138 of Artificial Intelligence Institution, Beijing, China) and PKU (Peking University, Beijing, China)<sup>1</sup>. In  
139 total, 448 participants took part in the competition, mainly from domestic competitive universities,  
140 research organizations, and top IT institutions.

141 There were two tracks in the competition (**Table 1**): Track 1 used original images with the size of

<sup>1</sup><https://www.biendata.xyz/competition/urisc/>



**Figure 3.** Examples of images with their annotations. (A) Organelle membranes were cautiously avoided to be annotated. (B) More time and patience were needed to annotate the image with low contrast.



**Figure 4.** Example of iterative human annotation. (A) Original image to be annotated. (B) Many errors were found out in the first round of annotation. (C) After correction, much fewer errors were detected in the second round of annotation, and the correction results were served as the final annotation. Red small triangles and boxes indicate false positive errors (enlargement in the bottom left), blue for false-negative errors (enlargement in the bottom right).

142 9958 pixel  $\times$  9959 pixel as training and testing datasets. In Track 2, images were downsampled to  
143 the size of 1024 pixel  $\times$  1024 pixel. The purpose of Track 2 was to allow researchers with limited  
144 computational resources to participate in the competition. The final round of human annotation  
145 was used as the ground truth to evaluate the algorithms, and an F1-score was applied as the  
146 evaluation metric (for details, please see Methods and Materials).

147 Surprisingly, from the competition, top 6 teams in each track gained F1-scores around 0.6 on  
148 U-RISC, which were far below the human levels (0.92 and 0.99, the first and second rounds of  
149 annotation). However, previous research has shown that the performance of the top teams in  
150 ISBI 2012 had already been reasonably close to the human level (*Arganda-Carreras et al., 2015*).  
151 To investigate causes of the performance gap between the methods and humans on U-RISC, we  
152 first surveyed the top 6 teams in our competition. It indicated that a variety of current popular  
153 approaches to segmentation were utilized (*Figure 5*). From the choice of models (*Figure 5(A)*), the  
154 participants used current popular image segmentation networks, such as U-Net (*Ronneberger*  
155 *et al., 2015*), Efficientnet (*Tan and Le, 2019*) and CASENet (*Yu et al., 2017*). For backbone selection,  
156 ResNet (*He et al., 2016*) and their variants were the most chosen architectures. Data augmentation  
157 was ubiquitously applied to improve the generalization of the models. About 13% of the participants

**Table 1.** Leaderboard of Track 1 and Track 2.

Track 1 (Original)			Track 2 (Downsample)		
Team Name	Institution	F1-score	Team Name	Institution	F1-score
Human 1st		0.92128± 0.012	Human 1st		0.96915± 0.014
Human 2nd		0.99334± 0.008	Human 2nd		0.99891± 0.003
SCP173	Tencent <sup>1</sup>	0.60704± 0.043	horch	UCAS <sup>2</sup>	0.56932± 0.053
yangsenwxy	SCU <sup>3</sup>	0.60701± 0.042	deadline	NJU <sup>4</sup>	0.56213± 0.055
SpongeBobbb	HDU <sup>5</sup>	0.60480± 0.042	SpongeBobbb	HDU <sup>5</sup>	0.56136± 0.049
VIDAR	USTC <sup>6</sup>	0.60303± 0.041	VIDAR	USTC <sup>6</sup>	0.55170± 0.046
deadline	NJU <sup>4</sup>	0.60066± 0.045	archer	THU <sup>7</sup>	0.55107± 0.047
Chasingstar	JLU <sup>8</sup>	0.59647± 0.044	scu_ws	SCU <sup>3</sup>	0.54847± 0.053

Mean and standard error are computed over test images.

158 used Hypercolumns (*Hariharan et al., 2015*) to improve the expressiveness of the model. From the  
 159 design of loss function, functions that can adjust penalty ratios according to sample distributions  
 160 were applied to reduce the effect of sample imbalance, such as Dice loss (*Dice, 1945*), Focal loss (*Lin*  
 161 *et al., 2017*), and BCE loss (*Cui et al., 2019*). And Adam (*Kingma and Ba, 2014*) was shown to be the  
 162 most chosen optimization method.

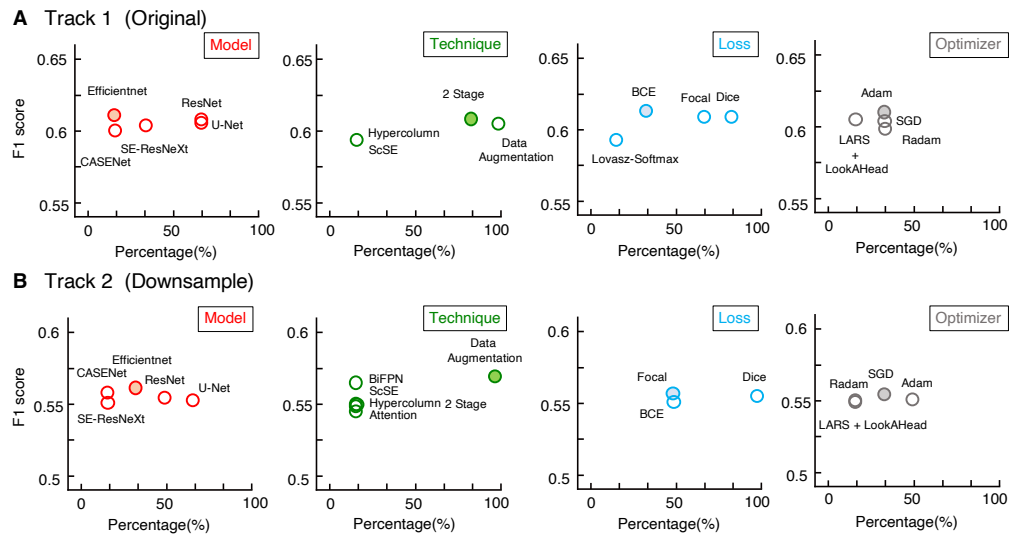
163 The analysis suggested that even though participants had considered many popular methods,  
 164 their performance was still not satisfactory and varied only slightly between each other. To identify  
 165 whether this was because of the challenges of U-RISC or the methods themselves, we picked out  
 166 three widely used methods: U-Net (*Ronneberger et al., 2015*), LinkNet (*Chaurasia and Culurciello,*  
 167 *2017*), and CASNet (*Yu et al., 2017*). We conducted a fair comparison between the performance of  
 168 each method on U-RISC (Track 1) and ISBI 2012. Results showed that these methods could reach  
 169 over 0.97 (F1-score) in ISBI 2012, but only between 0.57-0.61 in U-RISC (*Table 2*), which confirmed  
 170 that the performance gap in competition comes from the challenges of U-RISC.

171 What are the unique challenges brought by  
 172 **Table 2.** F1-scores in U-RISC and ISBI 2012. U-RISC to deep learning algorithms? Two types  
 173 of errors were analyzed first: false-positive errors, which led to incorrect membrane predic-  
 174 tions, and false-negative errors, which caused  
 175 discontinuity of cell membrane. According to our  
 176 analysis, both false-positive errors (pink boxes)  
 177 and false-negative errors (orange boxes) were  
 178 common in U-RISC, which were rare in ISBI 2012  
 179 (*Figure 6(B) (C)*). Further investigations for the networks are required to explore the reason and find  
 180 ways to reduce the errors.  
 181

### 182 Attribution analysis of the deep learning method on U-RISC and ISBI 2012

183 To acquire a deeper understanding of the different performances in U-RISC and ISBI 2012, we per-  
 184 formed an attribution analysis (*Ancona et al., 2019*) on the trained U-Net. We selected the gradient-

<sup>1</sup>Tencent Holdings Ltd (China)  
<sup>2</sup>University of Chinese Academy of Sciences (China)  
<sup>3</sup>Sichuan University (China)  
<sup>4</sup>Nanjing University (China)  
<sup>5</sup>Hangzhou Dianzi University (China)  
<sup>6</sup>University of Science and Technology of China  
<sup>7</sup>Tsinghua University (China)  
<sup>8</sup>Jilin University (China)



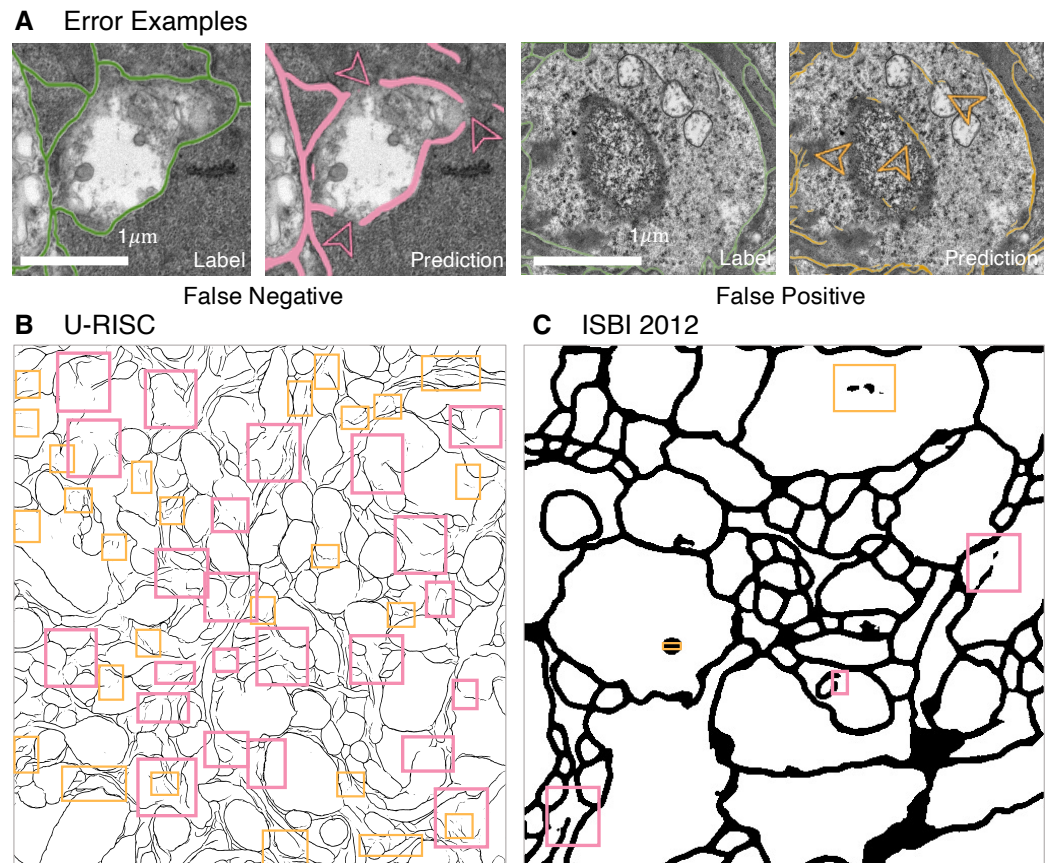
**Figure 5.** Mean F1 scores of teams with different methods used. (A and B) The statistics of Track 1 and Track 2, respectively. The X-axis represents the proportion of the team with the method, Y-axis represents the average of F1 scores.

185 based attribution method, **integral gradient (IG)** (Sundararajan et al., 2017), which is widely  
 186 applied on explainable artificial intelligence, such as understanding feature importance (Adadi and  
 187 Berrada, 2018), identifying data skew (Clark et al., 2019), and debugging model performance (Guidotti  
 188 et al., 2018). In brief, IG aims to explain the relationship between predictions and input features  
 189 based on gradients (Figure 7(A)). The IG output is plotted in **Attribution Fields** to reflect their con-  
 190 tribution to the final prediction. In the heatmap, each pixel was assigned with a normalized value  
 191 between [-1, 1]. With IG, we analyzed the attribution field of each predicted pixel of U-Net in U-RISC  
 192 and ISBI 2012. Color and shade were used to represent the normalized contribution values in  
 193 attribution fields (Figure 7(B)). For a fair comparison between U-RISC and ISBI 2012, areas of **Pixel**  
 194 **Attribution Fields**  $S_k$  were converted to physical size according to their respective resolutions.

195 **Figure 8** shows the examples of attribution fields, where bounding boxes with different colors  
 196 represented different pixel classifications, green for a correct predicted pixel, orange for a false  
 197 positive error, and pink for a false negative error. We noticed that the areas of attribution fields  
 198  $S_k$  of two datasets were both relatively minor to the whole images (Figure 7(B)). For example, at  
 199 the threshold of  $k > 0.01$ , the  $S_k$  of the correct cases accounted for only 5.1% and 0.8% relative to  
 200 the whole image (the green bounding boxes in Figure 8). This suggested that U-Net would focus  
 201 on local characteristics within small areas of the images when making predictions. In addition, we  
 202 found that the averaged  $S_k$  of each predicted pixel in U-RISC was significantly larger than that in  
 203 ISBI 2012, specifically  $46000 \text{ nm}^2$  in U-RISC and  $10300 \text{ nm}^2$  in ISBI 2012. Taken together, U-Net would  
 204 predict cell membrane according to local information around the pixel, and the average attribution  
 205 field was larger in U-RISC than that of ISBI 2012. All of these indicate that more information is  
 206 required for the segmentation in U-RISC.

### 207 **U-Net-transfer model achieve the state-of-the-art result on U-RISC benchmark**

208 Considering both the comprehensive analyses of competition and attribution analysis, we integrated  
 209 outperformed methods to develop our method (Figure 9(A)). For basic segmentation architecture,  
 210 we chose U-Net due to its better characteristic extraction ability. Many valuable techniques were  
 211 also considered, including cross-crop strategy for saving computational resources and data augmen-  
 212 tation to increase data diversity. We chose both focal loss and dice loss to deal with the imbalance  
 213 of samples for loss function design. Some parameters used for training were also optimized, such



**Figure 6.** Errors in segmentation predictions of U-RISC and ISBI 2012. (A) The examples of False Positive and False Negative errors. (B and C) The examples of two errors in the segmentations of U-RISC and ISBI 2012. Pink arrows and lines represent false-negative errors, and orange represents false positive errors.

**Figure 6-Figure supplement 1.** Supplements for segmentation predictions of U-RISC.

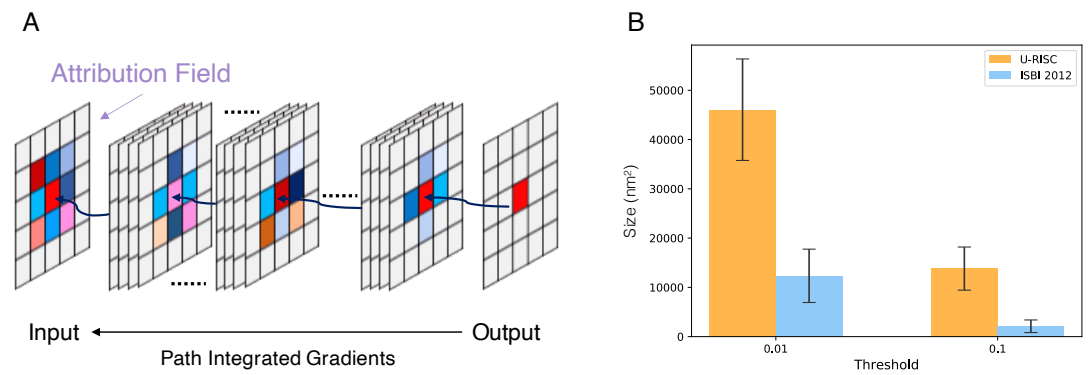
**Figure 6-Figure supplement 2.** Supplements for segmentation predictions of ISBI 2012.

214 as batch-size/GPU (4) and the number of GPUs (8). For more details, please refer to the part of  
215 Segmentation Networks in Methods and Materials. Especially, recent research has shown that trans-  
216 fer learning with domain-specific annotated datasets could be effective in elevating deep learning  
217 models' performance (Conrad and Narayan, 2021). Therefore, we introduced a pre-trained model,  
218 trained with MoCoV2 (He et al., 2020) on CEM500K (Conrad and Narayan, 2021). The segmentation  
219 result showed that the F1 scores of our method were 10% higher than the leader of competition  
220 (0.66 vs 0.61 in Table 2). Thus we provide a new benchmark on the cellular membrane segmentation  
221 of U-RISC.

## 222 Discussion

223 This paper first proposed the U-RISC, a cell membrane EM dataset created through intensive  
224 and elaborate annotation. The dataset is characterized by the highest resolution and the largest  
225 single image size compared with other current publicly available annotated EM datasets. Next,  
226 we organized a segmentation competition on U-RISC and proposed the benchmark. During the  
227 competition, we noticed that the performances of popular deep learning methods were far below  
228 that of humans, which motivated us to explore the causes. Thus, we carried out a comprehensive  
229 survey on the deep learning methods participants applied in the competition. To our surprise,  
230 methods such as U-Net, LinkNet, and CASENet exhibited a significant drop of F1-score on U-RISC





**Figure 7.** Attribution analysis. (A) Integrated gradients attribution method. (B) Statistics of attribution filed for U-RISC and ISBI 2012.

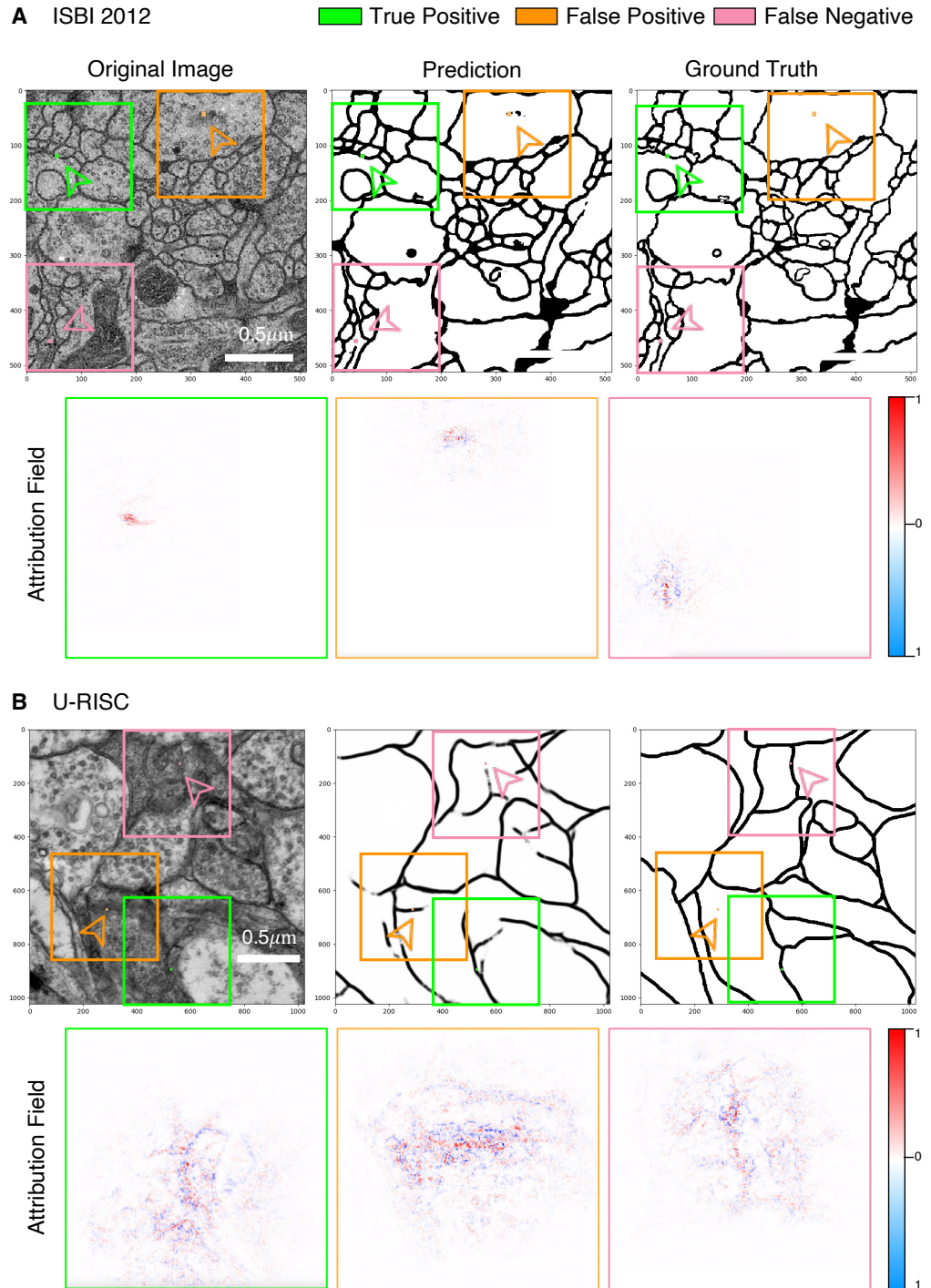
231 compared with ISBI 2012, from 0.9 to 0.6. To explore the mechanisms underlying this discrepant  
232 performance, we introduced a gradient-based attribution method, integrated gradient. Through  
233 attribution analysis of U-Net, we found the average pixel attribution field of U-RISC is larger than that  
234 of ISBI, corresponding to the size of cellular structure, and both of them are relatively small to the  
235 whole image size. By integrating currently available methods, we improve the benchmark to 0.67,  
236 about 10% higher than the top leader from the competition. Based on the analyses in this paper,  
237 here we raise some considerations about the challenges for deep learning-based segmentation  
238 algorithms brought by U-RISC and propose several suggestions for improving EM segmentation  
239 methods.

### 240 **Challenges for Deep Learning-based Segmentation**

241 Benchmark showed that the segmentation performance of deep learning algorithms on U-RISC was  
242 still far behind the human level. U-RISC poses challenges for deep learning-based segmentation  
243 in the following aspects: (1) high computational costs needed to deal with large images, (2) the  
244 extreme sample imbalance caused by low ratio of cellular membrane pixels in the whole image, (3)  
245 side effects of typical data processing methods.

246 Deep learning itself is already a computationally intensive method. It would require more  
247 computational resources to process the images with a much larger size in U-RISC. In practical  
248 terms, taking U-Net as an example, processing a 1024 pixel  $\times$  1024 pixel image requires a GPU  
249 with 12GB memory. This memory is enough to deal with the images in ISBI 2012, of which the size  
250 is 512 pixel  $\times$  512 pixel. But the size of a single image in U-RISC is 9958 pixel  $\times$  9959 pixel, which  
251 is far beyond the processing ability of the commonly used 12 GB memory GPU. Therefore, the  
252 additional computational burden brought by U-RISC raises the first challenge for deep learning-  
253 based segmentation.

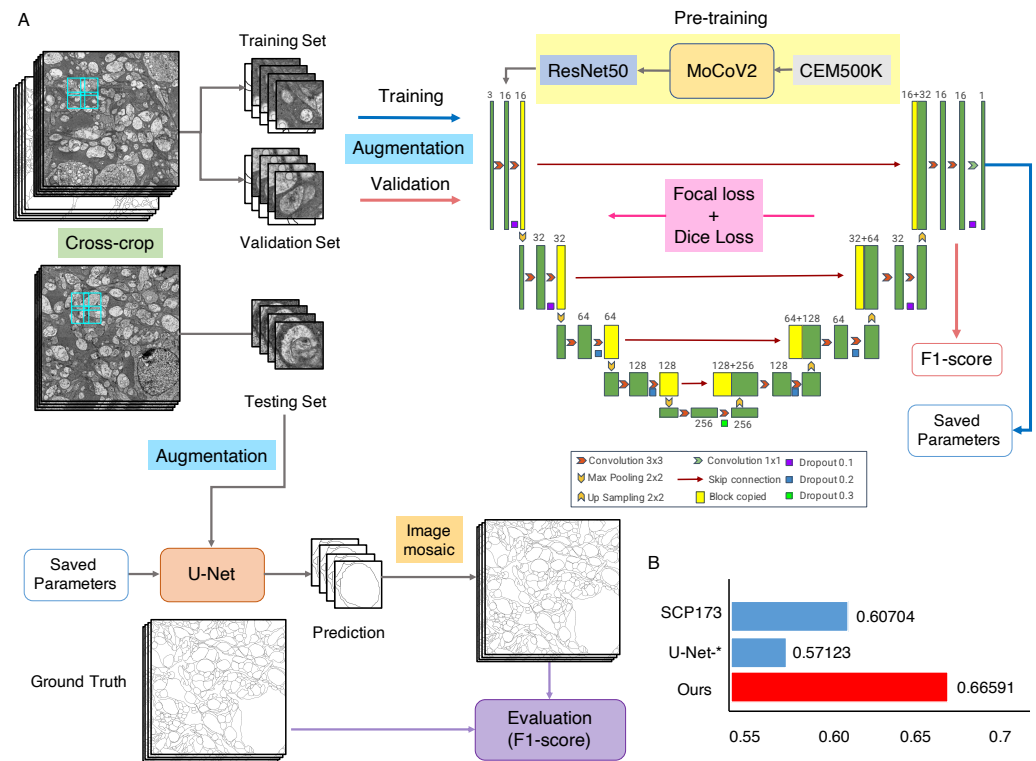
254 The problem of imbalanced samples widely exists in computational vision tasks (*Alejo et al.,*  
255 *2016; Li et al., 2010; Zhang et al., 2020*), which should be considered when designing algorithms.  
256 Cellular membrane segmentation is a typical situation of sample imbalance because cellular mem-  
257 brane only occupies a small proportion of the whole cell structure. According to statistics, the pixels  
258 belong to the cellular membrane account for 21.65% of the entire pixels of ISBI 2012. While the  
259 proportion in U-RISC is much smaller, 5.10%, making U-RISC an extremely imbalanced dataset.  
260 Pre-existing solutions were mainly proposed from several aspects: loss function design (*Lin et al.,*  
261 *2017; Cui et al., 2019*), data augmentation (*Yoo et al., 2020*), under/over-sampling (*Fernández et al.,*  
262 *2018; Yen and Lee, 2009*), and semantically multi-modal approaches (*Zhu et al., 2020*). However,  
263 even though the participants in the competition already used these approaches, the final results  
264 showed limited improvement of segmentation. So the imbalanced problem of U-RISC is yet to be



**Figure 8.** Attribution analysis. (A and B) Attribution fields of ISBI 2012 and U-RISC dataset. The first line represents the original image, network prediction result and annotation respectively. The pixels pointed by green (correct cell membrane pixel), orange (false positive predicted pixel) and pink (false negative predicted pixel) arrows are the prediction points used in attribution method. Images in the three color boxes with the same size in the second line represent the attribution field corresponding to the above three pixels. Blue means the network is likely to predict the pixels as cell membrane, while opposite for red.

**Figure 8–Figure supplement 1.** Supplements for attribution analysis on ISBI 2012.

**Figure 8–Figure supplement 2.** Supplements for attribution analysis on U-RISC.



**Figure 9.** The U-Net-transfer method achieves the best performance on U-RISC. (A) The pretraining, training, and testing processing for U-Net. (B) The comparison of the F1-scores. "SCP-173" represents the top performance in the competition. U-Net-\* represents the performance in *Table 2*. Ours represents the performance of the U-Net-transfer method in this section.

265 solved and becomes another challenge for deep learning-based segmentation.

266 Proper data processing is essential and helpful to deep learning algorithms. For example, a  
267 downsampling process on raw images with an enormous size is commonly adopted in segmentation  
268 tasks (*Marin et al., 2019; Chen et al., 2014*). And in the Track 2 of our competition, we used the  
269 downsampled dataset to reduce computational consumption as usual. Surprisingly, we found that  
270 the F1-score of the same method dropped as well as the overall performance decreased in Track 2  
271 compared with Track 1. We speculated a key reason might be the degradation of image quality from  
272 track 1 to Track 2. We confirmed the quality reduction through 4 representative indexes, including  
273 Brenner (*Subbarao and Tyan, 1998*), Variance (*Subbarao and Tyan, 1998*), SMD2 (*Thakkinstian*  
274 *et al., 2005*), and Vollath (*Vollath, 2008*) (shown in Appendix). More cautions should be paid when  
275 using traditional data processing methods, and more advanced data processing theories are  
276 expected from this point of view.

### 277 **Suggestions for the Improvement of Segmentation Methods**

278 To some degree, increasing computational resources is a possible way to cope with the challenges  
279 mentioned above. However, it might not be easy for all the community researchers to access  
280 sufficient computational power; therefore innovations in algorithms are still crucial for our future  
281 success. To improve the performance of deep learning in EM segmentation, we provide several  
282 suggestions for developing deep learning algorithms from the following perspectives: model design,  
283 training techniques, data processing, loss function design, and visualization tools.

284 **Model Design.** As shown in the attribution analysis, the current models for segmentation, such  
285 as U-Net (*Ronneberger et al., 2015*), Efficientnet (*Tan and Le, 2019*), and CASENet (*Yu et al., 2017*),  
286 are designed to focus on local information to make predictions. However, in a high-resolution image,  
287 other structures, like organelle membrane and synaptic vesicles, might share similar features with  
288 the cellular membrane in a local scale, which leads to false-positive results. And this constitutes one  
289 of the major error types in the competition. Therefore, it might not be enough for the classifiers of  
290 a model to make correct decisions with only local features. And studies have shown that models  
291 using global information could improve performance greatly (*Liu et al., 2018, 2020; Wang et al.,*  
292 *2019*). Therefore, more global information could also be considered in the future design of the  
293 segmentation network.

294 **Training Techniques.** Skillful training techniques also can be helpful to improve the segmenta-  
295 tion performance. According to our survey, a two-stage training strategy could be much better than  
296 a single-stage training strategy. A recent work also suggests that pretraining with domain-specific  
297 datasets can help network learning domain features (*Conrad and Narayan, 2021*). Besides that,  
298 much experience can be learnt from existing training methods. The Hypercolumns module (*Har-*  
299 *iharan et al., 2015*) is used to accelerate the convergence of training by combining features at  
300 different scales, and the combination of features from different scales can help bring in global  
301 information. ScSE (*Roy et al., 2018*) module introduces attention mechanism into network, thus  
302 bringing in global information. Hybrid architectures can also be considered because of its ability to  
303 expand the receptive field (*Goceri, 2019*). In a word, improvement can be made at the phase of  
304 training by utilizing advanced training techniques.

305 **Data Processing.** Data processing is commonly used in deep learning, while traditional down-  
306 sampling methods were shown to have side effects in the competition. To alleviate the side effects,  
307 some quality enhancing methods for downsampled images could be expected, such as edge and  
308 region-based image interpolation algorithms (*Asuni and Giachetti, 2008; Hwang and Lee, 2004*),  
309 low bit rate-based approaches (*Wu et al., 2009; Lin and Dong, 2006*), and quality assessment  
310 researches (*Vu et al., 2018; Wang and Bovik, 2006; Wang et al., 2003*). Meanwhile, other data  
311 processing methods can also be taken into account. Take data augmentation as an example; by  
312 augmenting the training data randomly, the dependence of the model on specific attributes can be  
313 reduced, which can be beneficial in EM segmentation with many imbalanced samples.

314 **Loss Function Design.** Loss function design is another important part of deep learning. But

315 many current loss functions have their own disadvantages in our competition. For example,  
316 Dice loss (*Dice, 1945*) was designed to optimize F1-score directly, without consideration of data  
317 imbalance. Focal loss (*Lin et al., 2017*) and BCE loss (*Cui et al., 2019*) were used in the competition  
318 to care more about data imbalance by giving different penalty according to sample difficulty, but the  
319 improvement was limited as the results have shown. A better design of loss function should take  
320 overall consideration of both sample imbalance and evaluation criteria. While the most common  
321 evaluation criterion, F1-score, a pixel-based statistic, is inconsistent with human subjective feeling  
322 to some extent. It might be a major cause of the performance gap between humans and algorithms.  
323 Some other structure-based criteria have appeared, such as V-Rand and V-info (*Arganda-Carreras*  
324 *et al., 2015*) integrating skeleton information of cell membrane, and ASSD (*Heimann et al., 2009*)  
325 considering the distance of point sets.

326 **Visualization tools.** Visualization tools can help us have a better understanding of the network.  
327 In this paper, with IG, we could learn the attribution fields of U-Net from the view of gradient,  
328 which inspires us to improve deep learning methods by paying more attention to global informa-  
329 tion. In comparison, many other visualization tools are starting from other characteristics of the  
330 network. Layer-wise relevance propagation (LRP) (*Bach et al., 2015*) and deep Taylor decompo-  
331 sition (DTD) (*Montavon et al., 2017*) get attribution distribution by modifying propagation rules.  
332 Information-based method IBA (*Schulz et al., 2020*) restricts the flow of information to accomplish  
333 attribution fields. Combining different visualization tools can help to promote much more insightful  
334 inspiration in improving deep learning methods.

335 Overall, we provide an annotated EM cellular membrane dataset, U-RISC, and its benchmark. It  
336 indeed brings many challenges for deep learning and promotes the development of deep learning  
337 methods for segmentation.

## 338 **Methods and Materials**

### 339 **Dataset**

340 The U-RISC dataset was annotated upon RC1, a large-scale retinal serial section transmission  
341 electron microscopic (ssTEM) dataset, publicly available upon request and detailedly described  
342 in the work of *Anderson et al. (2011)*. RC1 came from the retina of a light-adapted female Dutch  
343 Belted rabbit after in vivo excitation mapping. The imaged volume represents the retinal tissue with  
344 a diameter of 0.25 mm, spanning the inner nuclear, inner plexiform, and ganglion cell layers. Serial  
345 EM sections were cut at 70-90 nm with a Leica UC6 ultramicrotome and captured at the resolution  
346 of 2.18 nm/pixel across both axes using SerialEM (*Mastrorade, 2005*). In RC1, there are in total 341  
347 EM mosaics generated by the NCR Toolset (*Anderson et al., 2009*), and we clipped out 120 images  
348 in the size of 9958 pixel  $\times$  9959 pixel from randomly chosen sections.

349 To annotate cell membrane with high quality on 120 images, we launched an iterative annotation  
350 project that lasted for three months. All the annotators were trained to recognize and annotate  
351 cellular membrane in EM images, but only two-thirds of all, 53 annotators, were finally qualified  
352 to participate in the project according to their annotation results. In the iterative annotation  
353 procedure, each EM image would undergo three continuous rounds of annotation with the guidance  
354 of blind review. The final round of annotation was regarded as the “ground truth”. While since  
355 the first two rounds are valuable for analyzing the human learning process, we also reserved the  
356 intermediate results for public release. All of the U-RISC datasets are released at [https://brain.baai.  
357 ac.cn/biodb-rabbit-details.html](https://brain.baai.ac.cn/biodb-rabbit-details.html).

### 358 **Competition**

359 The goal of the competition was to predict cell membranes in the EM images of U-RISC. Participants  
360 were required to return images depicting the boundary of all neurons. F1-score was selected as  
361 the evaluation criterion for the accuracy of the results. There are two tracks in the competition,  
362 images in Track 1 were kept as the original size (9958 pixel  $\times$  9959 pixel), images in Track 2 were

363 downsampled to the size of 1024 pixel  $\times$  1024 pixel. Fifty images, 30 as the training dataset and 20  
364 as the test dataset, were released in Track 1. And Track 2 contained 70 images in total, 40 training  
365 images, and 30 testing images. The training dataset included EM images with their corresponding  
366 ground truth, while the ground truth of the test dataset was kept private. In both tracks, ten images  
367 from the training dataset served as the validation dataset for the participants to monitor and  
368 develop their models. No statistical methods were used to determine the assignment of images in  
369 the whole arrangement.

### 370 Segmentation Networks

371 We conducted experiments to compare the performance of the same methods on U-RISC (Track2)  
372 and ISBI 2012. Three representative deep learning networks were considered (**Table 2**), U-Net (**Ronneberger et al., 2015**), LinkNet (**Chaurasia and Culurciello, 2017**), and CASENet (**Yu et al., 2017**).  
373 The three networks are all pixel-based segmentation networks. To be specific, given the input image  
374  $x$ , the goal of the networks is to classify the corresponding semantic cell membrane pixel by pixel.  
375 For the input image  $x$  and classification function  $F(x)$ ,  $Y_{\{p|X, \Theta\}} \in [0, 1]$  is taken as the output of  
376 the network, which represents the edge probability of the semantic category of the pixel  $p$ .  $\Theta$  are  
377 the parameters in the network and are optimized in the training processes. Architectures of the  
378 three networks are described as follows.  
379

380 **U-Net.** U-Net (**Ronneberger et al., 2015**) is a classical fully convolutional network (that is, there  
381 is no fully connected operation in the network). The model is composed of two parts: contracting  
382 path and expansive path. The contracting path follows the typical architecture of a convolutional  
383 network. At each downsampling step, U-Net doubles the number of feature channels to gain a  
384 concatenation with the correspondingly cropped feature map from the contracting path. At the final  
385 layer a  $1 \times 1$  convolution is used to map each 64-component feature vector to the desired number  
386 of classes. In total the network has 23 convolutional layers. We use ResNet50 as its encoder.

387 **LinkNet.** The model structure of LinkNet (**Chaurasia and Culurciello, 2017**) is almost similar to  
388 U-Net, which is a typical encoder-decoder structure. The encoder starts with an initial block which  
389 performs convolution on input image with a kernel of size  $7 \times 7$  and a stride of 2. This block also  
390 performs spatial max-pooling in an area of  $3 \times 3$  with a stride of 2. The later portion of encoder  
391 consists of residual blocks and is represented as the encoder-block. To reduce parameters, LinkNet  
392 uses ResNet18 as its encoder.

393 **CASENet.** CASENet (**Yu et al., 2017**) is an end-to-end deep semantic edge learning architecture  
394 adopting ResNet-101 as backbone. The classification module here consists of a  $1 \times 1$  convolution  
395 and a bilinear interpolation up-sampling layer to generate  $M$  active images, each image size is  
396 the same as the original image. Each residual block is followed by a classification module to  
397 obtain five classification activation graphs. Then, a sliced concatenation layer is used to fuse the  
398  $M$  classification activation graphs, and finally a  $5 \times M$  channel activation graph is obtained. The  
399 activation graphs are used as the input of the fused classification layer to obtain a  $M$ -channel  
400 activation graph. The fusion classification layer is convolution of  $M$  group  $1 \times 1$ .

401 **Transfer learning.** The pre-trained model from (**Conrad and Narayan, 2021**) was used in our  
402 method, specifically MoCoV2 (**He et al., 2020**) and CEM500K (**Conrad and Narayan, 2021**) were  
403 respectively selected as the pretraining method and dataset.

404 **Training settings.** For each dataset, same training and testing data distribution was utilized on  
405 the three methods. For U-RISC, during the training, the original images were cut into  $1024 \times 1024$   
406 patches with overlaps. And the patches were randomly assigned into the training set and validation  
407 set according to the ratio of 50,000/20,000. For ISBI 2012, 20 images were used for training and 10  
408 images were used for testing.

409 **Loss function and optimization.** For each algorithm, we used the same loss function and  
410 optimization method. Specifically, focal loss and dice loss were chosen. Define  $y$  as the ground  
411 truth segmentation and  $\hat{y}$  as the predicted segmentation. The calculations of focal loss and dice

**Table 3.** Parameter settings.

Parameters	U-Net-*	CASENet-*	LinkNet-*	U-Net-transfer
Data augmentation	√	√	√	√
Pre-training	-	-	-	√
Learning Rate	1e-3	1e-7	5e-4	2e-5
Batch Size	4	2	1	4
GPUs	4	4	4	8
Epoch	100	300	300	50
Worker	16	16	8	32

412 loss were:

$$L_{\text{Focal}} = -(1 - \hat{y})^{\gamma} \log(\hat{y}) \quad (1)$$

413

$$L_{\text{Dice}} = 1 - \frac{2y\hat{y} + 1}{y + \hat{y} + 1} \quad (2)$$

414 And the final loss function is the summation of the two losses with the proportion of 1 :  $\lambda$ . That  
 415 is  $L = L_{\text{Focal}} + \lambda L_{\text{Dice}}$ . We set  $\lambda = 1$  in these experiments. When optimizing the parameters in the  
 416 network, we chose Adam (*Kingma and Ba, 2014*) as the optimizer.

417 **Parameter settings.** Data augmentation (random horizontal/vertical flip, random rotation,  
 418 random zoom, random cropping, random cropping, random translation, random contrast, and  
 419 random color jitter) were used. Four Nvidia V100 GPUs were used for training. In the testing stage,  
 420 the original images were cut into the same size as the training images, and the patches were tested.  
 421 These patches were eventually mosaiced back to the original size for evaluation. The parameters  
 422 settings are shown in **Table 3**. Mean value and standard error are computed over testing images of  
 423 each dataset. The methods with "-\*" in the table represent that they are implemented by ourselves.

#### 424 Image definition criteria

425 Four representative image definition criteria, **Brenner** (*Subbarao and Tyan, 1998*), **SMD2** (*Thakkins-*  
 426 *tian et al., 2005*), **Variance** (*Saltelli et al., 2010*), and **Vollath** (*Vollath, 2008*) were used for analyzing  
 427 the effects of downsampling on EM images. The former two consider the difference and variance  
 428 of gray values between adjacent pixels, while the latter two consider the whole image.

429 Brenner gradient function simply calculates the square of the gray difference between two  
 430 adjacent pixels.

$$D(f) = \sum_y \sum_x |f(x+2, y) - f(x, y)|^2. \quad (3)$$

431 where:  $f(x, y)$  represents the gray value of pixel  $(x, y)$  corresponding to image  $f$ , and  $D(f)$  is the  
 432 result of image definition calculation (the same below).

433 SMD2 multiplies two gray variances in each pixel field and then accumulates them one by one.

$$D(f) = \sum_y \sum_x |f(x, y) - f(x+1, y)| |f(x, y) - f(x, y+1)|. \quad (4)$$

434 The Variance function is defined as

$$D(f) = \sum_y \sum_x |f(x, y) - \mu|^2, \quad (5)$$

435 where:  $\mu$  is the average gray value of the whole image, which is sensitive to noise. The purer the  
 436 image, the smaller the function value.

437

438 The Vollath function is defined as

$$D(f) = \sum_y \sum_x f(x, y) \cdot f(x + 1, y) - M \cdot N \cdot \mu^2, \quad (6)$$

439 where:  $\mu$  is the average gray value of the whole image,  $M$  and  $N$  are the width and height of the  
440 image respectively.

#### 441 Attribution analysis

442 The purpose of the integral gradient (IG) method is to quantify the contribution of each part of the  
443 input feature to the decision. For a given input image  $x$  and model  $F(x)$ , the goal of the network is  
444 to find out which pixels or features in  $x$  have an important influence on the decision-making of the  
445 model or sort the importance of each pixel or feature in  $x$ . Such a process is defined as attribution.  
446 IG uses the integral value along the whole gradient line from input to output. In the cell membrane  
447 segmentation task, for the decision of a pixel of  $y$  (predicted as the cell membrane or not), we can  
448 get the contribution of each pixel of the input image. Put the contribution of each pixel together, we  
449 record it as an **Attribution Field**  $A$ , whose size is the same as the original image. And the value of  
450 each pixel is  $w_{i,j}$ , representing the contribution decision of pixel  $x_{i,j}$  to  $y$ .  $w_{i,j}$  is normalized to  $[-1, 1]$ .

451 In binary segmentation task, for the current input image  $x$ , if we know that the output  $y$  is a  
452 specific value, such as  $y = 0$ , and the corresponding reference image is  $x'$ , then we can take a linear  
453 interpolation, that is

$$x' + \alpha (x - x'). \quad (7)$$

454 If the constant  $\alpha = 0$ , then the input image is the base image as  $x'$ . And if  $\alpha = 1$ , then the input  
455 image is the current image, which is  $x$ . When  $0 \leq \alpha \leq 1$ , it can be other images.

456 For the output of the neural network  $F(x)$ , the formula of the IG method is shown as

$$A(F, x, x') = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\delta F(x' + \alpha(x - x'))}{\delta x_i} d\alpha. \quad (8)$$

457 In formula 8,  $\delta x_i$  on the denominator denotes variation. This design makes the whole partial  
458 derivative transform into the form of variation. The variation boundary is the reference image and  
459 the current image. The integral gradient method uses linear interpolation as the variational path.  
460 That is,

$$\gamma(\alpha) = x' + \alpha (x - x'). \quad (9)$$

461 Here we select the random noise image as the reference image.

462 As the resolution and image size of U-RISC and ISBI 2012 are different, for a fair comparison, we  
463 define the size of **Pixel Attribution Field** as  $S_k$ , which represents the physical size corresponding  
464 to the pixel area with fixed contribution value threshold  $k$ . If the contribution value  $w_{i,j}$  is greater  
465 than  $k$ , the pixel is the one with higher contribution to decision-making. Area of attribution field  $S_k$   
466 is obtained by multiplying the areas (attribution value  $w_{i,j} > k$ ) and the corresponding physical  
467 size of a pixel (square of resolution  $h$ ).

$$S_k = S(A_{w_{i,j}>k}) \times h^2, w_{i,j} \in A. \quad (10)$$

#### 468 Data analysis

469 All statistical tests used, including statistic values and sample sizes, are provided in the figure  
470 captions, including mean and standard. All analyses were performed using custom software  
471 developed using the following tools and software: MATLAB (R2018a), Python (3.6), PyTorch (1.6.0),  
472 NumPy (1.19.0), SciPy (1.5.1), and matplotlib (2.2.3).

#### 473 Acknowledgments

474 We thank Bryan W. Jones for kindly providing RC1 connectomic dataset. This project is supported  
475 by Cognitive Computing Grant No.62088102.



## References

- 476  
477 **Adadi A**, Berrada M. Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE*  
478 *Access*. 2018; 6:52138–52160.
- 479 **Alejo R**, Monroy de Jesús J, Pacheco Sánchez JH, López González E, Antonio Velázquez JA. A selective dynamic  
480 sampling back-propagation approach for handling the two-class imbalance problem. *Applied Sciences*. 2016;  
481 6(7):200.
- 482 **Ancona M**, Ceolini E, Öztireli C, Gross M. Gradient-based attribution methods. In: *Explainable AI: Interpreting,*  
483 *Explaining and Visualizing Deep Learning*; 2019.p. 169–191.
- 484 **Anderson JR**, Jones BW, Watt CB, Shaw MV, Yang JH, DeMill D, Lauritzen JS, Lin Y, Rapp KD, Mastronarde D, et al.  
485 Exploring the retinal connectome. *Molecular Vision*. 2011; 17:355–79.
- 486 **Anderson JR**, Jones BW, Yang JH, Shaw MV, Watt CB, Koshevoy P, Spaltenstein J, Jurrus E, Kannan U, Whitaker  
487 RT, et al. A computational framework for ultrastructural mapping of neural circuitry. *PLoS Biol*. 2009;  
488 7(3):e1000074.
- 489 **Arganda-Carreras I**, Turaga SC, Berger DR, Cireşan D, Giusti A, Gambardella LM, Schmidhuber J, Laptev D,  
490 Dwivedi S, Buhmann JM, et al. Crowdsourcing the creation of image segmentation algorithms for connectomics.  
491 *Frontiers in Neuroanatomy*. 2015; 9:142.
- 492 **Asuni N**, Giachetti A. Accuracy improvements and artifacts removal in edge based image interpolation. *VISAPP*.  
493 2008; 8:58–65.
- 494 **Bach S**, Binder A, Montavon G, Klauschen F, Müller KR, Samek W. On pixel-wise explanations for non-linear  
495 classifier decisions by layer-wise relevance propagation. *PLoS One*. 2015; 10(7):e0130140.
- 496 **y Cajal SR**. Estructura de los centros nerviosos de las aves; 1888.
- 497 **Chaurasia A**, Culurciello E. Linknet: Exploiting encoder representations for efficient semantic segmentation. In:  
498 *Proceedings of IEEE Visual Communications and Image*; 2017. p. 1–4.
- 499 **Chen LC**, Papandreou G, Kokkinos I, Murphy K, Yuille AL. Semantic image segmentation with deep convolutional  
500 nets and fully connected crfs. *arXiv preprint arXiv:14127062*. 2014; .
- 501 **Clark K**, Khandelwal U, Levy O, Manning CD. What does bert look at? an analysis of bert's attention. *arXiv*  
502 *preprint arXiv:190604341*. 2019; .
- 503 **Conrad R**, Narayan K. CEM500K, a large-scale heterogeneous unlabeled cellular electron microscopy image  
504 dataset for deep learning. *Elife*. 2021; 10:e65894.
- 505 **Cui Y**, Jia M, Lin TY, Song Y, Belongie S. Class-balanced loss based on effective number of samples. In: *Proceedings*  
506 *of the IEEE Conference on Computer Vision and Pattern Recognition*; 2019. p. 9268–9277.
- 507 **DeBello WM**, McBride TJ, Nichols GS, Pannoni KE, Sanculi D, Totten DJ. Input clustering and the microscale  
508 structure of local circuits. *Frontiers in neural circuits*. 2014; 8:112.
- 509 **Deng J**, Dong W, Socher R, Li LJ, Li K, Fei-Fei L. Imagenet: A large-scale hierarchical image database. In: *2009 IEEE*  
510 *conference on computer vision and pattern recognition*; 2009. p. 248–255.
- 511 **Dice LR**. Measures of the amount of ecologic association between species. *Ecology*. 1945; 26(3):297–302.
- 512 **Fei-Fei L**, Fergus R, Perona P. One-shot learning of object categories. *IEEE transactions on pattern analysis and*  
513 *machine intelligence*. 2006; 28(4):594–611.
- 514 **Fernández A**, García S, Herrera F, Chawla NV. SMOTE for learning from imbalanced data: progress and  
515 challenges, marking the 15-year anniversary. *Journal of Artificial Intelligence Research*. 2018; 61:863–905.
- 516 **Goceri E**. Challenges and recent solutions for image segmentation in the era of deep learning. In: *International*  
517 *Conference on Image Processing Theory, Tools and Applications*; 2019. p. 1–6.
- 518 **Golgi C**. Sulla fina anatomia degli organi centrali del sistema nervoso; 1885.
- 519 **Guidotti R**, Monreale A, Ruggieri S, Turini F, Giannotti F, Pedreschi D. A survey of methods for explaining black  
520 box models. *ACM computing surveys (CSUR)*. 2018; 51(5):1–42.

- 521 **Hariharan B**, Arbeláez P, Girshick R, Malik J. Hypercolumns for object segmentation and fine-grained localization.  
522 In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2015. p. 447–456.
- 523 **He K**, Fan H, Wu Y, Xie S, Girshick R. Momentum contrast for unsupervised visual representation learning. In:  
524 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2020. p. 9729–9738.
- 525 **He K**, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference*  
526 *on Computer Vision and Pattern Recognition*; 2016. p. 770–778.
- 527 **Heimann T**, Van Ginneken B, Styner MA, Arzhaeva Y, Aurich V, Bauer C, Beck A, Becker C, Beichel R, Bekes G,  
528 et al. Comparison and evaluation of methods for liver segmentation from CT datasets. *IEEE transactions on*  
529 *medical imaging*. 2009; 28(8):1251–1265.
- 530 **Hwang JW**, Lee HS. Adaptive image interpolation based on local gradient features. *IEEE Signal Processing*  
531 *Letters*. 2004; 11(3):359–362.
- 532 **Kingma DP**, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980. 2014; .
- 533 **Kornfeld J**, Denk W. Progress and remaining challenges in high-throughput volume electron microscopy. *Current*  
534 *Opinion in Neurobiology*. 2018; 50:261–267.
- 535 **Kosub S**. A note on the triangle inequality for the Jaccard distance. *Pattern Recognition Letters*. 2019; 120:36–38.
- 536 **Krizhevsky A**, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. *Advances*  
537 *in Neural Information Processing Systems*. 2012; 25:1097–1105.
- 538 **Lee K**, Zung J, Li P, Jain V, Seung HS. Superhuman accuracy on the SNEMI3D connectomics challenge. arXiv  
539 preprint arXiv:1706.00120. 2017; .
- 540 **Leitch B**. Ultrastructure of electrical synapses. *Electron microscopy reviews*. 1992; 5(2):311–339.
- 541 **Li DC**, Liu CW, Hu SC. A learning method for the class imbalance problem with medical data sets. *Computers in*  
542 *Biology and Medicine*. 2010; 40(5):509–518.
- 543 **Lin TY**, Goyal P, Girshick R, He K, Dollár P. Focal loss for dense object detection. In: *Proceedings of the IEEE*  
544 *International Conference on Computer Vision*; 2017. p. 2980–2988.
- 545 **Lin W**, Dong L. Adaptive downsampling to improve image compression at low bit rates. *IEEE Transactions on*  
546 *Image Processing*. 2006; 15(9):2513–2521.
- 547 **Liu L**, Wu FX, Wang YP, Wang J. Multi-Receptive-Field CNN for Semantic Segmentation of Medical Images. *IEEE*  
548 *Journal of Biomedical and Health Informatics*. 2020; 24(11):3215–3225.
- 549 **Liu Y**, Yu J, Han Y. Understanding the effective receptive field in semantic image segmentation. *Multimedia*  
550 *Tools and Applications*. 2018; 77(17):22159–22171.
- 551 **Marin D**, He Z, Vajda P, Chatterjee P, Tsai S, Yang F, Boykov Y. Efficient segmentation: Learning downsampling  
552 near semantic boundaries. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*; 2019. p.  
553 2131–2141.
- 554 **Mastrorade DN**. Automated electron microscope tomography using robust prediction of specimen move-  
555 ments. *Journal of structural biology*. 2005; 152(1):36–51.
- 556 **Montavon G**, Lapuschkin S, Binder A, Samek W, Müller KR. Explaining nonlinear classification decisions with  
557 deep Taylor decomposition. *Pattern Recognition*. 2017; 65:211–222.
- 558 **Ronneberger O**, Fischer P, Brox T. U-Net: Convolutional networks for biomedical image segmentation. In:  
559 *International Conference on Medical Image Computing and Computer-Assisted Intervention*; 2015. p. 234–241.
- 560 **Roy AG**, Navab N, Wachinger C. Concurrent spatial and channel ‘squeeze and excitation’ in fully convolutional  
561 networks. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*; 2018. p.  
562 421–429.
- 563 **Russakovsky O**, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, et al.  
564 Imagenet large scale visual recognition challenge. *International journal of computer vision*. 2015; 115(3):211–  
565 252.

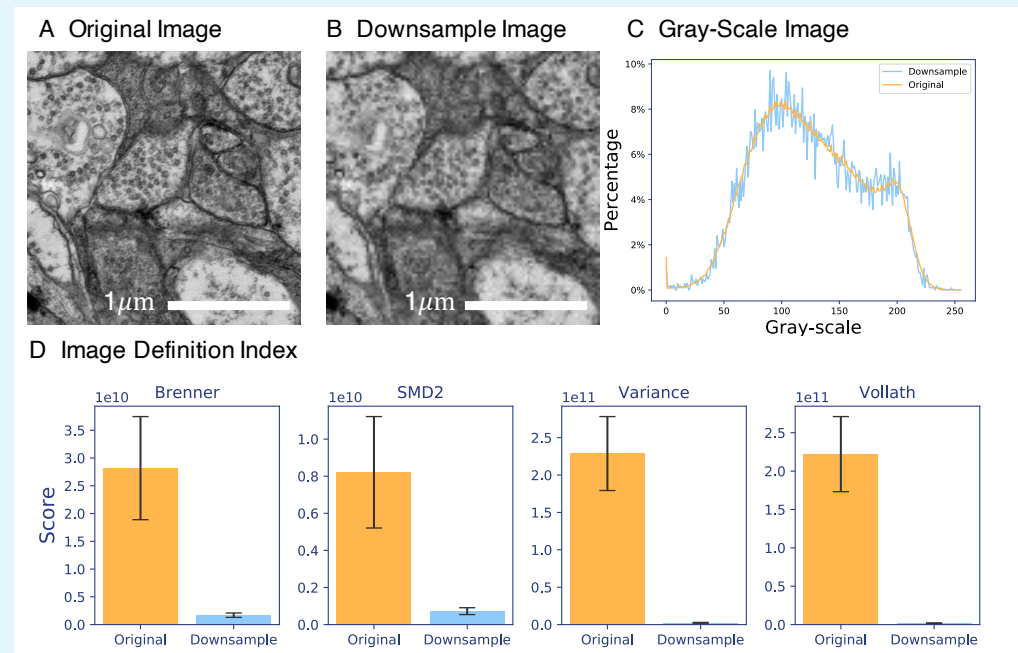
- 566 **Saltelli A**, Annoni P, Azzini I, Campolongo F, Ratto M, Tarantola S. Variance based sensitivity analysis of model  
567 output. Design and estimator for the total sensitivity index. *Computer Physics Communications*. 2010;  
568 181(2):259–270.
- 569 **Schulz K**, Sixt L, Tombari F, Landgraf T. Restricting the flow: Information bottlenecks for attribution. *arXiv*  
570 preprint arXiv:200100396. 2020; .
- 571 **Shapson-Coe A**, Januszewski M, Berger DR, Pope A, Wu Y, Blakely T, Schalek RL, Li P, Wang S, Maitin-Shepard J,  
572 et al. A connectomic study of a petascale fragment of human cerebral cortex. *bioRxiv*. 2021; .
- 573 **Shawn M**. Progress Towards Mammalian Whole-Brain Cellular Connectomics. *Frontiers in Neuroanatomy*. 2016;  
574 10.
- 575 **Simonyan K**, Zisserman A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint*  
576 arXiv:14091556. 2014; .
- 577 **Subbarao M**, Tyan JK. Selecting the optimal focus measure for autofocusing and depth-from-focus. *IEEE*  
578 *transactions on pattern analysis and machine intelligence*. 1998; 20(8):864–870.
- 579 **Sundararajan M**, Taly A, Yan Q. Axiomatic attribution for deep networks. In: *Proceedings of the 34th International*  
580 *Conference on Machine Learning*, vol. 70; 2017. p. 3319–3328.
- 581 **Szegedy C**, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A. Going deeper with  
582 convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2015. p. 1–9.
- 583 **Tan M**, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International*  
584 *Conference on Machine Learning*; 2019. p. 6105–6114.
- 585 **Thakkinstian A**, McElduff P, D'Este C, Duffy D, Attia J. A method for meta-analysis of molecular association  
586 studies. *Statistics in Medicine*. 2005; 24(9):1291–1306.
- 587 **Vollath D**. Nanomaterials an introduction to synthesis, properties and application. *Environmental Engineering*  
588 *and Management Journal*. 2008; 7(6):865–870.
- 589 **Vu T**, Van Nguyen C, Pham TX, Luu TM, Yoo CD. Fast and efficient image quality enhancement via desubpixel  
590 convolutional neural networks. In: *Proceedings of the European Conference on Computer Vision Workshops*; 2018.  
591 p. 0–0.
- 592 **Wang R**, Gong M, Tao D. Receptive field size versus model depth for single image super-resolution. *IEEE*  
593 *Transactions on Image Processing*. 2019; 29:1669–1682.
- 594 **Wang Z**, Bovik AC. Modern image quality assessment. *Synthesis Lectures on Image, Video, and Multimedia*  
595 *Processing*. 2006; 2(1):1–156.
- 596 **Wang Z**, Simoncelli EP, Bovik AC. Multiscale structural similarity for image quality assessment. In: *The Thirty-*  
597 *Seventh Asilomar Conference on Signals, Systems and Computers, 2003*, vol. 2; 2003. p. 1398–1402.
- 598 **White JG**, Southgate E, Thomson JN, Brenner S. The structure of the nervous system of the nematode *Caenorhab-*  
599 *ditis elegans*. *Philos Trans R Soc Lond B Biol Sci*. 1986; 314(1165):1–340.
- 600 **Wu X**, Zhang X, Wang X. Low bit-rate image compression via adaptive down-sampling and constrained least  
601 squares upconversion. *IEEE Transactions on Image Processing*. 2009; 18(3):552–561.
- 602 **Yen SJ**, Lee YS. Cluster-based under-sampling approaches for imbalanced data distributions. *Expert Systems*  
603 *with Applications*. 2009; 36(3):5718–5727.
- 604 **Yoo J**, Ahn N, Sohn KA. Rethinking data augmentation for image super-resolution: A comprehensive analysis  
605 and a new strategy. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2020.  
606 p. 8375–8384.
- 607 **Yu Z**, Feng C, Liu M, Srikumar R. Casenet: Deep category-aware semantic edge detection. In: *Proceedings of the*  
608 *IEEE Conference on Computer Vision and Pattern Recognition*; 2017. p. 5964–5973.
- 609 **Zhang K**, Wu Z, Yuan D, Luan J, Jia J, Meng H, Song B. Re-weighted interval loss for handling data imbalance  
610 problem of end-to-end keyword spotting. *Proceedings of Interspeech*. 2020; p. 2567–2571.
- 611 **Zhu Z**, Xu Z, You A, Bai X. Semantically multi-modal image synthesis. In: *Proceedings of the IEEE/CVF Conference*  
612 *on Computer Vision and Pattern Recognition*; 2020. p. 5467–5476.

## 613 Appendix 1

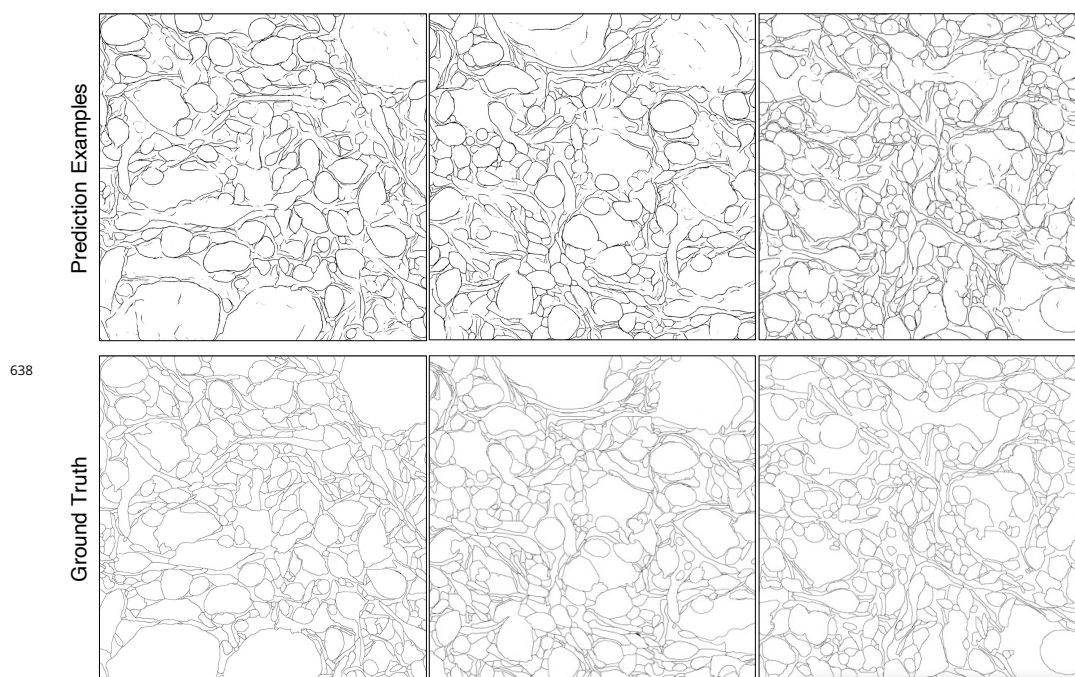
### 614 The Effects of Downsampling on U-RISC Dataset

615 To our best knowledge, downsampling is an effective approach to process images with  
616 enormous size in segmentation tasks (Marin et al., 2019; Chen et al., 2014). With the purpose  
617 of investigating the effects of downsampling on image definition quality, we analyzed the  
618 gray-scale histograms on a group of original and downsampled images in U-RISC and then  
619 calculated the definition values of images in Track 2.

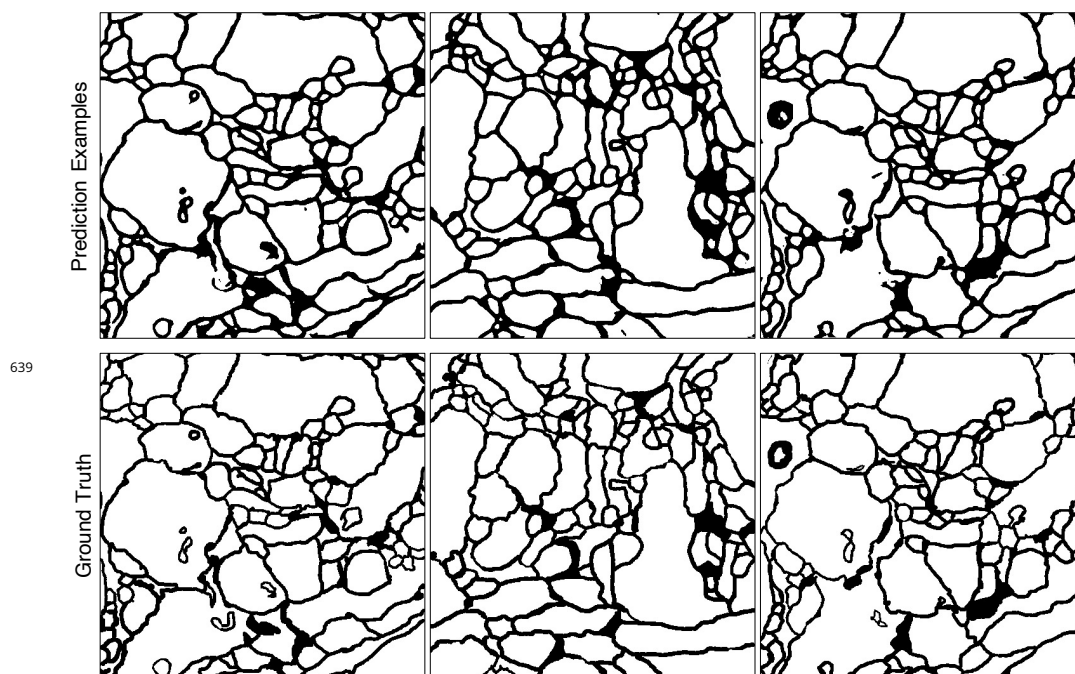
620 We first found that downsampling produced numerous sharp changes between adja-  
621 cent gray values (Figure 1(C)). In our analysis, the mutations might engender the texture  
622 information in the image to be blurred. For example, the bilayer structure of the membrane  
623 disappeared after downsampled (Figure 1 (A,B)), and some of the cell membranes which  
624 were hard to recognized became obscured. We then found that the defining quality of the  
625 images in Track 2 became lower after downsampling. Four definition indexes of all the  
626 images in Track 2 were calculated. Brenner (Subbarao and Tyan, 1998), Variance (Subbarao  
627 and Tyan, 1998), SMD2 (Thakkinstian et al., 2005), and Vollath (Vollath, 2008) are common  
628 indexes to show the gray value change between adjacent pixels. Results suggested that  
629 image indexes were significantly decreased after downsampling (Figure 1(D)), and all the  
630 four indexes dropped about ten times. Therefore, the analysis indicated that images were  
631 heavily blurred after the downsampling operation.



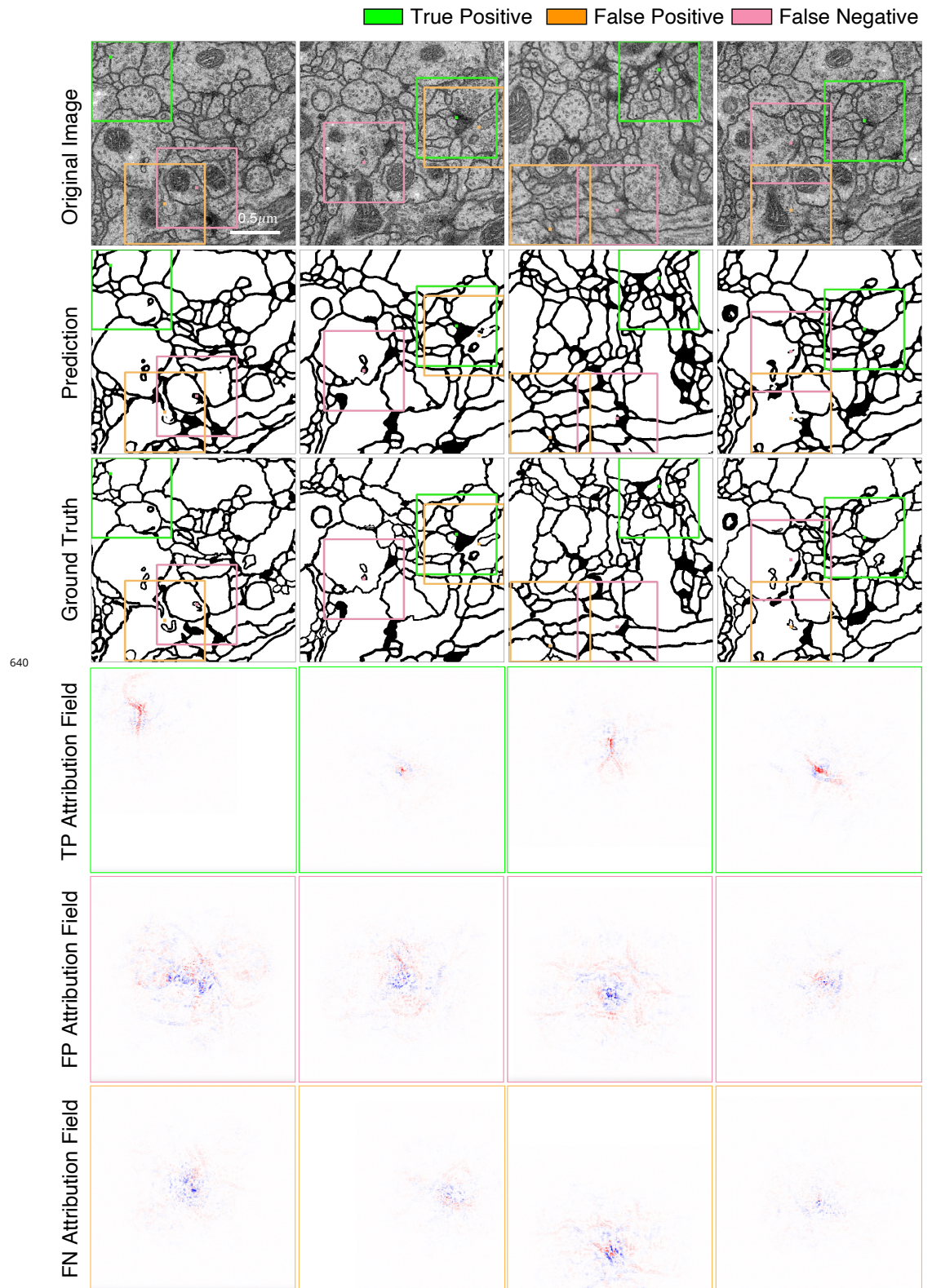
632 **Appendix 1 Figure 1.** Differences between original image and downsampled image.(A) The crop of  
633 original image. (B) The crop of downsampled image at the same position. The Gray-scale Histograms is  
634 calculated on A and B. (C) The scores of definition indexes calculated on the whole U-RISC dataset  
635 before and after downsampling. Details of indexes are described in Methods and Materials.  
636



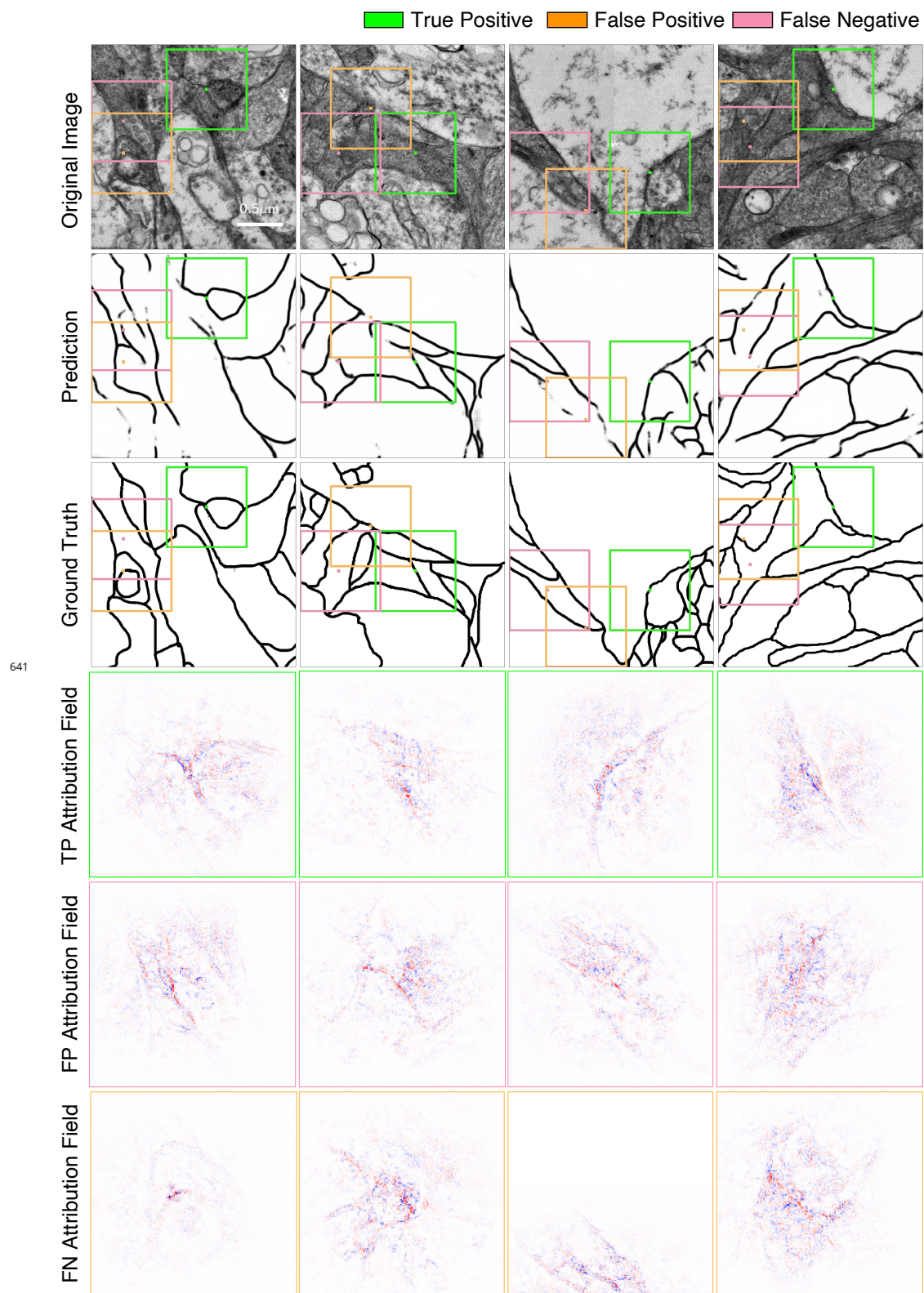
**Figure 6-Figure supplement 1.** Supplements for segmentation predictions of U-RISC.



**Figure 6-Figure supplement 2.** Supplements for segmentation predictions of ISBI 2012.



**Figure 8-Figure supplement 1.** Supplements for attribution analysis on ISBI 2012.



**Figure 8-Figure supplement 2.** Supplements for attribution analysis on U-RISC.